
Applied geostatistics

Tutorial: Regional mapping of climate variables from point samples

Ordinary Least Squares trend

Generalized Least Squares trend

Regression Kriging

Kriging with External Drift

Generalized Additive Models trend

Geographically-weighted regression

Data-driven methods: Regression trees, Random Forests, Cubist

Thin-plate splines

Local interpolators: Ordinary kriging, inverse-distance, Thiessen polygons

D G Rossiter

Cornell University, Section of Soil & Crop Sciences

ISRIC–World Soil Information

南京师范大学地理学学院

July 26, 2019

Contents

1	Introduction	1
1.1	Climate variables	2
2	Data manipulation	3
2.1	ESRI shapefiles for Northeast climate	3
2.2	Subsetting to four States	7
2.3	* Checking for duplicate locations	9
2.4	Projecting from geographic to metric CRS	10

Version 5.4 Copyright © 2014-2019 D G Rossiter. All rights reserved. Reproduction and dissemination of the work as a whole (not parts) freely permitted if this original copyright notice is included. Sale or placement on a web site where payment must be made to access this document is strictly prohibited. To adapt or translate please contact the author (d.g.rossiter@cornell.edu).

2.5	* State boundaries	13
2.6	Saving the dataset	16
3	Data exploration	16
3.1	Feature-space summary	16
3.2	Station locations	17
3.3	Postplots	17
3.4	* Viewing in geographic context	19
4	Naïve linear modelling	22
4.1	Exploring the relation between predictors and predictand	22
4.2	OLS fit to the linear model	24
4.3	Data cleaning	26
4.4	Spatial correlation of OLS model residuals	30
4.4.1	The empirical variogram	31
4.4.2	Fitting an authorized variogram model	33
4.5	* Close-range anomalies	34
5	Generalized least squares with REML	36
5.1	* GLS - theory	36
5.2	GLS - practice	38
5.3	Spatial correlation of GLS model residuals	41
5.4	* Fitting a GLS correlation structure with a nugget variance	43
5.5	Comparing OLS and GLS models	46
6	Creating a regional grid for mapping	49
6.1	* Creating a grid from global SRTM	49
6.2	Adjusting the grid for prediction	54
7	Prediction over the regional grid by OLS and GLS	55
8	Improving the trend prediction by regression kriging	59
8.1	* The Ordinary Kriging system	60
8.2	Predicting the residuals by Ordinary Kriging	62
8.3	The GLS-Regression Kriging prediction	64
9	Kriging with external drift (KED)	65
9.1	* The Universal Kriging system	66
9.2	Computing the empirical residual variogram	67
9.3	Fitting the residual variogram model	67
9.4	Predicting with KED	68
9.5	Accuracy assessment	72
9.6	KED in a local neighbourhood	73
9.7	* Demonstration that KED uses GLS to determine the trend	77
10	Generalized Additive Models	78
10.1	Fitting a Generalized Additive Model	81
10.2	GAM prediction over the study area	87
11	Geographically Weighted Regression	91
11.1	Theory	92

11.2 Dataset	94
11.3 Modelling approach	94
11.4 Determining the bandwidth	95
11.5 Computing the GWR	96
11.6 Interpretation	97
11.6.1 Feature space distribution	97
11.6.2 Statistical significance	101
11.7 Predictions	104
11.7.1 Predictions at known points	105
11.7.2 Predictions and model coefficients over the grid	105
11.7.3 Spatial distribution of coefficients	108
12 Data-driven models	111
12.1 Regression trees	112
12.1.1 Fitting a regression tree model	112
12.1.2 Regression tree prediction over the study area	121
12.2 Random forests	122
12.2.1 Fitting a Random Forest model	123
12.2.2 Random Forest prediction over the study area	129
12.3 Tuning data-driven models	132
12.4 Cubist	139
13 Thin-plate spline interpolation	147
13.1 * Theory	147
13.2 Practice	149
14 Local interpolators	152
14.1 Computing the empirical variogram	153
14.2 Fitting an authorized variogram model	155
14.3 Predicting by Ordinary Kriging	156
14.3.1 Accuracy assessment	159
14.4 Inverse-distance interpolation	162
14.5 Thiessen polygons	164
14.5.1 Accuracy assessment	169
15 Comparing the spatial patterns of two climate variables	172
15.1 Comparing variables with RK-GLS	175
15.2 Comparing variables with Random Forests	190
16 Answers	194
17 Challenge	201
References	203
A * Colour ramps with ggplot2	205
Index of R concepts	209

1 Introduction

This exercise presents various methods for regional mapping of climate variables from station information. The methods all relate to the **universal model of spatial distribution**:

$$Z(\mathbf{s}) = Z^*(\mathbf{s}) + \varepsilon(\mathbf{s}) + \varepsilon'(\mathbf{s}) \quad (1)$$

(\mathbf{s}) : a location in space, designated by a **vector** of coördinates;

$Z(\mathbf{s})$: **true** (unknown) value of some property at the location;

- when modelled, expressed as **most likely** value and some **uncertainty**, or as a **probability distribution**

$Z^*(\mathbf{s})$: **deterministic** component, due to some known or modelled **non-stochastic** process which operates over the entire region;

$\varepsilon(\mathbf{s})$: locally **spatially-autocorrelated stochastic** component;

$\varepsilon'(\mathbf{s})$: pure (“white”) **noise**, no structure.

In regional mapping, the deterministic component $Z^*(\mathbf{s})$ is called a **trend**, and models of these are called **trend surfaces**.

The focus is on prediction by Regression Kriging with Generalized Least Squares fitting of linear models (RK-GLS). This technique accounts for spatial correlation among the residuals $\varepsilon(\mathbf{s})$ of the trend surface $Z^*(\mathbf{s})$, fitting both together.

We contrast RK-GLS with predictions made by:

1. OLS trend surface;
2. GLS trend surface;
3. Kriging with External Drift (KED);
4. Data-driven (machine-learning) methods: Regression Trees, Random Forests, Cubist;
5. Generalized Additive Models (GAM);
6. Geographically-weighted regression (GWR);
7. Thin-plate splines;
8. Local interpolators: Ordinary Kriging (OK), Inverse Distance Weighting (IDW), Thiessen polygons.

This exercise also gives some practice in importing, manipulating, and verifying a large dataset, and also gives some practice with the `ggplot2`, `nlme`, `rgdal`, `sp`, `gstat`, `rpart`, `randomForest`, `ranger`, `Cubist`, `raster`, `plotKML` and `fields` R packages.

The exercise is organized as a set of **discussions**, **tasks**, **R code** to complete the tasks, self-study **questions** (with answers) to test your understanding, and a few **challenges**. §17 is an exercise to apply these techniques to a different study area and/or a different climate variable.

Note: The source for this document is a text file that includes ordinary \LaTeX source and “chunks” of R source code, using the Noweb¹ syntax. The formatted R source code, R text output, and R graphs in this document were automatically generated and incorporated into a \LaTeX source file by running the Noweb source document through R, using the `knitr` package [35]. The \LaTeX source was then compiled by \LaTeX into the PDF you are reading. The source code is provided in file `exRKGLS.R`.

1.1 Climate variables

We use agricultural climate as an example. The Northeast Regional Climate Center² has kindly provided a set of point ESRI shapefiles with various variables related to agricultural climate, for the entire USA and its possessions, in addition to the locations of weather stations and co-variables such as elevation. We use as an example *growing degree days, base 50°F*³ (GDD50);

Note: Temperatures T are expressed in °F in the USA, Bahamas, Belize, the Cayman Islands, and Palau. The conversion is $1^\circ\text{C} = (^\circ\text{F} - 32) \times (5/9)$.

These are defined for one day as [19]:

$$\text{GDD50} = \max([(T_{\max} + T_{\min})/2] - T_{\text{base}}, 0) \quad (2)$$

where $T_{\text{base}} = 50$.

For example, a day with maximum 86°F and minimum 60°F would account for $146/2 - 50 = 23$ GDD50. These are summed over some time period; we will use the sum over the year, i.e., Annual GDD50. These heat units are well-correlated with seasonal maize growth and are used to select maize varieties on the basis of the number of GDD50 required to reach maturity.

The aim is to **predict** GDD50 at an arbitrary location, using the records from the known stations, possibly also using covariables that are available over the whole region: latitude, longitude and elevation above sea level, as well as local spatial structure. While predicting, we also can examine the predictive models to **understand** the causes for Annual GDD50 spatial variation.

In §15 we compare models and predictions of this variable with models and predictions for others from the same dataset.

¹ <http://www.cs.tufts.edu/~nr/noweb/>

² <http://www.nrcc.cornell.edu>

³ 10°C

2 Data manipulation

In this section we load files with the weather information, limit the geographic area to a study area, make both a spatial and non-spatial versions of the dataset and re-project the spatial version. If you are not interested in this process, you can skip to §3, where you can load the dataset produced in this section.

2.1 ESRI shapefiles for Northeast climate

The ESRI shapefiles have been provided in a compressed file named `weather_stn_sums_1971_2000.zip`.

TASK 1 : Locate this file, unpack it in a working directory named NEweather, start R (preferably via RStudio) and connect to that directory. •

If using RStudio, you can navigate to the directory where you've stored the climate files with the "Files" window pane, and then select menu command Session | Set Working Directory... | To Files Pane Location.

If entering an R command at the command prompt, within the RStudio console window or directly in R, use the `setwd` function with the path name as its argument.

In this example – that is, on my system – I have set up a directory named NEweather, under a subdirectory for datasets that I call ds, this under my home directory on a Unix-like system (including Mac OS X), symbolized as `~`. You can choose any location on your file system.

For example, if the NEweather directory is inside a ds 'datasets' directory under my home directory, the following command will connect to it on a Unix-like system:

```
setwd("~/ds/NEweather")
```

On a Windows system, you will need to specify a full path, e.g.:

```
setwd("M://css6200/dgr2/datasets/NEweather")
```

The location of these files on your system will almost surely be different! So **adjust the command accordingly**.

TASK 2 : List the shapefiles in this directory. •

The `list.files` function lists the files to the console output; you can also look at the directory in a file manager.

- Its first argument is the directory in which to look.
- The optional `pattern` argument gives a *regular expression* to match the file names. Here we just want the base shape files; there are other "helper" files with the same name but different file extensions, so we specify a pattern of the `shp` extension, at the end of

the string, as symbolized by the special \$ regular expression character.

If the files are in the directory where you are connected⁴, you can just specify ".", which is a Unix and R abbreviation for "the current directory". If they are in a subdirectory, you need to name that, using the forward slash "/" character to show you are descending the directory tree.

So, I find the unpacked files like this:

```
list.files(".", pattern=".shp$")  
## [1] "cb_2014_us_state_500k.shp" "extmin_7100j.shp"  
## [3] "frz28_7100j.shp"          "frz32_7100j.shp"  
## [5] "gdd40_7100j.shp"          "gdd50_7100j.shp"  
## [7] "maat7100.shp"             "map7100.shp"
```

Each of these shapefiles has associated metadata, with extension .xml. Although this can be read into an R data frame with the read.metadata function of the plotKML package, I don't know any way to format it as human-readable within R.

We can guess the file contents from the file names: extreme minima, frost-free days at 28 and 32°F⁵, growing degree days based on 40 and 50°F⁶, mean annual air temperature and precipitation.

TASK 3 : Load the shapefile for the all-USA GDD50 and examine its structure. •

The readOGR function of the rgdal "R implementation of the Geographic Data Abstraction Library" package can read many formats into sp objects implemented by the sp package; one of the formats is the ESRI shapefile.

- The first argument dsn is the "data set name", which in ESRI terminology is not a data set, but rather the directory where the shapefile is located.
- The second argument layer is the layer name, *without* an extension – this is because there are several files, with different extensions, associated with one shapefile.
- There are many optional arguments, with defaults, that control the input, see ?readOGR for details. One that is needed on **64-bit operating systems**⁷, but which will not cause any problems on other systems, is integer64. By default it is set to convert integers into string, to avoid data loss of large integers (needing more than 32 bits). The numbers in most datasets, including this one, are much smaller. So we set this argument to allow.loss; this will read integers in correctly as numbers.

⁴ use getwd to find this

⁵ -2.2 and 0°C

⁶ 4.44 and 10°C

⁷ for example, Microsoft Windows 10

If you have set the working directory to the location of the dataset, the data set name is written as ".", meaning "the currently connected directory". The dsn argument should be the same as the directory argument to list.files, above.

First load the required packages:

```
library(sp)
library(rgdal)
```

And then the dataset:

```
usa <- readOGR(dsn=".", layer="gdd50_7100j",
               integer64="allow.loss")

## OGR data source with driver: ESRI Shapefile
## Source: "/Users/rossiter/data/edu/dgeostats/ex/ds/NEweather", layer: "gdd50_7100j"
## with 5556 features
## It has 26 fields
## Integer64 fields read as signed 32-bit integers: JAN_GDD50 FEB_GDD50 MAR_GDD50 APR_GDD50
```

Geographic objects defined in the sp package all have a coordinate reference system (CRS), which defines how the coordinates for the object relate to the Earth's surface. This can be NA "not available" or a list of CRS parameters.

TASK 4 : Show the CRS for the climate stations shapefile. •

We use the proj4string function of the rgdal package to show the CRS for the usa object. The strwrap "line-wrap a character string" prevents this long string from going outside the page margin.

```
bbox(usa)

##           min      max
## coords.x1 -177.35 174.10
## coords.x2  -14.33  71.28

names(usa)

## [1] "STATION_ID" "STATE"      "STATION_NA" "LATITUDE_D" "LONGITUDE_"
## [6] "ELEVATION_" "OID_"       "COOP_ID"    "STATE_1"    "STN_NAME"
## [11] "LAT_DD"     "LONG_DD"    "ELEV_FT"    "JAN_GDD50"  "FEB_GDD50"
## [16] "MAR_GDD50"  "APR_GDD50"  "MAY_GDD50"  "JUN_GDD50"  "JUL_GDD50"
## [21] "AUG_GDD50"  "SEP_GDD50"  "OCT_GDD50"  "NOV_GDD50"  "DEC_GDD50"
## [26] "ANN_GDD50"

strwrap(proj4string(usa))

## [1] "+proj=longlat +datum=NAD27 +no_defs +ellps=clrk66"
## [2] "+nadgrids=@conus,@alaska,@ntv2_0.gsb,@ntv1_can.dat"
```

For completeness, we change some fields that were read as numeric to character strings using the as.character function. These fields are not numbers in the mathematical sense. Rather, they are identification codes, with no analytical meaning other than to identify locations.

```
sapply(usa@data, "class")

## STATION_ID      STATE STATION_NA LATITUDE_D LONGITUDE_ ELEVATION_
## "integer"      "factor"  "factor"  "numeric"  "numeric"  "integer"
##      OID_      COOP_ID  STATE_1  STN_NAME  LAT_DD     LONG_DD
## "integer"      "integer"  "factor"  "factor"   "integer"  "integer"
```

```

##      ELEV_FT  JAN_GDD50  FEB_GDD50  MAR_GDD50  APR_GDD50  MAY_GDD50
##      "integer" "integer" "integer" "integer" "integer" "integer"
##      JUN_GDD50  JUL_GDD50  AUG_GDD50  SEP_GDD50  OCT_GDD50  NOV_GDD50
##      "integer" "integer" "integer" "integer" "integer" "integer"
##      DEC_GDD50  ANN_GDD50
##      "integer" "integer"

usa$STATION_ID <- as.character(usa$STATION_ID)
usa$STATION_NA <- as.character(usa$STATION_NA)
usa$STN_NAME <- as.character(usa$STN_NAME)
usa$OID_ <- as.character(usa$OID_)
usa$COOP_ID <- as.character(usa$COOP_ID)
summary(usa)

## Object of class SpatialPointsDataFrame
## Coordinates:
##           min      max
## coords.x1 -177.35 174.10
## coords.x2  -14.33  71.28
## Is projected: FALSE
## proj4string :
## [+proj=longlat +datum=NAD27 +no_defs +ellps=clrk66
## +nadgrids=@conus,@alaska,@ntv2_0.gsb,@ntv1_can.dat]
## Number of points: 5556
## Data attributes:
##      STATION_ID      STATE      STATION_NA
##      Length:5556      TX      : 366      Length:5556
##      Class :character  CA      : 298      Class :character
##      Mode  :character  MT      : 215      Mode  :character
##                      OR      : 176
##                      CO      : 166
##                      UT      : 160
##                      (Other):4175
##      LATITUDE_D      LONGITUDE_      ELEVATION_
##      Min.   : -14.33      Min.   : -177.35      Min.   : -194
##      1st Qu.: 35.17      1st Qu.: -110.09      1st Qu.: 465
##      Median : 39.56      Median : -96.52      Median : 1030
##      Mean   : 39.32      Mean   : -97.55      Mean   : 1866
##      3rd Qu.: 43.44      3rd Qu.: -85.37      3rd Qu.: 2710
##      Max.   : 71.28      Max.   : 174.10      Max.   : 11320
##
##      OID_      COOP_ID      STATE_1
##      Length:5556      Length:5556      TX      : 366
##      Class :character  Class :character  CA      : 298
##      Mode  :character  Mode  :character  MT      : 215
##                      OR      : 176
##                      CO      : 165
##                      (Other):4300
##                      NA's   : 36
##      STN_NAME      LAT_DD      LONG_DD      ELEV_FT
##      Length:5556      Min.   : -14.00      Min.   : -177.0      Min.   : -194
##      Class :character  1st Qu.: 35.00      1st Qu.: -110.0      1st Qu.: 460
##      Mode  :character  Median : 40.00      Median : -96.0      Median : 1025
##                      Mean   : 39.12      Mean   : -96.9      Mean   : 1863
##                      3rd Qu.: 43.00      3rd Qu.: -85.0      3rd Qu.: 2705
##                      Max.   : 71.00      Max.   : 174.0      Max.   : 11320
##
##      JAN_GDD50      FEB_GDD50      MAR_GDD50      APR_GDD50
##      Min.   : 0.0      Min.   : 0.00      Min.   : 0      Min.   : 0.0
##      1st Qu.: 0.0      1st Qu.: 0.00      1st Qu.: 7      1st Qu.: 44.0
##      Median : 9.0      Median : 13.00      Median : 23      Median : 99.0
##      Mean   : 43.3      Mean   : 49.09      Mean   : 91      Mean   : 176.5
##      3rd Qu.: 29.0      3rd Qu.: 42.00      3rd Qu.: 118      3rd Qu.: 269.0
##      Max.   :1013.0      Max.   :918.00      Max.   :1026      Max.   :1038.0
##
##      MAY_GDD50      JUN_GDD50      JUL_GDD50      AUG_GDD50
##      Min.   : 0.0      Min.   : 0.0      Min.   : 0.0      Min.   : 0.0
##      1st Qu.: 198.0      1st Qu.: 405.0      1st Qu.: 579.0      1st Qu.: 539.0
##      Median : 324.0      Median : 567.0      Median : 734.0      Median : 682.0

```

```
## Mean : 367.6 Mean : 568.4 Mean : 724.7 Mean : 685.3
## 3rd Qu.: 517.0 3rd Qu.: 744.0 3rd Qu.: 902.0 3rd Qu.: 855.0
## Max. :1078.0 Max. :1335.0 Max. :1568.0 Max. :1512.0
##
## SEP_GDD50 OCT_GDD50 NOV_GDD50 DEC_GDD50
## Min. : 0.0 Min. : 0.0 Min. : 0.0 Min. : 0.00
## 1st Qu.: 277.8 1st Qu.: 61.0 1st Qu.: 8.0 1st Qu.: 0.00
## Median : 430.0 Median : 154.0 Median : 24.0 Median : 10.00
## Mean : 461.0 Mean : 224.7 Mean : 87.3 Mean : 47.68
## 3rd Qu.: 636.0 3rd Qu.: 341.0 3rd Qu.:108.0 3rd Qu.: 35.00
## Max. :1203.0 Max. :1023.0 Max. :996.0 Max. :1023.00
##
## ANN_GDD50
## Min. : 0
## 1st Qu.: 2155
## Median : 3060
## Mean : 3527
## 3rd Qu.: 4584
## Max. :11985
##
```

The ESRI shapefile is correctly imported as a spatial object, with the correct coördinate reference system (CRS), here geographic coördinates on the NAD27 datum, which uses the Clarke 1866 ellipsoid.

The metadata lists the field names, but without any explanation; fortunately they are fairly self-explanatory, for example STATE for the US State responsible for the weather station, STATION_NAME for the station name, ELEVATION_ for the elevation, etc. We see fields for GDD50 for each of 12 months (e.g., JAN_GDD50), as well as annual (ANN_GDD50). The bounding box covers almost the whole world from W to E. This is because the USA controls possessions in the western Pacific ocean; also, the Aleutian Islands of Alaska cross the international date line.

The records are from the 48 USA states, several territories (e.g., PR for Puerto Rico), and one unspecified, code ??.

```
length(levels(usa$STATE))

## [1] 54

levels(usa$STATE)

## [1] "?? "AK" "AL" "AR" "AZ" "CA" "CO" "CT" "DE" "FL" "GA" "HI" "IA"
## [14] "ID" "IL" "IN" "KS" "KY" "LA" "MA" "MD" "ME" "MI" "MN" "MO" "MS"
## [27] "MT" "NC" "ND" "NE" "NH" "NJ" "NM" "NV" "NY" "OH" "OK" "OR" "PA"
## [40] "PI" "PR" "RI" "SC" "SD" "TN" "TX" "UT" "VA" "VI" "VT" "WA" "WI"
## [53] "WV" "WY"
```

Challenge: What does code ?? represent? Hint: use the selection technique of the next section to examine these.

2.2 Subsetting to four States

TASK 5 : Select the records reported by NY State and its neighbours VT, NJ and PA.

Both the %in% set operator and the [] matrix selection operator find matches between strings, and return a list of matrix row numbers. This

us the `ix` vector, whose entries are either TRUE or FALSE, i.e., a vector of **logical** type. The `sum` function takes the TRUE values as 1 and the FALSE as 0, so the sum is the number of TRUE values, i.e., those that met the logical condition. The `which` function is shorthand for “which of the elements in the logical vector are TRUE?” It returns a vector of row numbers for which the logical expression, in this case, that the record is from one of the four states, is TRUE. The syntax `[ix,]` then selects only the records (dataframe rows) where the value is TRUE.

```
dim(usa)[1] # number of rows in the full dataset
## [1] 5556

ix <- (usa$STATE %in% c("NY", "VT", "NJ", "PA")) # logical expression
str(ix) # each row is either TRUE or FALSE

## logi [1:5556] FALSE FALSE FALSE FALSE FALSE FALSE ...

sum(ix) # number of rows selected
## [1] 305

head(which(ix)) # first six row numbers of stations in these states
## [1] 2852 2853 2854 2855 2856 2857

ne <- usa[ix,] # select only these states' records
dim(ne)[1] # number of rows in the restricted dataset
## [1] 305
```

The `dim` function returns the dimensions of a data frame as a two-element vector: number of rows (records) and columns (fields). The first dimension is the number of rows; we see how this reduces from 5556 in the full dataset to 305 in the four-state dataset.

Even after selection the `STATE` factor has all the original levels for the whole USA (i.e. all state abbreviations), even though most have no records.

Therefore we renumber the levels of the `STATE` factor, with the `factor` function, since we only have four states in this dataset. Although `STATE` was already an R factor (categorical variable), re-applying the `factor` function discards the unused levels (states we did not select) and re-numbers the factor levels.

```
class(ne$STATE)
## [1] "factor"

levels(ne$STATE)
## [1] "???" "AK" "AL" "AR" "AZ" "CA" "CO" "CT" "DE" "FL" "GA" "HI" "IA"
## [14] "ID" "IL" "IN" "KS" "KY" "LA" "MA" "MD" "ME" "MI" "MN" "MO" "MS"
## [27] "MT" "NC" "ND" "NE" "NH" "NJ" "NM" "NV" "NY" "OH" "OK" "OR" "PA"
## [40] "PI" "PR" "RI" "SC" "SD" "TN" "TX" "UT" "VA" "VI" "VT" "WA" "WI"
## [53] "WV" "WY"

ne$STATE <- factor(ne$STATE)
levels(ne$STATE)
## [1] "NJ" "NY" "PA" "VT"
```

TASK 6 : Display the bounding box. •

```
bbox(ne)

##           min      max
## coords.x1 -80.47 -71.98
## coords.x2  38.95  44.93
```

The bounding box is now much smaller than for the original dataset; it is just big enough to enclose the stations from the four selected states.

2.3 * Checking for duplicate locations

Obviously there can not be two climate stations at the exact same location. However, since the database records station coordinates in geographic units with a limited resolution, it may be that two close-by stations may be recorded with the same coordinates. That will cause various geostatistical procedures in analysis to fail, because these procedures do not know how to understand co-located points with different attribute values, e.g., different GDD50.

TASK 7 : Check for co-located points. •

The `zerodist` function of the `sp` package checks for this:

```
length(ix <- zerodist(ne))

## [1] 0
```

In this case there are no co-located points, so for the purposes of this exercise you can **skip to the next section** (§2.4).

However, elsewhere in this dataset there are some duplicates. If you take up the challenge at the end of this exercise (§17) to work on other areas of the USA, you may need to deal with these.

```
length(ix <- zerodist(usa))

## [1] 6

usa@data[ix,1:5]

##      STATION_ID STATE      STATION_NA LATITUDE_D LONGITUDE_
## 1280      123777  IN HARTFORD CITY 4 ESE      40.43      -85.28
## 3201      314260  NC      HOT SPRINGS      35.90      -82.83
## 5538      914111  PI              CHUUK AP       7.45      151.83
## 1311      127069  IN      PORTLAND 1 SW      40.43      -85.28
## 3202      314265  NC      HOT SPRINGS 2      35.90      -82.83
## 5553      914851  PI              TRUK AP       7.45      151.83
```

In the GDD50 file `gdd50_7100j.shp` there are one pair of duplicates in Indiana, one in North Carolina, and one in the Pacific Islands.

In case there are co-located points, there are two solutions:

1. Examine the record and compare with other information, e.g., an internet map (Google, Bing, OpenStreet ...), to see if one named

station is obviously mis-located. If there is strong evidence for this, change the coördinates for one of the stations.

For example, of the duplicate points in Indiana, we confirm from Google Earth that the reported coördinates are indeed about 4 miles ESE of Hartford City for record 1280. The same coördinates are given for 1 mile SW of Portland, record 1311. But Portland is quite a bit to the E, the centre of the town is at (40.43, -84.98). The longitude differs by about 0.3° from the reported value. Looking at the Google Earth image, we measure 1 mile SW from the centre of Portland, which lies in a large cemetery. Looking closely, we see the administration building, which is the likely location of the weather station. Clicking on this point gives a corrected coördinates of (40.42, -84.99).

To adjust the incorrect point, we should change both the coördinates in the `coords` slot, and also the coördinates given as attributes in the `data` slot:

```
coordinates(usa)[ix[1],,]; coordinates(usa)[ix[4],,]

## coords.x1 coords.x2
## -85.28 40.43
## coords.x1 coords.x2
## -85.28 40.43

usa@data[ix[1],4:5]; usa@data[ix[4],4:5]

## LATITUDE_D LONGITUDE_
## 1280 40.43 -85.28
## LATITUDE_D LONGITUDE_
## 1311 40.43 -85.28

usa@coords[ix[4],] <- c(-84.99, 40.42)
usa@data[ix[4],4:5] <- c(40.42, -84.99) # this is the order in the dataframe
```

2. Delete one of the co-located points. The choice on which to delete could come from examining the whole record to see which seems more reliable, for example, which has a reported elevation which matches best the location given by the coördinates. Or, one record has a longer time-series than the other, so is presumably more reliable for long-term averages. For example, if the fourth point in the list of duplicates should be eliminated, the code would look like this:

```
usa <- (usa[-ix[4],])
```

The `-` operator applied to the first matrix subscript means all rows (records) except those in the list.

2.4 Projecting from geographic to metric CRS

This shapefile is in geographic coördinates. For geostatistical analysis it's necessary to convert it to metric (projected) coördinates, i.e., where distances between points can be accurately measured. The UTM CRS is not suitable, because the of this shapefile bounding box is about 8.5° wide, and a UTM zone only covers 6°. Further, the region is oriented more or less E-W and not N-S.

Note: Selecting an appropriate CRS (projection and coordinate system) is not simple. If the study area is in a single political entity, e.g., a USA State, it's common to select the CRS that is used for that entity's own mapping. For example, many USA States have a State Plane Coordinate System. Otherwise we have to consider the shape of the area: E-W areas often use conic projections, N-S areas often use transverse Mercator (Gauss-Krüger) projections. All projections distort; for geostatistics and trend surfaces we want minimal distortion of distances. Several comprehensive references [4, 25, 26] explain projections, their properties, and typical uses.

A reasonable choice for this study area is an Albers Equal-area conic projection, optimized for the northeastern USA, and on the WGS84 ellipsoid. Although this is an equal-area projection, scale is not distorted along the two standard parallels which define it, and if these are carefully chosen the error in measuring distances (which is what we want) is quite small [25, §14].

The `proj.4` “Cartographic Projections, version 4” project provides standard Unix functions to define and transform projections. This is used in QGIS and in the R `rgdal` package, which is used by the `sp` and other spatial packages. Classes defined in the `sp` package, such as the climate stations dataset which is an `SpatialPointsDataFrame` object, include a `proj4string` slot, where the CRS parameters are stored as a string constant, according to the `proj.4` definitions.

For example, the climate stations were imported from an ESRI shapefile, with the CRS defined; we can see this with the `proj4string` function of the `rgdal` package:

```
strwrap(proj4string(ne))  
  
## [1] "+proj=longlat +datum=NAD27 +no_defs +ellps=clrk66"  
## [2] "+nadgrids=@conus,@alaska,@ntv2_0.gsb,@ntv1_can.dat"
```

These are geographic coordinates, as we can see from the `" +proj=longlat"` substring.

We need to find the parameters of a suitable metric CRS, and defined a string for it in `proj.4` format; this can then be used in the `spTransform` method of the `sp` package to define the projection from geographic to metric coordinates.

The list of available projections and their codes can be found with the `projInfo` function of the `rgdal` package, using the `type` argument:

```
head(projInfo(type="proj"))  
  
##      name      description  
## 1   aea      Albers Equal Area  
## 2   aeqd     Azimuthal Equidistant  
## 3   airy      Airy  
## 4 aitoff     Aitoff  
## 5   alsk Mod. Stereographic of Alaska  
## 6   apian     Apian Globular I
```

We see that the Albers Equal Area projection has code `"aea"`.

Projections and their parameters are found on the proj.4 project webpage⁸, under the “Projections” category. Unfortunately as of now information on the Albers Equal Area projection is not included. Another source is the projections transform list at the GeoTiff site⁹. Here¹⁰ we find the required parameters and how they are represented in the proj.4 string:

```
+proj=aea +lat_1=Latitude of first standard parallel
          +lat_2=Latitude of second standard parallel
          +lat_0=Latitude of false origin
          +lon_0=Longitude of false origin
          +x_0=Easting of false origin
          +y_0=Northing of false origin
```

The “false origin” is the (0,0) of the coordinates. Since this is a metric CRS, where distances are measured, we also need to specify units of measure with the `+units=` substring.

We enter the parameters required for the selected projection as a text string argument to the CRS function, which builds a coordinate reference system (CRS) string; this is then the argument to the `proj4string` function which assigns a projection to a spatial object; this is then found in the `proj4string` slot of the spatial object.

We choose parallels (parameters `+lat_1` and `+lat_2`) about 1° of latitude inside the N and S edges of the four states, and a meridian (parameter `+lon_0`) near the middle. For convenience we use integer meridians and parallels, estimated with the `round`, `floor` and `ceiling` functions. There is no mathematical reason for this.

```
round(median(bbox(ne)[1,]))
## [1] -76

floor(bbox(ne)[2,1])+1
## [1] 39

ceiling(bbox(ne)[2,2])-1
## [1] 44
```

We set the false origin (0, 0) (parameters `+lat_0` and `+lon_0`) near the centre of the map, at (76°W, 42.5°N); so the projected coordinates will be negative to the south and west of this point. This is an arbitrary choice, any point could be used for the origin of the projected coordinates.

TASK 8 : Reproject the shapefile to an Albers equal-area projection with standard parallels 39 and 44°N, central meridian 76°W, origin of coordinate (in meters) at this meridian and 42.5°N (the centre of the two standard parallels), on the WGS84 ellipsoid.

⁸ <http://proj4.org/>

⁹ http://geotiff.maptools.org/proj_list/

¹⁰ http://geotiff.maptools.org/proj_list/albers_equal_area_conic.html

Display the bounding box in the transformed coördinates. •

Note: These choices minimize the scale distortion across the mapped area. The choice of origin near the centre keeps the coördinates as small as possible; this is just for convenience and has no effect on the calculations.

The `spTransform` method of the `sp` package performs the transformation, when given an `sp` object and a new CRS, which is specified with the `CRS` function.

```
ne.crs <-  
  CRS("+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76 +ellps=WGS84 +units=m")  
ne.m <- spTransform(ne, ne.crs)
```

We also change the coördinate names to match the concept of Easting and Northing, using the `coordnames` function.

```
coordnames(ne.m) <- c("E", "N")  
bbox(ne.m)
```

```
##           min           max  
## E -375745.2 318960.4  
## N -393922.9 276614.1
```

We also need a data frame version of this object – where the coördinates are just fields and not treated specially as they are in `sp` classes – for some functions that we will use in the exploration and analysis in §5.2.

TASK 9: Make a data frame from this `sp` object. •

The `as` method converts any objects for which a method has been written. The `sp` package has a method to convert `SpatialPointsDataFrame` objects to `data.frame` objects.

This method takes the coördinates from the `sp` object and adds them to the fields of the data slot of the `SpatialPointsDataFrame` object, and together makes these fields of the `data.frame` object.

```
ne.df <- as(ne.m, "data.frame")  
names(ne.df)  
  
## [1] "STATION_ID" "STATE" "STATION_NA" "LATITUDE_D" "LONGITUDE_"  
## [6] "ELEVATION_" "OID_" "COOP_ID" "STATE_1" "STN_NAME"  
## [11] "LAT_DD" "LONG_DD" "ELEV_FT" "JAN_GDD50" "FEB_GDD50"  
## [16] "MAR_GDD50" "APR_GDD50" "MAY_GDD50" "JUN_GDD50" "JUL_GDD50"  
## [21] "AUG_GDD50" "SEP_GDD50" "OCT_GDD50" "NOV_GDD50" "DEC_GDD50"  
## [26] "ANN_GDD50" "E" "N"
```

2.5 * State boundaries

For display purposes we'd like to see the State boundaries of the four selected states. These are not used in the analysis, only in map display, so this section can be omitted if you are only interested in the analysis.

The US Census Bureau provides various cartographic boundary shapefiles¹¹ at different generalized scales to be used in small-scale mapping.

¹¹ <https://www.census.gov/geo/maps-data/data/tiger-cart-boundary.html>

This includes State boundaries¹². Since we are only mapping four States, we select the largest available scale, namely 1:500 000.

TASK 10 : Download the State boundaries shapefile for the USA from the US Census Bureau, save in your working directory, and unpack the compressed file into a subdirectory. •

The downloaded compressed file is named `cb_2014_us_state_500k.zip`.

TASK 11 : Read the State boundaries into R as a `SpatialPolygonsDataFrame` object and select only the four states of interest in this analysis, •

Assuming you've put the State boundary file into the same directory as the NE weather stations and unpacked to that directory, the shapefile is read in using `readOGR`:

```
state <- readOGR(dsn=".", "cb_2014_us_state_500k")
summary(state)

## OGR data source with driver: ESRI Shapefile
## Source: "/Users/rossiter/data/edu/dgeostats/ex/ds/NEweather", layer: "cb_2014_us_state_500k"
## with 56 features
## It has 9 fields
## Integer64 fields read as strings:  ALAND AWATER

## Warning in readOGR(dsn = "./ds/NEweather", "cb_2014_us_state_500k"): Z-dimension
## discarded

## Object of class SpatialPolygonsDataFrame
## Coordinates:
##      min      max
## x -179.1489 179.77847
## y  -14.5487  71.36516
## Is projected: FALSE
## proj4string :
## [+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80
## +towgs84=0,0,0]
## Data attributes:
##      STATEFP      STATENS      AFFGEOID      GEOID
## 01      : 1      00068085: 1      0400000US01: 1      01      : 1
## 02      : 1      00294478: 1      0400000US02: 1      02      : 1
## 04      : 1      00448508: 1      0400000US04: 1      04      : 1
## 05      : 1      00481813: 1      0400000US05: 1      05      : 1
## 06      : 1      00606926: 1      0400000US06: 1      06      : 1
## 08      : 1      00662849: 1      0400000US08: 1      08      : 1
## (Other):50      (Other) :50      (Other)   :50      (Other) :50
##      STUSPS      NAME      LSAD      ALAND
## AK      : 1      Alabama      : 1      00:56      102262419204: 1
## AL      : 1      Alaska      : 1      :50      102283343474: 1
## AR      : 1      American Samoa: 1      :50      105831263791: 1
## AS      : 1      Arizona      : 1      :50      106800130794: 1
## AZ      : 1      Arkansas      : 1      :50      111901043977: 1
## CA      : 1      California   : 1      :50      115884023072: 1
## (Other):50      (Other)      :50      (Other)   :50
##      AWATER
## 1026252314 : 1
## 10266413579 : 1
## 1027790845 : 1
## 1034369068 : 1
## 104031344385: 1
## 1076856589 : 1
```

¹² https://www.census.gov/geo/maps-data/data/cbf/cbf_state.html

```
## (Other) :50
```

The summary shows that this file already has a CRS defined, it came with the shapefile. The bounding box covers pretty much the whole world E-W (note the USA crosses the international date line in the Aleutian Islands, Alaska), and from about 14°S to about 71°N.

We see that the STUSPS field contains the State abbreviations, so we use this to select the States of interest, using the %in% set membership operator.

```
state.ne <- state[state$STUSPS %in% c('NY', 'NJ', 'PA', 'VT'),]
summary(state.ne)

## Object of class SpatialPolygonsDataFrame
## Coordinates:
##      min      max
## x -80.51989 -71.46455
## y  38.92852  45.01666
## Is projected: FALSE
## proj4string :
## [+proj=longlat +datum=NAD83 +no_defs +ellps=GRS80
## +towgs84=0,0,0]
## Data attributes:
##      STATEFP      STATENS      AFFGEOID      GEOID      STUSPS
## 34      :1      01779795:1      0400000US34:1      34      :1      NJ      :1
## 36      :1      01779796:1      0400000US36:1      36      :1      NY      :1
## 42      :1      01779798:1      0400000US42:1      42      :1      PA      :1
## 50      :1      01779802:1      0400000US50:1      50      :1      VT      :1
## 01      :0      00068085:0      0400000US01:0      01      :0      AK      :0
## 02      :0      00294478:0      0400000US02:0      02      :0      AL      :0
## (Other):0      (Other):0      (Other):0      (Other):0      (Other):0
##      NAME      LSAD      ALAND      AWATER
## New Jersey :1      00:4      115884023072:1      1034369068 :1
## New York   :1      :1      122054577774:1      19242052501:1
## Pennsylvania:1      :1      19048769408 :1      3396010541 :1
## Vermont    :1      :1      23871896411 :1      3542752545 :1
## Alabama    :0      :1      102262419204:0      1026252314 :0
## Alaska     :0      :1      102283343474:0      10266413579:0
## (Other)    :0      :0      (Other):0      (Other):0
```

The summary now shows that the bounding box has been restricted to surround the extreme coordinates of these four States.

We then transform to this project's CRS using `spTransform`. To ensure consistency, we obtain the CRS from the `SpatialPointsDataFrame` we projected in §2.4, above, using the `proj4string` function to extract it from the `ne.m` object:

```
proj4string(ne.m)

## [1] "+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76 +ellps=WGS84 +units=m"

state.ne.m <- spTransform(state.ne, proj4string(ne.m))
proj4string(state.ne.m)

## [1] "+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76 +ellps=WGS84 +units=m"

coordnames(state.ne.m) <- c("E", "N")
bbox(state.ne.m)

##      min      max
## E -387135.2 357707.2
## N -396311.1 288684.4
```

2.6 Saving the dataset

We did a lot of work in this section to put the data into a proper form for exploration and analysis, namely, to extract the target variable and limit the stations to a study area, also to project to a metric CRS and create both spatial and dataframe objects. We can save these objects in the internal RData format, using the `save` function. These can be read into a workspace with the `load` function.

TASK 12 : Save the spatial and dataframe versions of the original and projected dataset as R objects. Also save the CRS we developed for this region. •

```
save(ne, ne.m, ne.df, state.ne, state.ne.m, ne.crs, file="./ne_stations.RData")
```

3 Data exploration

TASK 13 : If you did not load packages in the previous section, load them now. •

```
library(sp)
library(rgdal)
```

TASK 14 : If you did not create the dataset for exploration in analysis in the previous section, load it now from the prepared file. •

The `load` function is used to load R objects from an RData format file. The `save` function stores the object name(s) along with the object(s), so that after `load`, there are new objects in the workspace, with these names. Setting the optional `verbose` argument to `load` shows the objects that were loaded from the RData format file.

```
load(file="./ne_stations.RData", verbose=TRUE)
```

```
## Loading objects:
##   ne
##   ne.m
##   ne.df
##   state.ne
##   state.ne.m
##   ne.crs
```

3.1 Feature-space summary

TASK 15 : Summarize the elevations and annual GDD50, so we know the range of values we will use in the analysis. •

```
summary(ne.df[,c("ANN_GDD50", "ELEVATION_")])
```

```
##   ANN_GDD50      ELEVATION_
## Min.   : 795    Min.     :  5.0
## 1st Qu.:2100    1st Qu.: 330.0
## Median :2463    Median : 711.0
```

```
## Mean      :2518      Mean      : 799.7
## 3rd Qu.:2930      3rd Qu.:1160.0
## Max.      :4021      Max.      :3950.0
```

The stations go from almost at sea level up to almost 4 000' (1 219 m.a.s.l.). The maximum annual GDD50 is more than five times the minimum. So in both we have a good range to find statistical relations.

3.2 Station locations

TASK 16 : Display the locations of the weather stations. •

We use base graphics for this. The coordinates of the `sp` object are a two-dimensional array, so this is a scatterplot, also called “x-y plot”, with the axes being the coordinates. The generic `plot` method, when given a two-dimensional matrix, calls the internal function `plot.xy`. Colours are from the default palette¹³, according to the level number of the factor. These are in alphabetic sort order, so NJ has colour 1 (black), NY has colour 2 (red), etc. We also use the `labels` argument to the `text` function to display the first four characters of the station name.

We then can add the state boundaries to help visualize the study area. In this case the `plot` method specializes to an internal function¹⁴ which extracts the coordinates along the boundary lines and plots these. Here the `plot` function has an extra argument, `add`, specifying that we are adding to a plot, not starting a new one. The `border` argument specifies the border colour, and the `lwd` “line width” argument specifies how much the default line width should be expanded or contracted. The result is Fig. 1.

```
plot(coordinates(ne.m), asp=1, xlab="E", ylab="N", pch=20)
text(coordinates(ne.m),
      labels=substr(ne.m$STATION_NA, 1, 4), cex=0.5,
      pos=2, col=as.numeric(ne.m$STATE))
## add state boundaries if available
## otherwise omit the next command
plot(state.ne.m, add=TRUE, border="darkgray", lwd=2)
grid()
```

3.3 Postplots

We now display the locations, but also using the point symbol represent the data value. This is called a **postplot**, because it “posts” (puts into geographic position) the data values.

TASK 17 : Display a postplot of the annual growing degree days above 50° F (attribute `ANN_GDD50`). •

For this we use the “Grammar of graphics” [32] as implemented in the

¹³ See the help for the `palette` function

¹⁴ `plot.SpatialPolygons`

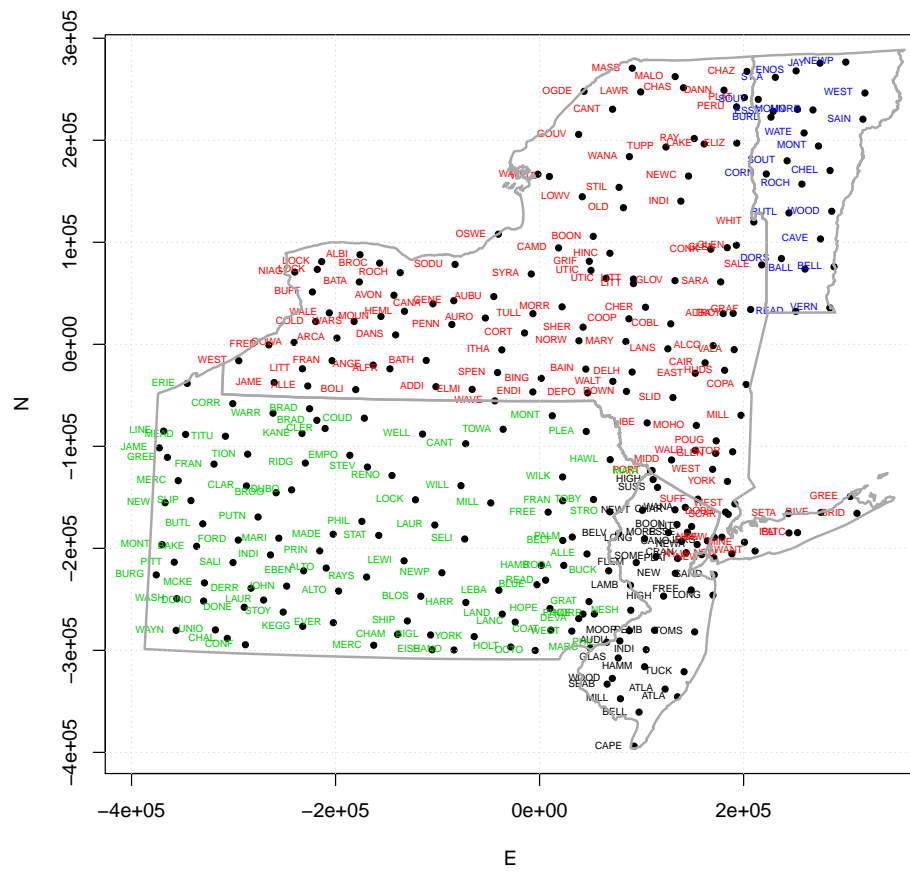


Figure 1: Location of weather stations

ggplot2 package [30, 31]. This is part of the so-called “tidyverse”¹⁵ set of packages from Hadley Wickham. The tidyverse web site has a comprehensive introduction to ggplot2¹⁶.

```
library(ggplot2)
```

The ggplot2 concept is that a graphic is *initialized* with ggplot and then elements are *added* to the graphic, each separated by a + operator, which in this context means “add to the graph”, not arithmetic addition.

In this example we first open the plot with the ggplot function, specifying the source of the data for the plot with the data argument, and then

1. specify how the graph should be set up on the page with the aes function: the x axis from the E coördinate and the y axis from the N coördinate;
2. add points with the the geom_point “geometry” function;
3. add axis labels with the xlab and ylab functions;
4. add a fixed-scale coördinate system for the graphic, with coord_fixed.

The points have an “aesthetic” (how they are displayed), specified with the aes function, and the name of the data frame where the names in the aesthetic can be found. We make the size of the points proportional by degree days, and the colour of the points by elevation, so we can visually assess if there is any relation both with coördinates and with elevation. We specify a printing character with the shape argument to geom_point.

```
ggplot(data=ne.df) +  
  aes(x=E, y=N) +  
  geom_point(aes(size=ANN_GDD50,  
                 colour=ELEVATION_),  
            shape=20) +  
  xlab("E") + ylab("N") + coord_fixed()
```

The result is Fig. 2.

Q1 : Does there appear to be any regional trend with North or East? with elevation? *Jump to A1 •*

3.4 * Viewing in geographic context

To check station locations, and to relate them to geographic features (land use, water bodies, cities ...) an easy method is to export the points as a KML file for display in Google Earth. Features in Google Earth must be in geographic coördinates (long/lat), on the WGS84 ellipsoid. As we saw above, these station records do have geographic coördinates, but on the NAD27 datum, which uses the Clarke 1866 ellipsoid. So we need to

¹⁵ <https://www.tidyverse.org>

¹⁶ <https://ggplot2.tidyverse.org>

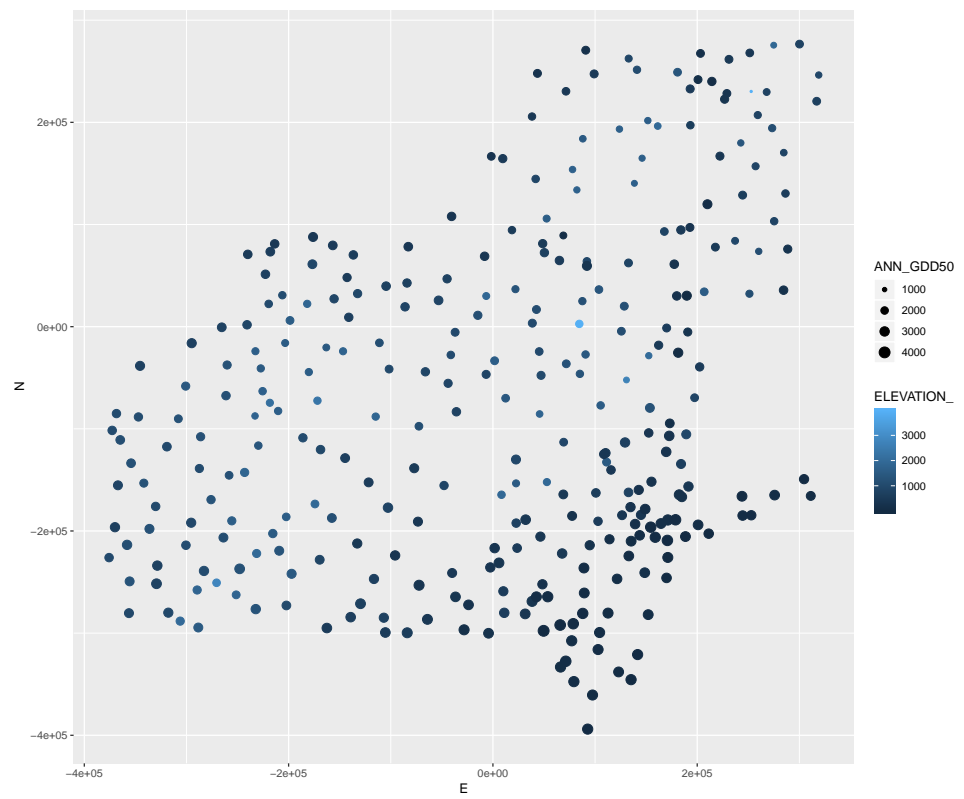


Figure 2: Annual Growing Degree days, 50F

4 Naïve linear modelling

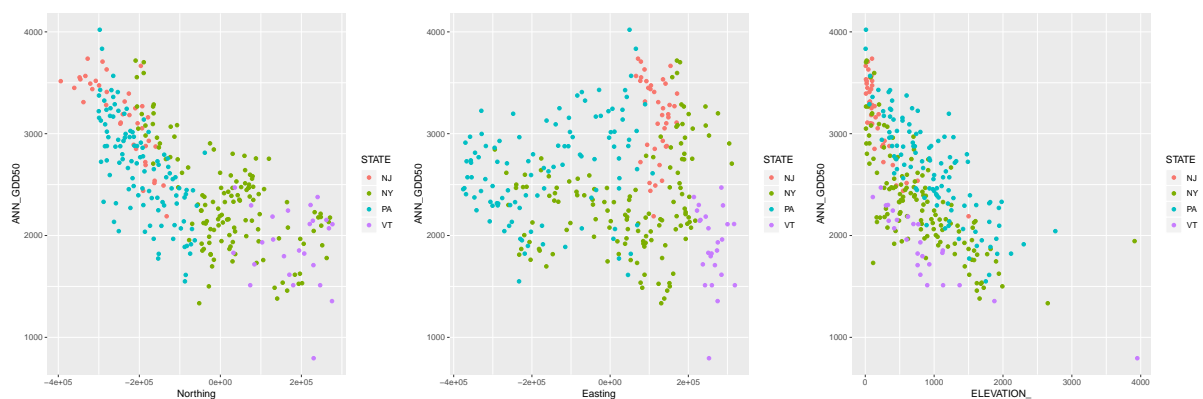
An obvious approach to prediction is to develop a model of GDD50 based on one or more of the available covariables: Easting, Northing coördinates and elevation. We know that in general higher elevations and more northerly latitudes are cooler; in this Atlantic region maybe the more easterly longitudes are warmer. We saw these relations spatially with the 2.5D postplots of the previous section, here we look at the relations in feature space for each possible covariable separately.

4.1 Exploring the relation between predictors and predictand

TASK 18 : Display the relation between ANN_GDD50 and covariates elevation and coordinates; colour the scatterplot points by State. •

```
p1 <- ggplot() +  
  geom_point(  
    aes(x=ELEVATION_, y=ANN_GDD50, colour=STATE),  
    data=ne.df  
  )  
p2 <- ggplot() +  
  geom_point(  
    aes(x=E, y=ANN_GDD50, colour=STATE),  
    data=ne.df  
  ) + xlab("Easting")  
p3 <- ggplot() +  
  geom_point(  
    aes(x=N, y=ANN_GDD50, colour=STATE),  
    data=ne.df  
  ) + xlab("Northing")
```

```
print(p3)  
print(p2)  
print(p1)
```



We now examine the scatterplots to see if these relations are useful predictors of GDD50, and if so, does the relation appear to be linear, or if a transformation to linearity is needed.

Q2 : Describe the relations between GDD50 and the three possible pre-

dictors.

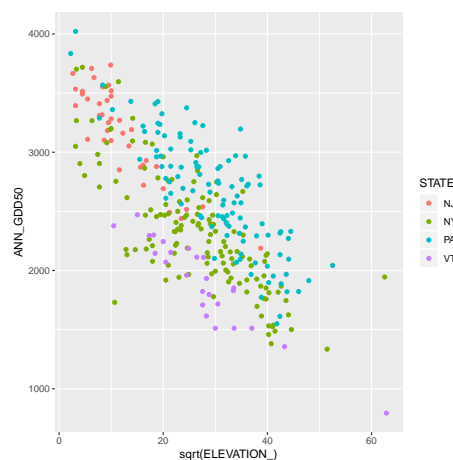
Jump to A2

-

The relation of GDD50 with elevation looks parabolic, so try a square-root transform.

TASK 19 : Re-display the relation between GDD50 and elevation, but with the elevation square-root transformed to correspond to the inverse parabolic shape perceived in the original graph. •

```
ggplot() +  
  geom_point(  
    aes(x=sqrt(ELEVATION_), y=ANN_GDD50, colour=STATE),  
    data=ne.df  
  )
```



A square-root transformation of elevation looks quite linear.

Note: This does not agree with theory, in that the adiabatic lapse rate (due to thinner air) of temperature with elevation is linear.

TASK 20 : Compute the linear correlations between each predictor and GDD50. •

The `cor` function computes bivariate correlations; the default method is Pearson product-moment (linear) correlation.

```
with(ne.df, cor(ANN_GDD50, N))  
## [1] -0.728741  
  
with(ne.df, cor(ANN_GDD50, E))  
## [1] -0.009425215  
  
with(ne.df, cor(ANN_GDD50, sqrt(ELEVATION_)))  
## [1] -0.7468479
```

Q3 : Which linear correlations are strongest? What does the sign (\pm) of

the correlation indicate?

[Jump to A3](#) •

4.2 OLS fit to the linear model

We first fit a linear model with the strongest single factor.

TASK 21 : Fit a linear model, using ordinary least squares (OLS), of annual GDD predicted by the square root of elevation. Display the model summary. •

The workhorse `lm` “linear modelling” function fits the model.

```
m.ols.elev <- lm(ANN_GDD50 ~ sqrt(ELEVATION_), data=ne.df)
summary(m.ols.elev)

##
## Call:
## lm(formula = ANN_GDD50 ~ sqrt(ELEVATION_), data = ne.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1346.24  -278.03    2.65   258.74  1012.87
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3472.287     53.520   64.88  <2e-16 ***
## sqrt(ELEVATION_) -37.000      1.893  -19.55  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 382.3 on 303 degrees of freedom
## Multiple R-squared:  0.5578, Adjusted R-squared:  0.5563
## F-statistic: 382.2 on 1 and 303 DF, p-value: < 2.2e-16
```

Q4 : How much of the variability of GDD50 over the four states is explained by this model? [Jump to A4](#) •

This is not so successful, so we try to use both predictors that showed a good correlation with GDD50.

TASK 22 : Fit a linear model, using ordinary least squares (OLS), of annual GDD predicted by the additive effects of Northing and square root of elevation. Display the model summary. •

Again we use the `lm` “linear modelling” function to fit the model:

```
m.ols <- lm(ANN_GDD50 ~ sqrt(ELEVATION_) + N, data=ne.df)
summary(m.ols)

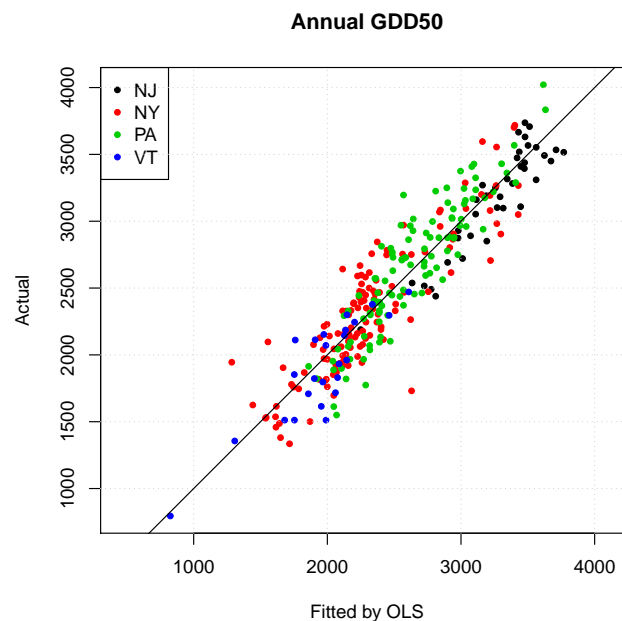
##
## Call:
## lm(formula = ANN_GDD50 ~ sqrt(ELEVATION_) + N, data = ne.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -899.15  -156.45   -9.78   153.12   660.08
##
## Coefficients:
```

```
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.119e+03  3.409e+01   91.49  <2e-16 ***
## sqrt(ELEVATION_) -2.924e+01  1.138e+00  -25.69  <2e-16 ***
## N              -1.981e-03  8.051e-05  -24.60  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 220.9 on 302 degrees of freedom
## Multiple R-squared:  0.8528, Adjusted R-squared:  0.8518
## F-statistic: 874.9 on 2 and 302 DF,  p-value: < 2.2e-16
```

Now the model explains much more (85.2%) of the variation in GDD50; we conclude that the North coördinate is an important predictor.

TASK 23 : Plot the actual vs. fits as a scatterplot, adding a 1:1 line. •

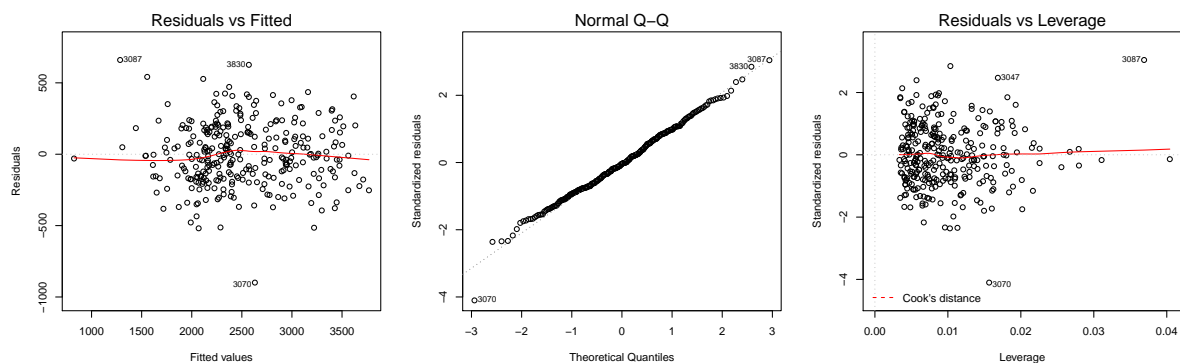
```
plot(ne.m$ANN_GDD50 ~ fitted(m.ols),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by OLS", ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)
```



Q5 : How well does the model fit the observations? Are any observations poorly-fit? Jump to A5 •

TASK 24 : Examine the feature-space model diagnostics. •


```
par(mfrow=c(1,3))
plot(m.ols, which=c(1,2,5))
par(mfrow=c(1,1))
```



Q6 : *Is there any relation between fits and residuals? Is the variability of the residuals approximately the same across the range of fits? Are the residuals normally-distributed? Are any high-leverage points inconsistent with the others?*

Jump to A6

We will look at the suspect observations in detail in the next §4.3.

We've now built a model, and we can use it to predict at any point for which we know the values of the predictors, i.e., the Northing and elevation. To make a predictive map of the whole study area we need to create a set of points that covers the area at some resolution, determine the predictor values at each point, and then predict at all those points using the model. These points are then considered to represent a small area, commonly called a grid cell. The predictions of each cell then form a complete surface covering the study area. For this we need a DEM of the study area to know the elevation at each prediction point; we will download and process one in §6, below.

4.3 Data cleaning

We now look at some poorly-fitted points, to determine if the poor fit is caused by an error in the database, or by an incorrect model.

TASK 25 : Find and display the points with the maximum and minimum model residuals.

The `which.min` and `which.max` functions find the index (position) in a vector of the minimum and maximum value, respectively.

```
ne.df$ols.resid <- residuals(m.ols)
(ix <- c(which.max(ne.df$ols.resid), which.min(ne.df$ols.resid)))
```

```
## [1] 113 96

ne.df[ix,]

##      STATION_ID STATE   STATION_NA LATITUDE_D LONGITUDE_ ELEVATION_
## 3087      305113   NY MARYLAND 6 SW      42.52      -74.97      3911
## 3070      303889   NY HINCKLEY 2 SW      43.30      -75.15      114
##      OID_ COOP_ID STATE_1   STN_NAME LAT_DD LONG_DD ELEV_FT
## 3087 3085  305113   NY MARYLAND 6 SW      43      -75      3911
## 3070 3068  303889   NY HINCKLEY 2 SW      43      -75      114
##      JAN_GDD50 FEB_GDD50 MAR_GDD50 APR_GDD50 MAY_GDD50 JUN_GDD50
## 3087          0          0          0          23          177          387
## 3070          0          0          0          22          155          349
##      JUL_GDD50 AUG_GDD50 SEP_GDD50 OCT_GDD50 NOV_GDD50 DEC_GDD50
## 3087          542          496          251          66          2          1
## 3070          505          443          210          46          1          0
##      ANN_GDD50 E N ols.resid
## 3087      1945 84568.04 2726.668 660.0801
## 3070      1731 68938.15 89269.546 -899.1468

ne.df[ix,c("LATITUDE_D", "LONGITUDE_")]

##      LATITUDE_D LONGITUDE_
## 3087      42.52      -74.97
## 3070      43.30      -75.15
```

The under-prediction (station 3087) is approximately six miles SW of Maryland, NY in Otsego County. This appears to be a mistake in the database; the elevation is listed as 3911' (feet), which is much higher than any terrain in the area.

Note: If you wish, find the approximate location on Google Earth by entering the geographic coördinates; note these are on NAD27 so will not match the WGS84 of Google Earth exactly. You can also find this on Google Maps and display the terrain.

Note that 3911' = 1192 m; the strong suspicion is that this height which should be in feet was assumed to be in meters and then multiplied by (3.28084 feet m⁻²), when in fact it was in feet. So, the incorrectly high elevation lead to an incorrectly low predicted value. After some research we find¹⁷ revised records for this station¹⁸:

NUM	DIV	ST	COUNTY	COOP	STATION	BEGINS	ENDS	LATITUDE	LONGITUDE	ELEV
						YEAR	YEAR	D M S	D M S	
305113	02	NY	OTSEGO	MARYLAND	6 SW	19831026	19881001	42 31 00	-074 58 00	1192
305113	02	NY	OTSEGO	MARYLAND	6 SW	19881001	20080702	42 31 00	-074 58 00	1192

Here the elevation is given as 1192', which matches well with the terrain; 6 miles SW from the village of Maryland is along Schenvus Creek , whose bluffs are at about 1200'. Further this matches the assumed mistake. So we feel justified in correcting this record accordingly.

The over-prediction is two miles SW of Hinckley, NY in Oneida County,

¹⁷ https://wrcc.dri.edu/Monitoring/Stations/station_inventory_show.php?snet=coop&sstate=NY

¹⁸ In fact this station was moved in July 2008, after our time period 1971-2000, so there is another record for it at a slightly different location

near Barneveld. Again this appears to be a mistake, the elevation listed as 114' but is about 1200'; so the model predicts as if it is at a much lower elevation, i.e., the prediction is too high. It seems here there was just a missing digit; the database record is:

```
303889 06 NY ONEIDA HINCKLEY 2 SW 19871014 19930301 43 18 00 -075 09 00 1141
```

We can correct these two points with their elevations from this (it appears correct) database. We make sure to make a record of these changes and our reasons for making them, in case other analysts want to check our work.

TASK 26 : Correct the elevation attribute of these two points in the dataset. •

We correct the ELEVATION_ field, which we are using for modelling. But we see that the dataset also has a ELEV_FT field, which duplicates this information – it is unclear why. At any rate, we do not want an inconsistent dataset, so we correct both elevation fields to the same value.

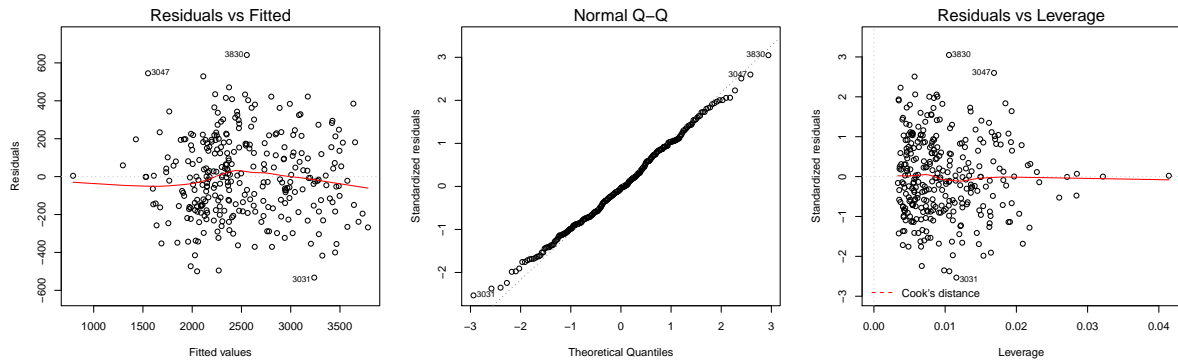
```
ne.m[ix[1], "ELEVATION_"] <- ne.m[ix[1], "ELEV_FT"] <- 1192
ne.df[ix[1], "ELEVATION_"] <- ne.df[ix[1], "ELEV_FT"] <- 1192
ne.m[ix[2], "ELEVATION_"] <- ne.m[ix[2], "ELEV_FT"] <- 1141
ne.df[ix[2], "ELEVATION_"] <- ne.df[ix[2], "ELEV_FT"] <- 1141
```

TASK 27 : Re-fit the linear model from the corrected database. •

```
m.ols <- lm(ANN_GDD50 ~ sqrt(ELEVATION_) + N, data=ne.df)
summary(m.ols)

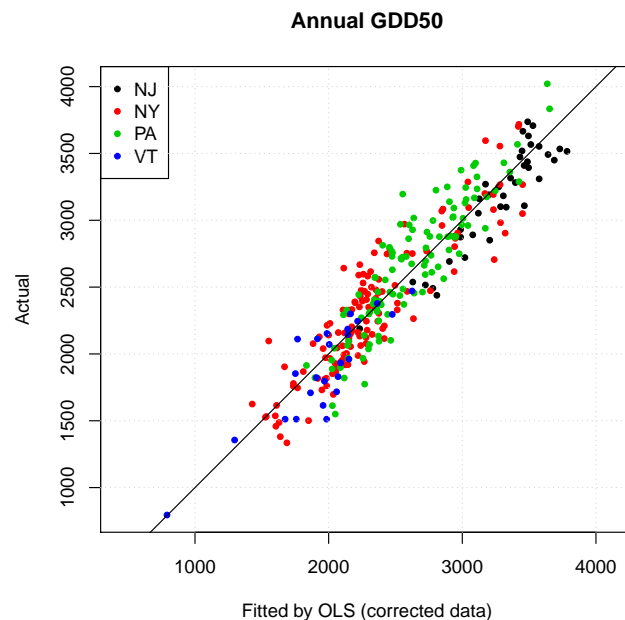
##
## Call:
## lm(formula = ANN_GDD50 ~ sqrt(ELEVATION_) + N, data = ne.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -532.76 -153.15   -7.84   155.44   641.76
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.151e+03  3.321e+01   94.87  <2e-16 ***
## sqrt(ELEVATION_) -3.040e+01  1.113e+00  -27.33  <2e-16 ***
## N              -1.952e-03  7.730e-05  -25.25  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 211.6 on 302 degrees of freedom
## Multiple R-squared:  0.865, Adjusted R-squared:  0.8641
## F-statistic: 967.3 on 2 and 302 DF,  p-value: < 2.2e-16
```

```
par(mfrow=c(1,3))
plot(m.ols, which=c(1,2,5))
par(mfrow=c(1,1))
```



These regression diagnostics are much better. The model coefficients have changed and the R^2 has increased about 1.5%, as a result of correcting just these two (out of a total 305) records.

```
plot(ne.m$ANN_GDD50 ~ fitted(m.ols),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by OLS (corrected data)",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)
```



TASK 28 : Summarize the OLS linear model residuals. •

```
summary(residuals(m.ols))
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	-532.764	-153.151	-7.838	0.000	155.435	641.758

To put these residuals in perspective, the difference between an early and medium maturity maize variety is about 300 GDD50. The inter-quartile range (IQR) is about ± 150 , which is about half of this on each side of perfect fit; i.e., the IQR is about the difference between maize maturities..

There are still some high residuals that could be checked for accuracy.

Challenge: Check the records for the largest residuals in the corrected model; if there is good evidence to do so, adjust the database accordingly and re-fit the model.

Challenge: Add Easting to the additive linear model and see if the model coefficient is significantly different from zero. You can try a pure trend surface (without elevation), a trend surface with elevation, and a trend surface with an elevation interaction – i.e., the effect of elevation is not the same across the region. Interpret the “best” model. Does this have a physical interpretation, in the same way we can relate GDD50 to Northing and elevation? Compare with ANOVA and/or AIC; confirm good linear model diagnostics. Identify the stations with the largest positive and negative residuals; try to explain why.

4.4 Spatial correlation of OLS model residuals

Are the model residuals spatially correlated? If so, that violates the assumption of independent residuals that is necessary for the OLS fit to be optimum.

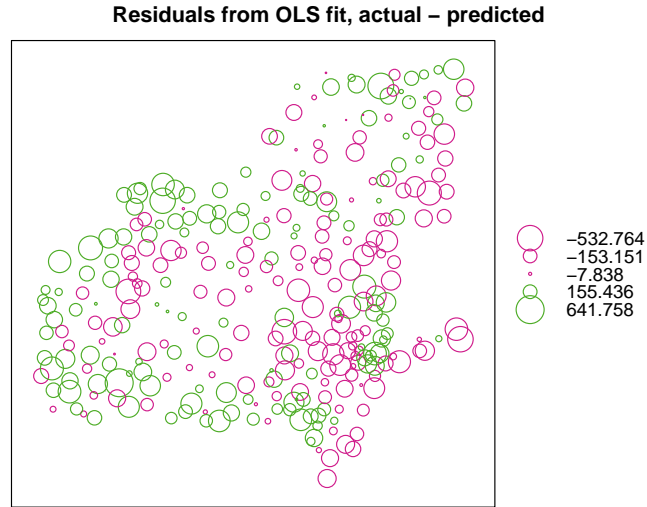
TASK 29: Add the linear model residuals to the `SpatialPointsDataFrame` object, so we can display them spatially. •

```
ne.m$ols.resid <- residuals(m.ols)
```

The so-called “bubble” plot is implemented by the `bubble` function of the `sp` package. It shows the absolute value of a variable by circle size (the “bubble”), and the sign (\pm) by colour.

TASK 30: Display a bubble plot of the model residuals. •

```
bubble(ne.m, zcol="ols.resid", pch=1,  
       main="Residuals from OLS fit, actual - predicted")
```



Q7 : *Is there evidence for spatial correlation of the OLS residuals? Describe some locations where this is most evident.* Jump to A7

•

We now quantify the spatial correlation with variogram analysis.

4.4.1 The empirical variogram

The empirical variogram shows the average **semivariance** as a function of **separation** between point pairs.

The semivariance between any two observations is defined as:

$$\gamma(\mathbf{s}_i, \mathbf{s}_j) = \frac{1}{2} [z(\mathbf{s}_i) - z(\mathbf{s}_j)]^2 \quad (3)$$

where \mathbf{s} is a geographic point and $z(\mathbf{s})$ is its **attribute value**, in this case the residual ANN_GDD50 from the OLS model. Each pair of points is separated by a vector \mathbf{h} , generally computed as the **Euclidean distance** between the points:

$$\mathbf{h} = ||\mathbf{x}_i, \mathbf{x}_j|| = \left(\sum_{k=1}^n (s_{i,k} - s_{j,k})^2 \right)^{\frac{1}{2}} \quad (4)$$

where n is the number of dimensions (in this example, 2). There are $(n \cdot (n - 1))/2$ point-pairs that can be compared this way; in our example this is $(305 \cdot 304)/2 = 46\,360$; clearly we need some way to summarize this.

The model of spatial dependence assumes **2nd-order stationarity**, i.e., the semivariance does not depend on the absolute location. Therefore an empirical variogram **averages** the individual semivariances in some **separation range** called a “bin”:

$$\gamma(\mathbf{h}) = \frac{1}{2N_{\mathbf{h}}} \sum_{i=1}^{N_{\mathbf{h}}} [z(\mathbf{s}_i) - z(\mathbf{s}_i + \mathbf{h})]^2 \quad (5)$$

where \mathbf{h} is a lag vector, i.e., a range of separations.

The analyst chooses the bin widths: wide enough to have enough point-pairs ($> \approx 150$) for reliable estimation, narrow enough to reveal the fine structure of spatial dependence.

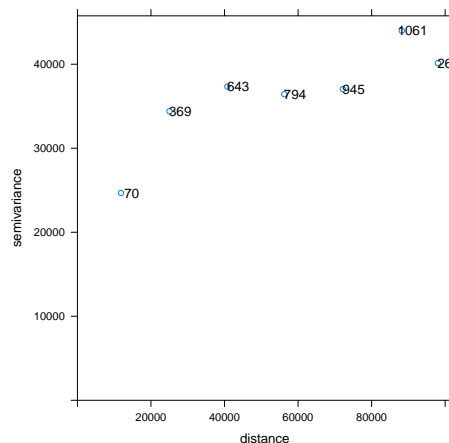
TASK 31 : Compute and display a variogram of the OLS residuals. Use a short cutoff, estimated from the bubble plot. •

The variogram can be computed with the `variogram` function of the `gstat` package. Load this package:

```
library(gstat)
```

Then compute and display the variogram. We use a cutoff of 100 km and a bin size of 16 km, to have enough points in the closest bin, and to avoid very local effects.

```
v.r.ols <- variogram(ols.resid ~ 1, locations=ne.m,
                     cutoff=100000, width=16000)
plot(v.r.ols, pl=T)
```



Note that the units of measure of the semivariances are $(\text{ANN_GDD50})^2$.

Q8 : What is the range of the spatial correlation? That is, the maximum separation distance at which there is lower semivariance than the maximum (total sill)? Jump to A8

•

4.4.2 Fitting an authorized variogram model

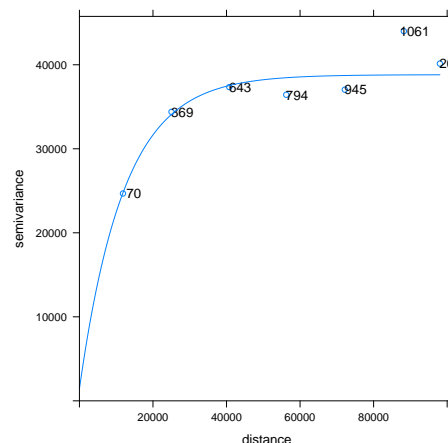
We require a continuous function of semivariance vs. separation, so that for any separation we can compute the semivariance to be used in the kriging equations. This must be an **authorized** variogram model that ensures that the kriging system will be positive semi-definite and thus invertible.

TASK 32 : Fit an exponential variogram model to this empirical variogram. •

The `fit.variogram` function of the `gstat` package uses a weighted least-squares (WLS) fit of a model form to the empirical variogram. The default weighting is proportional to the number of point-pairs per bin, and inverse-square proportional to the separation. So more weight is given to the short-range component of the model, and to bins with more reliable estimates of the semivariance.

An exponential variogram model is often used with residuals, as it is the simplest model form and corresponds to a diffusion process, which seems to fit our concept of air movement.

```
(vmf.r.ols <- fit.variogram(v.r.ols, vgm(15000, "Exp", 20000, 20000)))  
  
##      model      psill      range  
## 1  Nug 1435.842      0.00  
## 2  Exp 37382.603 12106.96  
  
plot(v.r.ols, pl=T, model=vmf.r.ols)
```



The effective range of the fitted model is 36.3 km, that is, pairs of points within this range partially duplicate each other's information.

This successful model fit shows that *there is spatial correlation among the residuals*, so the OLS fit of §4 is *not optimal*.

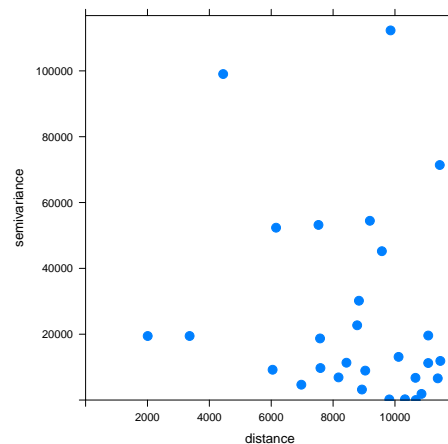
4.5 * Close-range anomalies

An interesting detour is to investigate anomalies which may show some features of micro-climate, using the **variogram cloud**. The results of this *optional* section are not used later in the exercise.

TASK 33 : Examine the short-range behaviour with the variogram cloud (all point-pairs): •

The `cloud` argument to the `variogram` function produces a variogram *cloud*, i.e., a variogram that shows each point-pair's semivariance vs. separation, rather than summarizing in (somewhat arbitrary) bins. This allows us to identify pairs of points with unusually high or low semivariations for their separation.

```
vc <- variogram(ols.resid ~ 1, locations=ne.m, cutoff=12000, cloud=T)
plot(vc, pch=20, cex=2)
```



We see some close point-pairs with fairly large semivariations (differences in their model residuals), especially the point-pair separated by about 4 000 m but with a semivariance of about 100 000 GDD50². Let's see which point-pairs are closest, and their individual semivariations.

The `as.data.frame` method applied to a variogram cloud creates two new fields, `left` and `right`, giving the identities of the two points in each pair.

Note: The order of the points is arbitrary, since the distance between them does not depend on which point is the origin and which the destination.

We then look for the highest semivariations at shortest distances (< 8 km).

We use the `order` function is used to sort the indices, here in order of increasing separation.

```
vc.df <- as.data.frame(vc)
vc.close <- subset(vc.df, vc.df$dist < 8000)
## sort by separation, look for anomalies.
```

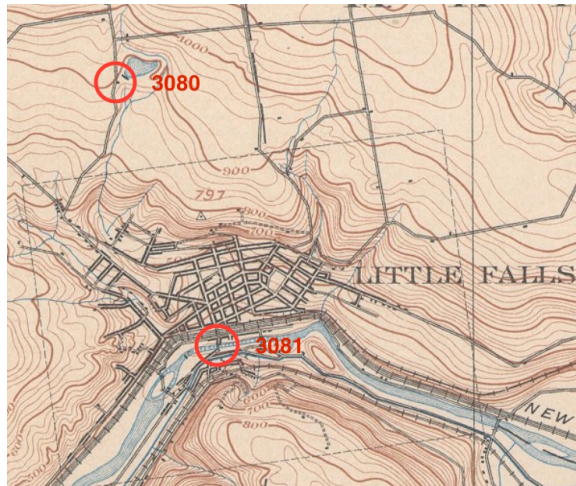


Figure 4: Little Falls NY weather station locations. Source: USGS 15 Minute Series, Little Falls NY quadrangle, 1900. Available from <http://nationalmap.gov/historical/>

```
vc.close[order(vc.close$dist),c("dist","gamma","left","right")]
```

```
##      dist      gamma left right
## 26 2007.455 19456.999 234  136
## 20 3362.049 19461.756 143   77
## 11 4446.408 99043.801 107  106
## 31 6044.983  9207.781 289  283
## 28 6156.757 52361.745 250  197
##  3 6973.391  4663.255  21   11
## 13 7527.628 53205.364 123   19
##  9 7578.549 18717.268  39   33
## 14 7590.756  9725.831 125  123
```

This list shows all the point-pairs separated by $< 8\text{km}$, see field `dist`. They have semivariances ranging from about 9000 to 55000 GDD50^2 , except for the pair of points 106 and 107, which have a very large semivariance, 9.9044×10^4 . This shows that the two have quite different residuals from the linear model fit. These two are separated by only about 4.5 km but are quite dissimilar. This is quite an anomaly, let's see the details for this two stations:

```
print(ne.m@data[c(106,107),c("STATE","STATION_NA","LATITUDE_D","LONGITUDE_",
                              "ELEVATION_", "ANN_GDD50","ols.resid")])
```

```
##      STATE      STATION_NA LATITUDE_D LONGITUDE_ ELEVATION_
## 3080    NY LITTLE FALLS CITY RSVR      43.07      -74.87      900
## 3081    NY LITTLE FALLS MILL ST      43.03      -74.87      360
##      ANN_GDD50 ols.resid
## 3080      1994 -119.8751
## 3081      2783  325.1952
```

These two stations are both near Little Falls (NY), one (3081) on Mill Street along the Mohawk River and one (3080) on Reservoir Road north of the village; see Figure 4.

Checking the topographic map, the elevations are correct; the discrepancy is because the Mill Street station is in a narrow river valley that

is much warmer than predicted by the model, hence the large positive residual (actual - predicted). The Reservoir Road station is somewhat over-predicted.

So, the data seems to be correct; the lesson is that close-by stations can have quite different micro-climates, and we have no factor in the model to account for this. Local interpolators such as kriging will also fail in this situation.

5 Generalized least squares with REML

To solve the problem of spatially-correlated OLS residuals, we turn to generalized least squares (GLS).

5.1 * GLS - theory

The key difference here is that in the linear model fit by OLS, the residuals ε are assumed to be *independently* and *identically* distributed with the same variance σ^2 :

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (6)$$

Whereas, now the residuals are themselves considered to be a random variable η that has a covariance structure:

$$\mathbf{y} = \mathbf{X}\beta + \eta, \quad \eta \sim \mathcal{N}(0, \mathbf{V}) \quad (7)$$

where \mathbf{V} is a positive-definite variance-covariance matrix of the model residuals. The covariances in this matrix (off-diagonals) are typically based on the distance between observations, using some model of spatial correlation.

Lark and Cullis [17, Appendix] point out that the error vectors can now not be assumed to be spherically distributed in feature space around the 0 expected value, but rather that error vectors in some directions are longer than in others. So, the measure of distance (the vector norm) is now a so-called “generalized” distance¹⁹, taking into account the covariance between error vectors:

$$S = (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) \quad (8)$$

The OLS equivalent is simpler:

$$S = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (9)$$

Comparing these equations, we see that the GLS formulation of Equation 8 includes the variance-covariance matrix of the residuals $\mathbf{V} = \sigma^2 \mathbf{C}$, where σ^2 is the variance of the residuals and \mathbf{C} is the correlation matrix. This reduces to the OLS formulation of Equation 9 if there is no covariance, i.e., $\mathbf{V} = \mathbf{I}$.

¹⁹ This is closely related to the Mahalanobis distance

Expanding Equation 8, taking the partial derivative with respect to the parameters, setting equal to zero and solving we obtain:

$$\begin{aligned}\frac{\partial}{\partial \beta} S &= -2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + 2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \beta \\ 0 &= -\mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \beta \\ \hat{\beta}_{\text{GLS}} &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}\end{aligned}\quad (10)$$

This reduces to the OLS estimate $\hat{\beta}_{\text{OLS}}$ if there is no covariance, i.e., $\mathbf{V} = \mathbf{I}$.

In the case of spatial correlation, we ensure positive-definiteness (i.e., always a real-valued solution) by using an **authorized covariance function** C and assuming that the entries are completely determined by the vector **distance** between points $\mathbf{x}_i - \mathbf{x}_j$:

$$\mathbf{C}_{i,j} = C(\mathbf{x}_i - \mathbf{x}_j) \quad (11)$$

In this formulation C has a three-parameter vector θ , as does the corresponding variogram model: the range a , the total sill σ^2 , and the proportion of total sill due to pure error, not spatial correlation s ²⁰.

In modelling terminology, the coefficients β are called **fixed** effects, because their effect on the response variable is fixed once the parameters are known. By contrast the covariance parameters η are called **random** effects, because their effect on the response variable is stochastic, depending on a random variable with these parameters.

Models with the form of Equation 7 are called **mixed** models: some effects are fixed (here, the relation between elevation or Northing and the GDD50) and others are random (here, the error variances) but follow a known structure; these models have many applications and are extensively discussed in Pinheiro and Bates [22]. Here the random effect η represents both the spatial structure of the residuals from the fixed-effects model, and the unexplainable (short-range) noise. This latter corresponds to the noise σ^2 of the linear model of Equation 6.

To solve Equation 10 we first need to compute \mathbf{V} , i.e., estimate the variance parameters $\theta = [\sigma^2, s, a]$, use these to compute \mathbf{C} with equation 11 and from this \mathbf{V} , after which we can use equation 10 to estimate the fixed effects β . But θ is estimated from the residuals of the fixed-effects regression, which has not yet been computed. How can this “chicken-and-egg”²¹ computation be solved?

The answer is to use **residual** (sometimes called “restricted”) **maximum likelihood** (REML) to maximize the likelihood of the random effects θ independently of the fixed effects β .

Here we fit the fixed effects (regression coefficients) at the same time as we estimate the spatial correlation.

²⁰ In variogram terms, this is the nugget variance c_0 as a proportion of the total sill $(c_0 + c_1)$.

²¹ from the question “which came first, the chicken or the egg?”

Lark and Cullis [17, Eq. 12] show that the likelihood of the parameters in Equation 6 can be expanded to include the spatial dependence implicit in the variance-covariance matrix \mathbf{V} , rather than a single residual variance σ^2 . The log-likelihood is then:

$$\ell(\beta, \theta | \mathbf{y}) = c - \frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) \quad (12)$$

where c is a constant (and so does not vary with the parameters) and \mathbf{V} is built from the variance parameters θ and the distances between the observations. By assuming **second-order stationarity**²², the structure can be summarized by the covariance parameters $\theta = [\sigma^2, s, a]$, i.e., the total sill, nugget proportion, and range.

However, maximizing this likelihood for the random-effects covariance parameters θ also requires maximizing in terms of the fixed-effects regression parameters β , which in this context are called *nuisance parameters* since at this point we don't care about their values; we will compute them after determining the covariance structure.

Both the covariance and the nuisance parameters β must be estimated, it seems at the same time ("chicken and egg" problem) but in fact the technique of REML can be used to first estimate θ without having to know the nuisance parameters. Then we can use these to compute \mathbf{C} with equation 11 and from this \mathbf{V} , after which we can use equation 10 to estimate the fixed effects β .

The maximum likelihood estimate of θ is thus called "restricted", because it only estimates the covariance parameters (random effects). Conceptually, REML estimation of the covariance parameters θ is ML estimation of both these and the nuisance parameters β , with the later integrated out [22, §2.2.5]:

$$\ell(\theta | \mathbf{y}) = \int \ell(\beta, \theta | \mathbf{y}) d\beta \quad (13)$$

Pinheiro and Bates [22, §2.2.5] show how this is achieved, given a likelihood function, by a change of variable to a statistic sufficient for β .

5.2 GLS – practice

The computations are performed with the `gls` function of the `nlme` 'Non-linear mixed effects models' package [1].

TASK 34: Set up and solve a GLS model, using the covariance structure estimated from the variogram of the OLS residuals. •

The linear model formulation is the same as for `lm`. However:

- It has an additional argument `correlation`, which specifies the correlation structure.

²² that is, the covariance structure is the same over the entire field, and only depends on the distance between pairs of points

- This is built with various correlation models; we use `corExp` for exponential spatial correlation, which is what we fit for the OLS residuals,
 - The `form` names the dimensions, here 2D with the Easting and Northing.
 - We initialize the search for the correlation structure parameters with the `value` argument, a list of the initial values. Here we only specify the range.
 - Our fitted variogram model for the OLS residuals showed a zero nugget, so here the function should not fit a nugget. So we set the `nugget` argument to `FALSE`.²³

Note: For a list of the predefined model forms see `?corClasses`. Users can also define their own `corStruct` classes.

```
library(nlme)
vmf.r.ols[1:2,]

##      model      psill      range
## 1   Nug 1435.842      0.00
## 2   Exp 37382.603 12106.96

m.gls <- gls(model=ANN_GDD50 ~ sqrt(ELEVATION_) + N,
             data=ne.df,
             correlation=corExp(
               value=c(vmf.r.ols[2,"range"]),
               form=~E + N,
               nugget=FALSE))
```

The `gls` function is not guaranteed to find a valid correlation structure. First, there may be no spatial correlation of the residuals. Second, we may have specified an inappropriate model form. Third, if the starting values are not close to good fits, the optimization method may not find the correct fit. Therefore it is crucial to check the results of the model fitting to see if they are reasonable.

TASK 35 : Display the model summary. •

This shows both the linear model coefficients and the spatial correlation structure.

```
summary(m.gls)

## Generalized least squares fit by REML
## Model: ANN_GDD50 ~ sqrt(ELEVATION_) + N
## Data: ne.df
##      AIC      BIC    logLik
## 4120.113 4138.665 -2055.056
##
## Correlation Structure: Exponential spatial correlation
## Formula: ~E + N
## Parameter estimate(s):
##      range
## 17456.99
##
```

²³ The nugget could also be fixed at a user-specified value.

```
## Coefficients:
##               Value Std.Error   t-value p-value
## (Intercept)    3136.3707  44.03754   71.22039     0
## sqrt(ELEVATION_) -30.0448   1.41085  -21.29546     0
## N              -0.0019   0.00011  -16.64114     0
##
## Correlation:
##               (Intr) s(ELEV
## sqrt(ELEVATION_) -0.888
## N              0.319 -0.183
##
## Standardized residuals:
##               Min      Q1      Med      Q3      Max
## -2.36510335 -0.68842210 -0.01831242  0.74010041  3.00789761
##
## Residual standard error: 217.344
## Degrees of freedom: 305 total; 302 residual
```

TASK 36 : Display the confidence intervals for the coefficients. •

The `intervals` function of the `nlme` gives approximate confidence intervals of the GLS fit.

```
intervals(m.gls, level=0.95)$coef
##               lower      est.      upper
## (Intercept)    3.049711e+03  3.136371e+03  3.223030e+03
## sqrt(ELEVATION_) -3.282113e+01 -3.004478e+01 -2.726843e+01
## N              -2.135757e-03 -1.909907e-03 -1.684056e-03
## attr("label")
## [1] "Coefficients:"
```

These intervals seem fairly wide, indicating that the model is perhaps not sufficiently specified to capture all the reasons for variation in GDD50 over this area..

TASK 37 : Display the correlation structure fit by `gls`. •

```
intervals(m.gls, level=0.95)$corStruct
##               lower      est.      upper
## range 12850.42 17456.99 23714.92
## attr("label")
## [1] "Correlation structure:"
```

Q9 : *What is the 95% confidence interval of the range parameter? Does this seem narrow or wide?* Jump to A9 •

TASK 38 : Compare with the correlation structure estimated from the OLS residuals. •

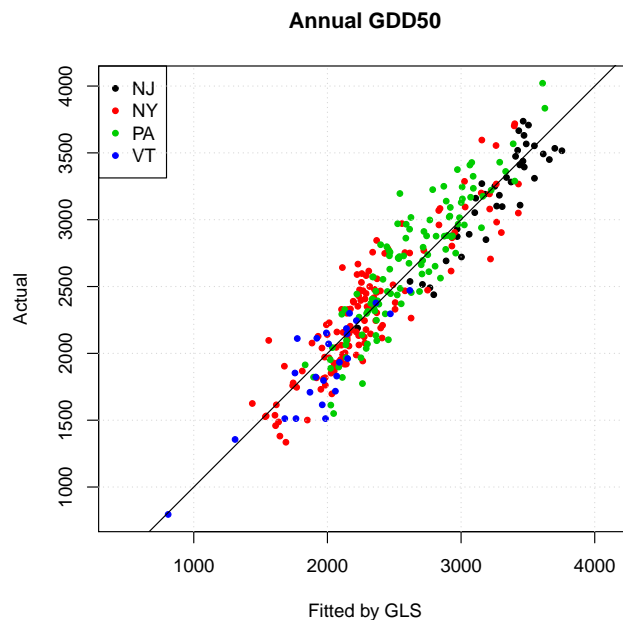
```
intervals(m.gls, level=0.95)$corStruct
##               lower      est.      upper
## range 12850.42 17456.99 23714.92
## attr("label")
## [1] "Correlation structure:"
```

```
vmf.r.ols[2,"range"]
## [1] 12106.96
```

Q10 : How closely does the correlation structure fitted by GLS match that estimated from the variogram of the OLS residuals? [Jump to A10](#) •

TASK 39 : Plot the actual vs. model fits on a 1:1 scatterplot. •

```
plot(ne.m$ANN_GDD50 ~ predict(m.gls),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by GLS",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)
```



The fit clusters well around the 1:1 line (good accuracy) but is diffuse (low precision).

5.3 Spatial correlation of GLS model residuals

As with the OLS model (§4.4), the GLS residuals may show spatial correlation. We examine this with an empirical variogram and then fit a variogram model.

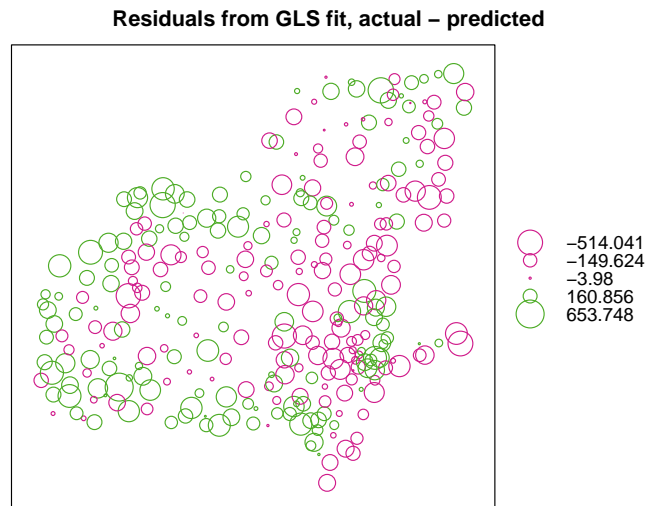
TASK 40 : Display a bubble plot of the model residuals. •

```
ne.m$gl.resid <- residuals(m.gls)
summary(ne.m$gl.resid)
```



```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -514.041 -149.624   -3.980    8.496  160.856  653.749

bubble(ne.m, zcol="gls.resid", pch=1,
       main="Residuals from GLS fit, actual - predicted")
```



Q11 : Does GLS remove the spatial correlation? Describe the spatial correlation structure of the GLS residuals. Jump to A11 •

These residuals should show the spatial correlation discovered in the REML fit.

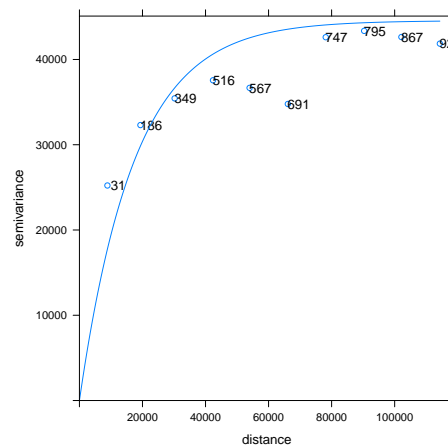
TASK 41 : Display a variogram of the GLS residuals, with the variogram model estimated as part of the GLS fit. •

This fit gives the range; we estimate the total sill from the overall variance. Knowing these variogram parameters, we construct a variogram model with the `vgm` function: We estimate the total sill as the variance of the residuals. Since there is no nugget variance in this model, we do not have to adjust it downward for the proportional nugget.

```
v.r.gls <- variogram(gls.resid ~ 1, ne.m, cutoff=12000, width=12000)
(vmf.r.gls <- vgm(psil1=var(ne.m$gls.resid),
                 model="Exp",
                 range=intervals(m.gls)$corStruct["range", "est."],
                 nugget=TRUE))

##      model      psil1      range
## 1      Nug       1.00       0.00
## 2      Exp 44558.48 17456.99

plot(v.r.gls, pl=T, model=vmf.r.gls)
```



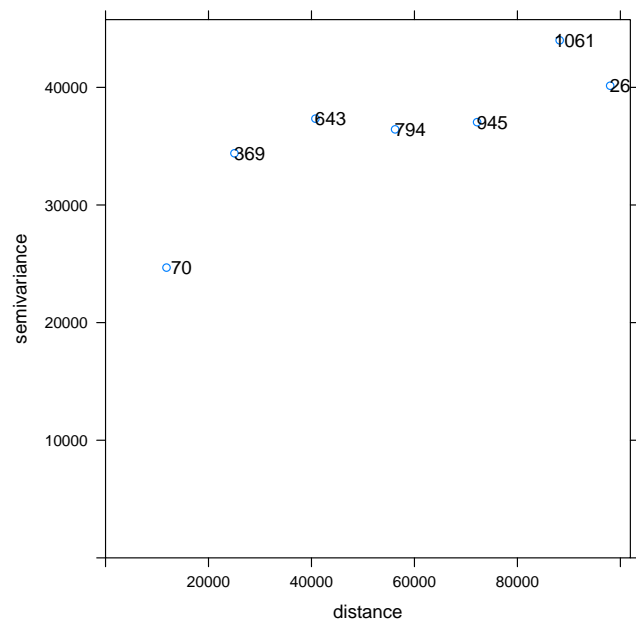
Q12 : *How well does this model fit the empirical variogram of the GLS residuals?* Jump to A12 •

5.4 * Fitting a GLS correlation structure with a nugget variance

In the previous § we built a model of spatial dependence of the residuals using the `corExp` model for exponential spatial correlation. Because our analysis of the OLS residual variogram in §4.4, we did not fit a nugget, by specifying the `nugget` argument to `FALSE`. In other analyses there may well be a nugget effect to be fit, so here we show how to do that, and the effect it has on the correlation structure fitted by `gls`.

In this example an exponential model has almost no nugget, so we select a spherical model that is nearly linear near the origin. We estimate the starting parameters from the empirical variogram.

```
plot(v.r.ols, pl=T)
```

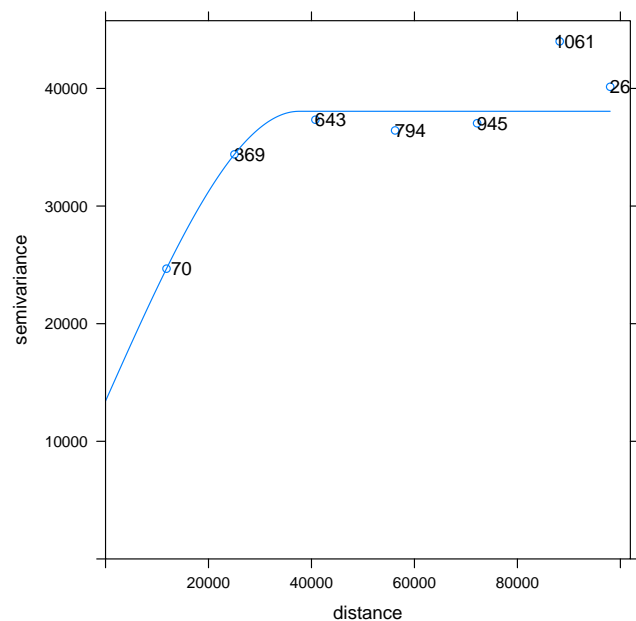


The total sill appears to be ≈ 35000 $(ANN_GDD50)^2$, the nugget ≈ 15000 $(ANN_GDD50)^2$, so, the partial sill is ≈ 20000 $(ANN_GDD50)^2$; the range (where the spherical model reaches the sill) ≈ 40000 m.

```
(vmf.r.ols.sph <- fit.variogram(v.r.ols,
                               vgm(psill=20000, model="Sph",
                                   range=40000, nugget=15000)))

##   model    psill   range
## 1  Nug 13426.14    0.0
## 2  Sph 24629.99 37585.4

plot(v.r.ols, pl=T, model=vmf.r.ols.sph)
```



This does not fit too badly, and clearly has a substantial nugget. Our initial estimates were close to the fitted values. This must be converted to a *proportional* nugget, i.e., the proportion of the total sill represented by the nugget.

```
(prop.nugget <- vmf.r.ols.sph[1,"psill"]/sum(vmf.r.ols.sph[, "psill"]))
## [1] 0.3527983
```

The nugget is about 35% of the total sill, with this model. We now use this in `gls`, substituting `corSpher` to specify a spherical model of spatial dependence.

To include a nugget variance, we must:

1. specify nugget argument as TRUE;
2. expand the `value` argument to a two-element list: (1) the starting value of the range, as before; and (2) the starting value of the nugget.

```
m.gls.2 <- gls(model=ANN_GDD50 ~ sqrt(ELEVATION_) + N,
  data=ne.df,
  correlation=corSpher(
    value=c(vmf.r.ols.sph[2,"range"], prop.nugget),
    form=~E + N,
    nugget=TRUE))
```

We compare the fitted correlation structures:

```
intervals(m.gls, level=0.95)$corStruct
##           lower      est.      upper
## range 12850.42 17456.99 23714.92
## attr("label")
## [1] "Correlation structure:"

intervals(m.gls.2, level=0.95)$corStruct
##           lower      est.      upper
## range 2.275362e+05 3.014994e+05 3.997193e+05
## nugget 2.661125e-01 4.622866e-01 6.708771e-01
## attr("label")
## [1] "Correlation structure:"
```

The nugget proportion has been fit as 0.462, somewhat higher than our original estimate 0.353.

Recall that the range parameter of the exponential model is 1/3 of the effective range, so to compare them:

```
intervals(m.gls)$corStruct["range", 2]*3
## [1] 52370.98

intervals(m.gls.2)$corStruct["range", 2]
## [1] 301499.4
```

The effective range is much longer, about 300 km for the spherical model with nugget, compared to about 50 km for the exponential model with no nugget. This seems unrealistic when we compare it with the residual variogram. The spherical model is thus not appropriate; we showed it

here in order to explain how to fit a proportional nugget, if the empirical residual variogram suggests that there is a nugget variance.

We can also compare the regression coefficients:

```
intervals(m.gls, level=0.95)$coef[,2]

##      (Intercept) sqrt(ELEVATION_)      N
## 3.136371e+03    -3.004478e+01    -1.909907e-03

intervals(m.gls.2, level=0.95)$coef[,2]

##      (Intercept) sqrt(ELEVATION_)      N
## 3.209772e+03    -3.241043e+01    -1.658746e-03
```

The changed spatial correlation structure has considerably changed the regression coefficients.

5.5 Comparing OLS and GLS models

How much has accounting for the spatial correlation of the model residuals affected the linear models?

TASK 42 : Compare the coefficients of the GLS and OLS models. Compute the relative change. •

```
coefficients(m.gls)

##      (Intercept) sqrt(ELEVATION_)      N
## 3.136371e+03    -3.004478e+01    -1.909907e-03

coefficients(m.ols)

##      (Intercept) sqrt(ELEVATION_)      N
## 3.150881e+03    -3.040452e+01    -1.952152e-03

round(100*(coefficients(m.gls) - coefficients(m.ols))/coefficients(m.ols),2)

##      (Intercept) sqrt(ELEVATION_)      N
##      -0.46      -1.18      -2.16
```

Q13 : *How much have the linear model coefficients changed from the OLS to the GLS fit? What explains the change in coefficients?* [Jump to A13](#) •

We can see this spatially by comparing the residuals.

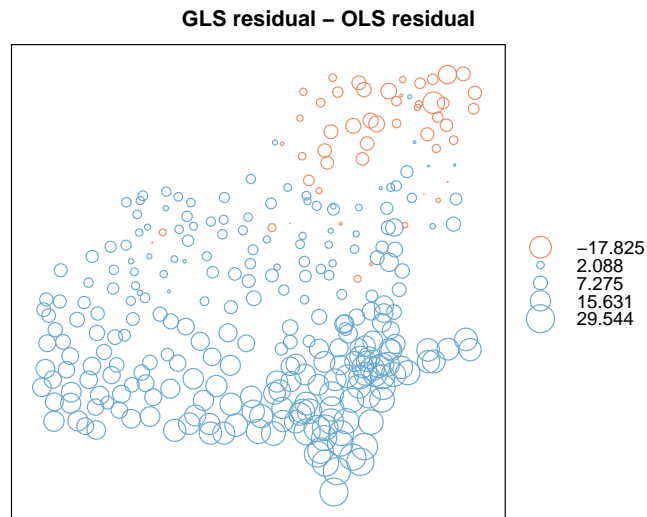
TASK 43 : Compute the difference between the GLS and OLS residuals, add them to the spatial points, and display as a bubble plot. •

We use a different colour scheme to emphasize that this is the *difference* between residuals, *not* the residuals themselves.

```
ne.m$diff.gls.ols.resid <- (ne.m$gl.s.resid - ne.m$ols.resid)
summary(ne.m$diff.gls.ols.resid)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -17.825   2.088   7.275   8.496  15.631  29.544
```

```
bubble(ne.m, zcol="diff.gls.ols.resid", pch=1,
       main="GLS residual - OLS residual",
       col=c("#EF8A62", "#67A9CF"))
```



Q14 : *Where are the largest differences between the OLS and GLS residuals?*

[Jump to A14](#)

•

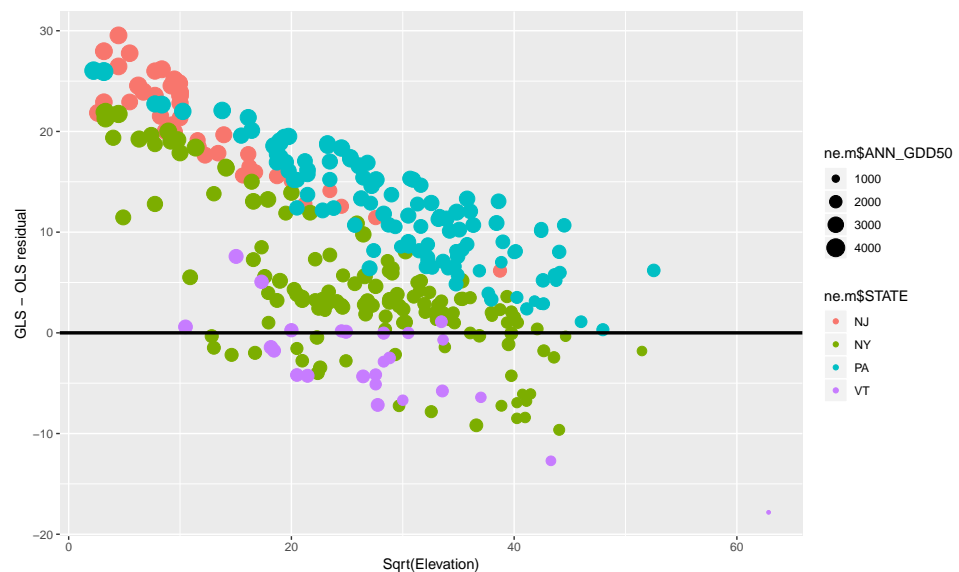
The coefficient for elevation was reduced by a smaller amount, and it is for the square root. To visualize this effect, we can use a scatterplot of the change in residuals vs. this marginal predictor.

TASK 44 : Display a scatterplot of change in residuals vs. the square root of elevation.

•

We can make the plot more informative by colouring the points by state and making their size proportional to the target variable

```
ggplot() +
  geom_point(aes(x=sqrt(ne.m$ELEVATION_), y=ne.m$diff.gls.ols.resid,
                size=ne.m$ANN_GDD50, colour=ne.m$STATE)) +
  xlab("Sqrt(Elevation)") +
  ylab("GLS - OLS residual") +
  geom_hline(yintercept=0, size=1.2)
```



This shows clearly that GLS residuals are larger at the lower elevations, and that the largest adjustments tend to be at the largest GDD50. Note that there would be a confounding effect if the two predictors (elevation and Northing) were not almost independent, as they are in this case.

So in combination, a southerly, low-lying station will have the largest positive GLS-OLS residuals, i.e., GLS predicts higher than OLS; a northerly, high-elevation station the largest negative residuals.

TASK 45 : Display the station name, elevation, and coördinates of the most positive and negative residuals. •

```
ix <- which.max(ne.m$diff.gls.ols.resid)
ne.df[ix,c("STATION_NA","STATE","ELEVATION_","N","E")]

##          STATION_NA STATE ELEVATION_      N      E
## 2859 CAPE MAY 2 NW    NJ          20 -393922.9 92756.51

ix <- which.min(ne.m$diff.gls.ols.resid)
ne.df[ix,c("STATION_NA","STATE","ELEVATION_","N","E")]

##          STATION_NA STATE ELEVATION_      N      E
## 4716 MOUNT MANSFIELD  VT          3950 230248.4 252835.5
```

The largest negative residual (GLS reduced the prediction the most) is for Mt. Mansfield (VT), the highest elevation station in the dataset, and near the N limit. The largest positive residual (GLS increased the prediction the most) is for Cape May (NJ), just above sea level and the southernmost station.

Conclusion: accounting for spatial correlation in the residuals significantly changed the linear model, resulting in differences up to 30 GDD50.

6 Creating a regional grid for mapping

A natural question is how target variable varies over the entire study area, not just at the observation points. For regional studies we want to predict and visualize over the entire area; i.e., we want to produce a **map** of the target variable. For this we need the predictors used in the models at a set of **grid cells** covering the **whole study area**. In this case these predictors are the metric coördinates and elevation. We can create an empty grid with coördinates with the `expand.grid` or `spsample` commands (see §13) but we also need the elevation.

To save you time, I've prepared the initial grid of elevations from the SRTM; if you want to use that skip to §6.2.

If you want to see how the grid was prepared, or prepare one for another study area, continue with the next subsection.

6.1 * Creating a grid from global SRTM

The Shuttle Radar Topography mission²⁴ has mapped most of the Earth at a nominal 3 arc second (≈ 90 m) resolution. This data has been greatly enhanced (e.g., gap filling, data cleaning) by the CGIAR Consortium for Spatial Information (CGIAR-CSI)²⁵; this is described on their SRTM products web page²⁶. The latest version is 4.1²⁷. We will use this to get the elevation over a grid covering the bounding box of the four selected States.

For regional applications the original resolution is too fine, so the dataset has been resampled to lower nominal resolutions (250 m, 500 m, and 1 km), conveniently available on Google Drive²⁸ (unless you are accessing the internet from a country where Google is blocked).²⁹ Two formats are provided: ESRI grid (.GRD, 152 Mb) and ASCII (.ASC, 185 Mb); the latter is easier to read into R. These are both provided as RAR archives, and must be decompressed before further use.

Although you can work with the original archives, I have processed the ASCII format file to restrict the area to just the lower 48 USA states. This is provided as a compressed file `srtm_1km_48.zip` (22.3 Mb), which contains both the .ASC file with the raster data (uncompressed 126.5 Mb) and a small .PRJ file with the coördinate reference system (CRS).

We process these with the `raster` package. A major advantage of `raster` is that it only reads data into memory when needed, so it can deal with large raster datasets.

²⁴ <http://www2.jpl.nasa.gov/srtm/>

²⁵ <http://www.cgiar-csi.org>

²⁶ <http://srtm.csi.cgiar.org>

²⁷ <http://www.cgiar-csi.org/data/srtm-90m-digital-elevation-database-v4-1>

²⁸ <https://goo.gl/6sFVwC>

²⁹ If you want to get the data yourself, it's easiest to copy to your own Google Drive and then sync to your computer's local drive, this is much faster than downloading.

TASK 46 : Load the nominal 1 km resolution DEM as a RasterLayer object. Display its R class with the `class` function. Examine its characteristics with the `projection`, `extent`, `res`, `nrow`, and `ncol` functions of the `raster` package. •

The `raster` function of the `raster` package can read many common formats directly, including ASCII rasters. Assuming you've put the downloaded `.asc` file into the same directory as the NE weather stations and unpacked to that directory, the ASCII grid is read in using `raster`:

```
library(raster)
dem.1km <- raster("./srtm_1km_48.asc")

class(dem.1km)

## [1] "RasterLayer"
## attr(,"package")
## [1] "raster"

extent(dem.1km)

## class      : Extent
## xmin       : -124.7667
## xmax       : -66.95
## ymin       : 24.525
## ymax       : 49.38333

projection(dem.1km)

## [1] "+proj=longlat +ellps=WGS84 +no_defs"

proj4string(dem.1km)

## [1] "+proj=longlat +ellps=WGS84 +no_defs"

res(dem.1km)

## [1] 0.008333333 0.008333333

nrow(dem.1km)

## [1] 2983

ncol(dem.1km)

## [1] 6938

ncell(dem.1km)

## [1] 20696054
```

We see that it covers the lower 48 States in a geographic projection on the WGS84 datum, with a resolution of 0.0083333° . These are not square grid cells, because degrees of longitude become smaller as the cosine of the latitude. For example, at 42°N one degree of longitude is only 0.743 the length of one degrees of latitude. One degree of latitude is about $10000/(90 * 60) \approx 111.11 \text{ km}$ ³⁰. So the actual resolution 0.0083333° corresponds to a ground resolution of about 0.926 km N-S and 0.688 km E-W at 42°N :

³⁰ The original definition of the meter was 1/10 000 of the distance from pole to equator.

```
round(res(dem.1km)[1]*(10000/90),3)

## [1] 0.926

round(res(dem.1km)[1]*(10000/90)*cos(42*pi/180),3)

## [1] 0.688
```

TASK 47 : Restrict the DEM to the four-state study area. •

The `crop` function of the `raster` package crops a raster to another extent; the two rasters must have the same CRS. The easiest way is if there is already another spatial object with the desired extent. In this case we have one, the four states, in the original geographic projection.

```
bbox(state.ne)

##           min           max
## x -80.51989 -71.46455
## y  38.92852  45.01666

dem.ne <- crop(dem.1km, state.ne)
bbox(dem.ne)

##           min           max
## s1 -80.51667 -71.46667
## s2  38.92500  45.01667

nrow(dem.ne)

## [1] 731

ncol(dem.ne)

## [1] 1086

ncell(dem.ne)

## [1] 793866
```

This is considerably smaller than the 48-State raster.

To apply our models we need the coördinates in the same metric system, and the heights in the same units (US feet).

TASK 48 : Convert the four-state DEM into the projection defined in §2.4 for the weather stations. •

The `projectRaster` command of the `raster` package re-projects a raster, keeping the same grid centres. In this case we get the CRS for the new raster from that of an existing object.

Note: The cells of reprojected raster must have the same dimensions in terms of the metric CRS, i.e., must have the same metric area (this is the definition of a raster). Since they were not the same area in the original raster (in geographic coördinates), there is not a 1-to-1 correspondence between them. The raster resolution is computed from the cell at the centre of the new raster, and then applied across the whole grid, to fill the bounding box. Thus the values of each cell must usually be computed

from those of several adjacent cells. A typical interpolation method for continuous variables such as elevation is bilinear.

```
dem.ne.m <- projectRaster(dem.ne, crs=proj4string(ne.m), method="bilinear")
nrow(dem.ne.m)

## [1] 751

ncol(dem.ne.m)

## [1] 1151

ncell(dem.ne.m)

## [1] 864401

res(dem.ne.m)

## [1] 690 926

summary(dem.ne.m)

##           srtm_1km_48
## Min.          -43.94964
## 1st Qu.       144.97001
## Median        277.67179
## 3rd Qu.       442.04596
## Max.          1503.89711
## NA's         160121.00000
```

The resolution is now in meters. Notice that the number of cells has increased.

TASK 49 : Convert the elevations to US feet. •

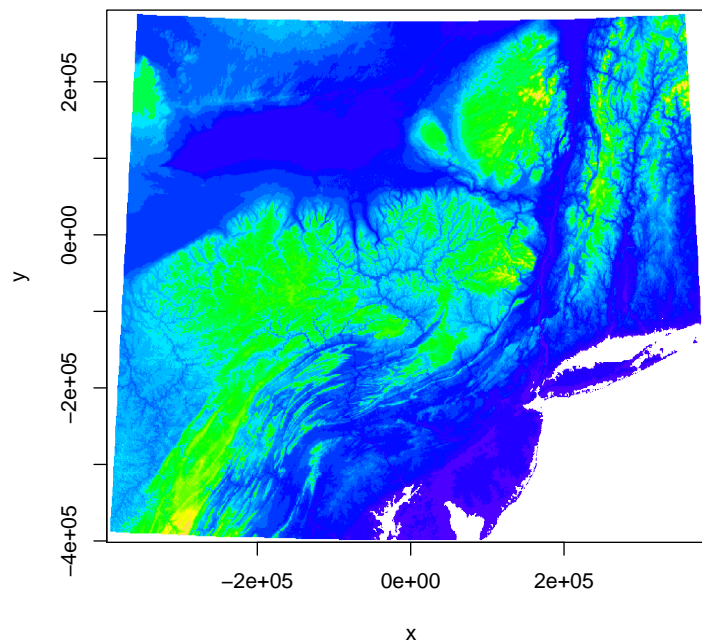
```
ft.to.m <- 0.3048
dem.ne.m <- dem.ne.m/ft.to.m
summary(dem.ne.m)

##           srtm_1km_48
## Min.          -144.1917
## 1st Qu.       475.6234
## Median        910.9967
## 3rd Qu.      1450.2820
## Max.          4934.0456
## NA's         160121.00000
```

TASK 50 : View the DEM as an image. •

The `image` function of the `raster` package provides a quick way to visualize a raster. We choose the `topo.colors` palette.

```
image(dem.ne.m, col=topo.colors(24))
```



The raster covers the whole area of the bounding box, including adjacent areas in Canada and surrounding United States.

The raster is still quite large for purposes of visualization at regional scale. Further, the grid cells are rectangular, not square.

TASK 51 : Reduce the resolution to approximately 4 x 4 km. •

The `aggregate` function of the `raster` package groups cells with a user-specified reduction factor, applying a function to the values of the cells to group; by default this is `mean`, which is appropriate for our application. The `fact` “reduction factor” argument can be one number, for equal reduction in both dimensions, or separate factors for two dimensions. We take advantage of this to get approximately square grid cells.

```
out.grid.res <- 4000
dem.ne.m <- aggregate(dem.ne.m,
  fact=c(floor(out.grid.res/res(dem.ne.m)[1]),
    floor(out.grid.res/res(dem.ne.m)[2])),
  fun=mean)
ncell(dem.ne.m)
## [1] 43428
```

TASK 52 : Save this grid as an R object. •

The `save` function is used to save R objects in the internal RData format.

```
save(dem.ne.m, file="./dem_ne_4km.RData")
```

6.2 Adjusting the grid for prediction

TASK 53 : If you did not create the grid object `dem.ne.m` in the previous section §6.1, load it from the provided R Data file. If you did create your own grid, jump to the next task. •

The `load` function is used to load R objects from an RData format file. The `save` function stores the object name(s) along with the object(s), so that after `load`, there are new objects in the workspace, with these names. Setting the optional `verbose` argument to `load` shows the objects that were loaded from the RData format file.

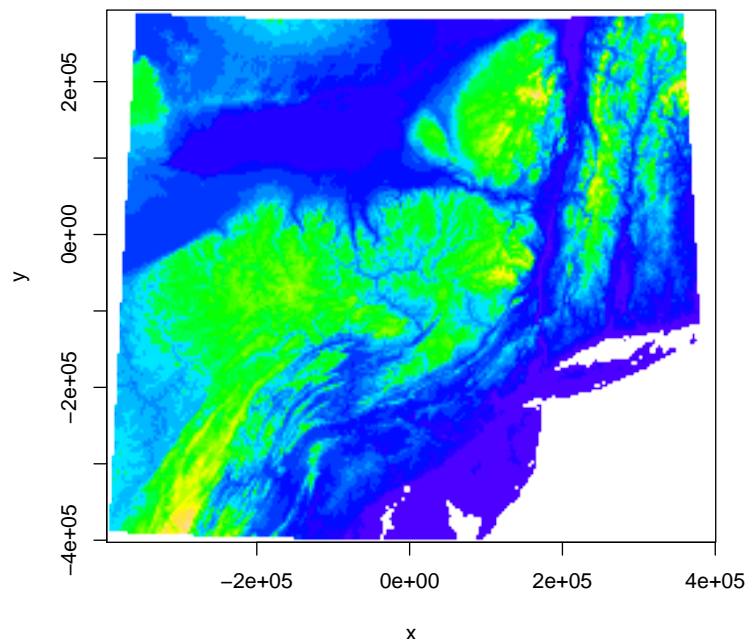
```
load(file="./dem_ne_4km.RData", verbose=TRUE)
```

```
## Loading objects:
##   dem.ne.m
```

TASK 54 : Display the prediction grid as an image. •

The `image` function of the `raster` package can be used to get a quick view of an image:

```
image(dem.ne.m, col=topo.colors(24))
```



To predict over this grid we need it to be an R data frame object, with fields for the factors used in the models, namely Northing, Easting and elevation. These must have the same names as in the model formulations.

TASK 55 : Convert the grid to a data frame with the appropriate names. •

This requires two steps: (1) convert to a `SpatialPixelsDataFrame` where the coordinates are explicit, and then (2) convert to a `data.frame` where the names are of fields, i.e., the coordinates are treated as ordinary variables.

```
class(dem.ne.m)

## [1] "RasterLayer"
## attr(,"package")
## [1] "raster"

dem.ne.m.sp <- as(dem.ne.m, "SpatialPixelsDataFrame")
coordnames(dem.ne.m.sp) <- c("E", "N")
names(dem.ne.m.sp) <- "ELEVATION_"
dem.ne.m.df <- as.data.frame(dem.ne.m.sp)
class(dem.ne.m.df)

## [1] "data.frame"

dim(dem.ne.m.df)

## [1] 35783      3

names(dem.ne.m.df)

## [1] "ELEVATION_" "E"          "N"
```

Some elevations are negative, which seems unlikely for the northeastern USA; these are likely the result of the global algorithm. These must be non-negative, otherwise the square root operation, as applied to the elevation in some models, fails.

TASK 56 : Change all negative elevations to zero. •

The `pmax` “parallel maximum” function applies `max` element-wise along a vector, such as one field of a data frame.

```
dem.ne.m.sp$ELEVATION_ <- pmax(0, dem.ne.m.sp$ELEVATION_)
dem.ne.m.df$ELEVATION_ <- pmax(0, dem.ne.m.df$ELEVATION_)
```

7 Prediction over the regional grid by OLS and GLS

Now we can predict over this grid with various models. We only use those meant for regional prediction, i.e., without local adjustments.

TASK 57 : Predict over the grid with the OLS and GLS models, add the results to the dataframe, and summarize them. •

```
dem.ne.m.df$pred.ols <- predict(m.ols, newdata=dem.ne.m.df)
dem.ne.m.df$pred.gls <- predict(m.gls, newdata=dem.ne.m.df)
summary(dem.ne.m.df[,-(1:3)])

##      pred.ols      pred.gls
## Min.   : 877    Min.   : 892.9
## 1st Qu.:1962    1st Qu.:1963.4
```

```
## Median :2263   Median :2257.7
## Mean   :2324   Mean     :2318.2
## 3rd Qu.:2592   3rd Qu.:2580.7
## Max.   :3912   Max.     :3880.8
```

user-defined
function

TASK 58: Create a function to display prediction maps using `ggplot`. •

The purpose of a user-defined **function** is to make a consistent approach to a repetitive task. In this tutorial we will make many maps showing the predictions of a mapping method over the study area – here we already have two, OLS and GLS. For consistency we define a function, using the `function` function (!). The function is assigned to a workspace name, and has a set of **arguments** with names that are used within the function.

Here since we know we always want to show the known points on the map (to show the discrepancy between predicted and known at the points) we hard-code that part into the function, and do not make it an argument. Also, we know the map data source will be a field in `dem.ne.m.df`, so that is also hard-coded.

The default colour palette ([§A](#)) is `YlGnBu`, a “sequential” palette from a “low” to a “high” value. This can be changed by the caller to show other kinds of maps.

We include an option `.plot.limits`, default `NULL`, to specify the limits of the legend scale, to allow side-by-side comparison of maps.

```
display.prediction.map <- function(.prediction.map.name,
                                   .plot.title,
                                   .legend.title,
                                   .plot.limits=NULL,
                                   .palette="YlGnBu") {
  ggplot(data=dem.ne.m.df) +
    geom_point(aes_(x=quote(E), y=quote(N),
                   color=as.name(.prediction.map.name))) +
    xlab("E") + ylab("N") + coord_fixed() +
    geom_point(aes(x=E, y=N, color=ANN_GDD50),
              data=ne.df, shape=I(20)) +
    ggtitle(.plot.title) +
    scale_colour_distiller(name=.legend.title,
                          space="Lab",
                          palette=.palette,
                          limits=.plot.limits)
}
```

Note: There are some tricky points in this function. The `aes_` function uses standard evaluation to capture the variable names, so the arguments in the call have to be quoted. If the argument is already a variable name, it needs to be quoted with the `quote` function. This has to do with how R evaluates expressions. Here it is necessary because `aes` uses non-standard evaluation rules.

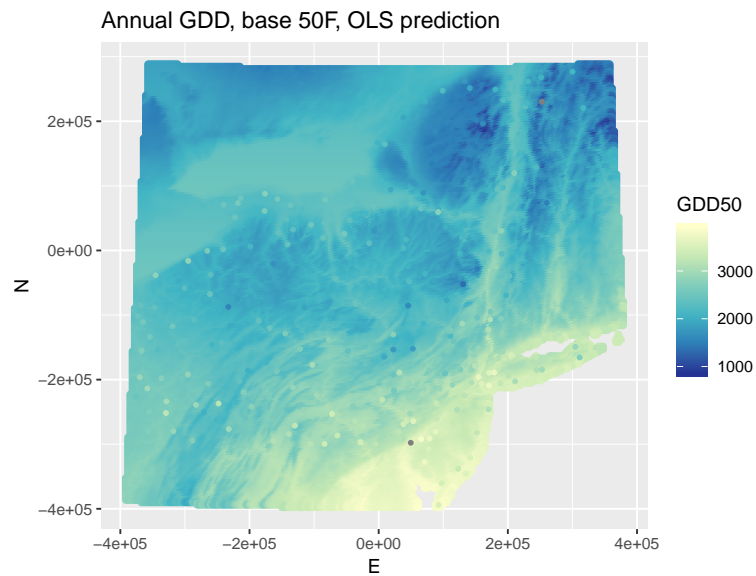
TASK 59: Plot these on the same visual scale. •

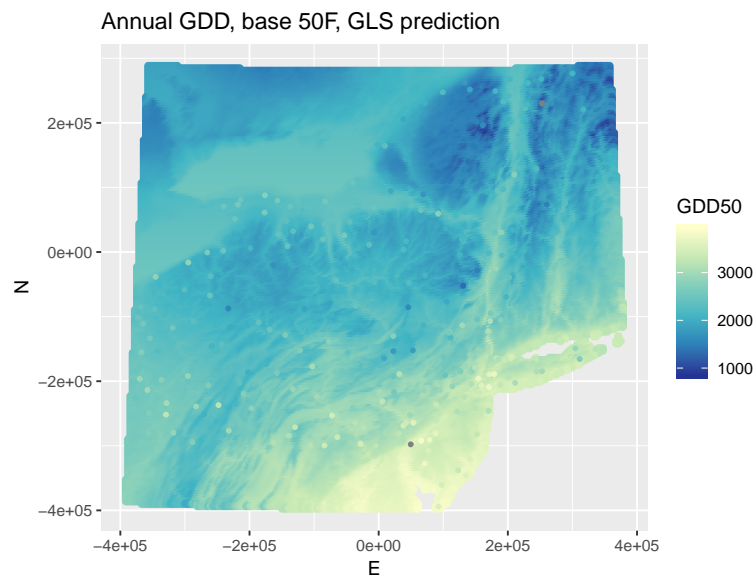
We now call the function for both predictions. The argument names used within the function are assigned to the names we give when calling the function. So the function operates on the map.

```
(gdd.pred.lim <- round(
  c(min(dem.ne.m.df[,4:5])-10, max(dem.ne.m.df[,4:5])+10)))

## [1] 867 3922

display.prediction.map("pred.ols", "Annual GDD, base 50F, OLS prediction",
  "GDD50", gdd.pred.lim)
display.prediction.map("pred.gls", "Annual GDD, base 50F, GLS prediction",
  "GDD50", gdd.pred.lim)
```





TASK 60 : Compute the differences between the OLS and GLS predictions, add them to the data frame, and display them. •

```
summary(dem.ne.m.df$diff.gls.ols <-
        dem.ne.m.df$pred.gls - dem.ne.m.df$pred.ols)
```

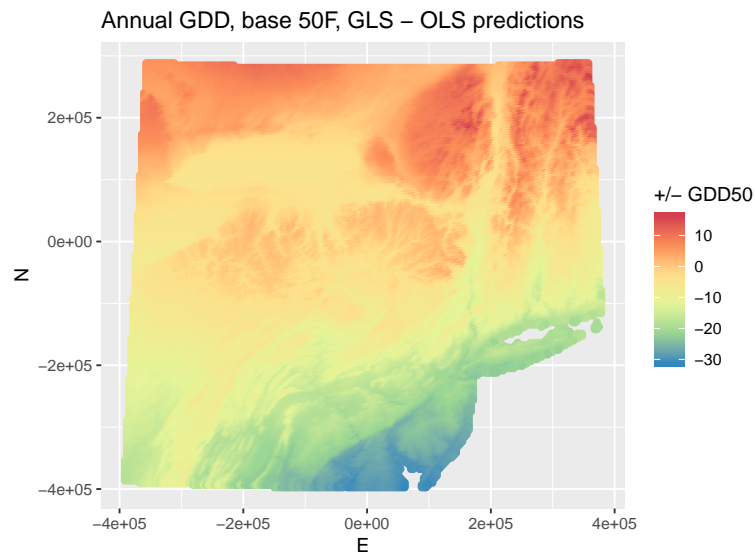
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -30.976 -11.493   -3.910   -5.371   1.034   16.096
```

The OLS and GLS predictions only vary slightly.

We also define a function to display difference maps. The default colour palette ([§A](#)) is `Spectral`, used for “diverging” palettes with a natural zero (so the two extremes are the most contrasting). This can be changed by the caller to show other kinds of maps.

```
display.difference.map <- function(.diff.map.name,
                                   .diff.map.title,
                                   .legend.name,
                                   .palette="Spectral") {
  ggplot(data=dem.ne.m.df) +
    geom_point(aes_(x=quote(E), y=quote(N),
                    color=as.name(.diff.map.name))) +
    xlab("E") + ylab("N") + coord_fixed() +
    ggtitle(.diff.map.title) +
    scale_colour_distiller(name=.legend.name,
                           space="Lab",
                           palette=.palette)
}

display.difference.map("diff.gls.ols",
                      "Annual GDD, base 50F, GLS - OLS predictions",
                      "+/- GDD50")
```



The GLS model predicts slightly higher to the north and at higher elevations.

8 Improving the trend prediction by regression kriging

In both GLS and OLS we have seen that the residuals are spatially-correlated. In §5.3 we fitted a variogram model to the GLS residuals:

```
print(vmf.r.gls)

##      model      psill      range
## 1    Nug       1.00       0.00
## 2    Exp 44558.48 17456.99
```

This exponential model has an effective spatial correlation range of about 52 km. So the residual at a prediction point is not simply the residual from the trend surface, but also depends on surrounding trend surface residuals, within this radius. Every location within the four-state area is within this effective range of one or more stations, and so the trend surface prediction can be locally adjusted in two steps:

1. krig the **residuals** to obtain the **local deviation** from the GLS trend surface. This deviation can be either “above” or “below” the surface;
2. **add** this to the trend surface prediction.

This procedure is called **Regression Kriging** [12].

The “kriging” part of the “regression kriging” term only applies to the residuals from the trend, and thus uses a variogram model and correla-

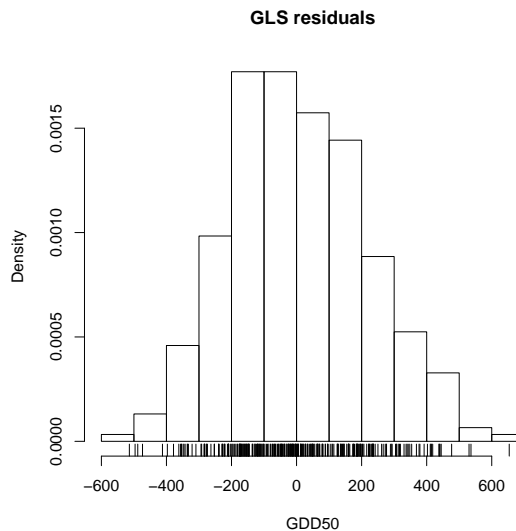
tion structure for these, not for the original values.. Most of the variation has been taken out by the trend surface (the “regression”, $Z^*(s)$). The kriging step is to adjust it locally (the “local spatial dependence” $\varepsilon(s)$), referring to Equation 1, the universal model of spatial distribution. Note that there will still be unexplained variation (“pure noise”, $\varepsilon'(s)$).

TASK 61 : Summarize the GLS trend surface residuals with a histogram and numerical summary. •

```
summary(ne.m$gls.resid)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -514.041 -149.624   -3.980    8.496 160.856  653.749
```

```
hist(ne.m$gls.resid, freq=F,
     xlab="GDD50",
     main="GLS residuals")
rug(ne.m$gls.resid)
```



Notice that the mean residual is not zero. This is because GLS trades unbiasedness for precision of the trend coefficients.

Now that we have (1) the known points; (2) a fitted authorized variogram model, we can predict at any location. The following optional section explains the mathematics of the OK system.

8.1 * The Ordinary Kriging system

OK predicts at an unknown point as a **weighted linear average** from the n known points with value z_i to obtain the value at each unknown point z_0 (Eqn.14).

$$z_0 = \sum_{i=1}^n \lambda_i z_i \quad (14)$$

The weights λ_i must sum to 1 (i.e., estimation of the mean is unbiased). Many weighting schemes can satisfy this equation, for example nearest neighbour (Thiessen polygons) (§14.5), inverse-distance weighting (§14.4), average of points within some radius, or average of some number of neighbours. The unique aspect of OK is that these weights are selected to **minimize the prediction variance** – this is the reason OK is called a “Best Unbiased Linear Predictor” (BLUP).

Note: Always remember, this “best” is with reference to the fitted variogram model. And there is no way to objectively know if that model is correct. So whether OK is “best” in the real world is not proveable.

In OK the weights λ_i are determined from the **Ordinary Kriging System**, which is derived from an expression for the prediction variance, which is minimized to derive these equations.

Weights are found by solving:

$$\mathbf{A}\lambda = \mathbf{b} \quad (15)$$

where:

$$\mathbf{A} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \gamma(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) & 1 \\ \gamma(\mathbf{x}_2, \mathbf{x}_1) & \gamma(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_2, \mathbf{x}_N) & 1 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \gamma(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) & 1 \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}$$

$$\lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_N \\ \psi \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_0) \\ \gamma(\mathbf{x}_2, \mathbf{x}_0) \\ \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_0) \\ 1 \end{bmatrix}$$

In the \mathbf{A} matrix the upper-left block $N \times N$ block is the **spatial correlation structure** of the observations; these are derived from the fitted variogram model. The last row ensures unbiasedness in estimating the spatial mean. The right-hand column is used to find the LaGrange multiplier that minimizes the variance. The kriging variance at a point is given by the scalar product of the weights (and multiplier) vector λ with the right-hand side of the kriging system.

$$\hat{\sigma}^2(\mathbf{x}_0) = \mathbf{b}^T \lambda \quad (16)$$

Note that λ includes as its last element ψ , which depends on covariance structure of the observation points. Note also that the prediction variance is computed only from semivariances, not from data values.

8.2 Predicting the residuals by Ordinary Kriging

TASK 62 : Predict the deviation from the trend surface at each location on the grid, using Ordinary Kriging (OK) of the GLS residuals, and display its summary. •

Note: We use OK instead of Simple Kriging (SK) because the *spatial* mean of the GLS residuals may not be zero. The non-spatial mean of the GLS mean is not required to be zero, as in OLS.

There are several R packages that implement kriging. Here we use the `krige` function of the `gstat` package, which uses the fitted variogram model as in `model` argument.

```
ok.gls.resid <- krige(gls.resid ~ 1,
                      loc=ne.m, newdata=dem.ne.m.sp,
                      model=vmf.r.gls)

## [using ordinary kriging]

summary(ok.gls.resid)

## Object of class SpatialPixelsDataFrame
## Coordinates:
##      min      max
## E -392899.3 379900.7
## N -399006.4 289937.6
## Is projected: TRUE
## proj4string :
## [+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76
## +ellps=WGS84 +units=m]
## Number of points: 35783
## Grid attributes:
##   cellcentre.offset cellsize cells.dim
## E          -394624.3     3450      231
## N          -400858.4     3704      188
## Data attributes:
##   var1.pred      var1.var
## Min.   :-465.7709   Min.    : 748
## 1st Qu.: -42.9739   1st Qu.:30020
## Median :  0.1104   Median :37006
## Mean   :  2.6124   Mean   :35517
## 3rd Qu.: 40.8275   3rd Qu.:44442
## Max.    :574.7911   Max.    :44900
```

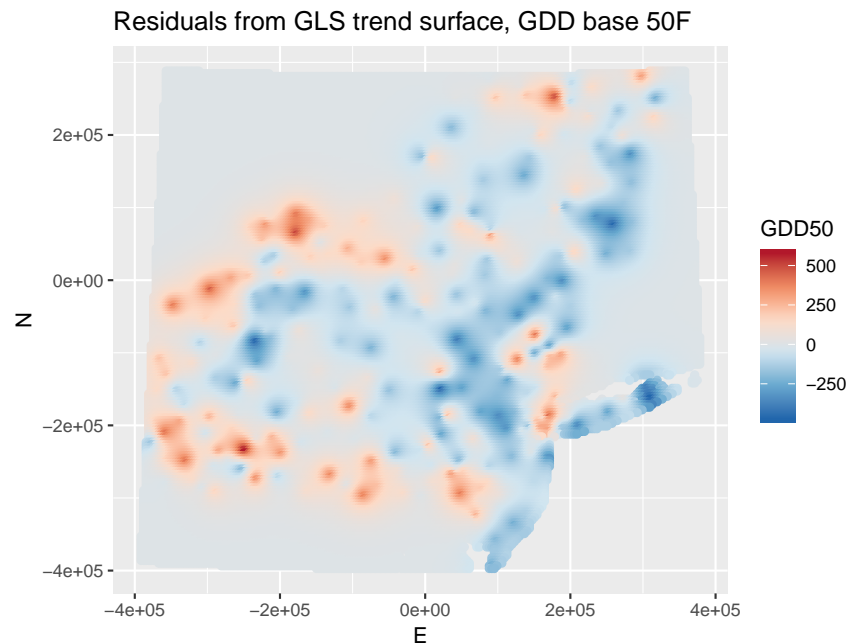
Half of the adjustments are between about ± 40 GDD50, so not very many; however there are some quite large adjustments at the extremes. The kriging prediction variance has a very low value at a grid cell centre that must be close to an observation, but otherwise is fairly large; the median prediction standard deviation is 188 GDD50.

TASK 63 : Add this deviation and its prediction variance to the GLS trend surface data frame. •

```
dem.ne.m.df$ok.gls.resid <- ok.gls.resid$var1.pred
dem.ne.m.df$ok.gls.resid.var <- ok.gls.resid$var1.var
```

TASK 64 : Display a map of the deviations. •

```
ggplot() +
  geom_point(aes(x=E, y=N, colour=ok.gls.resid), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Residuals from GLS trend surface, GDD base 50F") +
  scale_colour_distiller(name="GDD50", space="Lab", palette="RdBu")
```



Q15 : *Where are the largest adjustments to the GLS trend?* [Jump to A15](#) •

Q16 : *Notice there is no adjustment further than about 50 km from a station, why?* [Jump to A16](#) •

The predictions at these points are the **spatial** mean of the GLS residuals (*not* necessarily their arithmetic mean):

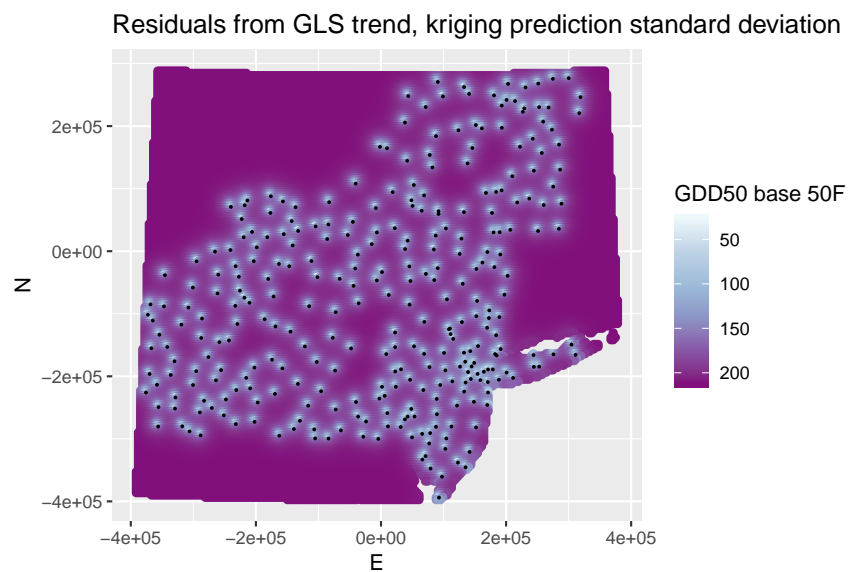
```
mean(dem.ne.m.df$ok.gls.resid) # spatial mean over the grid
## [1] 2.612386

mean(ne.m$gl.s.resid)          # arithmetic mean at observation points
## [1] 8.495601
```

TASK 65 : Display the kriging prediction standard deviation of the residuals. •

This is the square root of the prediction variance, so in the same units as the prediction. We also show the observation points, to show how the kriging prediction variance depends on the point configuration.

```
ggplot() +
  geom_point(aes(x=E, y=N, colour=sqrt(ok.gls.resid.var)),
             data=dem.ne.m.df) + xlab("E") + ylab("N") +
  coord_fixed() +
  geom_point(aes(x=E, y=N), data=ne.df, size=0.5,
             colour="black", shape=I(20)) +
  ggtitle("Residuals from GLS trend, kriging prediction standard deviation") +
  scale_colour_distiller(name="GDD50 base 50F", space="Lab",
                        palette="BuPu", trans="reverse")
```



8.3 The GLS-Regression Kriging prediction

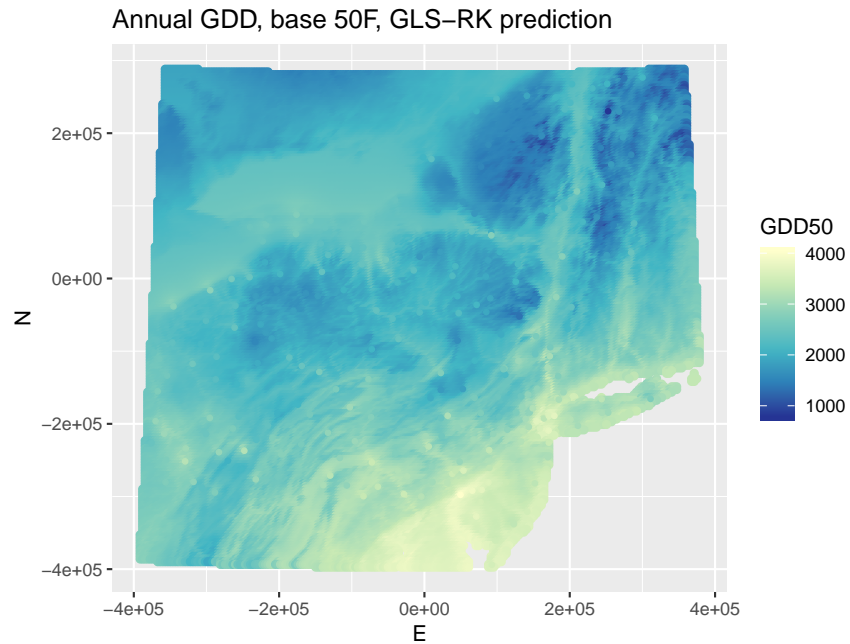
We have two predictions: the trend and the deviations from it. Adding these will give a final prediction.

TASK 66 : Add the kriged GLS residuals to the trend surfaces for a final GLS-RK prediction. •

```
summary(dem.ne.m.df$pred.rkgls <-
        dem.ne.m.df$pred.gls + dem.ne.m.df$ok.gls.resid)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      921.4  1935.0  2269.1  2320.8  2619.6  3977.5

display.prediction.map("pred.rkgls", "Annual GDD, base 50F, GLS-RK prediction",
                      "GDD50")
```



Q17 : *How does this compare to the GLS trend surface map?* [Jump to A17 •](#)

9 Kriging with external drift (KED)

Another way to predict using both the trend and local deviations from it is the one-step method Kriging with External Drift (KED)³¹. This is a form of **kriging**, i.e., a linear weighted average from the n known points with value z_i to obtain the value at each unknown point z_0 (Eqn.17). The weights λ_i must sum to 1 (i.e., estimation of the mean is unbiased), and are selected to **minimize the prediction variance** – this is the reason OK is called a “Best Unbiased Linear Predictor” (BLUP).

$$z_0 = \sum_{i=1}^n \lambda_i z_i \quad (17)$$

The difference between KED and OK (§8.1) is that KED also includes covariates in the kriging system, so that the linear trend with covariates and the local deviations at each prediction point are solved together to obtain the weights λ .

³¹ KED is mathematically equivalent to what is called Universal Kriging (UK); that term is often reserved for KED when only coördinates are used as covariables.

9.1 * The Universal Kriging system

The weights λ_i are determined from the **Universal Kriging System**, which is derived from an expression for the prediction variance, which is minimized to derive these equations.

Weights are found by solving:

$$\mathbf{A}_U \boldsymbol{\lambda}_U = \mathbf{b}_U \quad (18)$$

where

$$\mathbf{A}_U = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_1, \mathbf{x}_N) & 1 & f_1(\mathbf{x}_1) & \cdots & f_k(\mathbf{x}_1) \\ \vdots & \cdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_N, \mathbf{x}_N) & 1 & f_1(\mathbf{x}_N) & \cdots & f_k(\mathbf{x}_N) \\ 1 & \cdots & 1 & 0 & 0 & \cdots & 0 \\ f_1(\mathbf{x}_1) & \cdots & f_1(\mathbf{x}_N) & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ f_k(\mathbf{x}_1) & \cdots & f_k(\mathbf{x}_N) & 0 & 0 & \cdots & 0 \end{bmatrix}$$

In this matrix the upper-left block $N \times N$ block is the **spatial correlation structure** of the **residuals** from the trend; these are derived from the fitted variogram model. The lower-left $k \times n$ block (and its transpose in the upper-right) are the **trend** predictor values at sample points. In KED these are the covariate values at the prediction points; in UK these are the coördinate values at these points. The rest of the matrix fits with $\boldsymbol{\lambda}_U$ and \mathbf{b}_U to set up the solution:

$$\boldsymbol{\lambda}_U = \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_N \\ \psi_0 \\ \psi_1 \\ \vdots \\ \psi_k \end{bmatrix} \quad \mathbf{b}_U = \begin{bmatrix} \gamma(\mathbf{x}_1, \mathbf{x}_0) \\ \vdots \\ \gamma(\mathbf{x}_N, \mathbf{x}_0) \\ 1 \\ f_1(\mathbf{x}_0) \\ \vdots \\ f_k(\mathbf{x}_0) \end{bmatrix}$$

The $\boldsymbol{\lambda}_U$ vector contains the N weights for the sample points and the $k + 1$ LaGrange multipliers (1 for the overall mean and k for the trend model). The \mathbf{b}_U vector is structured like an additional column of \mathbf{A}_U , but referring to the point to be predicted. This contains the semivariances of the prediction point vs. the known points.

The kriging variance at a point is given by the scalar product of the weights (and multiplier) vector $\boldsymbol{\lambda}$ with the right-hand side of the kriging system.

$$\hat{\sigma}^2(\mathbf{x}_0) = \mathbf{b}_U^T \boldsymbol{\lambda}_U \quad (19)$$

Good explanations of KED are from Webster and Oliver [29] and Goovaerts [8]; in §14 we explain Ordinary Kriging (OK), where there is no trend, only local interpolation.

9.2 Computing the empirical residual variogram

KED uses the `krige` method of the `gstat` package directly with the residual variogram, and so does not require a separate regression prediction step. KED as implemented by `krige` uses GLS to compute the trend component, with a covariance structure specified by the analyst, generally from fitting a variogram model to the residual variogram. This differs from `gls`, which computes the covariance structure by REML.

In `gstat` the residuals are estimated from a linear model fit, using the `variogram` function with a formula for the trend. Since the covariance structure is not yet known, this must be by OLS.

TASK 67: Use the `variogram` function to compute the variogram of the residuals from a model of GDD50 predicted by Northing and square root of elevation, and fit an exponential model to it. Use a cutoff of 100 km and a bin size of 16 km, to have enough points in the closest bin, and to avoid very local effects. Compare to the variogram computed directly from the OLS residuals in §4.4. •

We use the `variogram` function, but instead of a null formula (right-hand side 1) to specify the spatial mean, we specify a formula for the trend. The left-hand side is the variable of interest, not the residuals from a previous step.

```
v.ked <- variogram(ANN_GDD50 ~ sqrt(ELEVATION_) + N, locations=ne.m,  
                  cutoff=100000, width=16000)
```

It is exactly the same variogram model as we found in §4.4, because these are exactly the same residuals. The difference is that we don't need to find the trend surface formula, we just need the trend surface residuals in order to compute the variogram for KED.

9.3 Fitting the residual variogram model

We require a continuous function of semivariance vs. separation, so that for any separation we can compute the semivariance to be used in the kriging equations. This must be an **authorized** variogram model.

Here, an exponential model appears to fit the empirical variogram of the residuals; this is a common and the simplest choice of variogram model.

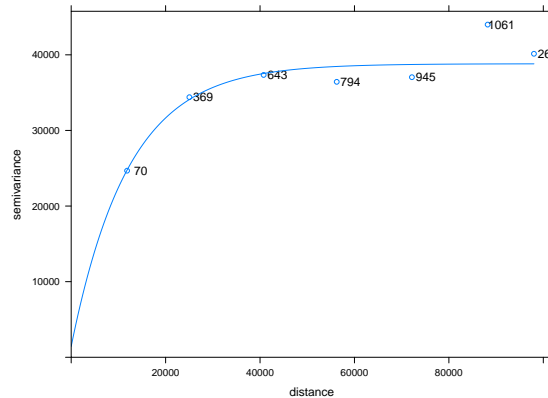
The `fit.variogram` function fits variogram models specified with the `vgm` “variogram model” function. We initialize the weighted least-squares (WLS) fit to the empirical variogram with our eyeball estimates.

```
(vmf.ked <- fit.variogram(v.ked, vgm(15000, "Exp", 20000, 20000)))
```

```
##      model      psill      range
```

```
## 1  Nug 1435.842    0.00
## 2  Exp 37382.603 12106.96

plot(v.ked, plot.numbers=TRUE, model=vmf.ked)
```



The effective range of the exponential model is three times the range parameter, so here about 363 km. This implies that there is local structure, not explained by the covariables, to this range.

With this covariance structure, we can now predict by KED, again specifying the dependence of the target variable on the covariables.

9.4 Predicting with KED

TASK 68 : Compute the KED prediction and its variance over the prediction grid. •

We call `krige` with a model formula that shows the linear dependence on covariables, exactly the same formula that we used to compute the empirical variogram of the residuals in the previous step. These two formulas must be identical.

```
k.ked <- krige(ANN_GDD50 ~ sqrt(ELEVATION_)+ N, locations=ne.m,
               newdata=dem.ne.m.sp, model=vmf.ked)
```

```
## [using universal kriging]
```

```
summary(k.ked)
```

```
## Object of class SpatialPixelsDataFrame
## Coordinates:
##      min      max
## E -392899.3 379900.7
## N -399006.4 289937.6
## Is projected: TRUE
## proj4string :
## [+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76
## +ellps=WGS84 +units=m]
## Number of points: 35783
## Grid attributes:
##   cellcentre.offset cellsize cells.dim
## E      -394624.3      3450        231
## N      -400858.4      3704        188
## Data attributes:
```

```
##      var1.pred      var1.var
## Min.   : 915.5    Min.    : 3689
## 1st Qu.:1943.9    1st Qu.:32454
## Median :2268.9    Median :36765
## Mean   :2321.1    Mean    :34754
## 3rd Qu.:2609.2    3rd Qu.:39185
## Max.   :3944.8    Max.    :42441
```

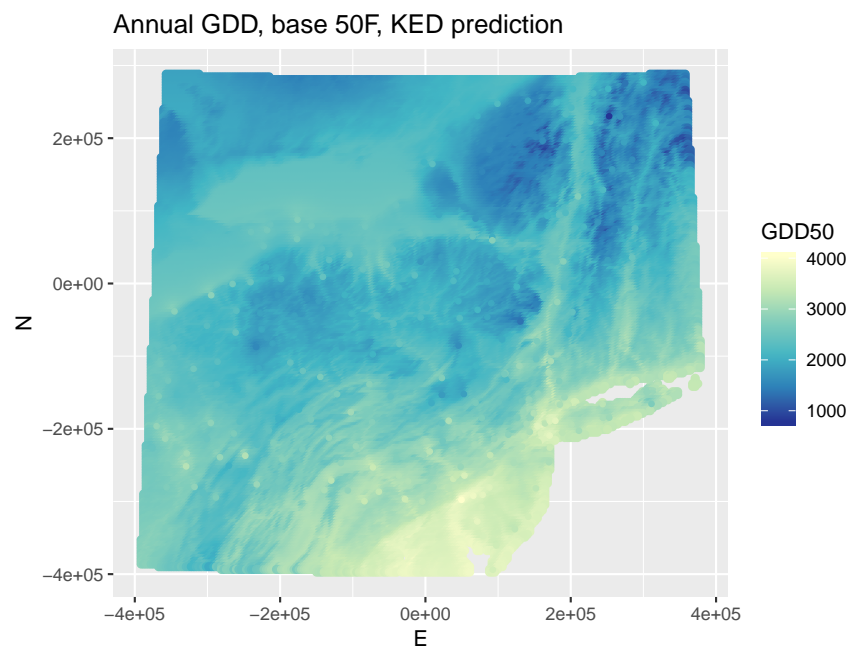
Note that the kriging prediction variance is also computed; this is known because by definition kriging minimizes it.

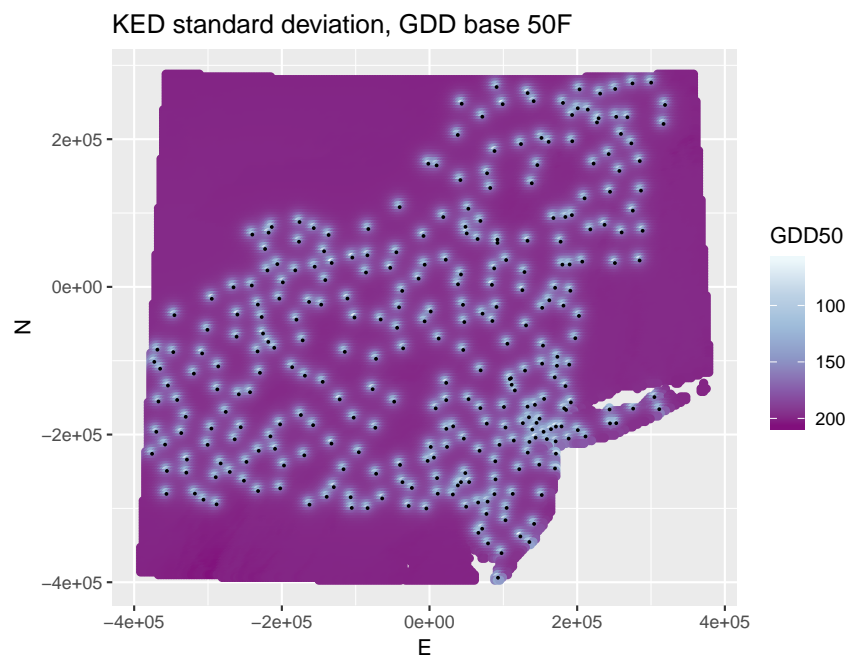
TASK 69 : Display the prediction and its standard deviation

To do this we add the prediction and its standard deviation to the grid data frame and then display these as maps.

```
dem.ne.m.df$pred.ked <- k.ked$var1.pred
dem.ne.m.df$pred.ked.sd <- sqrt(k.ked$var1.var)
display.prediction.map("pred.ked",
  "Annual GDD, base 50F, KED prediction",
  "GDD50")

#
ggplot() +
  geom_point(aes(x=E, y=N, colour=pred.ked.sd), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  geom_point(aes(x=E, y=N), data=ne.df, size=0.5,
    colour="black", shape=I(20)) +
  ggtitle("KED standard deviation, GDD base 50F") +
  scale_colour_distiller(name="GDD50", space="Lab",
    palette="BuPu", trans="reverse")
```





The predictions should be slightly different from the RK-GLS predictions, because the spatial correlation of the residuals was estimated from the OLS trend surface.

TASK 70 : Compute and display the differences between RK-GLS and KED over the grid. •

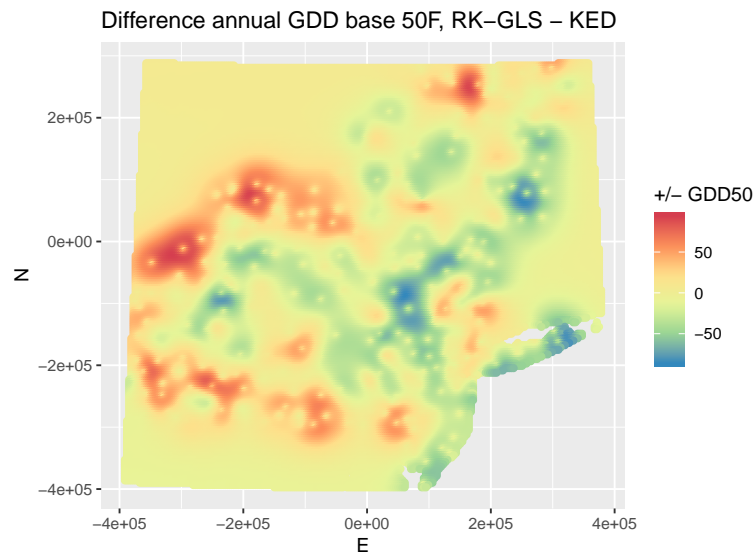
```
summary(dem.ne.m.df$diff.rkgls.ked <-
  dem.ne.m.df$pred.rkgls - dem.ne.m.df$pred.ked)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -85.1523 -15.6981   0.5366  -0.2692 12.5207   93.4232
```

Half of the differences are quite small, under about 15 GDD50. There are a few larger differences, but all less than 100 GDD50 (compare with the sample mean, 2518).

Display the locations of the differences:

```
display.difference.map("diff.rkgls.ked",
  "Difference annual GDD base 50F, RK-GLS - KED",
  "+/- GDD50")
```



The largest positive residuals (RK-GLS predicts higher) are along Lake Erie and western Lake Ontario. The cooling lake effect which we saw in the OLS residuals is increased in the GLS residuals, since the GLS trend predicts somewhat lower than the OLS trend in this area. So the RK is higher here. The largest negative residuals are in the Catskills and southern VT, where the GLS trend predicted somewhat higher than the OLS trend.

Challenge: There is a strong positive anomaly (RK-GLS predicts higher but just locally) near the Chazy station, on Lake Champlain in the centre North. There are several nearby stations also on the Lake. What is the reason for this anomaly?

In this section we have seen that KED has some practical advantages over RK-GLS:

1. It is easier to implement;
2. The covariance structure is estimated beforehand, and there is no risk that the procedure might not converge on a solution, as in REML;
3. It gives a prediction variance in the same step.

However, we can see from the results in this case study that there can be some fairly large differences in predictions.

9.5 Accuracy assessment

An objective way to evaluate the predictive power of JED is by **Leave-one-out cross-validation** (LOOCV). Here each point is removed from the dataset in turn, and predicted by the others, using the fitted variogram model. If the observation points well represent the total population, as they do here by design of the weather station network, this gives a good estimate of the prediction error.

TASK 71 : Compute and summarize the LOOCV for this KED prediction. •

The `krige.cv` function of the `gstat` package computes this:

```
kcv.ked <- krige.cv(ANN_GDD50 ~ sqrt(ELEVATION_)+ N, locations=ne.m, model=vmf.ked)

summary(kcv.ked$residual)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -477.630 -138.430  -11.306   -1.444   130.210   645.663
```

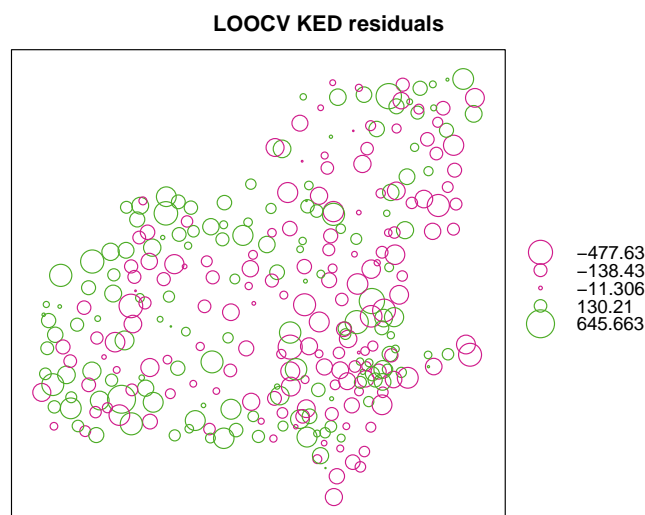
Overall the results are fairly good, but there are some large prediction errors at both extremes. An overall measure is the **root of the mean squared error**, RMSE:

```
(loocv.ked.rmse <- sqrt(sum(kcv.ked$residual^2)/length(kcv.ked$residual)))

## [1] 199.4083
```

TASK 72 : Display a bubble plot of the LOOCV residuals. •

```
bubble(kcv.ked, zcol="residual", pch=1, main="LOOCV KED residuals")
```



There are several regions with intermixed fairly large under- and over-

predictions; this means that in these regions there are local factors not accounted for. Other regions are consistently over- or under-predicted (Vermont mountains, Lake Ontario plain, respectively).

9.6 KED in a local neighbourhood

An advantage of KED over GLS-RK is that KED can be applied in some local neighbourhood, so that the relation with the covariables (here, Northing, Easting, square root of elevation) is re-fit at each prediction point. Besides the obvious computational advantage (fewer points → less computation), this allows a varying effect of the covariates over space. In our example, it may be that the effect on GDD50 of +100 km Northing may be more towards the south of the region than the north, or vice versa; or it may be a smaller effect in a north-south trending large valley such as the Hudson or Lake Champlain. The effect of elevation may be more, or less, in the Adirondacks in northern New York compared to the Allegheny Plateau in Pennsylvania.

Note: The linear model is *not* re-solved at each point; the UK system (§9.1) implicitly includes these in the solution. The $A_U \mathbf{a}$ matrix includes the covariance between the neighbourhood points, as well as their values of the covariates, and the \mathbf{b}_U vector includes the covariance between the neighbourhood points and the prediction point, as well as the prediction point's covariate values.

Whether the kriging is global or local, the \mathbf{b}_U vector must be computed at each prediction point. For local kriging a full $A_U \mathbf{a}$ matrix of all the observation points can be rapidly cut down to the set of local points closest to the prediction point.

Note: This has some relation to Geographically-Weighted Regression (§11), where the trend surface is explicitly re-computed at each prediction point, using only the points in some neighbourhood. But in local KED we also consider the spatial correlation between the residuals of the local trend.

TASK 73 : Recompute the KED prediction over the study area of §9.4 with a local neighbourhood. •

Q18 : *Why we not need or want to re-compute the residual variogram model (§9.3)?* *Jump to A18* •

The krig package has two optional arguments that can be used to implement this:

1. `nmax` “maximum number of neighbours to use”;
2. `maxdist` “maximum distance to a point to use”; this can be used along with `nmin` “minimum number of neighbours to use” to ensure that no predictions are made with few points. This latter will produce NA “not available” values if there are prediction points too far from the minimum number

We prefer `nmax` because here the points are well-distributed, and we can use the kriging prediction variance to find areas that are too poorly-predicted.

The obvious question is how to determine this number. One way is to try different numbers and compare their cross-validation statistics (see the Challenge at the end of this §). Here we choose to use 20% of the points, i.e., $305/5 \approx 61$ to show how the method works.

```
ked.n.max <- 61 # change this to try other numbers of neighbours
k.ked.nn <- krige(ANN_GDD50 ~ sqrt(ELEVATION_) + N, locations=ne.m,
                 newdata=dem.ne.m.sp, model=vmf.ked,
                 nmax=ked.n.max)

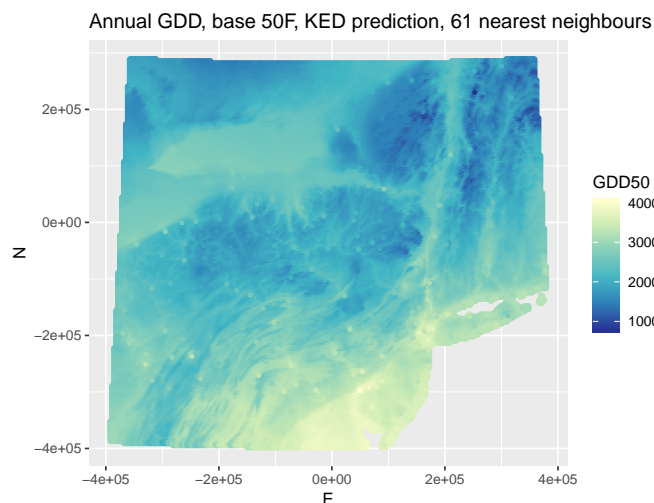
## [using universal kriging]

summary(k.ked.nn@data)

##      var1.pred      var1.var
## Min.   : 906.4    Min.   : 3829
## 1st Qu.:1922.7    1st Qu.:32979
## Median :2279.9    Median :37579
## Mean   :2333.2    Mean   :36625
## 3rd Qu.:2661.9    3rd Qu.:41240
## Max.   :3947.8    Max.   :69912
```

TASK 74 : Display the prediction map. •

```
dem.ne.m.df$pred.ked.nn <- k.ked.nn$var1.pred
dem.ne.m.df$pred.ked.nn.sd <- sqrt(k.ked.nn$var1.var)
display.prediction.map("pred.ked.nn",
                      paste("Annual GDD, base 50F, KED prediction,",
                             ked.n.max, "nearest neighbours"),
                      "GDD50")
```



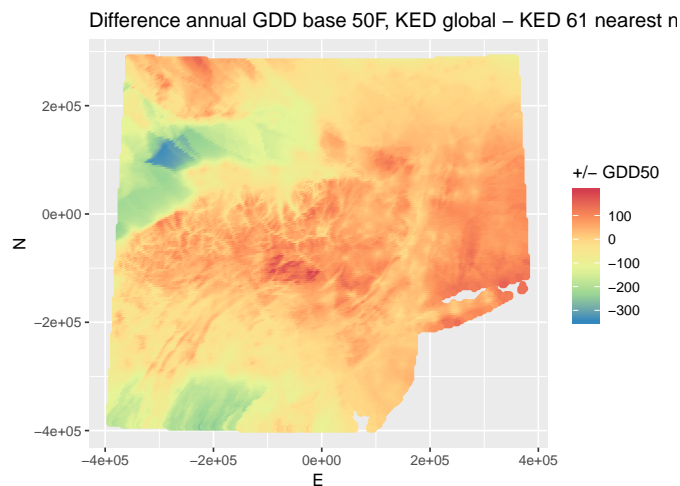
TASK 75 : Compute and display a map of the difference in predictions

between the global and local KED predictions.

```
summary(dem.ne.m.df$diff.ked <- dem.ne.m.df$pred.ked - dem.ne.m.df$pred.ked.nn)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -339.930 -45.461   3.778  -12.170  42.728  198.600

display.difference.map("diff.ked",
  paste("Difference annual GDD base 50F, KED global - KED",
    ked.n.max,"nearest neighbours"),
  "+/- GDD50")
```



Q19 : *Where are the largest differences? Explain.*

Jump to A19

TASK 76 : Compute and display the cross-validation statistics; compare them to global KED.

```
kcv.ked.nn <- krige.cv(ANN_GDD50 ~ sqrt(ELEVATION_)+ N,
  locations=ne.m, model=vmf.ked,
  nmax=ked.n.max)

summary(kcv.ked.nn$residual)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -430.555 -144.898  -18.085   -5.355  117.227  617.697

summary(kcv.ked$residual)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -477.630 -138.430  -11.306   -1.444  130.210  645.663

(loocv.ked.nn.rmse <- sqrt(sum(kcv.ked.nn$residual^2)/length(kcv.ked.nn$residual)))

## [1] 191.7587

(loocv.ked.rmse <- sqrt(sum(kcv.ked$residual^2)/length(kcv.ked$residual)))

## [1] 199.4083
```

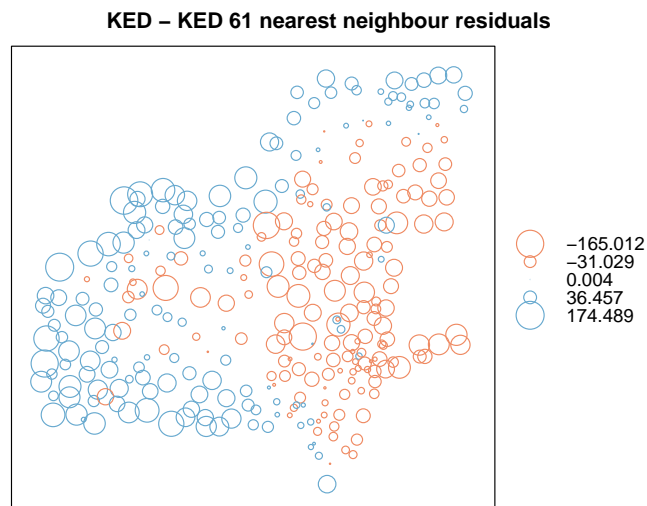
Q20 : Which KED, global or local, gives the best cross-validation results? What can you conclude about the strength of regional to local effects on GDD50? *Jump to A20 •*

TASK 77 : Display a bubble plot of the difference between the global and local cross-validation residuals. •

```
summary(kcv.ked$diff <- kcv.ked$residual - kcv.ked.nn$residual)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -165.0122 -31.0292   0.0038    3.9111   36.4572   174.4893

bubble(kcv.ked, zcol="diff",
       pch=1, main=paste("KED - KED",
                        ked.n.max, "nearest neighbour residuals"),
       col=c("#EF8A62", "#67A9CF"))
```



Q21 : Where are the largest differences? Does this seem to be geographically-consistent? Can you explain? *Jump to A21*

•

Challenge: Experiment with different numbers of neighbours to find the optimum for local KED.

Challenge: Since the trend surface is now fit locally, it may be that Easting is significant in some parts of the region. Refit the global residual variogram with this included in the linear model, and use it in the kriging prediction formula. How much and where does this change the prediction?

9.7 * Demonstration that KED uses GLS to determine the trend

Above we stated that KED as implemented by `krige` uses GLS to compute the trend component, with a covariance structure specified by the analyst, i.e., from modelling the residual variogram. This section shows the difference between this prediction and one where the trend is computed by OLS.

We first compute the OK trend surface using `krige` with a null model, that is, not using any local information; this is the OLS prediction of the trend. This is the equivalent of using the `lm` function, but is more convenient because it directly produces a gridded data structure, as in kriging.

```
k.ok <- krige(ANN_GDD50 ~ sqrt(ELEVATION_)+ N, locations=ne.m,
             newdata=dem.ne.m.sp, model=NULL)

## [ordinary or weighted least squares prediction]
```

We then krig the residuals from the OLS trend of the known points; we computed those in §4.3, and computed their variogram in §4.4.

```
k.okr <- krige(ols.resid ~ 1, locations=ne.m,
              newdata=dem.ne.m.sp, model=vmf.r.ols)

## [using ordinary kriging]
```

We then add these together to get a final prediction of the trend and the local deviations from it, which is what KED does in one step:

```
k.ok$rk.pred <- k.ok$var1.pred + k.okr$var1.pred
```

Finally, compare this to the KED prediction, and compute their differences:

```
k.ok$diff.pred <- k.ok$rk.pred - k.ked$var1.pred
summary(k.ok$rk.pred); summary(k.ked$var1.pred)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      901    1940    2267    2320    2609    3948
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      915.5 1943.9 2268.9 2321.1 2609.2 3944.8

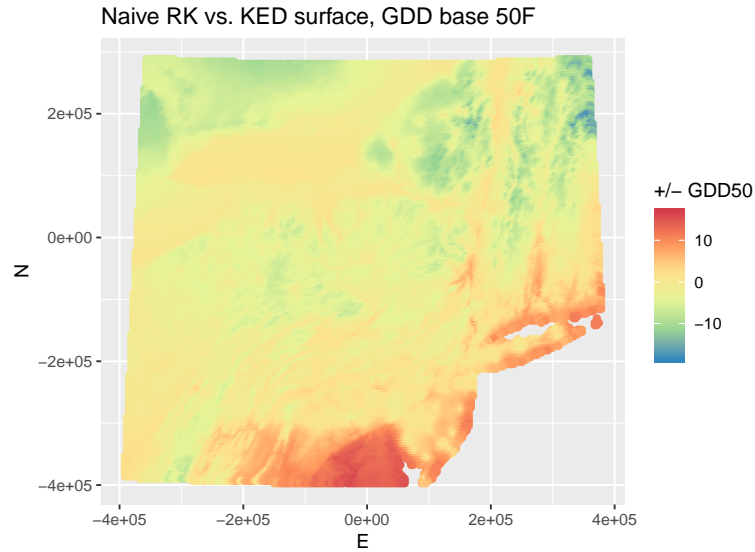
summary(k.ok$diff.pred)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -18.3806 -4.2696 -1.6590 -1.5411  0.6258 16.6615
```

We see a smaller magnitude difference here as when we compared KED to RK/GLS in the previous §.

The geographic distribution of the differences is because of the different trend surface:

```
dem.ne.m.df$diff.okrk.ked.pred <- k.ok$diff.pred
display.difference.map("diff.okrk.ked.pred",
                      "Naive RK vs.\ KED surface, GDD base 50F",
                      "+/- GDD50")
```



We also saw differences between GLS and OLS for the calibration points, in §5.5 and for the prediction grid in §7; here the differences are smaller because we also kriged the residuals from the two trend surfaces.

So this shows that KED as implemented by `krige` does use GLS, not OLS, to compute the trend component. The difference with RK/GLS is that the covariance structure is based on the OLS trend, not fit at the same time as the trend surface coefficients.

10 Generalized Additive Models

Generalized Additive Models (GAM) are similar to multiple linear regression, except that each term in the linear sum of predictors need not be the predictor variable itself, but can be an empirical smooth function of it. So instead of the linear model of k predictors:

$$y_i = \beta_0 + \sum_k \beta_k x_{k,i} + \varepsilon_i \quad (20)$$

we allow *functions* f_k of these:

$$y_i = \beta_0 + \sum_k f_k(x_{k,i}) + \varepsilon_i \quad (21)$$

The advantage is that non-linear relations in nature can be fit, without any need to try transformations or to fit piecewise regressions. If this is a better model fit, it should result in better predictions. The model is additive, so the marginal contribution of each predictor to the model

fit can be determined. The disadvantage is that it is just an empirical fit and can not be extrapolated beyond the range of calibration. A further disadvantage is that the choice of function is arbitrary; it is generally some smooth function of the predictor, with the degree of smoothness determined by cross-validation.

Note: The GAM should never be extrapolated (there is no data to support it), whereas a polynomial can, with caution, be extrapolated, on the theory that the data used to fit the model extends outside the range. This is of course very dangerous for higher-order polynomials, which are a main competitor to GAM.

Hastie et al. [11, §9.1] give a thorough explanation of GAM; a simplified explanation of the same material is given in James et al. [14, §7.7]. In a geostatistical setting, we can choose the coördinates as the predictors (as in a trend surface) but fit these with smooth functions, rather than polynomials. We can also fit any other predictor this way, e.g., in this example, the elevation.

The smooth functions can be chosen in many ways; the most common are cubic splines with knots at each value of the predictor. But we first examine whether a smooth curve, rather than one line (as in linear regression) better matches the dependence of the annual GDD50 on the three possible predictors.

TASK 78 : Display a scatterplot of the three predictors against the annual GDD50, with an empirical smoother. •

We use the `ggplot2` graphics package, introduced in §3, to produce the scatterplot and show a smoother with standard error. A simple way to visualize the trend is with a local polynomial regression, provided with the `loess` function, and incorporated into the scatterplot with the `geom_smooth` function.

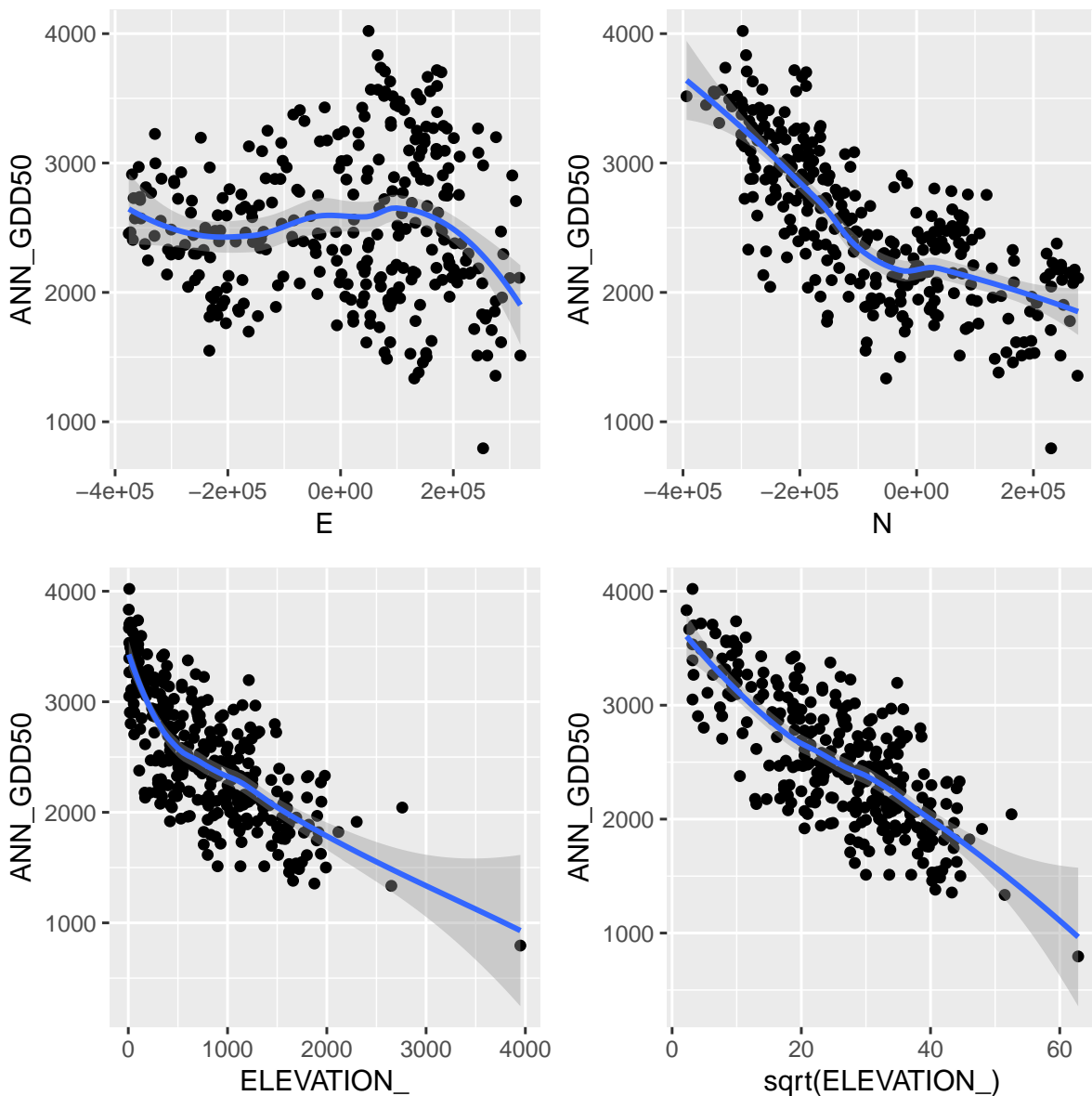
Note: The `loess` function has an `span` argument, which controls the degree of smoothing by setting the neighbourhood for the local fit as a proportion of the number of points. The default `span=0.75` thus uses the 3/4 of the total points closest each point. These are then weighted so that closer points have more weight; see `?loess` for details. The default works well in most situations, and here we only want a visual impression, not a “best fit” in a statistical sense.

We use the `gridExtra` package, which includes a function `grid.arrange` to arrange saved plots in a grid.

```

g1 <- ggplot(ne.df, aes(x=E, y=ANN_GDD50)) +
  geom_point() +
  geom_smooth(method="loess")
g2 <- ggplot(ne.df, aes(x=N, y=ANN_GDD50)) +
  geom_point() +
  geom_smooth(method="loess")
g3 <- ggplot(ne.df, aes(x=ELEVATION_, y=ANN_GDD50)) +
  geom_point() +
  geom_smooth(method="loess")
g4 <- ggplot(ne.df, aes(x=sqrt(ELEVATION_), y=ANN_GDD50)) +
  geom_point() +
  geom_smooth(method="loess")
require(gridExtra)
grid.arrange(g1, g2, g3, g4, ncol = 2)

```



Q22 : Do these marginal relations appear to be linear in the predictors?

These marginal plots motivate us to try a GAM.

10.1 Fitting a Generalized Additive Model

GAM can be fit in R with the `gam` function of the `mgcv` “Mixed GAM Computation Vehicle” package. This specifies the model with a formula, as with `lm`, but terms can now be arbitrary functions of predictor variables, not just the variables themselves or simple transformations that apply to the whole range of the variable, e.g. `sqrt` or `log`. Smooth functions of one or more variables are specified with the `s` function of the `mgcv` package.

TASK 79 : Load the `mgcv` package into the workspace. •

```
require(mgcv)
```

Common practice in GAM for models using coördinates is to smooth them together with a bivariate smoother, by default a thin plate regression spline; see §13 below for details. Other covariates, here the elevation, are smoothed with penalized regression splines. These control the degree of smoothness by penalizing increasingly complex models, i.e., those with more curvature; see the help text `?s` “defining smooths in GAM formulae” for details. In practice the default parameters work well.

TASK 80 : Fit a GAM to the annual GDD50 at the observation stations, with the predictors being a two-dimensional thin-plate spline of the coördinates and a one-dimensional penalized regression spline of the elevation. •

```
m.g.xy <- gam(ANN_GDD50 ~ s(E, N) + s(ELEVATION_), data=ne.df)
summary(m.g.xy)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## ANN_GDD50 ~ s(E, N) + s(ELEVATION_)
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 2517.518      9.986   252.1   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df    F p-value
## s(E,N)         23.529 27.300 37.8  <2e-16 ***
## s(ELEVATION_)   8.521  8.922 51.6  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.908   Deviance explained = 91.7%
## GCV = 34111   Scale est. = 30415        n = 305
```



```
summary(residuals(m.g.xy))
```

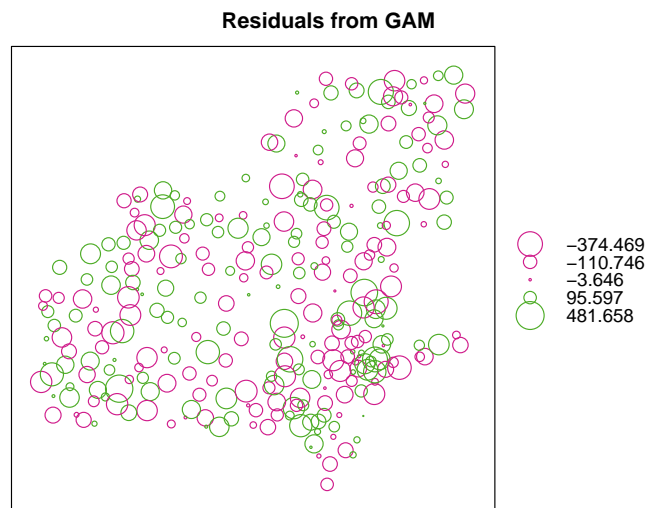
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	-374.469	-110.746	-3.646	0.000	95.597	481.658

Q23 : *How well does this model fit the calibration observations?* [Jump to A23](#) •

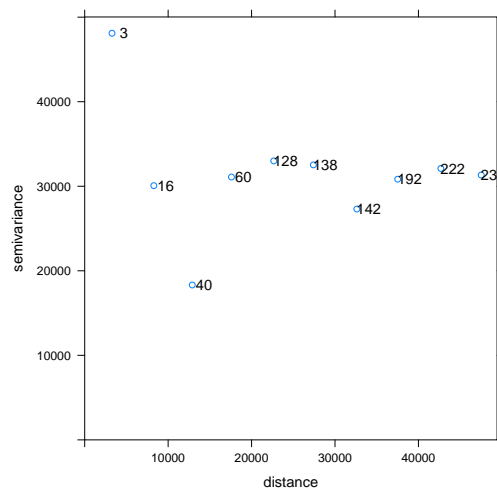
An important consideration is whether the residuals have any spatial structure; recall this is why we replaced OLS with GLS.

TASK 81 : Plot the residuals as a bubble plot, and examine their spatial structure with a variogram. •

```
ne.m@data$resid.m.g.xy <- residuals(m.g.xy)
bubble(ne.m, zcol="resid.m.g.xy", pch=1, main="Residuals from GAM")
```



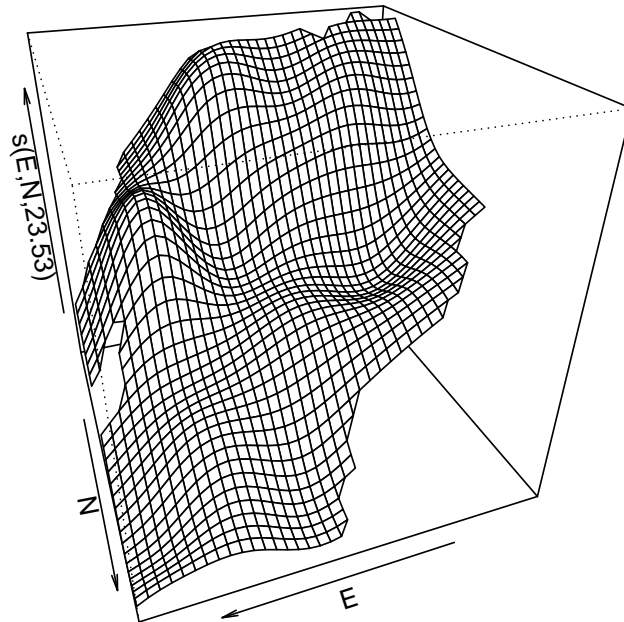
```
vr <- variogram(resid.m.g.xy ~ 1, loc=ne.m, cutoff=50000, width=5000)
plot(vr, pl=T)
```



Q24 : *Does there appear to be any local spatial correlation of the residuals? Does the empirical variogram support your conclusion?* [Jump to A24](#) •

The `plot.gam` function of the `mgcv` package displays the marginal smooth fit. For the 2D surface (model term $s(E, N)$), this is shown as a wireframe plot if the optional `scheme` argument is set to 1. The `select` argument selects which model term to display. We orient it to see lowest GDD towards viewer, using the `theta` argument:

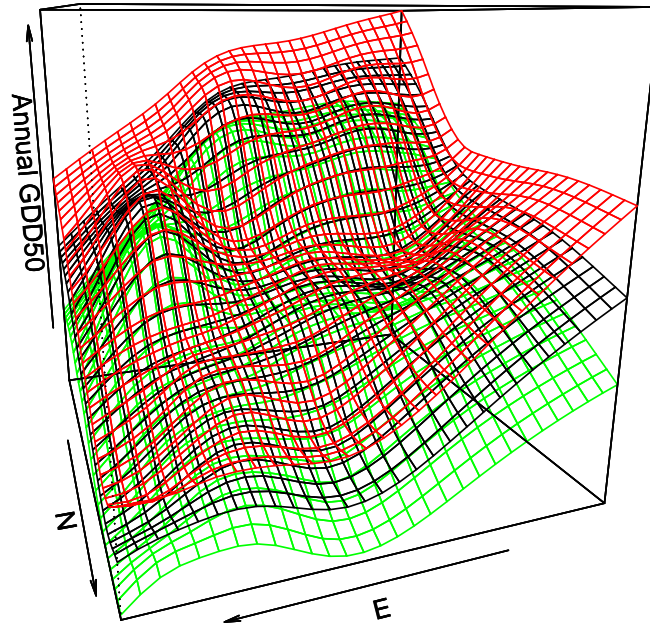
```
plot.gam(m.g.xy, rug=T, se=T, select=1,
         scheme=1, theta=30+130, phi=30)
```



Q25 : Does the GAM 2D trend differ from a linear trend surface? [Jump to A25](#) •

This surface can also be shown with the `vis.gam` function of the `mgcv` package, also showing ± 1 standard error of fit:

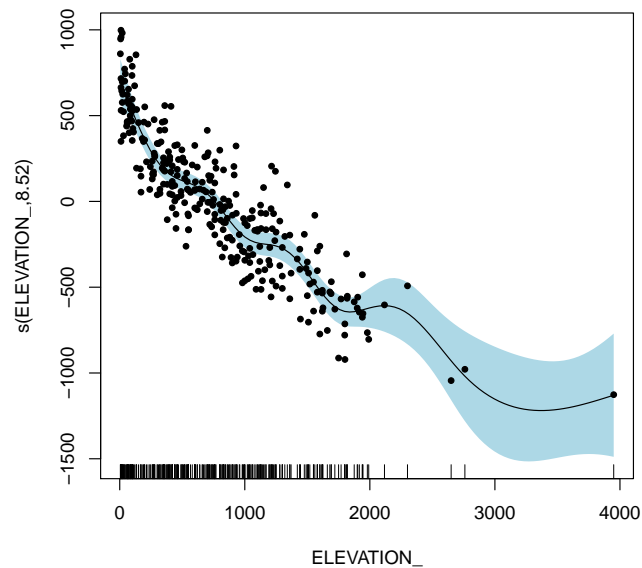
```
vis.gam(m.g.xy, plot.type="persp", color="terrain",
        theta=160, zlab="Annual GDD50", se=1.96)
```



red/green are ± 1.96 s.e.

The marginal relation with elevation can be presented as a scatterplot, with the confidence intervals and residuals from the fit:

```
plot.gam(m.g.xy, select=2, rug=T, se=T, residuals=T, pch=20,
          shade=T, seWithMean=T, shade.col="lightblue")
```



Q26 : Does the fitted marginal relation with elevation appear to be linear? Jump to A26 •

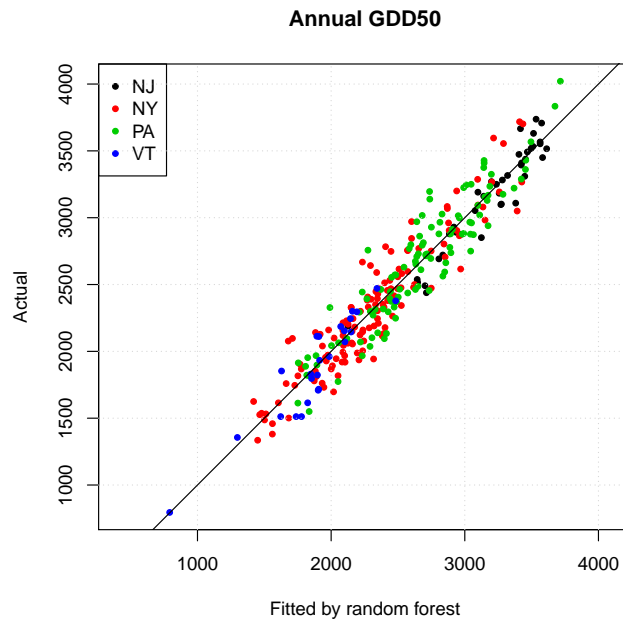
Notice the very large confidence interval at the high elevations – we have so few points there that the smoother is quite uncertain in this range.

TASK 82 : Compare the GAM model fits with the actual values. •

```
(rmse.gam <- sqrt(sum(residuals(m.g.xy)^2)/length(residuals(m.g.xy))))

## [1] 164.6781

plot(ne.m$ANN_GDD50 ~ predict(m.g.xy, newdata=ne.df),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest", ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.df$STATE), pch=20, col=1:4)
grid(); abline(0,1)
```



The RMSE is 164.7; there are no observations that are particularly badly-fit.

10.2 GAM prediction over the study area

Since we now have a model which uses the covariables known across the prediction grid, we can use the model to predict.

TASK 83 : Predict the annual GDD50, and the standard error of prediction, across the prediction grid, using the fitted GAM, and display the predictions. •

The `predict.gam` function predicts from a fitted GAM. The `se.fit` optional argument specifies that the standard error of prediction should also be computed.

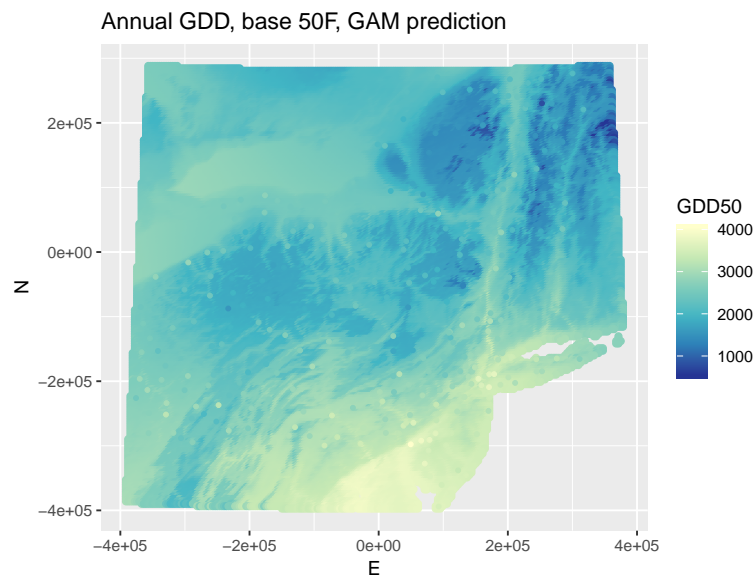
```
tmp <- predict.gam(object=m.g.xy, newdata=dem.ne.m.df, se.fit=TRUE)
summary(tmp$fit)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      563   1972   2326   2357   2683   3868

summary(tmp$se.fit)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      36.01  55.34  63.30   73.09  84.93  250.57

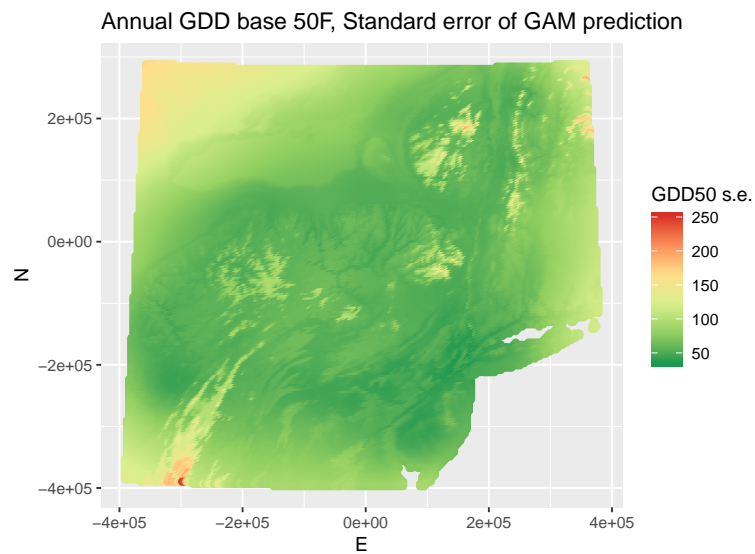
dem.ne.m.df$pred.gam <- tmp$fit
dem.ne.m.df$pred.gam.se <- tmp$se.fit
display.prediction.map("pred.gam", "Annual GDD, base 50F, GAM prediction",
                       "GDD50")
```



This map shows more detail than the OLS and GLS maps, especially in the high elevations and along the Atlantic coast.

TASK 84 : Display the map of the standard errors of prediction. •

```
ggplot() +
  geom_point(aes(x=E, y=N, colour=pred.gam.se), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Annual GDD base 50F, Standard error of GAM prediction") +
  scale_colour_distiller(name="GDD50 s.e.", space="Lab", palette="RdYlGn",
    direction=-1)
```



Consistent with the marginal plots, we see that the standard error is much higher at the highest elevations, where there are few observations to support the GAM.

An obvious question is where this map differs from the GLS and GLS-RK maps.

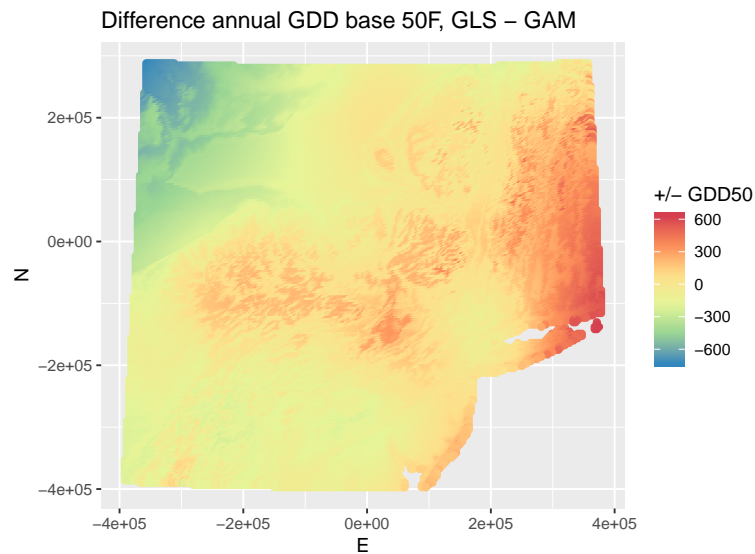
TASK 85 : Compute the differences in GLS and GAM model predictions over the grid, summarize numerically, and display as a difference map.

```
summary(dem.ne.m.df$diff.gls.gam <-
  dem.ne.m.df$pred.gls - dem.ne.m.df$pred.gam)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -721.75 -151.86  -34.99  -38.67   92.93   624.67
```

There are some large differences. See where these are located:

```
display.difference.map("diff.gls.gam",
  "Difference annual GDD base 50F, GLS - GAM",
  "+/- GDD50")
```

Q27 : *Where are the largest differences between the GAM and GLS predictions?*

[Jump to A27](#)

•

The RK component of GLS-RK allowed for local adjustment to the overall trend surface, based on local spatial dependence. GAM also adjusts locally, via its smoothers. How close are these predictions?

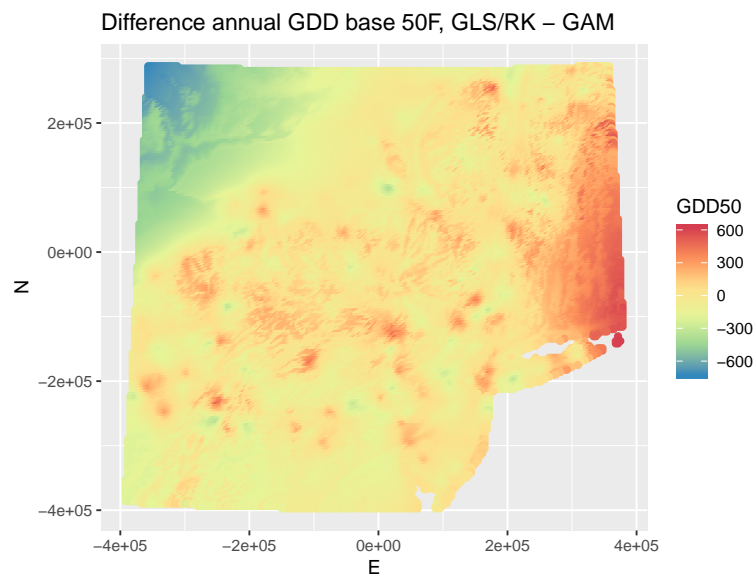
TASK 86 : Compute the differences in GLS-RK and GAM model predictions over the grid, summarize numerically, and display as a difference map.

```
dem.ne.m.df$diff.rkgls.gam <- dem.ne.m.df$pred.rkgls - dem.ne.m.df$pred.gam
summary(dem.ne.m.df$diff.rkgls.gam)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -721.75 -120.97  -25.17  -36.06   72.78   612.40
```

There are some large differences. See where these are located:

```
ggplot() +
  geom_point(aes(x=E, y=N, colour=diff.rkgls.gam), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Difference annual GDD base 50F, GLS/RK - GAM") +
  scale_colour_distiller(name="GDD50", space="Lab", palette="Spectral")
```



Q28 : *Where are the largest differences between the GAM and GLS-RK predictions?* Jump to A28 •

11 Geographically Weighted Regression

Geographically Weighted Regression (GWR) is an extension of linear or generalized linear regression, which re-fits the regression equation at each data point, based on some neighbourhood and weighting scheme. This is indicated if the model is **spatially non-stationary**, i.e., the relation is not the same over the whole map, so that a single global model, although representing the overall relation, misses important local variations. A main use of GWR is to detect if this is the case. Prediction is a secondary objective. It should have a physical or social basis, i.e., some reason to think this might occur, and why, and over what spatial extent.

GWR was first developed and explained in the text of Fotheringham et al. [6]; since then there have been developments to extend to multiple scales [7]. See also simulations and applications in [9, 10, 28, 36].

GWR may be compared with Regression Kriging (RK) (§8), where the global trend is fit and then adjusted locally by kriging the residuals. This assumes that the global trend is correct, but affected by local factors. By contrast, GWR considers that any fitted global trend is not correct, there is no overall effect of the covariates, only local effects. A related competing method for mapping is Kriging with External Drift (KED) (§9) in a restricted neighbourhood, i.e., the feature-space trend is re-fit at each

prediction point according to some restricted radius, and the residuals from this local trend, in the same neighbourhood are at the same time kriged.

There are other local regression methods, aimed mainly at mapping, e.g., minimum-curvature (thin-plate splines), see §13. GWR has several advantage over these: It gives explicit values of (1) the bandwidth within which a local regression should be fit; (2) the regression coefficients at each point; therefore (3) the variability and spatial pattern of these.

11.1 Theory

GWR is a locally-adaptive version of ordinary linear regression applied over the entire point set, where the aim is to fit:

$$y_i = \beta_0 + \sum_k \beta_k x_{ik} + \varepsilon_i \quad (22)$$

where the errors ε_i are I.I.D. This is the classic OLS formulation; the solution is the well-known:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (23)$$

By contrast, in GWR the equation:

$$y_i = \beta_0(u_i, v_i) + \sum_k \beta_k(u_i, v_i) x_{ik} + \varepsilon_i \quad (24)$$

is fit separately at every point to be evaluated, but only using some subset of the x_{ik} .³² Thus the coefficients vary with the coordinates (u_i, v_i) .

The solution to Equation 24 is a weighted version of the OLS solution, i.e. weighted least squares (WLS):

$$\hat{\beta}(u_i, v_i) = (\mathbf{X}^T \mathbf{W}_{(u_i, v_i)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_{(u_i, v_i)}^{-1} \mathbf{y} \quad (25)$$

where the $\mathbf{W}_{(u_i, v_i)}$ are the diagonals of a matrix containing the weights of the points in the neighbourhood to be used to fit the regression for point (u_i, v_i) ; the \mathbf{X} matrix includes all points (as in a global model) but some will have very low weight (Gaussian kernel) or even zero weight (bisphere kernel), see below.

Note: This same formulation can be used globally, with the weights based on some criterion such as the precision of the observation, with more precise observation receiving higher weights.

The key issue is how to select the weights and neighbourhood. These are from a **spatial kernel**, of which there are several forms, considering the distance d_{ij} between the target point i and a distant point j and the

³² If all x_{ik} were used, the solution would be the same at every point.

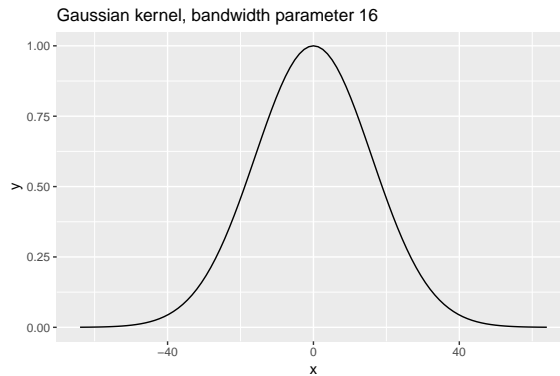


Figure 5: Gaussian kernel

bandwidth h which determines the range of influence of points in the regression. There are two common kernels:

Gaussian : $w_{ij} = e^{-\frac{1}{2}(d_{ij}/h)^2}$

Bisquare : $w_{ij} = (1 - (d_{ij}^2/h^2))^2$ if $d_{ij} \leq h$, else $w_{ij} = 0$

Thus the Gaussian kernel considers all the points, with exponentially decaying weight, whereas the bisquare only considers points within some radius. In this section we continue with the Gaussian kernel, because of its smooth decay properties. See Fig. 5.

```
require(ggplot2)
h <- 16
ggplot(data=data.frame(x=c(-64, 64)), aes(x)) +
  stat_function(fun=function(x) { exp(-1/2 * (x/h)^2) }) +
  labs(title = paste("Gaussian kernel, bandwidth parameter", h))
g.kernel <- function(d, h) { exp(-1/2 * (d/h)^2) }
```

It reaches a weight of 0.5 at $d = h\sqrt{-2\log(.5)} \approx 1.18h$, 0.1 at $d = h\sqrt{-2\log(.1)} \approx 2.15h$, and 0.02 at $d = h\sqrt{-2\log(.02)} \approx 2.80h$

There are two ways to define the bandwidth:

- fixed : Use the distance parameter h in the above formulation, and only consider points within that distance (bisquare) or decay according to it (Gaussian). This is appropriate if points are more or less evenly spread, although it can produce artefacts at the edges of the study area.
- adaptive : Always consider a proportion of the points to use for each local fit, and then weight them according to the kernel function. This is appropriate if points irregularly distributed, as it ensures that there are enough points with sufficient weight to calibrate the regression.

Selecting the bandwidth is a critical step in a GWR analysis. We do not want to fit too locally because (1) not enough sample points to reliably calibrate the regression model; (2) we may miss broad-scale trends; (3) we may fit artificial local variability, not corresponding to the process. We do not want to fit too broadly, because we might miss true local

spatial variability in the process being modelled by the regression. The bandwidth can be specified by the analyst based on prior knowledge of the spatial scale of the process, but normally it is estimated by cross-validation, see below §11.4.

11.2 Dataset

We continue with the regional climate dataset `ne.m`. However, to make the coefficients more interpretable, we convert the coordinates to km.

TASK 87 : Make a version with the coordinates in km, to make the regression coefficients easier to read. •

To do this, we create a `proj4string` string by editing the current one which uses meters, and then apply `spTransform`.

```
require(rgdal)
(p4str.km <- paste0(substr(proj4string(ne.m), 1,
                           nchar(proj4string(ne.m))-1), "km"))

## [1] "+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76 +ellps=WGS84 +units=km"

ne.km <- spTransform(ne.m, p4str.km)
bbox(ne.km)

##           min           max
## E -375.7452 318.9604
## N -393.9229 276.6141
```

The study area is about 695 km E-W and 670 km N-S.

TASK 88 : Add the metric E and N in km to the data frame, to be used in the linear models. These fields can not have the same names as the coordinates, so change those to x and y. •

```
coordnames(ne.km) <- c("x", "y")
ne.km$E <- coordinates(ne.km)[,1]
ne.km$N <- coordinates(ne.km)[,2]
```

11.3 Modelling approach

We model the annual GDD50 from the coordinates and elevation. However, we now do not assume that the relation is constant over the entire area, but might vary. There should be a physical basis for this. For example:

- Although there is a physical basis for the adiabatic lapse rate based on decreasing air density with altitude, this could be muted by increased vertical air mixing, e.g., near oceans or lakes, or in mountains with strong thermal winds.
- Although there is a physical basis for the decreasing annual solar energy with increasing North latitude, this could be masked by lake effects or local topography.

- Certainly if we only use a local neighbourhood, the average value (intercept) will vary – this is an artefact of the GWR method and substitutes for much of the effect of Northing.

And in any case we can *explore* to see if the coefficient change, and then try to determine the physical (geographic) cause.

Several packages implement GWR, we use `spgwr` by the authors of `sp`³³. Another choice is `GWmodel`[18].

TASK 89 : Load the `spgwr` package.

```
require(sp)
require(spgwr)
```

11.4 Determining the bandwidth

The `gwr.sel` function selects an appropriate bandwidth, depending on the type of kernel by trying a set of bandwidths and comparing these by **cross-validation**. This compares the root mean square prediction error at the known points using GWR with the proposed bandwidth, and then selects the bandwidth that minimizes this objective function. In other words, it computes the GWR with the selected bandwidth using all the points, but then applies it at known points without considering that point, and collects the errors. This is the default method argument `cv`.

TASK 90 : Determine the optimum **fixed** bandwidth.

```
(bw.Gauss.f <- gwr.sel(ANN_GDD50 ~ E + N + sqrt(ELEVATION_), data = ne.km,
                      adapt = F, verbose = F,
                      gweight=gwr.Gauss))

## [1] 60.50139
```

The fixed bandwidth parameter h is selected as 60.5 km. At this separation the weight of a distant point is $e^{1/2} = 0.6065$.

TASK 91 : Determine the optimum **adaptive** bandwidth.

```
(bw.Gauss.a <- gwr.sel(ANN_GDD50 ~ E + N + sqrt(ELEVATION_), data = ne.km,
                      adapt = T, verbose = F,
                      gweight=gwr.Gauss, method="cv"))

## [1] 0.02308575
```

The adaptive bandwidth parameter h is selected as 0.0231, i.e., about 3.3% of the data will be used for each regression, so in this case about 10 of the 305 known points. This is the number of points within one “radius”, i.e., where points get at least $e^{1/2} = 0.6065$ weight. In the Gaussian

³³ Although the authors provide the disclaimer: “NOTE: This package does not constitute approval of GWR as a method of spatial analysis; see `example(gwr)`”.

kernel all points are included; however after about $2.8h \approx 0.093 \approx 28$ points the more distant points get less than 2% weight.

11.5 Computing the GWR

We now compute the GWR at the observation points, with both the fixed and adaptive Gaussian kernels, using the `gwr` function. Note the use of the `adapt` argument for the adaptive kernel, but the `bandwidth` argument for the fixed kernel.

TASK 92 : Compute the GWR, i.e., the locally-varying regression equations, at each known point. •

```
(gwr.Gauss.f <- gwr(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
  data = ne.km,
  bandwidth = bw.Gauss.f,
  gweight=gwr.Gauss))

## Call:
## gwr(formula = ANN_GDD50 ~ E + N + sqrt(ELEVATION_), data = ne.km,
##     bandwidth = bw.Gauss.f, gweight = gwr.Gauss)
## Kernel function: gwr.Gauss
## Fixed bandwidth: 60.50139
## Summary of GWR coefficient estimates at data points:
##               Min.      1st Qu.      Median      3rd Qu.
## X.Intercept. 2421.53760 3103.43217 3292.95917 3426.26408
## E            -3.42389  -0.95632  -0.55876  -0.17502
## N            -3.69416  -2.30000  -1.48158  -0.81772
## sqrt.ELEVATION_ -51.85338 -40.02365 -36.39290 -32.67106
##               Max.      Global
## X.Intercept. 4366.31917 3306.2372
## E            1.42417  -0.5387
## N            1.64766  -1.6614
## sqrt.ELEVATION_ -21.99504 -35.4905

(gwr.Gauss.a <- gwr(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
  data = ne.km,
  adapt=bw.Gauss.a,
  gweight=gwr.Gauss))

## Call:
## gwr(formula = ANN_GDD50 ~ E + N + sqrt(ELEVATION_), data = ne.km,
##     gweight = gwr.Gauss, adapt = bw.Gauss.a)
## Kernel function: gwr.Gauss
## Adaptive quantile: 0.02308575 (about 7 of 305 data points)
## Summary of GWR coefficient estimates at data points:
##               Min.      1st Qu.      Median      3rd Qu.
## X.Intercept. 2252.497817 3002.119640 3258.560180 3416.371004
## E            -4.573084  -1.120115  -0.614117  0.011682
## N            -4.862175  -2.512653  -1.341797  -0.578654
## sqrt.ELEVATION_ -50.329280 -40.173975 -36.251682 -31.731118
##               Max.      Global
## X.Intercept. 5387.760774 3306.2372
## E            8.380249  -0.5387
## N            4.585209  -1.6614
## sqrt.ELEVATION_ -20.941457 -35.4905
```

Note that the summary compares the range of the GWR coefficients with the single global OLS fit.

TASK 93 : Extract just the results from the `gwr` object for easier coding.

The main results are in a `SpatialPointsDataFrame` `SDF`; the global coefficients in a vector `lm`, within the `gwr` object – these of course are the same for both kernels, since the kernel is not used in the global regression.

```
gwr.f <- gwr.Gauss.f$SDF@data
gwr.a <- gwr.Gauss.a$SDF@data
(gwr.global <- gwr.Gauss.f$lm$coefficients)
```

##	(Intercept)	E	N	sqrt(ELEVATION_)
##	3306.2372354	-0.5387108	-1.6613608	-35.4905240

11.6 Interpretation

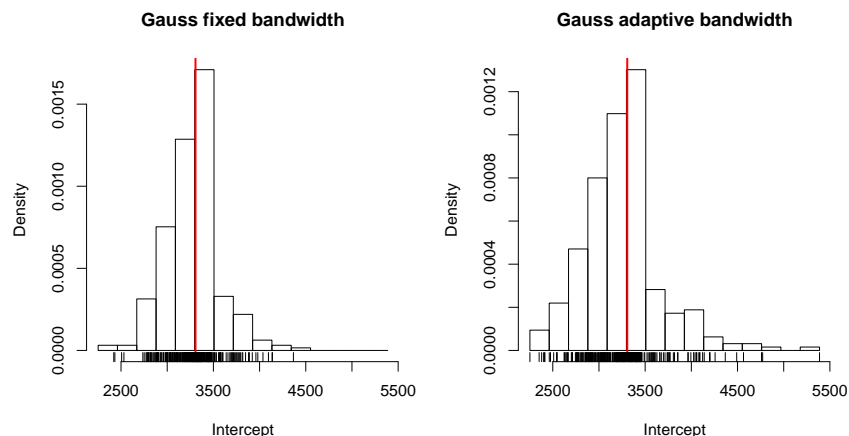
The GWR results are **interpreted** to determine if there is non-stationarity in the regression, and from that try to determine the causes. We can examine these in both feature and geographic space.

11.6.1 Feature space distribution

TASK 94 : Display histograms of the feature-space distribution of the coefficients and goodness-of-fit for the two methods, each shown with a common scale, and with the value from the global regression as a thicker red line.

First, the intercepts:

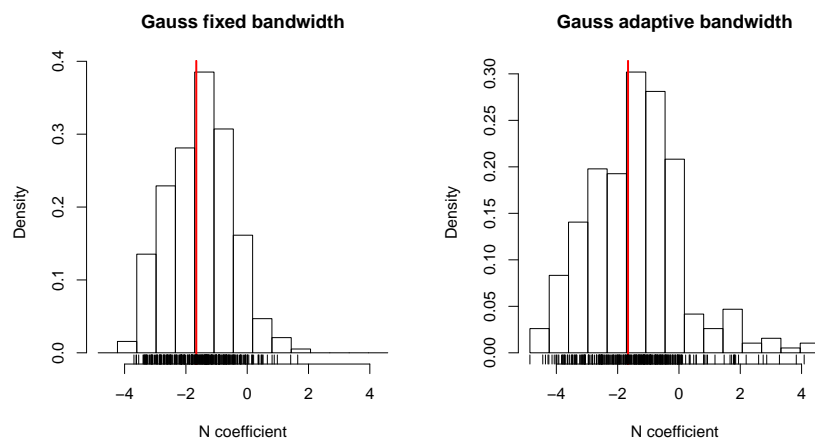
```
xlim.b0 <- range(gwr.f$'(Intercept)', gwr.a$'(Intercept)')
breaks.b0 <- seq(xlim.b0[1], xlim.b0[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$'(Intercept)', main="Gauss fixed bandwidth",
     xlab="Intercept", xlim=xlim.b0, freq=F, breaks=breaks.b0)
abline(v=gwr.global['(Intercept)'], col="red", lwd=2)
rug(gwr.f$'(Intercept)')
hist(gwr.a$'(Intercept)', main="Gauss adaptive bandwidth",
     xlab="Intercept", xlim=xlim.b0, freq=F, breaks=breaks.b0)
abline(v=gwr.global['(Intercept)'], col="red", lwd=2)
rug(gwr.a$'(Intercept)')
par(mfrow=c(1,1))
```



There is quite a range, corresponding to the four-state study area. The local regression is centred at the intercept in its smaller region, so the intercept is a spatial average of the stations in that region, of course weighted by the kernel function.

Second, the North coefficient:

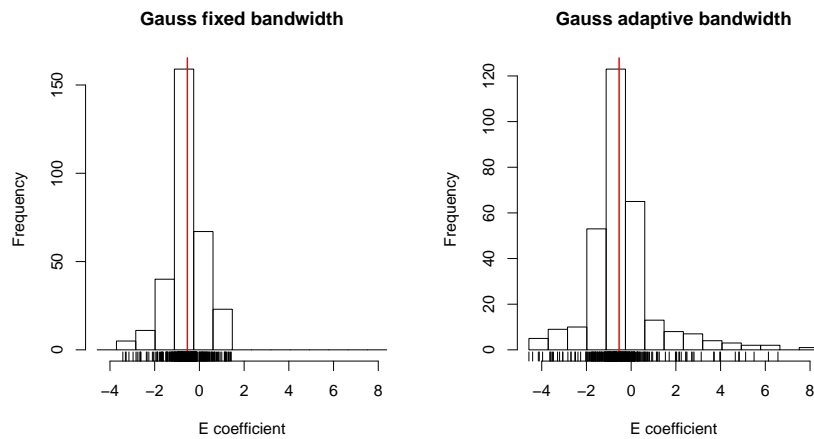
```
xlim.N <- range(gwr.f$N, gwr.a$N)
breaks.N <- seq(xlim.N[1], xlim.N[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$N, main="Gauss fixed bandwidth",
     xlab="N coefficient", xlim=xlim.N, freq=F, breaks=breaks.N)
abline(v=gwr.global["N"], col="red", lwd=2)
rug(gwr.f$N)
hist(gwr.a$N, main="Gauss adaptive bandwidth",
     xlab="N coefficient", xlim=xlim.N, freq=F, breaks=breaks.N)
abline(v=gwr.global["N"], col="red", lwd=2)
rug(gwr.a$N)
par(mfrow=c(1,1))
```



The adaptive bandwidth finds some positive coefficients, i.e., GDD50 increasing northward! This is clearly a local effect.

Third, the East coefficient:

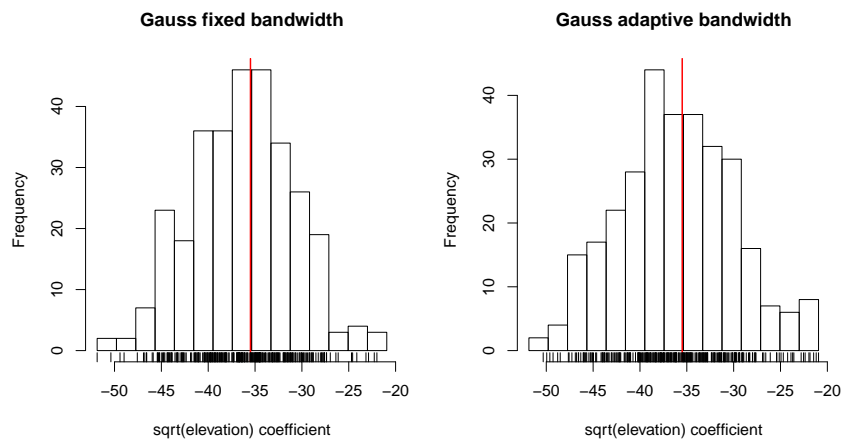
```
xlim.E <- range(gwr.f$E, gwr.a$E)
breaks.E <- seq(xlim.E[1], xlim.E[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$E, main="Gauss fixed bandwidth",
     xlab="E coefficient", xlim=xlim.E, freq=T, breaks=breaks.E)
abline(v=gwr.global["E"], col="red", lwd=1.5)
rug(gwr.f$E)
hist(gwr.a$E, main="Gauss adaptive bandwidth",
     xlab="E coefficient", xlim=xlim.E, freq=T, breaks=breaks.E)
abline(v=gwr.global["E"], col="red", lwd=1.5)
rug(gwr.a$E)
par(mfrow=c(1,1))
```



Again the adaptive bandwidth gives a wider range of coefficients.

Finally, the elevation coefficient:

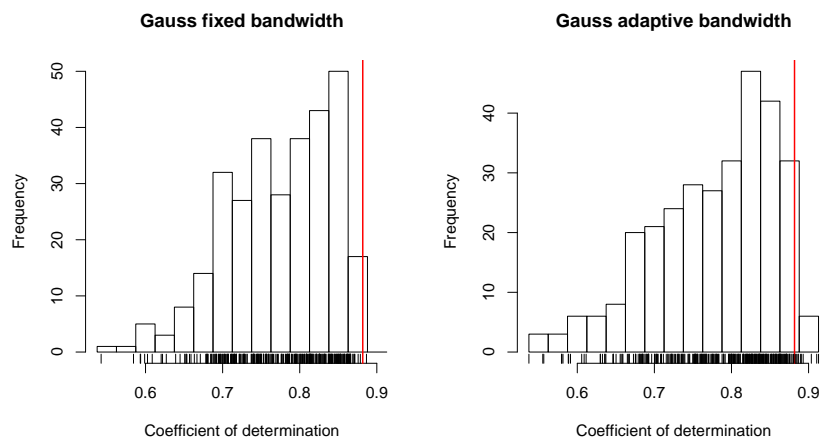
```
xlim.elev <- range(gwr.f$'sqrt(ELEVATION_)', gwr.a$'sqrt(ELEVATION_)')
breaks.elev <- seq(xlim.elev[1], xlim.elev[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$'sqrt(ELEVATION_)', main="Gauss fixed bandwidth",
     xlab="sqrt(elevation) coefficient", xlim=xlim.elev, freq=T,
     breaks=breaks.elev)
abline(v=gwr.global['sqrt(ELEVATION_)'], col="red", lwd=1.5)
rug(gwr.f$'sqrt(ELEVATION_)')
hist(gwr.a$'sqrt(ELEVATION_)', main="Gauss adaptive bandwidth",
     xlab="sqrt(elevation) coefficient", xlim=xlim.elev, freq=T,
     breaks=breaks.elev)
abline(v=gwr.global['sqrt(ELEVATION_)'], col="red", lwd=1.5)
rug(gwr.a$'sqrt(ELEVATION_)')
par(mfrow=c(1,1))
```



There is a wide range of elevation coefficients, although all are negative, as expected. This shows that in some neighbourhoods there is not much effect of elevation. We don't expect this from the basic physics of the adiabatic lapse rate (less dense atmosphere as elevation increases), so this must reflect local effects of e.g., winds mixing air across elevations, more in some neighbourhoods than in others.

TASK 95 : Show the feature-space distribution of the local fits (adjusted R^2). •

```
# global fit
global.r2 <- summary(lm(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
                        data=ne.km@data))$adj.r.squared
xlim.r2 <- range(gwr.f$localR2, gwr.a$localR2)
breaks.r2 <- seq(xlim.r2[1], xlim.r2[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$localR2, main="Gauss fixed bandwidth",
     xlab="Coefficient of determination", xlim=xlim.r2, freq=T, breaks=breaks.r2)
abline(v=global.r2, col="red", lwd=1.5)
rug(gwr.f$localR2)
hist(gwr.a$localR2, main="Gauss adaptive bandwidth",
     xlab="Coefficient of determination", xlim=xlim.r2, freq=T, breaks=breaks.r2)
abline(v=global.r2, col="red", lwd=1.5)
rug(gwr.a$localR2)
par(mfrow=c(1,1))
```

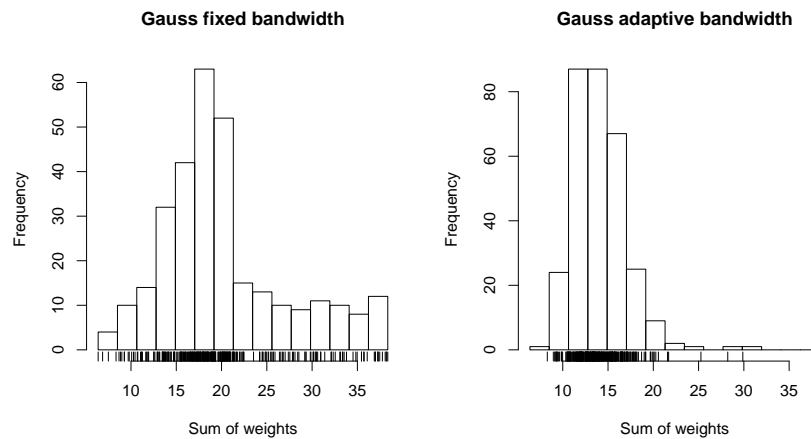


The local R^2 are almost all less than the global R^2 , this makes sense because of the smaller dataset in each fit.

An interesting diagnostic is the sum of weights. This is the sum of the number of points that participated in the fit multiplied by their weights.

TASK 96 : Show the feature-space distribution of the sum of weights. •

```
xlim.sum.w <- range(gwr.f$sum.w, gwr.a$sum.w)
breaks.sum.w <- seq(xlim.sum.w[1], xlim.sum.w[2], length=16)
par(mfrow=c(1,2))
hist(gwr.f$sum.w, main="Gauss fixed bandwidth",
     xlab="Sum of weights", xlim=xlim.sum.w, freq=T, breaks=breaks.sum.w)
rug(gwr.f$sum.w)
hist(gwr.a$sum.w, main="Gauss adaptive bandwidth",
     xlab="Sum of weights", xlim=xlim.sum.w, freq=T, breaks=breaks.sum.w)
rug(gwr.a$sum.w)
par(mfrow=c(1,1))
```



The fixed kernel uses more information, as measured by the sum of weights; it is also more variable. The adaptive bandwidth aims at equal numbers of points; if their distribution is fairly uniform the sum of weights is similar.

11.6.2 Statistical significance

With the global model we tested for statistical significance with the F-statistic computed from the ratio of the total and residual sums of squares, corrected for degrees of freedom. We can do the same for the GWR model, but this requires extra computation, specifying the `se.fit` and `hatmatrix` arguments. The “hat” matrix is used to compute normalized residual sums of squares.

TASK 97 : Re-compute the GWR at the observation points, for the adaptive kernel, also computing the standard errors of the coefficients and their effective degrees of freedom. •

```
gwr.with.se <- gwr(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
  data = ne.km,
  adapt = bw.Gauss.a,
  gweight=gwr.Gauss,
  hatmatrix=TRUE,
  se.fit=TRUE)
names(gwr.with.se$SDF@data)

## [1] "sum.w"          "(Intercept)"
## [3] "E"              "N"
## [5] "sqrt(ELEVATION_)" "(Intercept)_se"
## [7] "E_se"           "N_se"
## [9] "sqrt(ELEVATION_)_se" "gwr.e"
## [11] "pred"           "pred.se"
## [13] "localR2"        "(Intercept)_se_EDF"
## [15] "E_se_EDF"       "N_se_EDF"
## [17] "sqrt(ELEVATION_)_se_EDF" "pred.se"
```

We see that `gwr` has computed standard errors for all coefficients; these vary considerably among points, as expected because we already saw that the adjusted R^2 varies.

TASK 98 : Test for significance, using a two-tailed t-test with the effective degrees of freedom used to compute the standard error. •

The significant coefficients are those with large absolute t-values, i.e., low standard errors relative to their values.

```
summary(t.n <- gwr.with.se$SDF$N / gwr.with.se$SDF$N_se_EDF)

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
## -8.4179 -3.7729 -2.1408 -2.4625 -0.6526  1.9715

summary(t.e <- gwr.with.se$SDF$E / gwr.with.se$SDF$E_se_EDF)

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
## -4.05157 -1.85139 -0.99541 -0.79138  0.01712  3.63942

summary(t.h <- gwr.with.se$SDF$'sqrt(ELEVATION_)' /
         gwr.with.se$SDF$'sqrt(ELEVATION_)_se_EDF')

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
## -13.008  -8.330  -6.661  -6.802  -4.888  -1.940
```

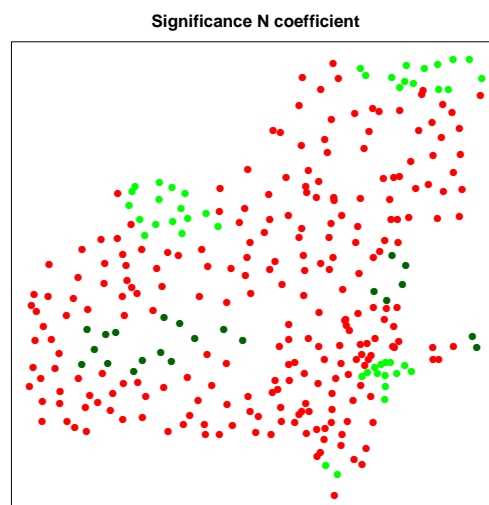
Note: There is no point in computing significance for the intercept, since they are all greatly different from 0. They are all significant, they set the local mean, which is not zero GDD50.

TASK 99 : Show a map of the non-significant coefficients in red, the negative significant in dark green, and the positive in light green. •

We select $(-4, 4)$ as the limits.

First, the N coefficient:

```
sig.map <- SpatialPointsDataFrame(ne.km, data.frame(t.n, t.e, t.h))
sig.colours <- c("darkgreen", "red", "green")
breaks <- c(min(t.n), -4, 4, max(t.n))
spplot(sig.map, col.regions=sig.colours, breaks=breaks,
       zcol="t.n",
       main="Significance N coefficient",
       auto.key=FALSE)
```



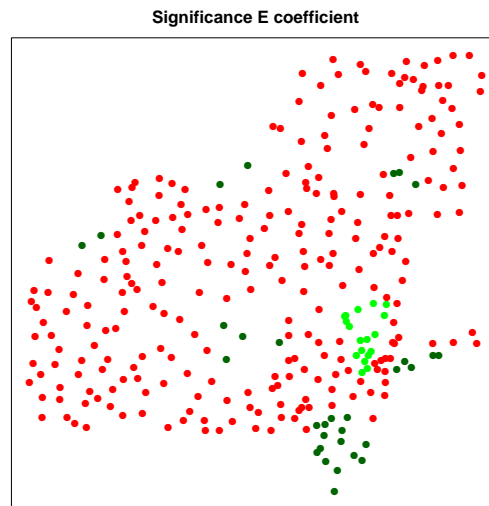
With the local neighbourhood many N coefficients are not significant -

this is because of the limited size of the area. The negative relation with N (expected by physical principles) is only found in central PA and as artefacts at a few points in the east. A positive relation is found near Lake Ontario and northern VT.

Does this mean that Northing does not affect GDD50? Obviously not, because the global model (as well as physical principles) shows it to be highly significant. What we see here is its effect in a local neighbourhood on the order of 70–100 km radius. In hilly areas the elevation may be more important locally.

Second, the E coefficient:

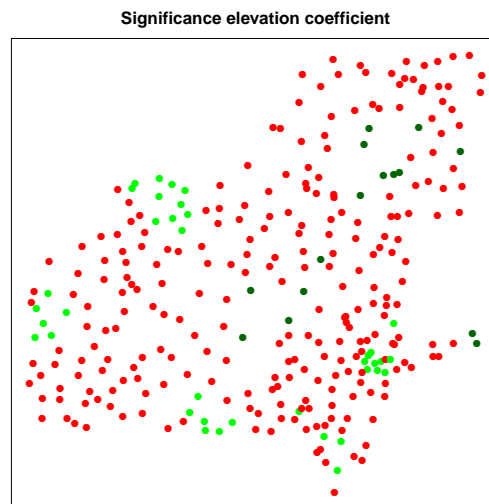
```
breaks <- c(min(t.e), -4, 4, max(t.e))
sppplot(sig.map, col.regions=sig.colours, breaks=breaks,
        zcol="t.e",
        main="Significance E coefficient",
        auto.key=FALSE)
```



This is important in the lower Hudson valley (positive relation, i.e., warmer towards the East), and the southeast corner of Lake Ontario and southern NJ (negative relation, i.e., warmer towards the West), These likely reflect the influence of water bodies.

Finally, the elevation coefficient:

```
breaks <- c(min(t.h), -3, 3, max(t.h))
sppplot(sig.map, col.regions=sig.colours, breaks=breaks,
        zcol="t.h",
        main="Significance elevation coefficient",
        auto.key=FALSE)
```



In low-relief areas (NYC, Rochester) the elevation is a significant positive predictor of GDD50, likely because elevation is correlated with distance from the water body.

TASK 100 : Determine if any locations have *no* significant coefficients, other than the intercept. •

```
table(is.sig <- (abs(t.n) > 3) | (abs(t.e) > 3) | (abs(t.h) > 3))

##
## FALSE  TRUE
##      6    299

ne.m@data[which(!is.sig), "STATION_NA"]

## [1] "CANOE BROOK"          "NEWARK INTL AP"
## [3] "AVON"                  "BATAVIA"
## [5] "ROCHESTER MONROE CO AP" "GREENVILLE 2 NE"
```

At this location, only the weighted average values of other stations are used. This is inverse distance with a Gaussian kernel.

11.7 Predictions

Although GWR is not primarily intended for mapping, it can be used for that purpose. A local model is fit at *each* prediction point, and then used to predict from the known points, as weighted by the bandwidth.

“Please also be aware that using GWR for prediction has no good basis anywhere for anything - and the standard errors should not be given any credibility. This is not what GWR is for at all.” – Roger Bivand³⁴

Note also that GWR does *not* account for local spatial correlation of the residuals, i.e., it can not estimate coefficients by GLS, because residuals in each window are too few to model their spatial structure.

³⁴ <https://grokbase.com/t/r/r-sig-geo/106bgxksy4/gwr-analysis>

11.7.1 Predictions at known points

TASK 101 : Compare the predictions at the known points with those made by other methods. •

We can not compare with kriging methods, because they predict exactly at a known point³⁵. So we compare with the two global regressions computed previously: OLS and GLS trend surfaces, with the same predictors.

First, compute and compare the RMSE:

```
n <- length(residuals(m.ols))
tmp.rmse <- data.frame(cbind(sqrt(sum(residuals(m.ols)^2)/n),
  sqrt(sum(residuals(m.gls)^2)/n),
  sqrt(sum((gwr.Gauss.f$SDF$pred - ne.m$ANN_GDD50)^2)/n),
  sqrt(sum((gwr.Gauss.a$SDF$pred - ne.m$ANN_GDD50)^2)/n)))
row.names(tmp.rmse) <- "RMSE"

tmp <- data.frame(round(
  cbind(summary(residuals(m.ols)),
    summary(residuals(m.gls)),
    summary(gwr.f$gwr.e),
    summary(gwr.a$gwr.e)), 1))
tmp <- rbind(tmp.rmse, tmp)
names(tmp) <- c("OLS", "GLS", "GWR-fixed", "GWR-adaptive")
print(tmp)

##              OLS              GLS GWR-fixed GWR-adaptive
## RMSE      210.5514  210.9136  165.2229    155.9703
## Min.     -532.8000 -514.0000 -396.4000   -384.2000
## 1st Qu.  -153.2000 -149.6000 -107.8000   -101.1000
## Median    -7.8000  -4.0000  -7.5000     -7.1000
## Mean         0.0000   8.5000  -2.0000     -6.3000
## 3rd Qu.   155.4000 160.9000  95.4000     88.6000
## Max.      641.8000 653.7000 514.3000    467.9000

rm(tmp, tmp.rmse)
```

GWR clearly has a narrower IQR and lower RMSE at the fitted points. So the local regression has found structure that was ignored in the global regression.

11.7.2 Predictions and model coefficients over the grid

Here we want to predict over the grid to map the four-state area. We have the raster DEM covering the study area, `dem.ne.m`. However, GWR should not be extrapolated, since it depends on the local points. Therefore we clip this to the four-state study area. We have the map of the four states in object `state.ne.m`.

TASK 102 : Clip the grid to the four-state area. •

```
require(raster)
dem.ne4.m <- raster::mask(dem.ne.m, state.ne.m)
```

TASK 103 : Convert to a `SpatialPointsDataFrame` and make a version

³⁵ Although, we could compare their leave-one-out cross-validation (LOOCV) results.

with the distances in km, to make the regression coefficients easier to read, and to match the GWR points.

```
dem.ne4.m.spdf <- as(dem.ne4.m, "SpatialPointsDataFrame")
## match the coordinate names with the calibration points
coordnames(dem.ne4.m.spdf)

## [1] "x" "y"

dem.ne4.km.spdf <- spTransform(dem.ne4.m.spdf, p4str.km)
## remove negative elevations so sqrt() is real
names(dem.ne4.km.spdf) <- "ELEVATION_"
dem.ne4.km.spdf[(dem.ne4.km.spdf$ELEVATION_ < 0), "ELEVATION_"] <- 0
## add the coordinates to the data frame
dem.ne4.km.spdf$E <- coordinates(dem.ne4.km.spdf)[,1]
dem.ne4.km.spdf$N <- coordinates(dem.ne4.km.spdf)[,2]
## match the names with the calibration points
summary(dem.ne4.km.spdf)

## Object of class SpatialPointsDataFrame
## Coordinates:
##      min      max
## x -384.2743 354.0257
## y -393.4504 288.0856
## Is projected: TRUE
## proj4string :
## [+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76
## +ellps=WGS84 +units=km]
## Number of points: 22649
## Data attributes:
##      ELEVATION_      E      N
## Min.   : 0.0      Min. :-384.27 Min. :-393.45
## 1st Qu.: 602.5    1st Qu.: -160.02 1st Qu.: -189.73
## Median :1173.9    Median : 12.48  Median : -67.50
## Mean   :1125.1    Mean  : -12.93  Mean  : -53.66
## 3rd Qu.:1587.4    3rd Qu.: 129.78 3rd Qu.: 69.55
## Max.   :3992.5    Max.   : 354.03  Max.   : 288.09
```

TASK 104: Predict over the grid, and save the results in the DEM object.

This calls the `gwr` function, with the `predict` argument set to `TRUE`.

```
gwr.Gauss.f.pred <- gwr(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
  data=ne.km,
  predictions=TRUE, # for prediction
  fit.points=dem.ne4.km.spdf,
  bandwidth=bw.Gauss.f,
  gweight=gwr.Gauss)
gwr.Gauss.a.pred <- gwr(ANN_GDD50 ~ E + N + sqrt(ELEVATION_),
  data=ne.km,
  predictions=TRUE, # for prediction
  fit.points=dem.ne4.km.spdf,
  adapt=bw.Gauss.a,
  gweight=gwr.Gauss)
summary(gwr.Gauss.f.pred$SDF$pred)

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##      860.3 1883.4 2203.7 2268.6 2609.8 4006.5

summary(gwr.Gauss.a.pred$SDF$pred)

##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##      860.4 1884.6 2206.7 2270.9 2611.2 4012.1

dem.ne4.km.spdf$pred.f <- gwr.Gauss.f.pred$SDF$pred
dem.ne4.km.spdf$pred.a <- gwr.Gauss.a.pred$SDF$pred
```

TASK 105 : Also save the fitted coefficients at each prediction grid cell in the DEM object; we will map these in §11.7.3 , below. •

```
dem.ne4.km.spdf$b0.f <- gwr.Gauss.f.pred$SDF$(Intercept)"
dem.ne4.km.spdf$b0.a <- gwr.Gauss.a.pred$SDF$(Intercept)"
dem.ne4.km.spdf$N.f <- gwr.Gauss.f.pred$SDF$N
dem.ne4.km.spdf$N.a <- gwr.Gauss.a.pred$SDF$N
dem.ne4.km.spdf$E.f <- gwr.Gauss.f.pred$SDF$E
dem.ne4.km.spdf$E.a <- gwr.Gauss.a.pred$SDF$E
dem.ne4.km.spdf$H.f <- gwr.Gauss.f.pred$SDF$"sqrt(ELEVATION_)"
dem.ne4.km.spdf$H.a <- gwr.Gauss.a.pred$SDF$"sqrt(ELEVATION_)"
```

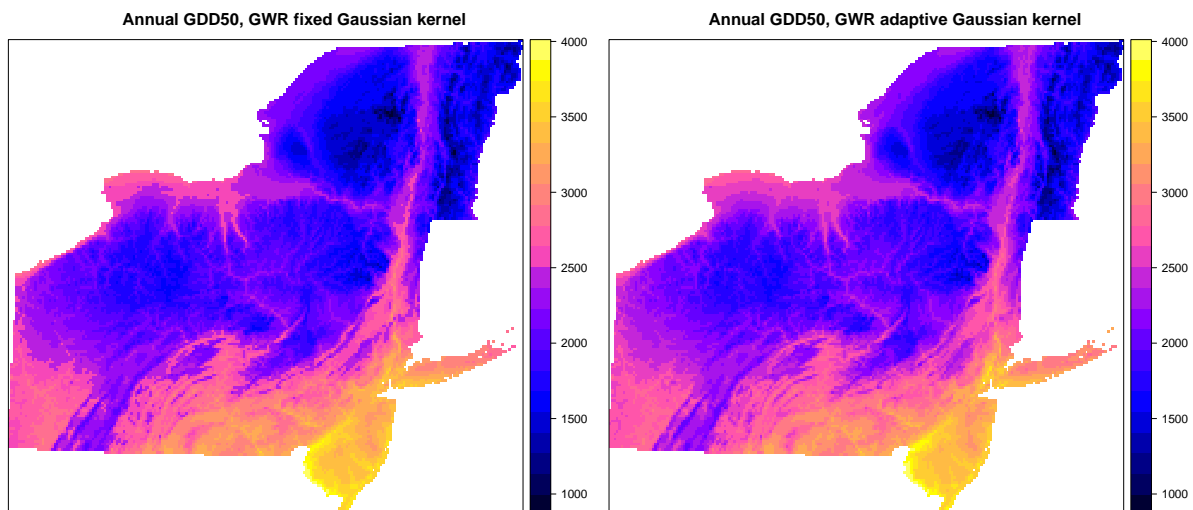
TASK 106 : Make a gridded version for display of results. •

```
dem.ne4.km.grid <- dem.ne4.km.spdf
# convert to a gridded object
gridded(dem.ne4.km.grid) <- TRUE; fullgrid(dem.ne4.km.grid) <- TRUE
```

TASK 107 : Display the prediction maps. •

```
pred.range <- range(na.omit(dem.ne4.km.grid$pred.f),
                    na.omit(dem.ne4.km.grid$pred.a))
cuts <- seq(pred.range[1], pred.range[2], length=24)
```

```
p1 <- spplot(dem.ne4.km.grid, zcol="pred.f",
             main="Annual GDD50, GWR fixed Gaussian kernel",
             at=cuts, key.space="right", col.regions=bpy.colors(24))
p2 <- spplot(dem.ne4.km.grid, zcol="pred.a",
             main="Annual GDD50, GWR adaptive Gaussian kernel",
             at=cuts, key.space="right")
print(p1, split=c(1,1,2,1), more=T); print(p2, split=c(2,1,2,1), more=F)
```



Challenge: Compare the prediction map with others produced by locally-adaptive methods, e.g., OK or KED.

Challenge: Recompute the GWR with just an intercept. This is a weighted moving average. Determine the bandwidth and prediction errors. How much worse is this than using the full GWR?

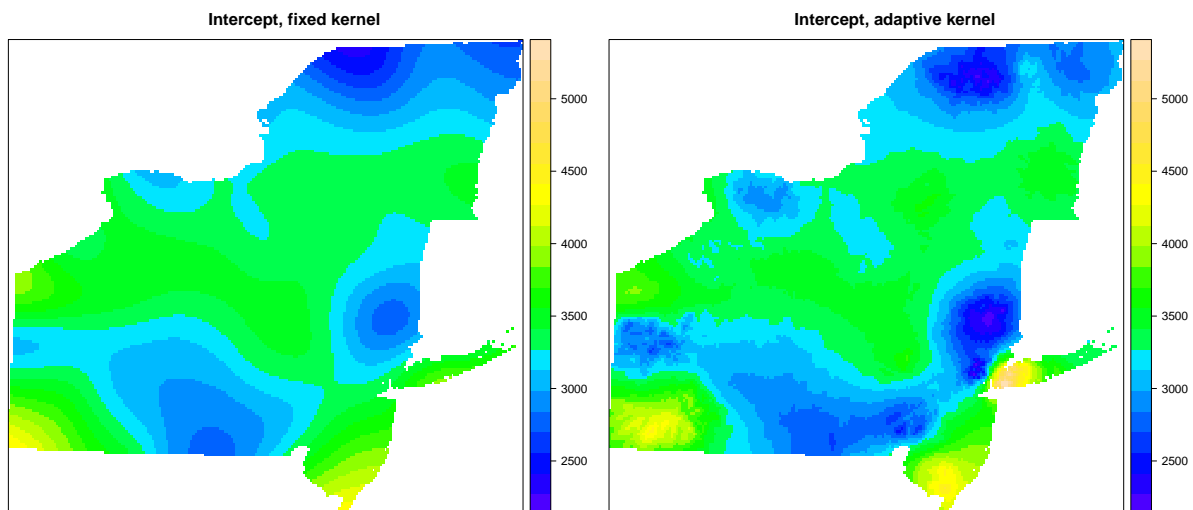
11.7.3 Spatial distribution of coefficients

The main interest in GWR is the spatial heterogeneity of the regression model. We can visualize this by mapping the coefficient values over the prediction grid.

TASK 108: Display the coefficient maps: (1) intercept. •

```
pred.range <- range(na.omit(dem.ne4.km.grid$b0.f),  
                    na.omit(dem.ne4.km.grid$b0.a))  
cuts <- seq(pred.range[1], pred.range[2], length=24)
```

```
p1 <- spplot(dem.ne4.km.grid, zcol="b0.f",  
             main="Intercept, fixed kernel",  
             at=cuts, key.space="right", col.regions=topo.colors(24))  
p2 <- spplot(dem.ne4.km.grid, zcol="b0.a",  
             main="Intercept, adaptive kernel",  
             at=cuts, key.space="right", col.regions=topo.colors(24))  
print(p1, split=c(1,1,2,1), more=T); print(p2, split=c(2,1,2,1), more=F)
```



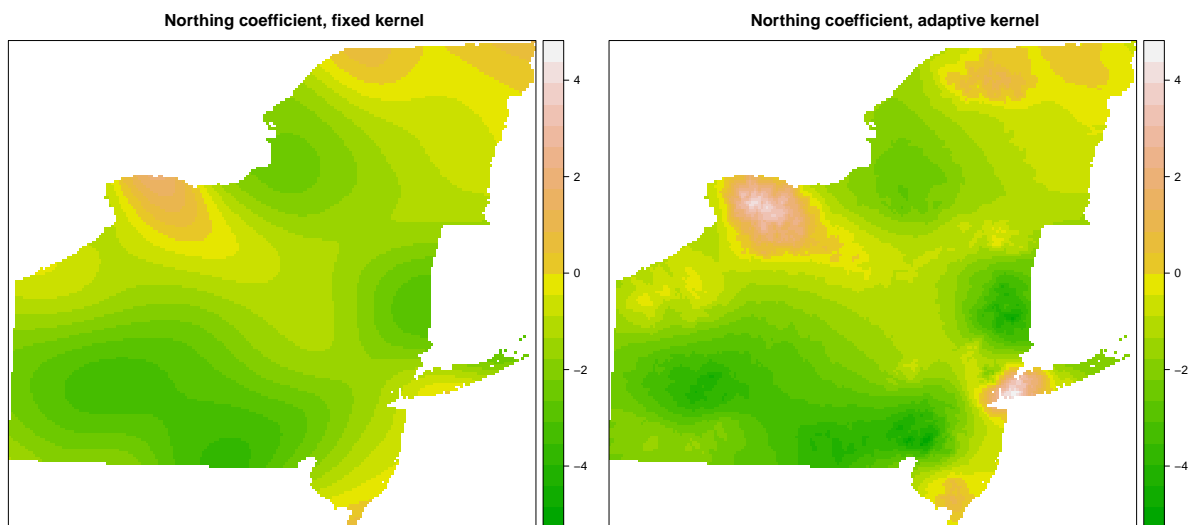
This shows the overall pattern of the GDD, adjusted for regression in a “local” neighbourhood. The adaptive kernel has more extreme values, both high (yellow, near NYC) and low (dark blue, Dutchess County).

This map shows clearly that most of the local prediction is in the intercept, i.e., a moving window of the station values, without considering the other coefficients. Below we will see that the range of these is much narrower.

TASK 109 : Display the coefficient maps: (2) Northing:

```
pred.range <- range(na.omit(dem.ne4.km.grid$N.f),  
                    na.omit(dem.ne4.km.grid$N.a))  
cuts <- seq(pred.range[1], pred.range[2], length=24)
```

```
p1 <- spplot(dem.ne4.km.grid, zcol="N.f",  
             main="Northing coefficient, fixed kernel",  
             at=cuts, key.space="right", col.regions=terrain.colors(24))  
p2 <- spplot(dem.ne4.km.grid, zcol="N.a",  
             main="Northing coefficient, adaptive kernel",  
             at=cuts, key.space="right", col.regions=terrain.colors(24))  
print(p1, split=c(1,1,2,1), more=T); print(p2, split=c(2,1,2,1), more=F)
```

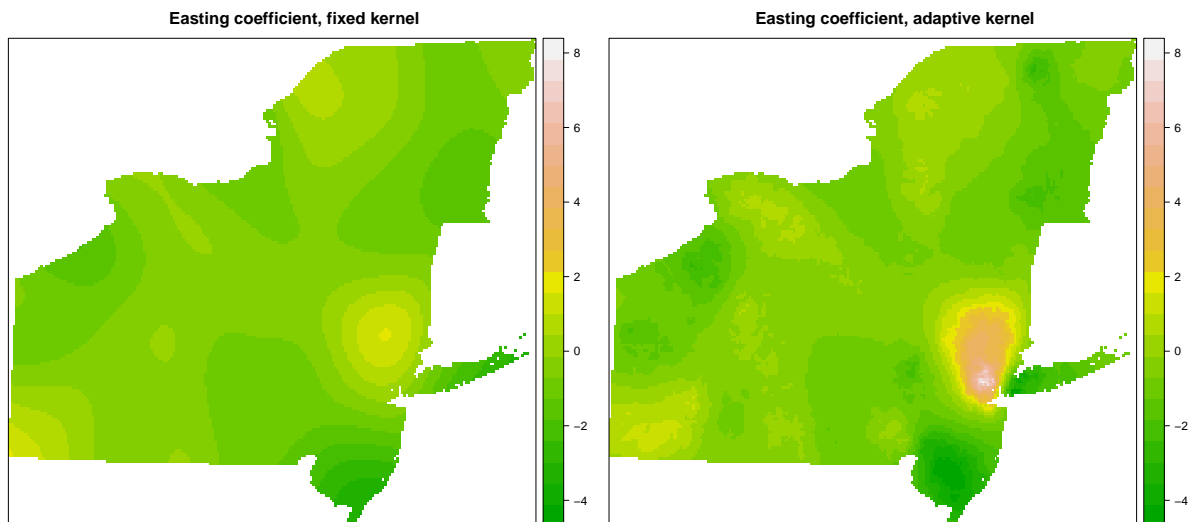


The North coefficient is much more negative, i.e., has more effect, in southern PA than in the Great Lakes margin and northern NY/VT. In some local areas the Northing coefficient is positive, which does not agree with the overall model. Again, the adaptive coefficient has more extreme values.

TASK 110 : Display the coefficient maps: (3) Easting:

```
pred.range <- range(na.omit(dem.ne4.km.grid$E.f),  
                    na.omit(dem.ne4.km.grid$E.a))  
cuts <- seq(pred.range[1], pred.range[2], length=24)
```

```
p1 <- spplot(dem.ne4.km.grid, zcol="E.f",
  main="Easting coefficient, fixed kernel",
  at=cuts, key.space="right", col.regions=terrain.colors(24))
p2 <- spplot(dem.ne4.km.grid, zcol="E.a",
  main="Easting coefficient, adaptive kernel",
  at=cuts, key.space="right", col.regions=terrain.colors(24))
print(p1, split=c(1,1,2,1), more=T); print(p2, split=c(2,1,2,1), more=F)
```

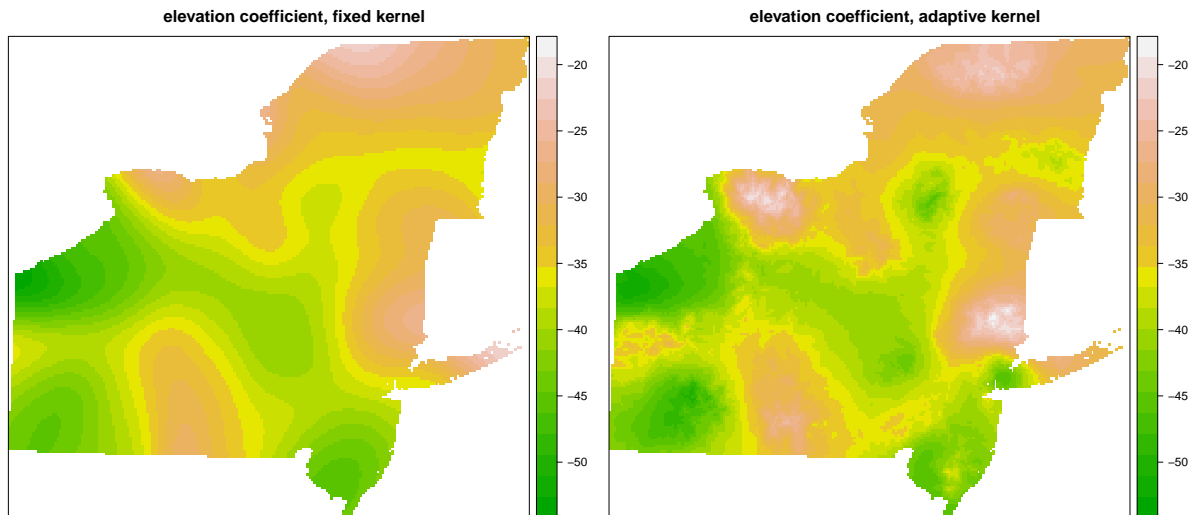


This coefficient has some effect in the lower Hudson Valley/ NYC (positive towards the East) and southern NJ (negative towards the East), as an adjustment to the other coefficients.

TASK 111 : Display the coefficient maps: (4) Square root of elevation: •

```
pred.range <- range(na.omit(dem.ne4.km.grid$H.f),
  na.omit(dem.ne4.km.grid$H.a))
cuts <- seq(pred.range[1], pred.range[2], length=24)
```

```
p1 <- spplot(dem.ne4.km.grid, zcol="H.f",
  main="elevation coefficient, fixed kernel",
  at=cuts, key.space="right", col.regions=terrain.colors(24))
p2 <- spplot(dem.ne4.km.grid, zcol="H.a",
  main="elevation coefficient, adaptive kernel",
  at=cuts, key.space="right", col.regions=terrain.colors(24))
print(p1, split=c(1,1,2,1), more=T); print(p2, split=c(2,1,2,1), more=F)
```



Higher elevations always lead to lower GDD50; this effect is greatest in windows with high relief. We see a much more pronounced effect of elevation at the Lake Erie margin, the western Allegheny plateau, and southern NJ. The coefficient is much less in the Tug Hill plateau and Adirondacks.

12 Data-driven models

A data-driven model is an alternative to linear modelling. It makes no assumptions about linearity; rather, it uses a set of **regression trees**. These partition the feature space of predictors into a set of rectangles, i.e., defined by limits of each predictor in feature space. These rectangles each then have a simple prediction model, in the simplest case just a constant, which is a single predicted value of the response variable for all combinations of predictor variables in that feature-space rectangle. The **advantages** of this approach are:

1. no assumption that the functional form is the same throughout the range of the predictors;
2. over-fitting can be avoided by specifying large enough rectangles; their optimum size can be calculated by cost-complexity pruning.

This is a high-variance, low-bias method. This means that different fits of the model with different subsets of the data may result in quite different

fitted models, but the predictions not be systematically different from the expected values.

A **disadvantage** of this approach is that it can not not extrapolate outside of its range of calibration. But that might equally be considered an advantage³⁶.

Note: Hastie et al. [11, §9.2] give a thorough explanation of a tree-based regression method known as CART (“Classification and Regression Trees”) [2]; these are implemented in R by the `rpart` “Recursive Partitioning” package. A simplified explanation of the same material is given in James et al. [14, §8.1].

12.1 Regression trees

We first start with the simplest data-driven approach: the **regression tree**. This replaces a regression equation, as developed in the previous section, with a **decision tree** based only on optimal splitting of the response variable by the predictors.

12.1.1 Fitting a regression tree model

The procedure is as follows:

1. We first specify the response variable and the possible predictors.
2. We specify a calibration (“training”) dataset, as for the linear model.
3. The algorithm then looks for one predictor variable that “best” splits the data into two groups, and the value of that predictor on which to split. Intuitively, “best” refers to the maximum reduction in sum of within-group sums of squares in the response variable, compared to its overall sum of squares with no split; this is the same measure as used in Analysis of Variance (ANOVA); symbolically the reduction is $SS_T - (SS_L + SS_R)$, where L, R represent the “left” and “right” branches of the tree.
4. Following the split, this process is applied separately to both subgroups; this continues recursively until the subgroups either reach a minimum size (specified by us) or until no improvement can be made; that is the sum of the within-groups sum of squares can not be further reduced.
5. This model is then **pruned**, i.e., some branches are combined, by cross-validation, to avoid over-fitting.

Regression trees are implemented in the `rpart` function of the `rpart` package.

TASK 112 : Load the `rpart` package. •

³⁶ “Elke naadeel heeft zijn voordeel” – Johann Cruijff

```
library(rpart)
```

TASK 113 : Compute a regression tree for the response GDD50, from the predictors N, E, and ELEVATION_. Note that there is no need to transform any predictor. •

The `rpart` function has several control options: (1) the minimum number of observations which can be considered for a split (using the `minsplit` argument); and (2) the minimum value of a complexity parameter (using the `cp` argument). This corresponds to the improvement in R^2 with each split. A small complexity parameter (close to 0) grows a larger tree, which may be over-fitting.

We set these to allow maximum splitting: split even if only two cases, using the `minsplit` optional argument. Also specify a small complexity parameter with the `cp` optional argument: keep splitting until there is less than 0.3% improvement in (unadjusted) R^2 .

The model formulation is the same as for linear modelling: specify the predictand (dependent variable) on the left side of the `~` formula operator and the predictors on the right side, separated by the `+` formula operator. Note there is no interaction possible in tree models; the predictors are considered separately when determining which to use for a split.

```
m.rt <- rpart(ANN_GDD50 ~ N + E + ELEVATION_,
              data=ne.df,
              minsplit=2,
              cp=0.003)
```

TASK 114 : Display the fitted tree in text form. •

```
print(m.rt)

## n= 305
##
## node), split, n, deviance, yval
##      * denotes terminal node
##
##  1) root 305 100137700.0 2517.518
##    2) N>=-155967.2 183 27412290.0 2182.536
##      4) ELEVATION_>=1275 45 3699167.0 1778.200
##        8) N>=119836.3 12 921737.0 1493.500
##          16) ELEVATION_>=2945 1 0.0 795.000 *
##          17) ELEVATION_< 2945 11 389480.0 1557.000
##            34) ELEVATION_>=1355 10 68720.0 1503.000 *
##            35) ELEVATION_< 1355 1 0.0 2097.000 *
##        9) N< 119836.3 33 1451091.0 1881.727
##          18) E>=-237967.7 31 1059064.0 1854.065
##            36) ELEVATION_>=1721 11 354595.6 1715.818 *
##            37) ELEVATION_< 1721 20 378607.8 1930.100 *
##          19) E< -237967.7 2 612.5 2310.500 *
##    5) ELEVATION_< 1275 138 13957200.0 2314.384
##      10) N>=82678.87 40 2854040.0 2076.125
##        20) ELEVATION_>=520 19 568665.7 1878.263 *
##        21) ELEVATION_< 520 21 868542.6 2255.143
##          42) N>=142190.2 16 195405.8 2173.375 *
##          43) N< 142190.2 5 223838.8 2516.800 *
```



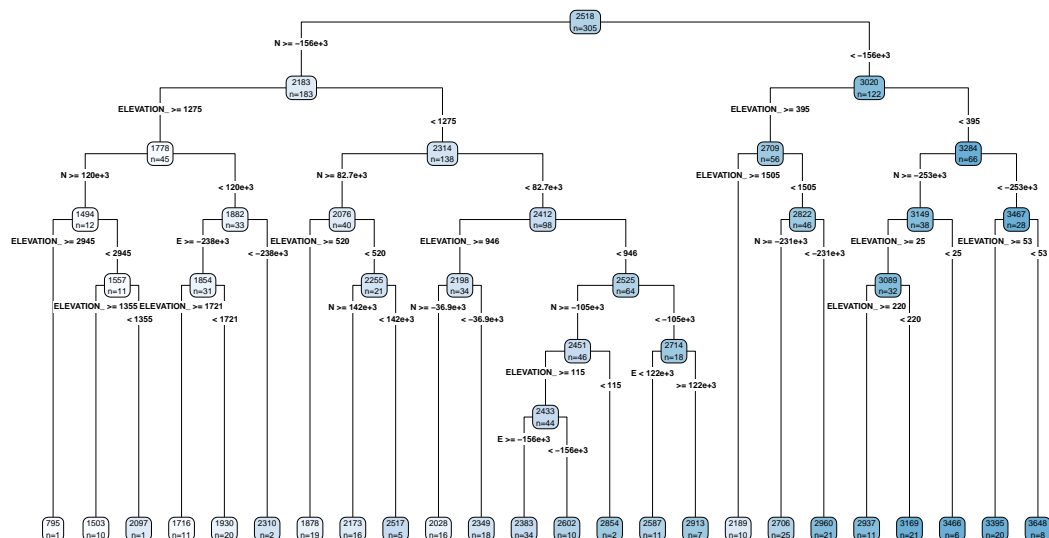
```
##      11) N< 82678.87 98    7905649.0 2411.633
##      22) ELEVATION_>=946 34    1936208.0 2197.941
##      44) N>=-36918.74 16    529437.0 2028.250 *
##      45) N< -36918.74 18    536519.1 2348.778 *
##      23) ELEVATION_< 946 64    3592056.0 2525.156
##      46) N>=-104679.8 46    1681030.0 2451.326
##      92) ELEVATION_>=115 44    1336795.0 2433.023
##      184) E>=-156122.8 34    803000.0 2383.382 *
##      185) E< -156122.8 10    165155.6 2601.800 *
##      93) ELEVATION_< 115 2    5202.0 2854.000 *
##      47) N< -104679.8 18    1019503.0 2713.833
##      94) E< 122469.2 11    451684.0 2587.000 *
##      95) E>=122469.2 7    112794.9 2913.143 *
##      3) N< -155967.2 122    21387870.0 3019.992
##      6) ELEVATION_>=395 56    6252615.0 2709.107
##      12) ELEVATION_>=1505 10    202742.9 2189.100 *
##      13) ELEVATION_< 1505 46    2757956.0 2822.152
##      26) N>=-231009.3 25    1153309.0 2706.320 *
##      27) N< -231009.3 21    869901.0 2960.048 *
##      7) ELEVATION_< 395 66    5130590.0 3283.773
##      14) N>=-252628.4 38    2433245.0 3148.526
##      28) ELEVATION_>=25 32    1337325.0 3088.969
##      56) ELEVATION_>=220 11    161948.7 2936.545 *
##      57) ELEVATION_< 220 21    785949.2 3168.810 *
##      29) ELEVATION_< 25 6    377040.8 3466.167 *
##      15) N< -252628.4 28    1058940.0 3467.321
##      30) ELEVATION_>=52.5 20    420733.0 3394.950 *
##      31) ELEVATION_< 52.5 8    271573.5 3648.250 *
```

TASK 115 : Plot the regression tree. •

The tree is graphed with the `rpart.plot` function of the `rpart.plot` package.

The `rpart.plot` function has several options to control the display; we choose to show the values of the response variable at the interior nodes as well as at the leaves, and to show the number of observations in each split and leaf. The information is the same as given with a printout of the model object, but easier to visualize.

```
library(rpart.plot)
rpart.plot(m.rt, digits=3, type=4, extra=1)
```



Q29 : What is the first (root) splitting variable? At what value is the split? What is the mean value of GDD50 of the whole dataset, and of the two branches? How many observations in each branch? [Jump to A29](#) •

Although there is no model with coefficients, we can still see which predictor variables had the most influence on the tree.

TASK 116 : Display the variable importance as a proportion. •

This is the proportion of the variance explained by the tree which is due to each variable.

```
x <- m.rt$variable.importance
data.frame(variableImportance = 100 * x / sum(x))

##           variableImportance
## N                        48.90337
## ELEVATION_                38.60227
## E                        12.49436
```

Q30 : Which variables are most important?

[Jump to A30](#) •

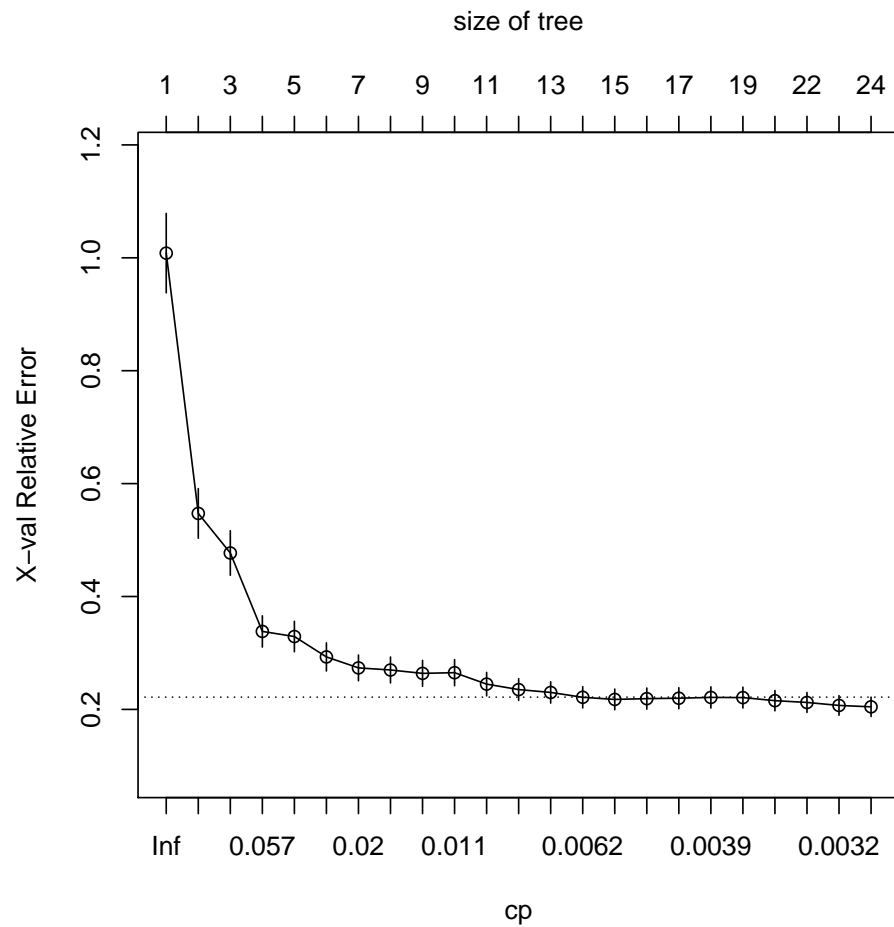
We now examine the reduction in fitting and cross-validation error with the `printcp` “print the complexity parameter” function.

TASK 117 : Print and plot the error rate vs. the complexity parameter and tree size. •

```
printcp(m.rt)

##
## Regression tree:
## rpart(formula = ANN_GDD50 ~ N + E + ELEVATION_, data = ne.df,
##       minsplit = 2, cp = 0.003)
##
## Variables actually used in tree construction:
## [1] E           ELEVATION_ N
##
## Root node error: 100137734/305 = 328320
##
## n= 305
##
##           CP nsplit rel error  xerror   xstd
## 1  0.5126697    0  1.000000  1.00827  0.070337
## 2  0.0999090    1  0.487330  0.54720  0.044003
## 3  0.0974250    2  0.387421  0.47709  0.039461
## 4  0.0328739    3  0.289996  0.33815  0.027674
## 5  0.0319311    4  0.257122  0.32911  0.026982
## 6  0.0237411    5  0.225191  0.29299  0.024940
## 7  0.0163615    6  0.201450  0.27357  0.022835
## 8  0.0141488    7  0.185089  0.26986  0.023029
## 9  0.0132452    8  0.170940  0.26385  0.023009
## 10 0.0089030    9  0.157695  0.26509  0.023120
## 11 0.0086905   10  0.148792  0.24484  0.020820
## 12 0.0073373   11  0.140101  0.23515  0.019238
## 13 0.0071789   12  0.132764  0.23006  0.018855
## 14 0.0053152   13  0.125585  0.22150  0.018797
## 15 0.0045440   14  0.120270  0.21770  0.018348
## 16 0.0044868   15  0.115726  0.21912  0.018778
## 17 0.0039088   16  0.111239  0.21982  0.018672
## 18 0.0038889   17  0.107330  0.22125  0.018801
## 19 0.0036613   18  0.103441  0.22096  0.018469
## 20 0.0035335   19  0.099780  0.21551  0.018011
## 21 0.0032541   21  0.092713  0.21233  0.017656
## 22 0.0032032   22  0.089459  0.20707  0.017381
## 23 0.0030000   23  0.086256  0.20459  0.017167

plotcp(m.rt)
```



Note: Your results will likely be different. This is because the cross-validation makes a *random* split of the full dataset into a number of subsets for model building and evaluation. Each run gives a different random split.

The xerror field in the summary shows the **cross-validation error**; that is, applying the model to the original data split K -fold, each time excluding some observations. If the model is over-fitted, the cross-validation error increases; note that the fitting error, given in the **error** field, always decreases. By default, the split is 10-fold; this can be modified by the `control` argument to the `rpart` function.³⁷

Q31 : Does this model appear to be overfit? Why or why not? What appears to be the optimum complexity parameter to avoid over-fitting?

Jump to A31 •

TASK 118 : Prune the tree back to complexity level estimated from the previous answer. •

³⁷ See the help for `rpart.control`.

We do this with the `prune` function, specifying the `cp` “complexity parameter” argument.

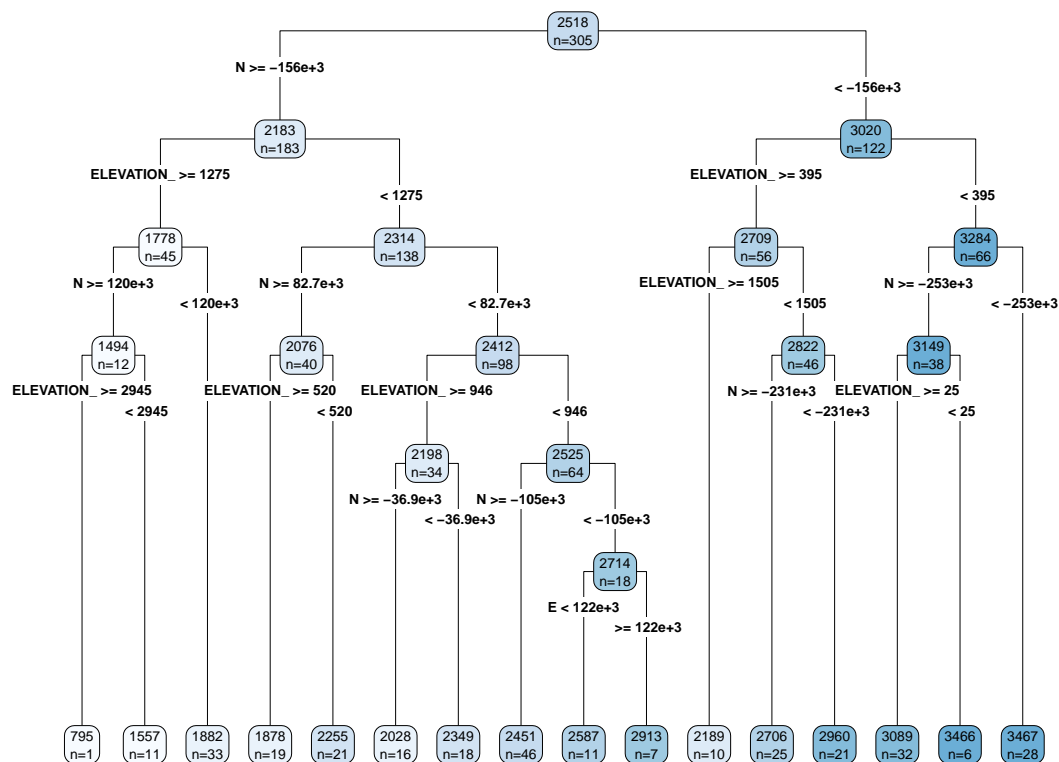
```
(m.rt.p <- prune(m.rt, cp=0.0045))

## n= 305
##
## node), split, n, deviance, yval
##      * denotes terminal node
##
## 1) root 305 100137700.0 2517.518
##    2) N>=-155967.2 183 27412290.0 2182.536
##      4) ELEVATION_>=1275 45 3699167.0 1778.200
##        8) N>=119836.3 12 921737.0 1493.500
##          16) ELEVATION_>=2945 1 0.0 795.000 *
##            17) ELEVATION_< 2945 11 389480.0 1557.000 *
##          9) N< 119836.3 33 1451091.0 1881.727 *
##        5) ELEVATION_< 1275 138 13957200.0 2314.384
##          10) N>=82678.87 40 2854040.0 2076.125
##            20) ELEVATION_>=520 19 568665.7 1878.263 *
##              21) ELEVATION_< 520 21 868542.6 2255.143 *
##            11) N< 82678.87 98 7905649.0 2411.633
##              22) ELEVATION_>=946 34 1936208.0 2197.941
##                44) N>=-36918.74 16 529437.0 2028.250 *
##                45) N< -36918.74 18 536519.1 2348.778 *
##              23) ELEVATION_< 946 64 3592056.0 2525.156
##                46) N>=-104679.8 46 1681030.0 2451.326 *
##                47) N< -104679.8 18 1019503.0 2713.833
##                  94) E< 122469.2 11 451684.0 2587.000 *
##                  95) E>=122469.2 7 112794.9 2913.143 *
##            3) N< -155967.2 122 21387870.0 3019.992
##              6) ELEVATION_>=395 56 6252615.0 2709.107
##                12) ELEVATION_>=1505 10 202742.9 2189.100 *
##                13) ELEVATION_< 1505 46 2757956.0 2822.152
##                  26) N>=-231009.3 25 1153309.0 2706.320 *
##                  27) N< -231009.3 21 869901.0 2960.048 *
##                7) ELEVATION_< 395 66 5130590.0 3283.773
##                  14) N>=-252628.4 38 2433245.0 3148.526
##                    28) ELEVATION_>=25 32 1337325.0 3088.969 *
##                    29) ELEVATION_< 25 6 377040.8 3466.167 *
##                  15) N< -252628.4 28 1058940.0 3467.321 *
```

Task 119: Plot the pruned regression tree.

•

```
rpart.plot(m.rt.p, digits=3, type=4, extra=1)
```



Q32 : How does this tree differ from the original regression tree? [Jump to A32](#) •

We can compare the fitted vs. actual values.

TASK 120 : Use the pruned regression tree to predict at the calibration points. •

We do this with the `predict` method applied to a `rpart` object; this automatically calls function `predict.rpart`. The points to predict and

the values of the predictor variables at those points are supplied in a dataframe as argument `newdata`. We count the number of predicted values with the `unique` function; there is only one value per “box” in the feature space defined by the predictor variables.

```
summary(p.rt.p <- predict(m.rt.p, newdata=ne.df))
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      795   2028   2451   2518   2960   3467
```

TASK 121 : Count the unique predicted values. •

```
length(unique(p.rt.p))
```

```
## [1] 16
```

TASK 122 : Compute the residuals: (actual - predicted) and the RMSE. Summarize them and present as a histogram. •

```
summary(residuals.rt.p <- ne.df$ANN_GDD50 - p.rt.p)
```

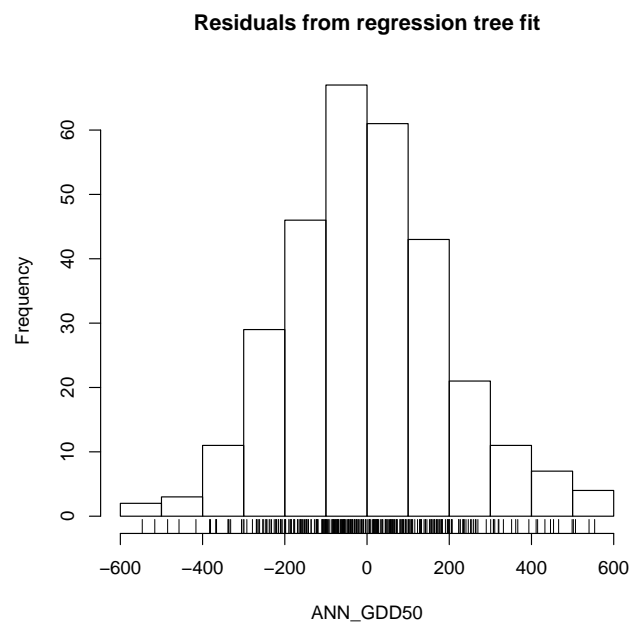
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -546.73 -122.73  -10.14    0.00  127.90   553.68
```

```
hist(residuals.rt.p, main="Residuals from regression tree fit",
     xlab="ANN_GDD50")
```

```
rug(residuals.rt.p)
```

```
sqrt(mean(residuals.rt.p^2)/length(residuals.rt.p))
```

```
## [1] 11.16128
```



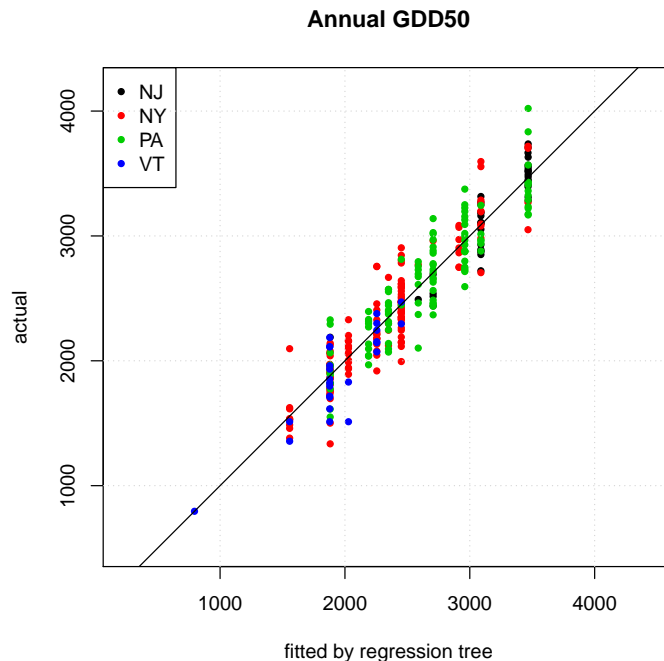
TASK 123 : Plot the actual vs. fitted values. •

```
plot(ne.df$ANN_GDD50 ~ p.rt.p, asp=1, pch=20,
     xlab="fitted by regression tree", ylab="actual",
```

```

xlim=c(500,4200), ylim=c(500,4200),
col=ne.df$STATE,
main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)

```



Q33 : How many unique values are predicted by the pruned regression tree? How close is the fit to the actual values, compared to the OLS or GLS models? Explain. [Jump to A33](#) •

12.1.2 Regression tree prediction over the study area

TASK 124 : Predict over the grid with the regression tree model, add the results to the dataframe, and summarize them. •

TASK 125 : Display the regression tree surface. •

```

dem.ne.m.df$pred.rt <- predict(m.rt.p, newdata=dem.ne.m.df)
summary(dem.ne.m.df$pred.rt)

```

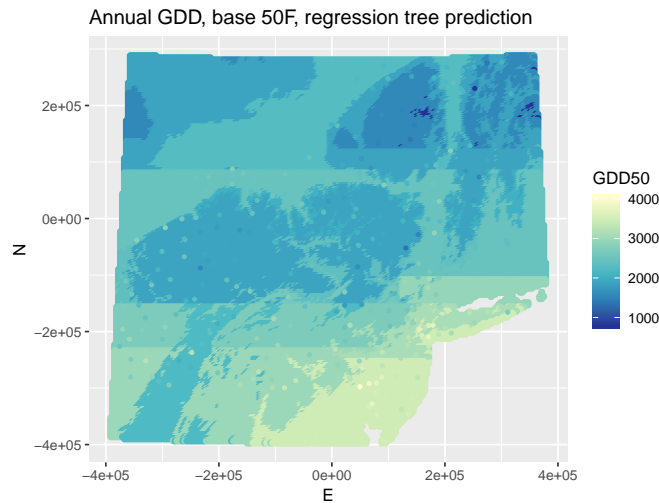
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	795	1882	2255	2311	2706	3467

We also add the point set in the same colour scheme; if the point is visible it means the residual (lack of fit) is large.

```

display.prediction.map("pred.rt",
  "Annual GDD, base 50F, regression tree prediction",
  "GDD50")

```

Q34 : What is the spatial pattern of the regression tree prediction?
 Explain why. Jump to A34 •

12.2 Random forests

A problem with regression trees is that a small change in the sample set, for example a missing or erroneous observation, can radically change the tree. Pruning the tree to avoid overfitting is also somewhat subjective. Also, correlated predictors can appear in the tree as surrogates for each other, depending on the details of the calibration set; this makes interpretation difficult. Finally, trees do not give smooth predictions; rather they only predict a single value in each “box”. To solve these problems, a method known as “random forests” was developed; see Hastie et al. [11, §15] (advanced) or James et al. [14, §8] (simplified). There are a lot of details to this method, but the basic idea is straightforward:

1. Build a **large number of regression trees**, independently, using *different* sets of observations; these are built by *sampling with replacement* from the actual observations; this is a technique known as *bagging* or *bootstrap aggregation*.
2. Save all these trees; when predicting, use all of them and **average their predictions**.
3. In addition we can summarize the whole set of trees to see how different they are, thus how robust is the final model.
4. Also, for each tree we can use observations that were not used to construct it for true **validation**, called *out-of-bag* validation. This gives a good idea of the true prediction error.

The first step may seem suspicious. The underlying idea is that what we observe is the best sample we have of reality; if we sampled again we'd expect to get similar values. So we simulate re-sampling by sampling from this same set, *with replacement*, to get a sample of the same size but a different sample. If this idea bothers you, read Efron and Gong [5], Shalizi [24] or Hastie et al. [11, §8.2].

So, how well does the random forest method work for a spatially-distributed variable?

We know that by definition a linear model, whether OLS or GLS fit, will vary smoothly with the predictors. That is, if these predictors vary smoothly in space the map of the linear model predictions will also look smooth. In this example the predictor Northing is by definition smooth, and the predictor elevation often changes smoothly, so we can expect a smooth prediction map. However, random forests are not linear. They are based on an ensemble of regression trees, with no requirement for smooth cut points.

Another attractive feature of random forests is that there is no need for predictor variable selection to avoid colinearity. Since the predictors to compare are randomly chosen at each split, it is possible for “minor” predictors which would not be included in a parametric approach to contribute to some of the trees, and thus to the ensemble prediction. In this example Easting was not used in the regression models, but should be used here. The relative importance of the predictors is reported by the RF function.

12.2.1 Fitting a Random Forest model

TASK 126 : Fit a random forest model of GDD50 based on all three possible covariates: elevation, Northing and Easting. •

Since this does not assume linearity we can use the untransformed station elevation.

The `randomForest` function of the `randomForest` package fits this model. This model requires two parameters:

1. the number of trees in the forest, optional argument `ntree`, default value 500;
2. the number of predictors to compare at each split, optional argument `mtry`, default value is 1/3 of the predictors.

Here we accept the default for the number of predictors, which in this case will be one of the three predictors, but require more than the default number of trees

We also specify the optional `importance` argument as `TRUE` to see two measures of predictor importance.

```
library(randomForest)
m.rf <- randomForest(ANN_GDD50 ~ ELEVATION_ + N + E,
                      data=ne.df, ntree=1200,
                      importance=TRUE)
randomForest::importance(m.rf)[,1]

## ELEVATION_      N      E
## 69.19142    82.44646  58.05839

randomForest::importance(m.rf)[,2]/max(randomForest::importance(m.rf)[,2])

## ELEVATION_      N      E
## 0.9342154    1.0000000  0.3543592
```

The first measure of predictor importance is the increase in mean squared error (MSE) in the out-of-bag (OOB) validations if the predictor is removed from the set. From the help text³⁸:

For each tree, the prediction error on the out-of-bag portion of the data is recorded [i.e., MSE]. Then the same is done after permuting each predictor variable. The difference between the two are then averaged over all trees, and normalized by the standard deviation of the differences.

Q35 : Which predictor is most important? How much would assigning elevations randomly to stations increase the MSE? [Jump to A35](#) •

The second measure is the increase in node purity, which we normalized to the maximum over all variables.

[This] ... measure is the total decrease in node impurities from splitting on the variable, averaged over all trees ... measured by residual sum of squares.

As in the first measure, Northing is the most important, but elevation is almost as important. Easting is much less important.

TASK 127 : Plot the actual vs. fitted values from the random forest model. Compute the mean error (bias, ME) and the root mean squared error (RMSE). •

The `predict` function gives predicted values based on a model. To predict back at the calibration points, the `newdata` argument to the `predict` function must specify this point set.

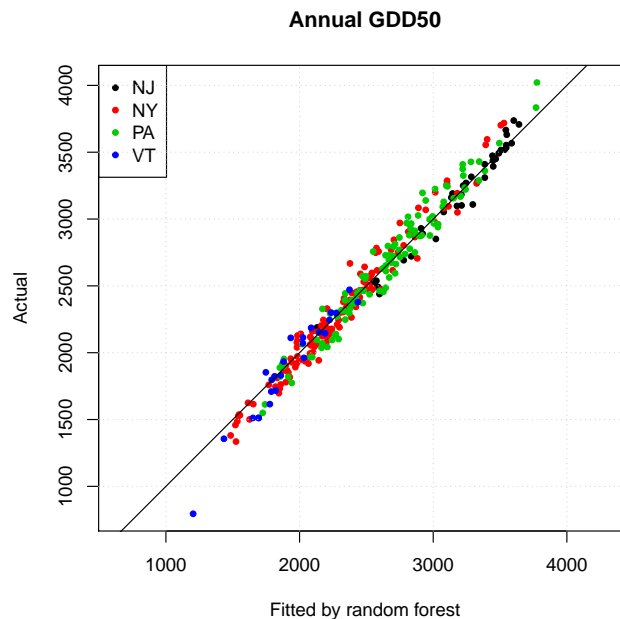
Note: Each run of the random forest will produce different fits, so your graph may not look exactly like this one.

```
summary(predict(m.rf, newdata=ne.m))

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    1204    2125    2470    2518    2895    3776
```

³⁸ ?importance

```
plot(ne.m$ANN_GDD50 ~ predict(m.rf, newdata=ne.m),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)
```



```
summary(rf.resid <- ne.m$ANN_GDD50 - predict(m.rf, newdata=ne.m))

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -408.6057 -64.6425  -8.2177  -0.2273  55.5859  290.7625

(rf.me <- mean(rf.resid))

## [1] -0.2273227

(rf.rmse <- sqrt(mean(rf.resid^2)/length(rf.resid)))

## [1] 5.675547
```

The fits are very close. However, this is not a good measure of the prediction accuracy. For that we use out-of-bag (OOB) **cross-validation**.

OOB predictions are also computed with the `predict` function, but with the `newdata` argument omitted. Each point is predicted as the average of the RF predictions for those regression trees where that point was *not* included in the “bag”.

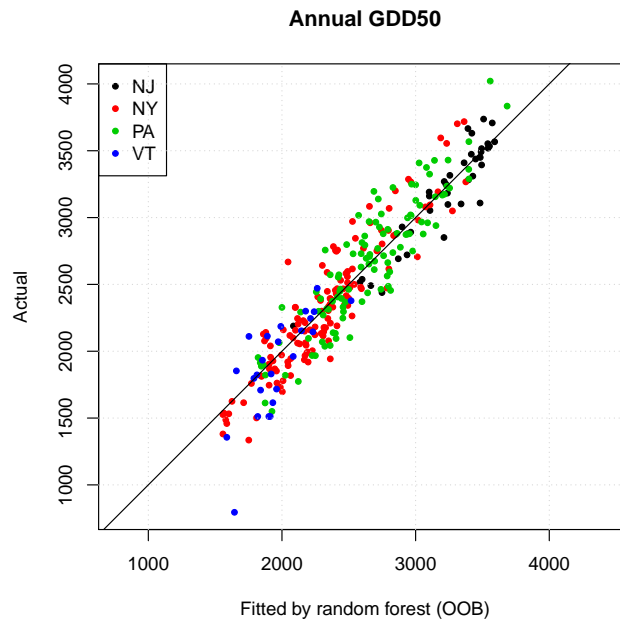
TASK 128 : Plot the actual vs. out-of-bag validation values from the random forest model. Compute the mean error (bias, ME) and the root mean squared error (RMSE). •

```
summary(predict(m.rf))

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
```

```
##      1558      2145      2457      2520      2870      3685

plot(ne.m$ANN_GDD50 ~ predict(m.rf),
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest (OOB)",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid()
abline(0,1)
```



```
summary(rf.oob.resid <- ne.m$ANN_GDD50 - predict(m.rf))

##      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
## -849.554 -128.116    -11.079     -2.063    108.612    622.538

(rf.oob.me <- mean(rf.oob.resid))

## [1] -2.06279

(rf.oob.rmse <- sqrt(mean(rf.oob.resid^2)/length(rf.oob.resid)))

## [1] 11.58191
```

This is a realistic view of the prediction accuracy. It is much more scattered around the 1:1 line than the calibration fit. Note in particular the high negative residual for Mt. Mansfield (VT); this is the blue point at the bottom of the plot.

TASK 129: Compare the ME and RMSE for the model fits and out-of-bag residuals. •

```
(rf.oob.me/rf.me)

## [1] 9.074284

(rf.oob.rmse/rf.rmse)
```

```
## [1] 2.040668
```

Q36 : *How does the OOB RMSE compare to the RMSE from the model fits at the training points? Which is more realistic as a measure of prediction accuracy?* [Jump to A36](#) •

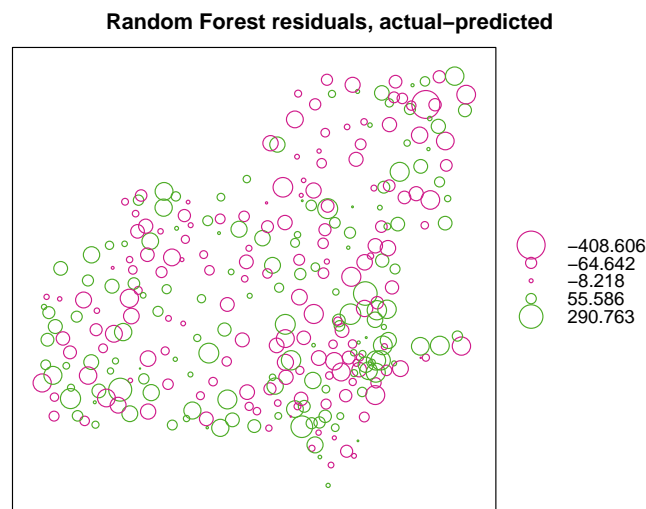
TASK 130 : Extract the RF model residuals, summarize them, and show them as a bubble plot. •

The model residuals are computed from the fits at each point and the actual values; for this we use the fits, not the out-of-bag estimates.

```
ne.m$rf.resid <- (ne.df$ANN_GDD50 - predict(m.rf, newdata=ne.m))
summary(ne.m$rf.resid)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -408.6057 -64.6425  -8.2177  -0.2273  55.5859  290.7625

bubble(ne.m, zcol="rf.resid", pch=1,
       main="Random Forest residuals, actual-predicted")
```



The residuals are fairly similar in range to the linear model. There are some very poorly-fit points. Notice that the mean residual is *not* zero, as in the least-squares fit of the linear model.

TASK 131 : Compute the RF out-of-bag residuals and compare them to the RF model residuals. •

```
ne.m$rf.resid.oob <- (ne.df$ANN_GDD50 - predict(m.rf))
summary(ne.m$rf.resid.oob)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
```

```
## -849.554 -128.116 -11.079 -2.063 108.612 622.538

summary(ne.m$rf.resid)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -408.6057 -64.6425  -8.2177  -0.2273  55.5859  290.7625

summary(ne.m$rf.resid.oob/ne.m$rf.resid)

##      Min. 1st Qu.  Median     Mean 3rd Qu.    Max.
## -59.676  1.843   2.027   1.767  2.235   15.311
```

As seen in the 1:1 plots, the out-of-bag residuals are much larger: The mean ratio of the two is 1.77.

TASK 132 : List the eight worst-fit points, sorted by their absolute residuals, along with the residual from the GLS fit linear model. •

```
(ix <- order(abs(ne.m$rf.resid), decreasing=TRUE)[1:8])

## [1] 293 118 220 232 115 199 107 278

ne.m@data[ix,c("STATE","STATION_NA","ELEVATION_",
               "gls.resid", "rf.resid")]

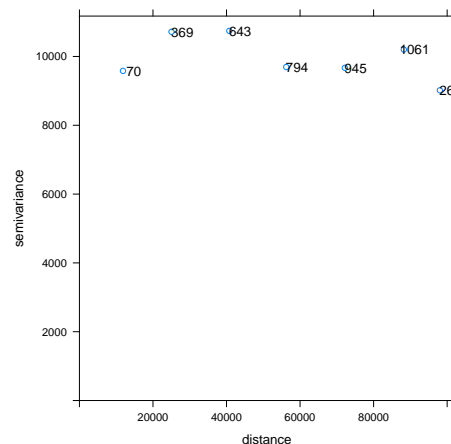
##      STATE      STATION_NA ELEVATION_ gls.resid rf.resid
## 4716    VT      MOUNT MANSFIELD    3950 -13.33249 -408.6057
## 3092    NY      MOHONK LAKE    1245 439.87959 290.7625
## 3830    PA      JOHNSTOWN    1214 653.74848 276.2142
## 3842    PA      MARCUS HOOK     10 411.03014 244.9626
## 3089    NY      MIDDLETOWN 2 NW    700 413.05043 220.6852
## 3809    PA      DONORA 1 SW     762 437.40026 211.0848
## 3081    NY      LITTLE FALLS MILL ST 360 330.36606 207.5222
## 3888    PA WILKES BRE SCTN AP AVOCA 930 288.51349 206.1419
```

The RF almost always comes closer to these worse-fit points, because it is only using fairly similar points, in terms of the predictors (N, E, elevation) to predict, and does not try to fit a regional trend, which must consider all the calibration points.

However, in most runs of the random forest model the lowest actual GDD value (Mt. Mansfield in VT) is badly under-predicted. This is because it is so unlike any other point, because it is so much higher than the others. In the linear model this has strong leverage (highest elevation by far, well to the North) and is thus closely fit. The other poorly-predicted stations are also “unusual” in their covariate-space neighbourhood.

TASK 133 : Compute and display an empirical variogram of the RF model residuals. •

```
v.rf <- variogram(rf.resid ~ 1, locations=ne.m, cutoff=100000, width=16000)
plot(v.rf, pl=T)
```



Q37 : *Do the residuals have spatial structure? This depends again on each run of the model; your results will look different from the ones presented here.* [Jump to A37](#) •

TASK 134 : Compare the statistics of the regression model residuals. •

```
summary(ne.m$ols.resid)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## -532.764 -153.151   -7.838    0.000   155.435   641.758
```

```
summary(ne.m$rf.resid)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## -408.6057 -64.6425   -8.2177   -0.2273   55.5859   290.7625
```

```
summary(ne.m$gls.resid)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## -514.041 -149.624   -3.980    8.496   160.856   653.749
```

```
sd(ne.m$ols.resid)
```

```
## [1] 210.8974
```

```
sd(ne.m$rf.resid)
```

```
## [1] 99.2818
```

```
sd(ne.m$gls.resid)
```

```
## [1] 211.0888
```

The statistics are similar; there is no clear “winner” Note that the GLS and RF models are not unbiased – the mean residual is not zero.

12.2.2 Random Forest prediction over the study area

We can use the RF model to predict over the study area, as we did in §6 for the OLS and GLS models.

TASK 135 : Predict over the grid with the RF model, add the results to the dataframe, and summarize them. •

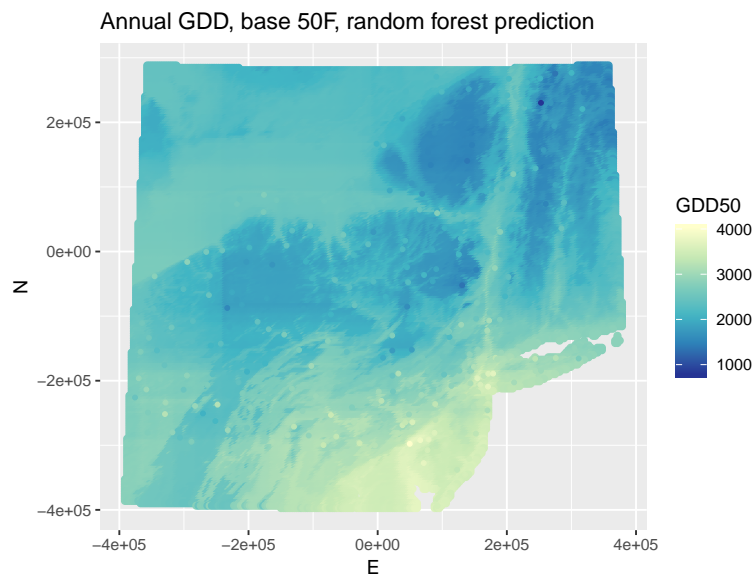
```
dem.ne.m.df$pred.rf <- predict(m.rf, newdata=dem.ne.m.df)
summary(dem.ne.m.df$pred.rf)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1356	2024	2320	2363	2625	3764

TASK 136 : Display the RF surface. •

We also add the point set in the same colour scheme; if the point is visible it means the residual (lack of fit) is large.

```
display.prediction.map("pred.rf",
  "Annual GDD, base 50F, random forest prediction",
  "GDD50")
```



Q38 : What is the difference in the spatial pattern between the RF surface and the RK-GLS surface? [Jump to A38](#) •

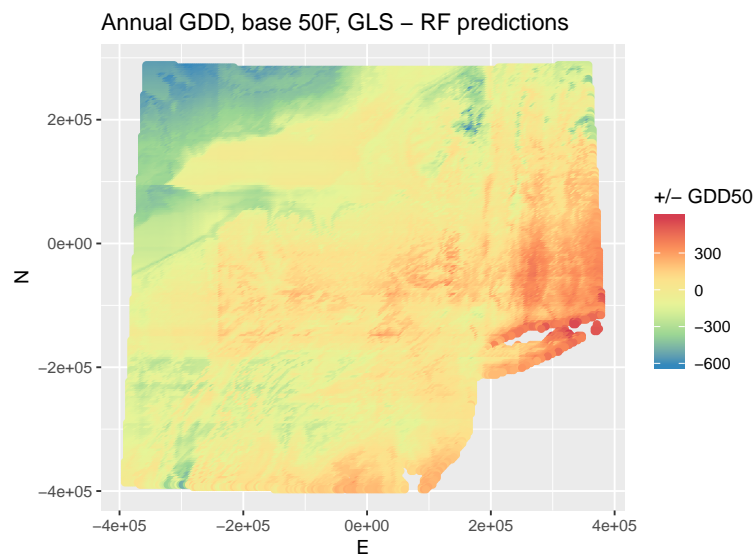
Q39 : Why is there only one predicted value in most of Lake Erie, another in the north-most part of Lake Erie, Lake Ontario, and the adjacent areas of Ontario (Canada)? [Jump to A39](#) •

TASK 137 : Compute the differences between the GLS and RF predictions, add them to the data frame, and display them. •

```
summary(dem.ne.m.df$diff.gls.rf <-
  dem.ne.m.df$pred.gls - dem.ne.m.df$pred.rf)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -608.36 -137.40  -26.59  -44.56   71.69   581.10

display.difference.map("diff.gls.rf",
  "Annual GDD, base 50F, GLS - RF predictions",
  "+/- GDD50")
```



Q40 : *Why are the largest discrepancies in the east, especially in Connecticut and Massachusetts, and in the northwest (Ontario)?* [Jump to A40](#) •

Q41 : *Why is the border of PA with OH on the west visible in this difference map?* [Jump to A41](#) •

Q42 : *Why are there positive differences (GLS-RK greater than RF) in the PA mountains and in the Catskills and Taconics in NY? Why does this not occur in the Adirondacks?* [Jump to A42](#) •

Q43 : *What do you conclude about the geographical area of applicability of the random forest prediction?* [Jump to A43](#) •

12.3 Tuning data-driven models

Data-driven models have parameters that control their behaviour and can significantly affect their predictive power.

For example, regression trees (§12.1) can be adjusted by the minimum number of observations which can be considered for a split (using the `minsplit` argument) and the minimum value of a complexity parameter (using the `cp` argument).

Random forests (§12.2) can also be controlled by the minimum number of observations in a terminal node (optional argument `nodesize`), as well as the number of predictors to compare at each split (optional argument `mtry`). In the `randomForest` function these have default values of 5 and 1/3 of the number of predictors, respectively. These can have a large influence on the resulting forest. Too small terminal nodes will result in over-fit trees, too large in poorer fits. Too many predictors tested at each split will not allow less powerful predictors into the forest; too few will result in many poorly-fitted trees.

Note: The number of trees `ntree` also has an influence on the random forest model. Too few trees will cause repeated model fits to be too variable, too many wastes computing time. This parameter is not optimized as such, a large value is used and the graph of out-of-bag RMSE vs. number of trees is examined to select an appropriate value.

So how to decide on these parameter values? A powerful method is to examine a space of possible values, and select the best values by the performance of each model. The performance is often evaluated by repeated **cross-validation**:

1. For each combination of parameters to be optimized:
 - (a) Split the dataset into some disjunct subsets, for example 10, by random sampling.
 - (b) For each subset:
 - i. Fit the model with the selected parameters on all but one of the subsets.
 - ii. Predict at the remaining subset, i.e., the one not used for model building, with the fitted model.
 - iii. Compute the goodness-of-fit statistics of fitting, typically the root mean square error (RMSE) of prediction and the squared correlation coefficient between the actual and fitted values, i.e., R^2 against a 1:1 line.
 - (c) Average the statistics for the disjunct subsets.
2. Search the table of results for the best results: lowest RMSE and highest R^2 .

These values are then the recommended ones to fit a final model on all data.

The caret “Classification And REgression Training” package [15] implements this procedure.

```
library(caret)
```

This package is highly flexible and can be used to optimize 238 kinds of data-driven models, including the ones we use in this tutorial. To see the list of possible models:

```
length(names(getModelInfo()))

## [1] 238

head(names(getModelInfo()), 24)

## [1] "ada"          "AdaBag"       "AdaBoost.M1"
## [4] "adaboost"     "amdai"       "ANFIS"
## [7] "avNNet"      "awnb"        "awtan"
## [10] "bag"         "bagEarth"    "bagEarthGCV"
## [13] "bagFDA"      "bagFDAGCV"   "bam"
## [16] "bartMachine" "bayesglm"    "binda"
## [19] "blackboost"  "blasso"      "blassoAveraged"
## [22] "bridge"     "brnn"        "BstLm"
```

Many more details are given in the caret package on-line book³⁹. The package is highly adaptable, and before applying it in a production environment, carefully read the documentation. Here we just present a simple case.

The principal function of the caret package is `train`, which implements the cross-validation procedure and reports the optimum combination of parameters. To use this, we have to set up the following arguments:

1. a design matrix `x` of predictor values for each observation;
2. a response vector `y` of known values for each observation;
3. the `tuneGrid` argument with the names and range of values for the parameters to be tuned;
4. the `method` argument, which says which data-driven method to optimize;
5. the `trControl` argument, which specifies the optimization criterion, using the `trainControl` function.

TASK 138 : Optimize a random forest model. •

The caret package does not support `randomForest` models, instead it supports random forest models implemented in the more efficient and newer `ranger` package [34].

We create a matrix with all combinations of the two parameters, with the `expand.grid` function. The columns are named the same as the tuning parameters, which we obtain with the `getModelInfo` function:

³⁹ <http://topepo.github.io/caret/index.html>

```
getModelInfo("ranger")$ranger$parameters

##      parameter      class      label
## 1      mtry      numeric #Randomly Selected Predictors
## 2    splitrule character      Splitting Rule
## 3 min.node.size      numeric      Minimal Node Size
```

For the meaning of each of these, see `?ranger` and the explanations in the journal article. Here we only have three predictors, so `mtry` can vary from 1 to 3. The splitting rule is set to "variance", the normal criterion for regression trees: the split is where the between-class variance is maximized and the within-class variance is minimized. The minimum node size `min.node.size` defaults to 5, we try smaller (more complex trees in the forest) and larger (less complex) values.

```
library(ranger)
dim(preds <- ne.df[, c("E", "N", "ELEVATION_")])

## [1] 305    3

length(response <- ne.df[, "ANN_GDD50"])

## [1] 305

system.time(
  ranger.tune <- train(x = preds, y = response, method="ranger",
    tuneGrid = expand.grid(.mtry = 1:3,
      .splitrule = "variance",
      .min.node.size = 1:10),
    trControl = trainControl(method = 'cv'))
)

##      user  system elapsed
## 55.200   2.478  13.157

print(ranger.tune)

## Random Forest
##
## 305 samples
## 3 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 274, 275, 275, 274, 275, 274, ...
## Resampling results across tuning parameters:
##
##      mtry min.node.size RMSE      Rsquared    MAE
## 1      1      1      199.7651  0.8862826  156.1662
## 1      1      2      200.5215  0.8851154  156.3225
## 1      1      3      200.6421  0.8854146  156.2801
## 1      1      4      201.4257  0.8851293  156.8953
## 1      1      5      200.3734  0.8870097  156.2643
## 1      1      6      202.2373  0.8851730  157.4440
## 1      1      7      203.7624  0.8832940  158.9196
## 1      1      8      202.9067  0.8844527  158.6500
## 1      1      9      203.2735  0.8845912  158.8762
## 1      1     10      205.4798  0.8824785  160.3347
## 2      2      1      198.5917  0.8842978  155.7959
## 2      2      2      199.2312  0.8838005  156.2332
## 2      2      3      199.8000  0.8829345  156.5254
## 2      2      4      198.8858  0.8836315  156.4780
## 2      2      5      199.7606  0.8829220  156.5060
## 2      2      6      198.2823  0.8844972  155.7579
## 2      2      7      199.0948  0.8838903  156.1358
## 2      2      8      199.7685  0.8830005  156.6918
```

```
##      2      9      199.6730 0.8832496 156.7499
##      2     10      199.8913 0.8825836 157.1165
##      3      1      202.8704 0.8791472 160.5457
##      3      2      202.8491 0.8788279 159.8376
##      3      3      202.2480 0.8794636 159.1250
##      3      4      202.9873 0.8787474 159.7276
##      3      5      203.4839 0.8782558 160.4607
##      3      6      201.8906 0.8796622 159.6206
##      3      7      202.3331 0.8789335 159.0573
##      3      8      201.9809 0.8793349 158.7097
##      3      9      202.9065 0.8781754 159.7739
##      3     10      202.5687 0.8788200 159.5980
##
## Tuning parameter 'splitrule' was held constant at a value of variance
## RMSE was used to select the optimal model using the smallest value.
## The final values used for the model were mtry = 2, splitrule
## = variance and min.node.size = 6.

names(ranger.tune$result)

## [1] "mtry"          "splitrule"      "min.node.size" "RMSE"
## [5] "Rsquared"      "MAE"           "RMSESD"        "RsquaredSD"
## [9] "MAESD"

ix <- which.min(ranger.tune$result$RMSE)
ranger.tune$result[ix, c(1,3,4)]

##      mtry min.node.size      RMSE
## 16      2              6 198.2823

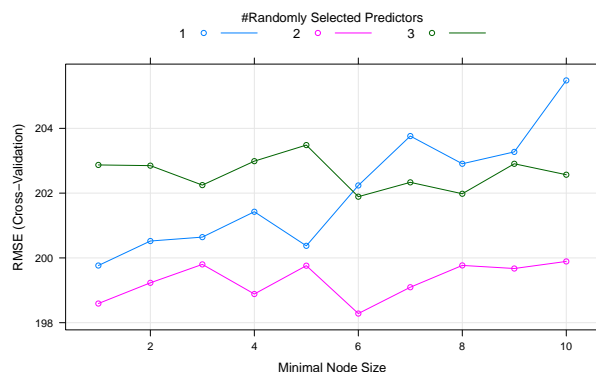
ix <- which.max(ranger.tune$result$Rsquared)
ranger.tune$result[ix, c(1,3,5)]

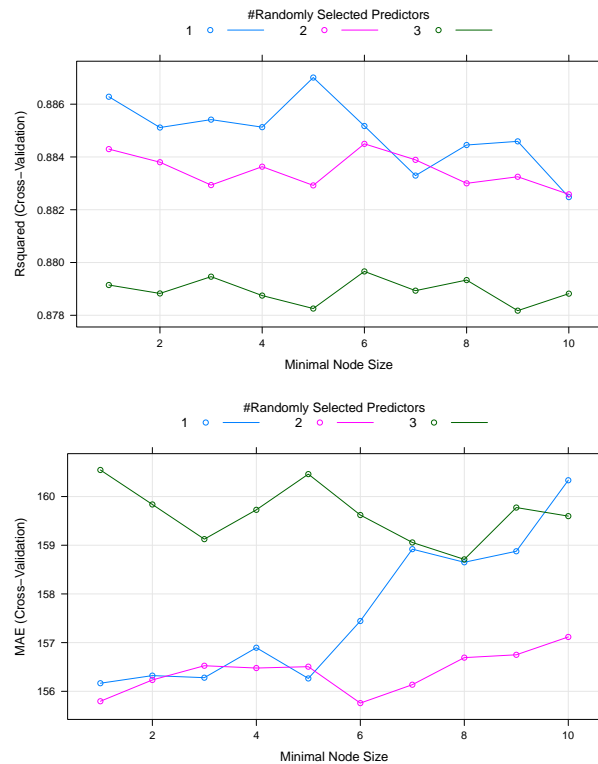
##      mtry min.node.size Rsquared
## 5      1              5 0.8870097

ix <- which.min(ranger.tune$result$MAE)
ranger.tune$result[ix, c(1,3,6)]

##      mtry min.node.size      MAE
## 16      2              6 155.7579

plot.train(ranger.tune, metric="RMSE")
plot.train(ranger.tune, metric="Rsquared")
plot.train(ranger.tune, metric="MAE")
```





In this case the three optimization criteria give different recommendations. Here it's clear that `mtry=2` is optimal in terms of RMSE and MAE. This is different from the default $3/3 = 1$. For both of these `min.node.size=6` is best. However in terms of R^2 `mtry=1` is best, and with node size `min.node.size=5`.

Each run of `train` will give different results, because of the different random splits into test and train sets. Rather than accept the recommendation, look at the graph and select the combination that is almost optimal (lowest RMSE) but less complex.

Once we've selected an optimal combination we fit a final model.

TASK 139: Build an optimal model and display its fit to known points.

The `ranger` function requires a formula argument (as does `randomForest`):

```
(ranger.rf <- ranger(ANN_GDD50 ~ N + E + ELEVATION_, data=ne.df,
                     mtry=2, min.node.size=9))

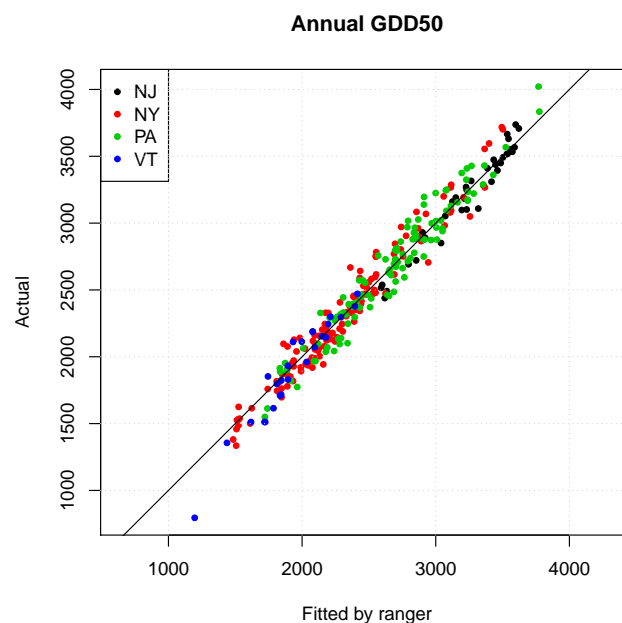
## Ranger result
##
## Call:
## ranger(ANN_GDD50 ~ N + E + ELEVATION_, data = ne.df, mtry = 2, min.node.size = 9)
##
## Type: Regression
## Number of trees: 500
## Sample size: 305
## Number of independent variables: 3
## Mtry: 2
## Target node size: 9
```

```
## Variable importance mode:      none
## Splitrule:                    variance
## OOB prediction error (MSE):    40654.16
## R squared (OOB):              0.8765813
```

```
summary(ranger.fits <- predict(ranger.rf, data=ne.df)$predictions)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1196    2119    2444    2519    2907    3775
```

```
plot(ne.df$ANN_GDD50 ~ ranger.fits,
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by ranger",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
```



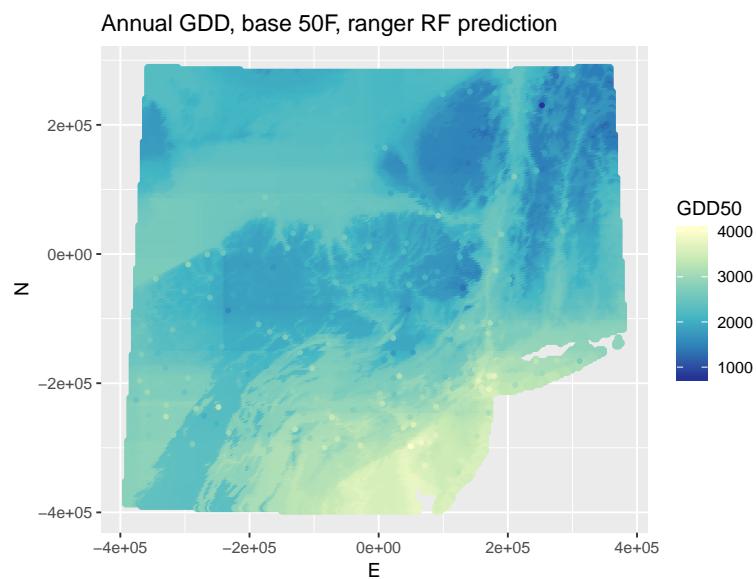
TASK 140 : Predict over the grid with the optimized ranger model, add the results to the dataframe, and summarize them. •

```
summary(dem.ne.m.df$pred.ranger <-
  predict(ranger.rf, data=dem.ne.m.df)$predictions)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1235    1978    2262    2339    2638    3777
```

TASK 141 : Display the map. •

```
display.prediction.map("pred.ranger",
  "Annual GDD, base 50F, ranger RF prediction",
  "GDD50")
```

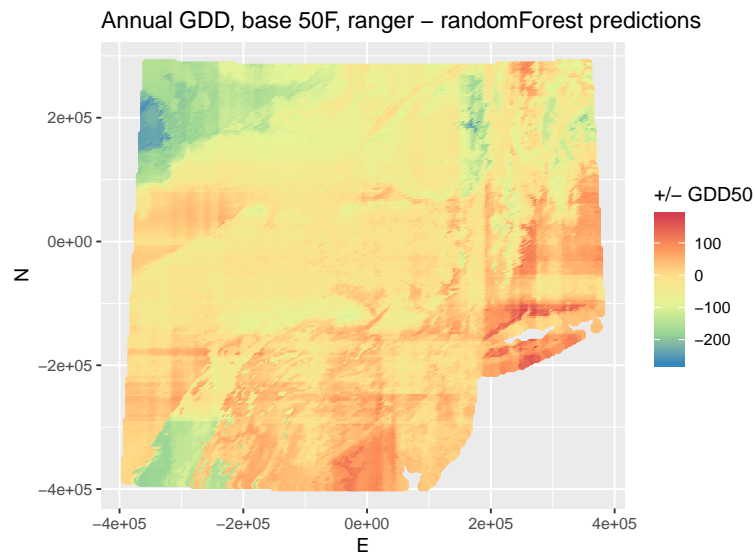



TASK 142 : Compare this “optimal” ranger random forest map with that made by randomForest (§12.2.2). •

```
summary(dem.ne.m.df$diff.ranger.rf <-
  dem.ne.m.df$pred.ranger - dem.ne.m.df$pred.rf)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -270.85  -54.11   -15.57   -23.24   18.97   181.90
```

```
display.difference.map("diff.ranger.rf",
  "Annual GDD, base 50F, ranger - randomForest predictions",
  "+/- GDD50")
```



Although these are two implementations of random forests, for `ranger` we used `mtry=2`, whereas for `randomForest` we accepted the default $1/3$ of the number of predictors, i.e., $3/3 = 1$.

12.4 Cubist

A popular data-driven model is Cubist, derived from the C4.5 models [23] but extensively modified and implemented as R package `Cubist`, an R port of the Cubist GPL C code released by RuleQuest⁴⁰.

```
library(Cubist)
```

A good introduction to Cubist is by Kuhn and Johnson [16, §8.7] as well as the vignette distributed with the `Cubist` package:

```
vignette("cubist")
```

Cubist is similar to regression trees, but instead of single values at leaves it creates a multivariate linear regression for the cases in the leaf. As the vignette explains:

“A tree is grown where the terminal leaves contain linear regression models. These models are based on the predictors used in previous splits. Also, there are intermediate linear models at each step of the tree. A prediction is made using the linear regression model at the terminal node of the tree, but is ‘smoothed’ by taking into account the prediction from the linear model in the previous node of the tree (which also occurs recursively up the tree). The tree is reduced to a set of

⁴⁰ <http://rulequest.com/cubist-info.html>

rules, which initially are paths from the top of the tree to the bottom. Rules are eliminated via pruning and/or combined for simplification.”

The advantage over regression trees is that the predictions are continuous, not discrete values equal to the number of leaves in the regression tree. The advantage over random forests is that the model can be interpreted, to a certain extent. A disadvantage of Cubist is that its algorithm is not easy to understand; however its results are generally quite good.

Cubist models can be improved in two ways: (1) with “committees” of models and (2) by adjusting predictions based on nearest neighbours in feature (predictor) space.

committees

Committees are a form of boosting. A set of model trees are built in sequence. The first tree is the standard Cubist best tree, using the original data in the training set. Subsequent trees are built from adjusted versions to the training set. If the previous Cubist tree over(under)-predicted a value, the response is adjusted down(up)ward for the next model, before it is fit. The final prediction is the average of the predictions from each model tree. The idea here is that the predictions by the sequence of trees vary around the “true” value.

nearest
neighbours

This prediction from the set of trees can then be adjusted using the values of some number of **nearest neighbours** in feature space. The idea here is that the overall model fits all the training data, but locally we may have some unknown factor that operates only in a local region of feature space, so if we have data from that region, we should give it more weight. Specifically, if the single model or committee predict a value \hat{y} , the adjusted prediction based on the K nearest-neighbours in feature space is:

$$\hat{y}' = \frac{1}{K} \sum_{i=1}^K w_i [t_i + (\hat{y} - \hat{t}_i)] \quad (26)$$

where t_i is the actual value of the neighbour, \hat{t}_i is its value predicted by the model tree(s), and w_i is the weight given to this neighbour for the adjustment, based on its distance D_i from the target point. These are computed as $w_i = 1/(D_i + 0.5)$ and normalized to sum to one.

In addition, to guard against using neighbours that are too far from the target point, the proposed set of neighbors are filtered based on the average pairwise distance of data points in the training set. In our case study, it is likely that Mt. Mansfield (VT) is too far from any other observation points to have any neighbours.

Obviously, the question is how many committees and neighbours to use, if any. The `caret` package can be used to tune these parameters (S12.3).

TASK 143 : Determine the optimum number of committees and neigh-

bours for a Cubist model to predict the growing degree days from Northing, Easting, and Elevation. •

We do this with the `train` function, specifying the method argument as 'cubist'.

```
require(caret)
all.preds <- ne.df[, c("N", "E", "ELEVATION_")]
all.resp <- ne.df[, "ANN_GDD50"]
system.time(
  cubist.tune <- train(x = all.preds, y = all.resp, "cubist",
    tuneGrid = expand.grid(.committees = 1:12,
      .neighbors = 0:8),
    trControl = trainControl(method = 'cv'))
)

##    user  system elapsed
##   8.664   0.029   8.705
```

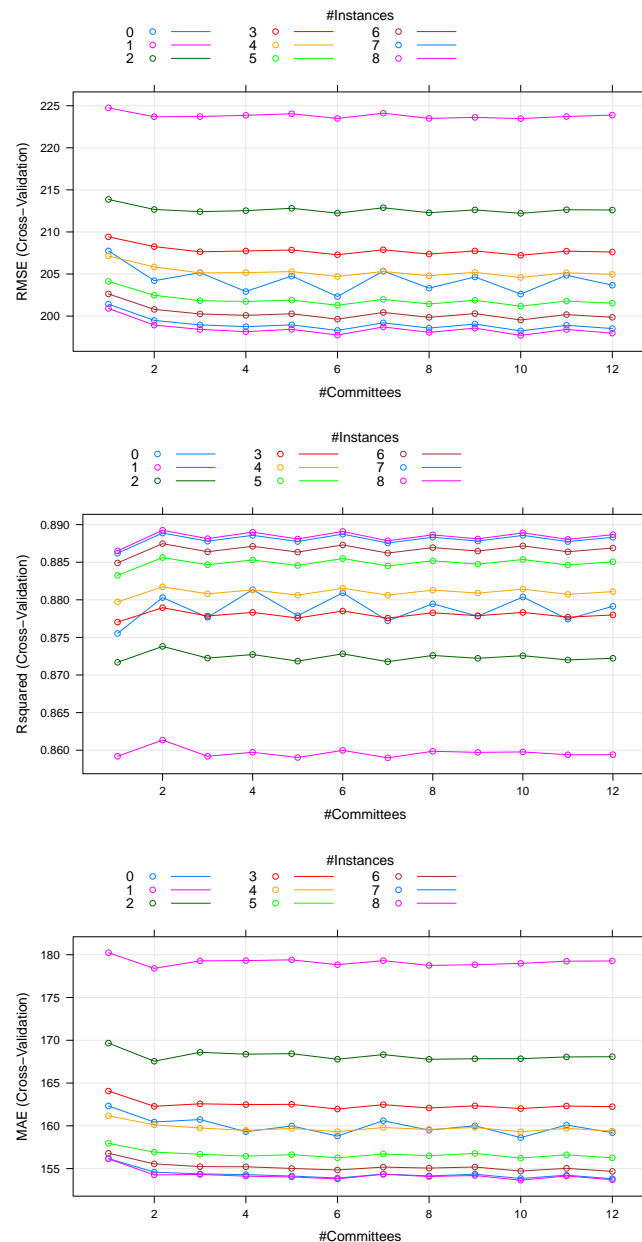
```
print(cubist.tune)

## Cubist
##
## 305 samples
##   3 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 273, 275, 275, 275, 274, 275, ...
## Resampling results across tuning parameters:
##
##  committees neighbors RMSE      Rsquared  MAE
##    1          0      207.7487  0.8755315 162.3228
##    1          1      224.7484  0.8591894 180.2378
##    1          2      213.8634  0.8716873 169.6693
##    1          3      209.4158  0.8770407 164.0687
##    1          4      207.1419  0.8797171 161.1776
##    1          5      204.1171  0.8832538 157.9565
##    1          6      202.6330  0.8849017 156.7901
##    1          7      201.4227  0.8861899 156.1779
##    1          8      200.9056  0.8865034 156.1447
##    2          0      204.2075  0.8803000 160.4451
##    2          1      223.6898  0.8613473 178.4033
##    2          2      212.6706  0.8738085 167.5547
##    2          3      208.2557  0.8789518 162.2772
##    2          4      205.8434  0.8817342 160.1162
##    2          5      202.4716  0.8856152 156.9241
##    2          6      200.7948  0.8874701 155.5540
##    2          7      199.4826  0.8888780 154.5913
##    2          8      198.9309  0.8892507 154.2609
##    3          0      205.1556  0.8776715 160.7356
##    3          1      223.7132  0.8591908 179.2812
##    3          2      212.4087  0.8722457 168.5927
##    3          3      207.6415  0.8778458 162.5732
##    3          4      205.1335  0.8807990 159.7426
##    3          5      201.8301  0.8846648 156.6787
##    3          6      200.2557  0.8863954 155.2295
##    3          7      198.9532  0.8878005 154.3835
##    3          8      198.4194  0.8881376 154.3195
##    4          0      202.9172  0.8813563 159.3132
##    4          1      223.8600  0.8597248 179.3147
##    4          2      212.5275  0.8727231 168.3751
##    4          3      207.7342  0.8783242 162.4805
##    4          4      205.1646  0.8813156 159.4704
##    4          5      201.7388  0.8852968 156.4552
##    4          6      200.0868  0.8871201 155.2128
##    4          7      198.7423  0.8885780 154.2961
##    4          8      198.1466  0.8890008 154.1155
##    5          0      204.7672  0.8778889 159.9803
```

##	5	1	224.0468	0.8590284	179.4008
##	5	2	212.8105	0.8718464	168.4341
##	5	3	207.8558	0.8775811	162.5104
##	5	4	205.2882	0.8806216	159.6961
##	5	5	201.8975	0.8845709	156.6291
##	5	6	200.2722	0.8863484	155.0193
##	5	7	198.9553	0.8877546	154.1300
##	5	8	198.4225	0.8880925	154.0242
##	6	0	202.3220	0.8809207	158.8087
##	6	1	223.4878	0.8599804	178.8349
##	6	2	212.2288	0.8728150	167.7835
##	6	3	207.2906	0.8785115	161.9613
##	6	4	204.7057	0.8815347	159.3266
##	6	5	201.2862	0.8854901	156.2590
##	6	6	199.6319	0.8873003	154.8363
##	6	7	198.2966	0.8887348	153.8970
##	6	8	197.7459	0.8890948	153.7957
##	7	0	205.3520	0.8772038	160.5959
##	7	1	224.1112	0.8589880	179.3086
##	7	2	212.8728	0.8717725	168.3237
##	7	3	207.8752	0.8775669	162.4743
##	7	4	205.2987	0.8806245	159.8012
##	7	5	201.9814	0.8845065	156.7110
##	7	6	200.4333	0.8862200	155.1742
##	7	7	199.1963	0.8875543	154.3628
##	7	8	198.7111	0.8878471	154.3406
##	8	0	203.3202	0.8794670	159.4943
##	8	1	223.4954	0.8598596	178.7526
##	8	2	212.2826	0.8726029	167.7763
##	8	3	207.3700	0.8782757	162.0853
##	8	4	204.7955	0.8812939	159.5780
##	8	5	201.4404	0.8851888	156.5082
##	8	6	199.8493	0.8869440	155.0606
##	8	7	198.5758	0.8883208	154.1437
##	8	8	198.0656	0.8886398	154.0757
##	9	0	204.6509	0.8778234	160.0073
##	9	1	223.6139	0.8597099	178.8303
##	9	2	212.6175	0.8722231	167.8424
##	9	3	207.7421	0.8778849	162.3367
##	9	4	205.1965	0.8808933	159.8104
##	9	5	201.8795	0.8847465	156.7730
##	9	6	200.2982	0.8864854	155.1845
##	9	7	199.0518	0.8878149	154.3446
##	9	8	198.5693	0.8880959	154.1941
##	10	0	202.6142	0.8803737	158.6199
##	10	1	223.4708	0.8597699	178.9933
##	10	2	212.2136	0.8725701	167.8452
##	10	3	207.2250	0.8783267	162.0196
##	10	4	204.5850	0.8814157	159.3126
##	10	5	201.1794	0.8853635	156.2361
##	10	6	199.5381	0.8871732	154.7109
##	10	7	198.2367	0.8885710	153.8447
##	10	8	197.7055	0.8889094	153.6400
##	11	0	204.8540	0.8774217	160.0968
##	11	1	223.7206	0.8593948	179.2459
##	11	2	212.6357	0.8720026	168.0535
##	11	3	207.7186	0.8776943	162.3113
##	11	4	205.1433	0.8807237	159.7070
##	11	5	201.7819	0.8846289	156.6128
##	11	6	200.1767	0.8863936	155.0306
##	11	7	198.9069	0.8877460	154.2308
##	11	8	198.4075	0.8880418	154.1380
##	12	0	203.6696	0.8791484	159.1982
##	12	1	223.8857	0.8594049	179.2720
##	12	2	212.6048	0.8722285	168.0769
##	12	3	207.6146	0.8779872	162.2367
##	12	4	204.9556	0.8810915	159.3881
##	12	5	201.5261	0.8850577	156.2663
##	12	6	199.8546	0.8868914	154.6737

```
## 12      7      198.5165 0.8883212 153.8216
## 12      8      197.9749 0.8886656 153.6956
##
## RMSE was used to select the optimal model using the smallest value.
## The final values used for the model were committees = 10
## and neighbors = 8.
```

```
plot(cubist.tune, metric="RMSE")
plot(cubist.tune, metric="Rsquared")
plot(cubist.tune, metric="MAE")
```



Your results will be different, because of the randomness in the splits from method. So, your choice of optimal values may also be different than those presented here.

In this case committees improve the model. Although the best result

is with 10 committees, looking at the graphs we can see that 6 gives almost equally good results, so we prefer the simpler model. Neighbours definitely do improve the model. Using one neighbour (the closest in feature space) makes the model much worse – too much fine adjustment to the training set. Using two to seven neighbours gives improvement, using eight is only slightly better. This shows that the overall model can benefit by local adjustment.

TASK 144 : Built an “optimum” Cubist model. •

The model is fit with the `cubist` function; the optimal number of committees is specified at model building with the `committees` argument. Adjustment by neighbours is done at prediction, because only then do we know which point we are predicting, hence which are its nearest neighbours. This is specified by the `neighbors` to the `predict` function of the `Cubist` package.

```
c.model <- cubist(x = all.preds, y = all.resp, committees=6)
summary(c.model)

##
## Call:
## cubist.default(x = all.preds, y = all.resp, committees = 6)
##
## Cubist [Release 2.07 GPL Edition] Fri Jul 26 16:40:11 2019
## -----
##      Target attribute `outcome'
##
## Read 305 cases (4 attributes) from undefined.data
##
## Model 1:
##
## Rule 1/1: [109 cases, mean 2114.8, range 795 to 2845, est err 151.9]
##
##   if
##   N > -20368.92
##   then
##   outcome = 2715.4 - 0.496 ELEVATION_ - 0.00105 N - 0.00068 E
##
## Rule 1/2: [196 cases, mean 2741.5, range 1335 to 4021, est err 168.7]
##
##   if
##   N <= -20368.92
##   then
##   outcome = 2704 - 0.00264 N - 0.606 ELEVATION_ - 0.00021 E
##
## Model 2:
##
## Rule 2/1: [305 cases, mean 2517.5, range 795 to 4021, est err 171.2]
##
## outcome = 2826.9 - 0.599 ELEVATION_ - 0.00205 N - 0.00023 E
##
## Model 3:
##
## Rule 3/1: [174 cases, mean 2171.9, range 795 to 3084, est err 172.4]
##
##   if
##   N > -145485.8
##   then
##   outcome = 2642.4 - 0.474 ELEVATION_ - 0.00077 E - 0.00045 N
##
```

```

## Rule 3/2: [74 cases, mean 2729.1, range 1774 to 3428, est err 178.6]
##
##   if
## N <= -145485.8
## ELEVATION_ > 330
##   then
## outcome = 1956.9 - 0.00513 N - 0.457 ELEVATION_ - 9e-05 E
##
## Rule 3/3: [76 cases, mean 3090.3, range 2078 to 4021, est err 182.4]
##
##   if
## ELEVATION_ <= 330
##   then
## outcome = 3106.6 - 2.148 ELEVATION_ - 0.00188 N
##
## Model 4:
##
## Rule 4/1: [305 cases, mean 2517.5, range 795 to 4021, est err 170.7]
##
## outcome = 2842.7 - 0.613 ELEVATION_ - 0.00206 N - 0.0003 E
##
## Model 5:
##
## Rule 5/1: [174 cases, mean 2171.9, range 795 to 3084, est err 174.1]
##
##   if
## N > -145485.8
##   then
## outcome = 2629.8 - 0.462 ELEVATION_ - 0.00071 E - 0.00043 N
##
## Rule 5/2: [74 cases, mean 2729.1, range 1774 to 3428, est err 183.9]
##
##   if
## N <= -145485.8
## ELEVATION_ > 330
##   then
## outcome = 1925.4 - 0.0052 N - 0.45 ELEVATION_ - 6e-05 E
##
## Rule 5/3: [76 cases, mean 3090.3, range 2078 to 4021, est err 182.7]
##
##   if
## ELEVATION_ <= 330
##   then
## outcome = 3112.7 - 2.176 ELEVATION_ - 0.00185 N
##
## Model 6:
##
## Rule 6/1: [305 cases, mean 2517.5, range 795 to 4021, est err 171.1]
##
## outcome = 2854.2 - 0.623 ELEVATION_ - 0.00208 N - 0.00036 E
##
## Evaluation on training data (305 cases):
##
##   Average |error|           169.9
##   Relative |error|         0.36
##   Correlation coefficient    0.93
##
##
## Attribute usage:
##   Conds  Model
##
##   43%    100%    N
##   16%    100%    ELEVATION_
##           92%    E
##
##
## Time: 0.0 secs

```


Rule 1 splits at $N = -20368.92$, near the centre of the map. There is then a slightly different linear regression for the two halves. The elevation and Northing coefficients are both larger for the southern half. Rule 2 has no split, just a single linear model. It is not the same as the linear model fit in §4, because it is fit based on the values adjusted by Rule 1 predictions. Rule 3 again splits on Northing, but much further south, at $N = -145485.8$, and on low elevations, 330.

TASK 145 : Examine the fit of the Cubist model to the known points. •

```
# predictive accuracy
cubist.fits <- predict(c.model, newdata=all.preds,
                      neighbors=cubist.tune$bestTune$neighbors)

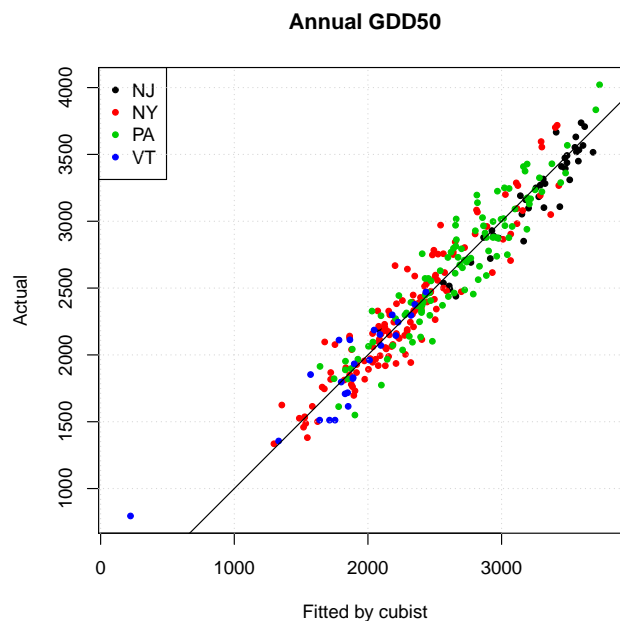
## Test set RMSE
sqrt(mean((cubist.fits - all.resp)^2))

## [1] 164.9617

cor(cubist.fits, all.resp)^2 # R^2

## [1] 0.918133

plot(ne.df$ANN_GDD50 ~ cubist.fits,
     col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by cubist",
     ylab="Actual",
     main="Annual GDD50")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
```



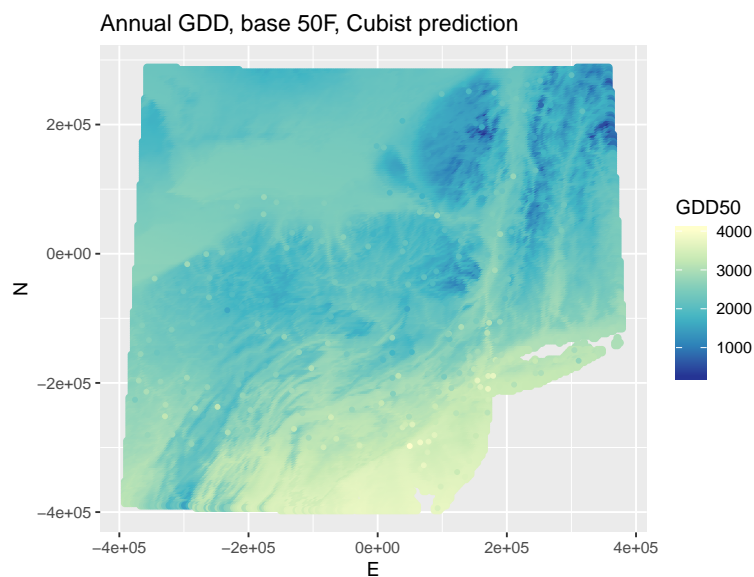
Notice that this Cubist model severely underpredicts the lowest GDD50 at Mount Mansfield (VT). This is because it is too far in elevation space from any other observation points to have any neighbours that could modify the linear regressions in the rule set, which reduce GDD50 with elevation.

TASK 146 : Predict over the study area with the Cubist model and display the resulting map. •

```
summary(dem.ne.m.df$pred.cubist <-
  predict(c.model, newdata=dem.ne.m.df,
    neighbours=cubist.tune$bestTune$neighbors))

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 265.7 1999.6 2278.4 2343.1 2609.1 3733.5

display.prediction.map("pred.cubist",
  "Annual GDD, base 50F, Cubist prediction",
  "GDD50")
```



Although several of the rules split on Northing, that is not visible in this map, because the results of each rule are averaged.

13 Thin-plate spline interpolation

A quick way to see the distribution of a variable in space as a surface is with an empirical method that adjusts locally to the data. A common empirical method is **thin-plate splines** (TPS), also referred to as “minimum curvature”, which are implemented in the `fields` package.

13.1 * Theory

Hastie et al. [11, §5.7] explains the mathematics of multi-dimensional smoothing splines. A more thorough mathematical treatment is given by Wood [33] and Mitsova and Mitsova [21]; these are developments from

the “minimum curvature” methods of Briggs [3]. Applications include Hutchinson [13] and Mitasova and Hofierka [20].

TPS is the mathematical equivalent of a thin (so, flexible) plate that is warped to fit the data. This can range from very “rigid”, i.e., just a single surface (the usual least-squares plane of a first-order trend surface) to very “flexible”, i.e., perfectly fitting every observation. In general we want something in between: if we think there is an overall surface we just fit it as one polynomial (first, second ... order polynomials on the coördinates), but if we want to fit more locally, we must expect local noise which should be somehow locally averaged-out.

Fitting a TPS depends on the k data points with known coördinates and attribute values. They can be described by $2(k + 3)$ parameters, six of which are overall affine transformation parameters (to center the function in 2D) and $2k$ of which link to the control points.

The general method is to minimize the residual sum of squares (RSS) of the fitted function, subject to a constraint that the function be “smooth” in some sense; this is expressed by a **roughness penalty** which balances the fit to the observations with smoothness. This is a minimization problem. If \mathbf{x}_i is one point in 2D space (i.e., it has two coördinates) and y_i is the attribute value at the same points, the aim is to minimize:

$$\min_f \sum_{i=1}^N \{y_i - f(\mathbf{x}_i)\}^2 + \lambda J[f] \quad (27)$$

where J is the penalty function and λ controls how important it is; $\lambda = 0$ means there is no roughness penalty and the data will be fit exactly; as $\lambda \rightarrow \infty$ the solution approximates the least-squares plane, i.e., the trend surface averaged over all the points.

In 2D an appropriate penalty is:

$$J[f] = \int_{\mathbb{R}} \int_{\mathbb{R}} \left[\left(\frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} \right)^2 + 2 \left(\frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} \right)^2 \right] dx_1 dx_2 \quad (28)$$

where (x_1, x_2) are the two coördinates of the vector \mathbf{x} . In practice the double integral is discretized over some grid known as **knots**; these may be defined by the observations or may be a different set, maybe an evenly-spaced grid.

This penalty can be interpreted as the “bending energy” of a thin plate represented by the function $f(\mathbf{x})$; by minimizing this energy the spline function in over the 2D plane is a thin (flexible) plate which, according to the first term of Equation 27 would be forced to pass through data points, with minimum bending. However the second term of Equation 27 allows some smoothing: the plate does not have to bend so much, since it is allowed to pass “close to” but not necessarily through the data points. The higher the λ , the less exact is the fit. This has two purposes:

(1) it allows for measurement error; the data points are not taken as exact; (2) it results in a smoother surface. So cross-validation is used to determine the degree of smoothness.

The solution to Equation 28 is a linear function:

$$f(\mathbf{x}) = \beta_0 + \beta^T \mathbf{x} + \sum_{j=1}^N \alpha_j h_j(\mathbf{x}) \quad (29)$$

where the β account for the overall trend and the α are the coefficients of the warping.

The set of functions $h_j(\mathbf{x})$ is the **basis kernel**, also called a **radial basis function** (RBF), for thin-plate splines:

$$h_j(\mathbf{x}) = \|\mathbf{x} - \mathbf{x}_j\|^2 \log \|\mathbf{x} - \mathbf{x}_j\| \quad (30)$$

where the norm distance $r = \|\mathbf{x} - \mathbf{x}_j\|$ is also called the **radius** of the basis function. The norm is usually the Euclidean (straight-line) distance.

13.2 Practice

Here we only uses the Northing and Easting as predictors.

TASK 147 : Set up for thin-plate splines and compute the minimum-curvature spline, subject to roughness constraint determined by generalized cross-validation. •

The `Tps` function of the `fields` package compute this; however the co-ordinates must be formatted as a matrix field in the dataframe, using the `matrix` function.

```
require(fields)
ne.tps <- ne.df
ne.tps$coords <- matrix(c(ne.tps$E, ne.tps$N), byrow=F, ncol=2)
str(ne.tps$coords)

##  num [1:305, 1:2] 123045 135231 78656 97523 77534 ...

surf.1 <-Tps(ne.tps$coords, ne.tps$ANN_GDD50)
summary(surf.1)

## CALL:
## Tps(x = ne.tps$coords, Y = ne.tps$ANN_GDD50)
##
## Number of Observations:      305
## Number of unique points:     305
## Number of parameters in the null space 3
## Parameters for fixed spatial drift    3
## Effective degrees of freedom:      88.5
## Residual degrees of freedom:      216.5
## MLE sigma                          220.7
## GCV sigma                          228.6
## MLE rho                           129700000
## Scale passed for covariance (rho)    <NA>
## Scale passed for nugget (sigma^2)    <NA>
## Smoothing parameter lambda          0.0003755
##
```

```
## Residual Summary:
##      min      1st Q      median      3rd Q      max
## -882.40000 -113.60000  -0.01774  127.30000  574.90000
##
## Covariance Model: Rad.cov
## Names of non-default covariance arguments:
##      p
##
## DETAILS ON SMOOTHING PARAMETER:
## Method used: GCV Cost: 1
##      lambda      trA      GCV      GCV.one GCV.model      shat
## 3.755e-04 8.846e+01 7.359e+04 7.359e+04      NA 2.286e+02
##
## Summary of all estimates found for lambda
##      lambda      trA      GCV      shat -lnLike Prof converge
## GCV      0.0003755 88.46 73587 228.6      2149      13
## GCV.model      NA      NA      NA      NA      NA      NA
## GCV.one      0.0003755 88.46 73587 228.6      NA      13
## RMSE      NA      NA      NA      NA      NA      NA
## pure error      NA      NA      NA      NA      NA      NA
## REML      0.0009690 59.35 74532 245.0      2146      7
```

TASK 148 : Set up a grid covering the four States at approximately 6 x 6 km resolution, and convert to a dataframe with the coordinates as a matrix field. •

The `spsample` function of the `sp` package can make various sampling plans, including a regular grid, within a study area.

We compute the approximate area of the four states, in km, from its bounding box in the US Census state shapefile; these are in m², and so must be converted to km². We then ask for a grid with each cell covering about (9 km)².

Note: Recall, this is not for exact prediction, just to get an overview of the regional distribution of the variable of interest.

```
resolution <- 9
bbox(state.ne.m)

##      min      max
## E -387135.2 357707.2
## N -396311.1 288684.4

(ne.sq.km <- diff(bbox(state.ne.m)[1,]) * diff(bbox(state.ne.m)[2,])/10^6)

##      max
## 510213.8

(approx.n.grid.cells <- ne.sq.km/(resolution^2))

##      max
## 6298.935

states.grid <- spsample(state.ne.m, n=approx.n.grid.cells, type="regular")
gridded(states.grid) <- TRUE
states.grid.df <- as.data.frame(states.grid)
states.grid.df$coords <-
  matrix(c(states.grid.df$x1, states.grid.df$x2), byrow=F, ncol=2)
str(states.grid.df)

## 'data.frame': 6294 obs. of 3 variables:
## $ x1 : num 91201 97979 97979 104757 97979 ...
```

```
## $ x2      : num  -389755 -389755 -382977 -382977 -376199 ...
## $ coords: num [1:6294, 1:2] 91201 97979 97979 104757 97979 ...
```

TASK 149 : Predict over the four-states grid using the fitted thin-plate spline. •

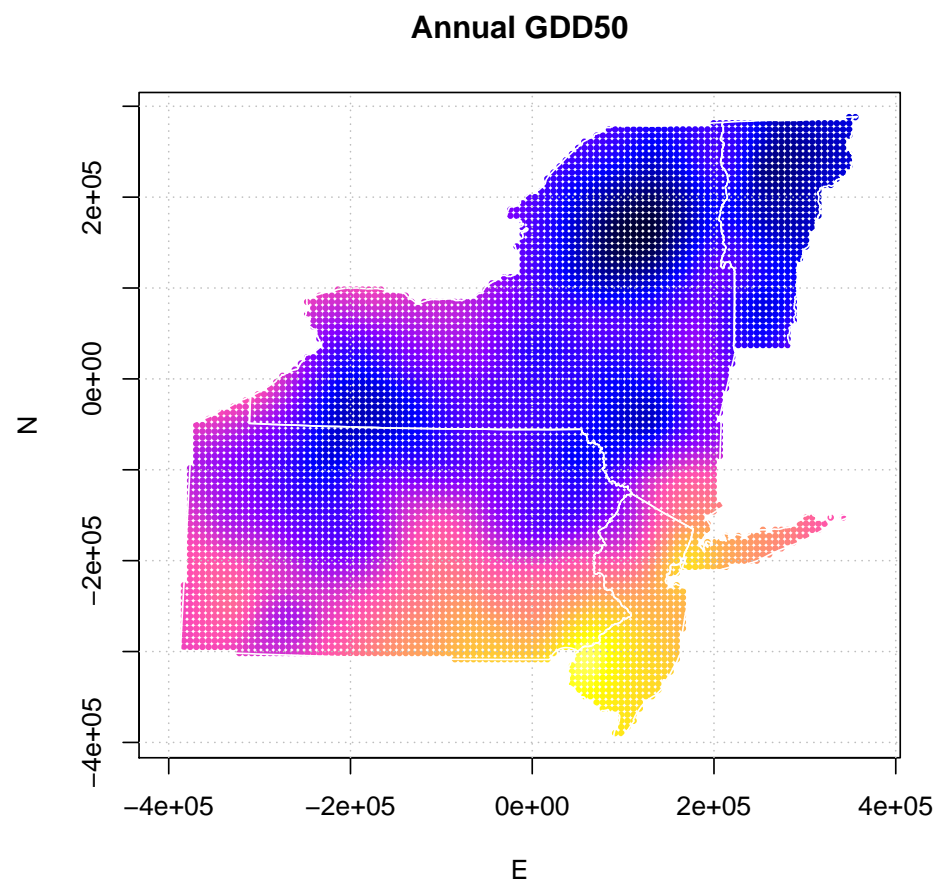
The `predict.Krig` method of the `fields` package computes the prediction.

```
surf.1.pred <- predict.Krig(surf.1, states.grid.df$coords)
summary(as.vector(surf.1.pred))
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1459   2057   2288   2390   2675   3694
```

TASK 150 : Display the gridded prediction. •

```
plot(states.grid.df$coords, pch=20, asp=1, cex=.6,
     col=bpy.colors(256)[cut(surf.1.pred, 256)],
     xlab="E", ylab="N",
     main="Annual GDD50")
lines(state.ne.m, col="white")
grid(col="gray")
```



This map captures the main features of the annual GDD50 fairly well,

even though elevation was not used in the thin-plate spline empirical model. In particular, it captures the high-GDD areas along the Lake Ontario and Lake Erie plains, the low-GDD cold spots in the Adirondack, Catskill and Green Mountains, as well as in the Allegheny State Park area of SW NY/NE PA, and the very high GDD-area around the Delaware Bay. It does not account for local variations in GDD because of elevation.

Note: It is also possible to include the elevation in the thin-plate spline model, but but because it varies at such short ranges, the result is non-sense.

To determine the predicted value at any location, the kriged surface must be appended as a data frame to the `SpatialPixels` object, creating a `SpatialPixelsDataFrame` object.

```
states.grid.spdf <- SpatialPixelsDataFrame(
  states.grid,
  data=data.frame(tps.pred=as.vector(surf.1.pred))
)
```

We can then determine the value at a location with the over “spatial overlay” method of the `sp` package. This queries its second argument (layer from which attributes are queried) and assigns attribute values to the geometries in the first argument.

```
pt <- SpatialPoints(coords=cbind(-76.402175, 42.453271))
proj4string(pt) <- CRS("+init=EPSG:4326")
pt <- spTransform(pt, CRS(proj4string(states.grid.spdf)))
coordinates(pt)

##      coords.x1 coords.x2
## [1,] -33055.67 -5118.213

pt.sp <- SpatialPointsDataFrame(pt, over(pt, states.grid.spdf))
pt.sp$tps.pred

## [1] 2202.954
```

The thin-plate spline interpolation predicts 2203 GDD50 for this location.

14 Local interpolators

A purely **local** approach to prediction is to ignore any causitive factors (in this case, northing and elevation) and just use “nearby” known observations to predict at any location. This is an operational realization of the well-known Tobler’s First Law of Geography: “everything is related to everything else, but near things are more related than distant things” [27]⁴¹. In the current case this is not adviseable, because of the strong and useful relation of the target variable `ANN_GDD50` with the covariables, which we have seen in earlier sections. However, for completeness we illustrate this method.

Local approaches can be **model-based** or **model-free**.

⁴¹ Tobler goes on to point out that “the ...model used is thus very parochial, and ignores most of the world”

The best-known model-based method is Ordinary Kriging, which relies on a model of local spatial dependence of the target variable. In this case the universal model of spatial distribution shown in 1 is simplified to:

$$Z(\mathbf{s}) = \varepsilon(\mathbf{s}) + \varepsilon'(\mathbf{s}) \quad (31)$$

\mathbf{s} : a location in space, designated by a **vector** of coördinates;

$Z(\mathbf{s})$: **true** (unknown) value of some property at the location;

- when modelled, expressed as **most likely** value and some **uncertainty**, or as a **probability distribution**

$\varepsilon(\mathbf{s})$: locally **spatially-autocorrelated stochastic** component;

$\varepsilon'(\mathbf{s})$: pure (“white”) **noise**, no structure.

In the purely local model there is no deterministic component, i.e., global coördinates or covariables.

We model $\varepsilon(\mathbf{s})$ with an **authorized model of spatial dependence**, usually with an authorized variogram model. This model is then used in Ordinary Kriging (OK). This also reveals the magnitude of $\varepsilon'(\mathbf{s})$, i.e., the pure noise that can not be modelled nor predicted.

Good explanations of Ordinary Kriging are from Webster and Oliver [29] and Goovaerts [8].

14.1 Computing the empirical variogram

The definition of the empirical variogram was explained in §4.4.1. Here we apply it to the original values, not the trend residuals as in that section.

TASK 151 : Display an empirical variogram to 400 km maximum separation. •

We use the `variogram` function of the `gstat` package. The default bin width for this function is 1/15 the cutoff. By default the cutoff is 1/3 of the diagonal across the bounding box of the point-set. These are both ‘rules of thumb’ and should be adjusted by the analyst until (1) the range of spatial dependence (if any) is found; (2) each bin has enough point-pairs; (3) the structure is clear.

First with the defaults. The cutoff is then 322 km, and the bin widths about 21 km.

```
library(gstat)
bbox(ne.m)

##           min           max
## E -375745.2 318960.4
## N -393922.9 276614.1

(cutoff.default <- (sqrt(diff(bbox(ne.m)["E",])^2 + diff(bbox(ne.m)["N",])^2))/3)
```

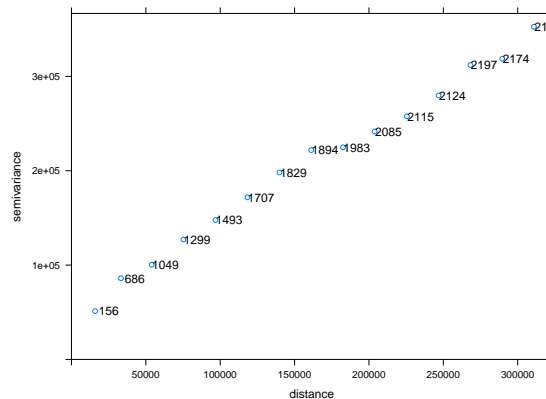


```
##      max
## 321841.2

(binwidth.default <- cutoff.default/15)

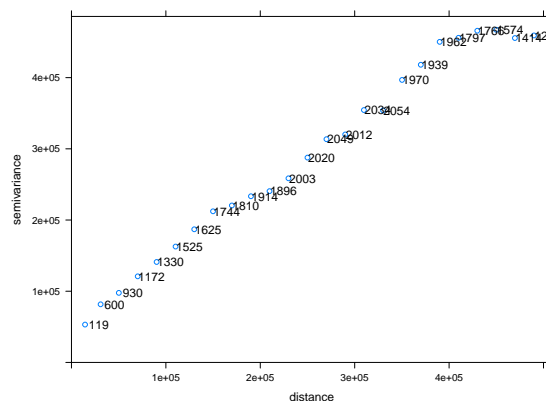
##      max
## 21456.08

plot(v.o <- variogram(ANN_GDD50 ~ 1, locations=ne.m),
     plot.numbers=TRUE)
```



Within this range the variogram is unbounded, which means the maximum variation has not been reached even between observations at the largest separations. This is because of the strong regional effect of Northing. Extend the cutoff to see if we can find a bound:

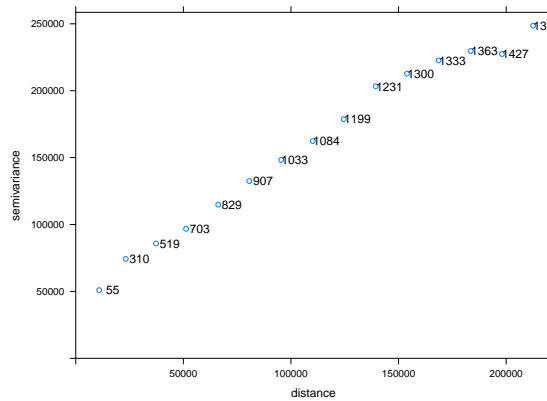
```
plot(v.o <- variogram(ANN_GDD50 ~ 1, locations=ne.m, cutoff=500000, width=20000),
     plot.numbers=TRUE)
```



Here we finally do get a bound, at about 430 km. There appears to be two structures: one to about 180 km and one to about 430 km. Since OK is a local predictor, and the nearer points receive most of the weight, this short-range variation is what we want to model for local prediction.

TASK 152 : Recompute and display the empirical variogram to 220 km cutoff.

```
v.o <- variogram(ANN_GDD50 ~ 1, locations=ne.m,
                 cutoff=220000)
plot(v.o, plot.numbers=TRUE)
```



14.2 Fitting an authorized variogram model

We require a continuous function of semivariance vs. separation, so that for any separation we can compute the semivariance to be used in the kriging equations. This must be an **authorized** variogram model that ensures that the kriging system will be positive semi-definite and thus invertible.

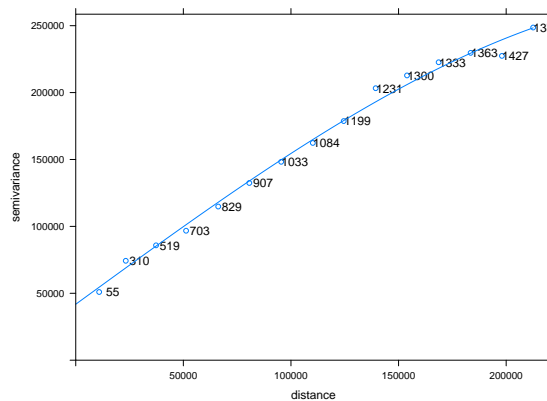
This empirical variogram can be well-modelled with a spherical variogram model; at short ranges this is almost linear. The nugget appears to be about 40 000 GDD², the partial sill at the maximum semivariance of about 250 000 GDD² less the nugget, and the range can be approximated by the maximum separation.

The `fit.variogram` function fits variogram models specified with the `vgm` “variogram model” function. We initialize the weighted least-squares (WLS) fit to the empirical variogram with our eyeball estimates.

```
(vmf.o <- fit.variogram(v.o,
                        vgm(psill=210000, model="Sph",
                           range=220000, nugget=40000)))

##   model    psill    range
## 1  Nug  41816.39     0.0
## 2  Sph 232353.41 297706.6

plot(v.o, plot.numbers=TRUE, model=vmf.o)
```



This fits fairly well; our original estimates were not too far off. The range is fitted to be about 35% longer than our estimate.

14.3 Predicting by Ordinary Kriging

Once we have the fitted model and the data points, we can predict at any location, by solving the OK system (§8.1) for the weights λ_i to be used in the weighted average (Eqn. 14). The mathematics of OK were presented in (§8.1).

TASK 153 : Predict over this grid with OK. •

For Ordinary Kriging, we can use stations that we know are in the “local” region, for example the nearest “few” stations, or we can use all points within the range of the empirical variogram. The latter will have a strong averaging effect; the former will show artefacts where some stations come into or leave the range from a prediction point. Theoretically we should use all the stations; but if we have reason to suspect first-order non-stationarity, we may prefer to just use the closest stations.

For the local OK we specify the optional `nmax` to use only the twenty-four closest stations.

We also specify the optional `block` argument, to average the prediction over a 1 km² block.

```
# dem.ne.m.sp was set up in \S6.2 "Adjusting the grid for prediction"
k.ok <- krige(ANN_GDD50 ~ 1, locations=ne.m, newdata=dem.ne.m.sp,
             model=vmf.o, nmax=24, block=c(1000,1000))

## [using ordinary kriging]

summary(k.ok)

## Object of class SpatialPixelsDataFrame
## Coordinates:
##      min      max
## E -392899.3 379900.7
## N -399006.4 289937.6
## Is projected: TRUE
## proj4string :
## [+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76]
```

```
## +ellps=WGS84 +units=m]
## Number of points: 35783
## Grid attributes:
##   cellcentre.offset cellsize cells.dim
## E           -394624.3    3450      231
## N           -400858.4    3704      188
## Data attributes:
##   var1.pred      var1.var
## Min.   :1483    Min.   : 12201
## 1st Qu.:2157    1st Qu.: 27576
## Median :2456    Median : 34580
## Mean   :2478    Mean   : 68018
## 3rd Qu.:2715    3rd Qu.: 85242
## Max.   :3740    Max.   :337793
```

TASK 154 : List the 24 stations that were used for the local prediction at the block containing the Ithaca weather station, in order of their separations from Ithaca. •

```
e.ith.pt <- subset(ne.m, (ne.m$STATION_NAME=="ITHACA CORNELL UNIV"))
dist.pt <- spDists(ne.m, e.ith.pt)
(ix <- order(dist.pt)[1:24])

## [1] 100 148 72 49 154 54 81 167 82 48 133 85 119 40 53 128
## [17] 271 152 62 241 145 51 76 147

print(cbind(ne.df[ix,c('STN_NAME', 'STATE', 'ELEVATION_', 'ANN_GDD50')],
            dist=round(dist.pt[ix,]/1000,1)))

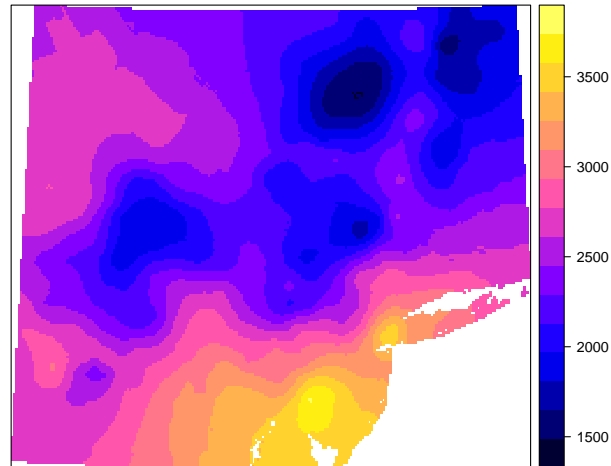
##           STN_NAME STATE ELEVATION_ ANN_GDD50 dist
## 3074      ITHACA CORNELL UNIV NY          960      2160 0.0
## 3122      SPENCER 2 N NY         1050      1988 22.6
## 3046      CORTLAND NY         1129      2329 27.7
## 3023      AURORA RESEARCH FARM NY          830      2590 35.2
## 3128      TULLY HEIBERG FOREST NY         1899      1746 46.8
## 3028      BINGHAMTON BROOME CO AP NY         1600      2141 47.7
## 3055      ELMIRA NY          844      2395 48.4
## 3141      WAVERLY NY          845      2304 50.5
## 3056      ENDICOTT NY          827      2245 51.2
## 3022      AUBURN NY          744      2352 52.9
## 3107      PENN YAN NY          830      2430 55.0
## 3059      GENEVA RESEARCH FARM NY          718      2396 67.4
## 3093      MORRISVILLE 6 SW NY         1300      1818 72.6
## 3014      ADDISON NY          980      2123 74.0
## 3027      BATH NY         1120      2055 74.9
## 3102      NORWICH NY         1020      2203 76.1
## 3881      TOWANDA 1 ESE PA          750      2447 77.8
## 3126      SYRACUSE HANCOCK INTL AP NY          410      2467 79.8
## 3036      CANANDAIGUA 3 S NY          720      2531 81.3
## 3851      MONTROSE PA         1420      2064 81.3
## 3119      SHERBURNE 2 S NY         1080      2157 82.7
## 3025      BAINBRIDGE 2 E NY          994      2063 84.4
## 3050      DEPOSIT NY         1000      2245 94.1
## 3121      SODUS CENTER NY          420      2558 95.5
```

Notice the wide range of elevations in this group of stations, from 410' to almost 1900'.

TASK 155 : Display a map of the ANN_GDD50 predicted by OK. •

A simple way to present a map of an attribute stored as a field in a `SpatialGridDataFrame` is with the `splot` function of the `sp` package:

```
spplot(k.ok, zcol="var1.pred")
```

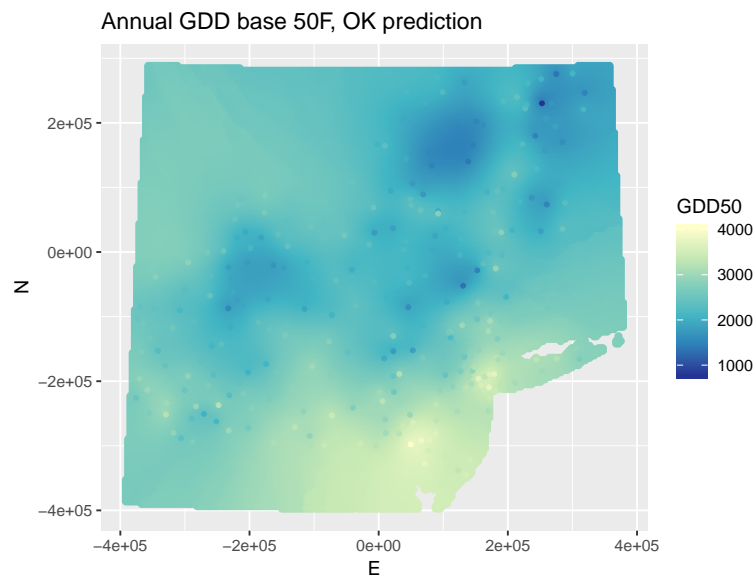


For comparison with the maps made with other techniques, we also add the results to the `data.frame` covering the same area and then display it with `ggplot`:

```
dim(dem.ne.m.df)

## [1] 35783    26

dem.ne.m.df$pred.ok <- k.ok$var1.pred
display.prediction.map("pred.ok",
  "Annual GDD base 50F, OK prediction",
  "GDD50")
```



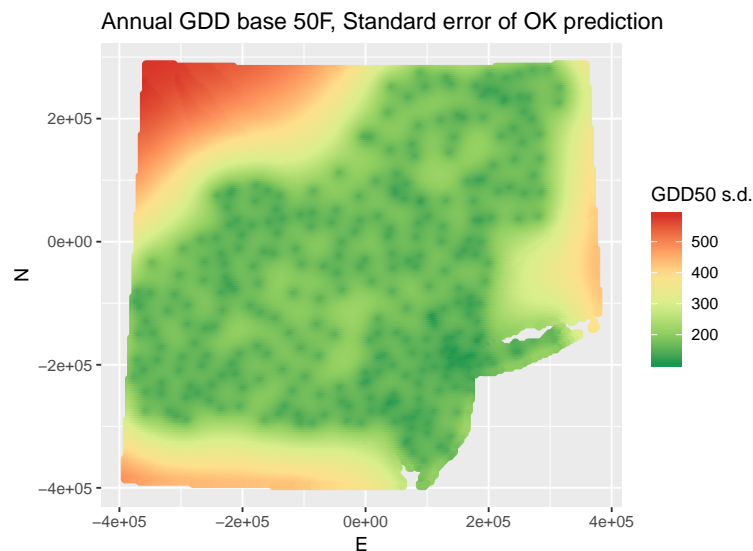
We see this is much smoother than other methods, mainly because it does not take elevation into account.

14.3.1 Accuracy assessment

An advantage of OK is that it gives a direct estimate of its prediction variance at each prediction location – this is because the variance must be minimized, so it must be computed.

TASK 156 : Display a map of the OK prediction standard deviations. •

```
dem.ne.m.df$sd.ok <- sqrt(k.ok$var1.var)
ggplot() +
  geom_point(aes(x=E, y=N, colour=sd.ok), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Annual GDD base 50F, Standard error of OK prediction") +
  scale_colour_distiller(name="GDD50 s.d.", space="Lab", palette="RdYlGn",
    direction=-1)
```



Clearly, the further from local information, the more the prediction is uncertain, e.g., northwest of Lake Ontario. The prediction uncertainty is least near stations, especially near a cluster of stations, e.g., the NYC area. We could choose to not predict at areas with too much uncertainty, based on user requirements.

The minimum and mean prediction standard deviations:

```
summary(dem.ne.m.df$ok.sd)

## Length Class Mode
##      0  NULL  NULL

ix <- which.min(dem.ne.m.df$ok.sd)
k.ok@data[ix,"var1.pred"]

## numeric(0)

round(100*dem.ne.m.df[ix,"ok.sd"]/k.ok$var1.pred[ix],1)

## numeric(0)
```

The minimum is quite small, only about 3% of the prediction at that point, and the mean is only about 2.5 times as poor.

A better way to evaluate the predictive power OK is by **Leave-one-out cross-validation** (LOOCV). Here each point is removed from the dataset in turn, and predicted by the others, using the fitted variogram model. If the observation points well represent the total population, as they do here by design of the weather station network, this gives a good estimate of the prediction error.

TASK 157 : Compute and summarize the LOOCV for this OK prediction.

The `krige.cv` function of the `gstat` package computes this:

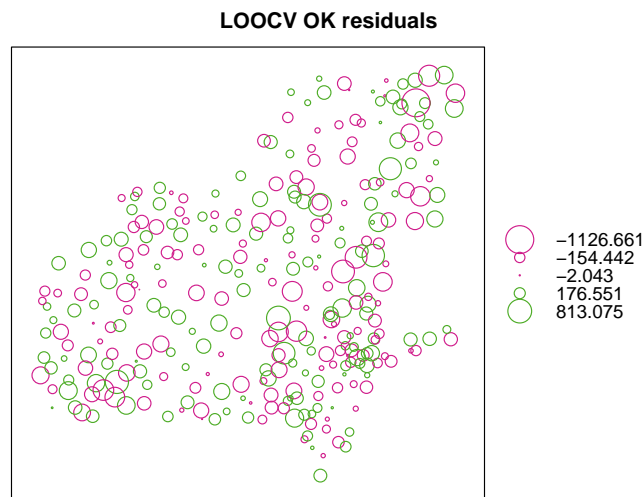
```
kcv.ok <- krige.cv(ANN_GDD50 ~ 1, locations=ne.m, model=vmf.o)
summary(kcv.ok$residual)
```

Overall the results are fairly good, but there are some extremely bad predictions. An overall measure is the **root of the mean squared error**, RMSE:

```
(loocv.ok.rmse <- sqrt(sum(kcv.ok$residual^2)/length(kcv.ok$residual)))
## [1] 268.1796
```

TASK 158 : Display a bubble plot of the LOOCV residuals.

```
bubble(kcv.ok, zcol="residual", pch=1, main="LOOCV OK residuals")
```



There are several regions with intermixed fairly large under- and over-predictions; this means that in these regions there are local factors, most notable elevation.

TASK 159 : Find the worst predictions, try to explain in geographic terms.

```
ne.m@data[which.min(kcv.ok$residual),2:6]

##      STATE      STATION_NA LATITUDE_D LONGITUDE_ ELEVATION_
## 4716    VT MOUNT MANSFIELD      44.53      -72.82      3950

ne.m@data[which.max(kcv.ok$residual),2:6]

##      STATE      STATION_NA LATITUDE_D LONGITUDE_ ELEVATION_
## 3888    PA WILKES BRE SCTN AP AVOCA      41.33      -75.73      930
```


The largest under-predictions (i.e., most negative LOOCV residual) is Mt. Mansfield (VT). This station is the only one on the higher parts of the Green Mountains, and is geographically near to Lake Champlain to its west. The largest over-prediction is the AVP airport near Wilkes-Barre (PA). This station is on an exposed plateau above the nearby Lackawanna Valley stations.

Challenge: Find the optimal local neighbourhood for OK by testing several numbers of nearest neighbours and/or maximum distance, and comparing their cross-validation statistics.

Challenge: Much of the LOOCV errors from OK could be because local elevation is not taken into account. There can be a substantial elevation variation within the nearest 24 stations. Apply KED (§9) with elevation (or its square root) as the covariable and see how much this improves the LOOCV error.

14.4 Inverse-distance interpolation

In case a model of spatial dependence can not be built, a fallback is a model-free empirical interpolation. There are many such, none of which have any theoretical basis, but which can be tested by cross-validation:

- inverse distance weighting, to some power (generally one or two);
- average of observations within some radius;
- average of some number of nearest observations.

The `krige` function with a null model computes inverse-distance interpolation; another way is with the convenience function `idw` function, which has an optional `idp` “inverse distance power” argument gives the decay power; the default is two, i.e., quadratic decay. Since there is no model, and hence no optimization of the prediction variance, there is here no internal estimate of uncertainty.

```
k.idw <- idw(ANN_GDD50 ~ 1, locations=ne.m, newdata=dem.ne.m.sp,
             idp=2, nmax=24)

## [inverse distance weighted interpolation]

summary(k.idw$var1.pred)

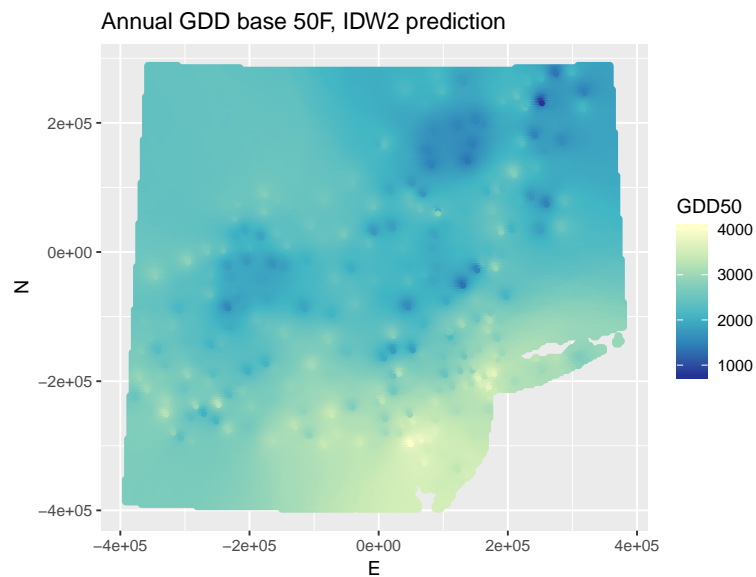
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  851.8  2126.7  2394.1  2441.2  2692.3  4013.0

summary(k.ok$var1.pred)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  1483    2157    2456    2478    2715    3740
```

Plot the result:

```
dem.ne.m.df$pred.idw <- k.idw$var1.pred
display.prediction.map("pred.idw",
                       "Annual GDD base 50F, IDW2 prediction",
                       "GDD50")
```



LOOCV also can be applied to IDW:

```
kcv.idw <- krige.cv(ANN_GDD50 ~ 1, locations=ne.m, model=NULL)

summary(kcv.idw$residual)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -1183.393 -204.227   -2.284   -23.102   176.446   745.680

summary(kcv.ok$residual)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -1126.6606 -154.4425   -2.0428   -0.8718   176.5507   813.0746

(loocv.idw.rmse <- sqrt(sum(kcv.idw$residual^2)/length(kcv.idw$residual)))

## [1] 300.457

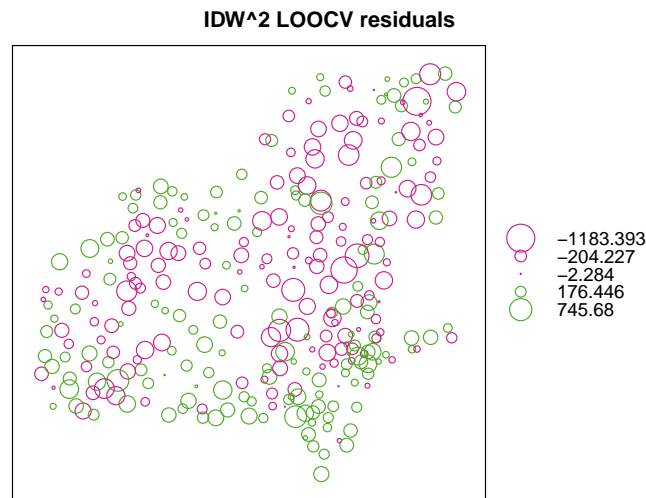
(loocv.ok.rmse)

## [1] 268.1796
```

The OK cross-validation is somewhat better than that from IDW.

TASK 160 : Display a bubble plot of the LOOCV residuals. •

```
bubble(kcv.idw, zcol="residual", pch=1, main="IDW^2 LOOCV residuals")
```



Challenge: Find the optimal power for IDW by testing several inverse-distance powers, and comparing their cross-validation statistics.

14.5 Thiessen polygons

The simplest way to predict a variable at one position is to use the value of that variable at the **nearest neighbour** in **geographic** space. For climate variables with a fairly dense network, as in this example, this is a common procedure: look at the climate record for the nearest station and consider that the local climate can not differ “too much”.

A spatial expression of this is to divide the prediction area into **Thiessen polygons**⁴², where each location is in a polygon whose centroid is its nearest neighbour. The advantage of this approach is that it requires no statistical model; in particular, there is no assumption of second-order stationarity as required by kriging.

TASK 161 : Compute and display the Thiessen polygons over the study area. •

The computation is with the `voronoi` “Voronoi tessellation” function of the `dismo` “Species Distribution Modeling” package. This in turn uses the `deldir` “Delaunay Triangulation and Dirichlet (Voronoi) Tessellation” function of the `deldir` package.

We specify the optional `ext` “extent” argument, to limit the polygons to the bounding box of the study area.

```
library(dismo)
v <- voronoi(ne.m, ext=t(bbox(state.ne.m)))
```

⁴² also known as a Voronoi tessellation or a Dirichlet tessellation

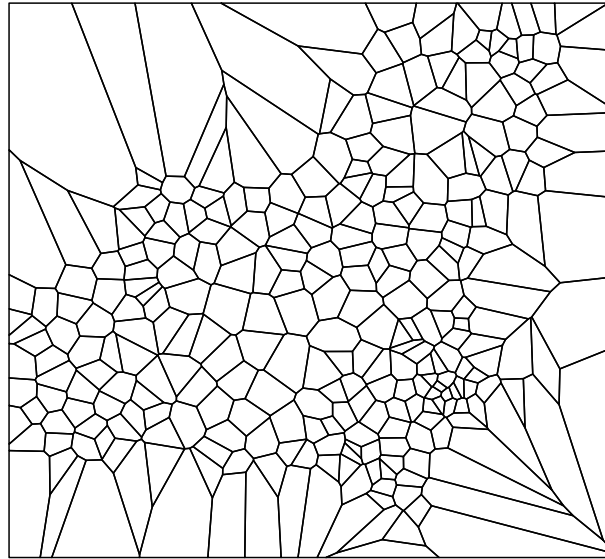
```

class(v)

## [1] "SpatialPolygonsDataFrame"
## attr("package")
## [1] "sp"

plot(v)

```



TASK 162 : Clip the tessellation to the four states. •

Various polygon geometry operations are implemented by the rgeos “Interface to Geometry Engine - Open Source (GEOS)” package.

The first operation is the **union** of the polygons from the four states, using the `gUnaryUnion` function. This is because the state boundaries should not come between an area and its nearest climate station.

```

require(rgeos)
area.ne.m <- gUnaryUnion(state.ne.m)
class(area.ne.m)

## [1] "SpatialPolygons"
## attr("package")
## [1] "sp"

plot(area.ne.m)

```



Note that there are still several polygons, because of the islands (Staten, Plum etc.), but all have the same identification.

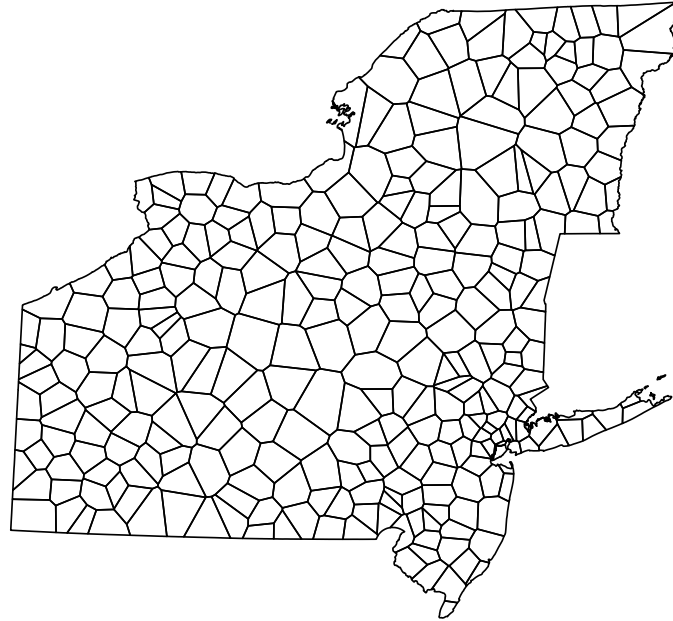
Now we intersect this area with the Thiessen polygons, using the `gIntersection` function.

Note: The optional argument `byid` seems to be necessary in this case, perhaps because of the complicated topology of the island polygons. If `byid=TRUE` the function is applied separately across polygons with different identifications, rather than the entire object.

```
v2 <- gIntersection(v, area.ne.m, byid=TRUE)
class(v2)

## [1] "SpatialPolygons"
## attr("package")
## [1] "sp"

plot(v2)
```



This map now can be used for prediction. Simply, the entire area of each polygon is predicted with the value from the nearest station, i.e., its centroid.

TASK 163 : Predict GDD50 over the study area by assigning the GDD50 from the centroid weather station to the entire polygon. •

The `over` “spatial overlay” method of the `sp` package queries its second argument (layer from which attributes are queried) and assigns attribute values to the geometries in the first argument. In this case each polygon covers an area; within that area is only one climate station in the `SpatialPointsDataFrame` object `ne.m`, so the value of the attributes at that point will be assigned to the polygon.

Note: For overlays where several points are in the same polygon, a user-specified function must be applied to return a single value.

```
v3 <- over(v2, ne.m)
class(v3)

## [1] "data.frame"
```

Now we convert this dataframe to a `SpatialPolygonsDataFrame`, using

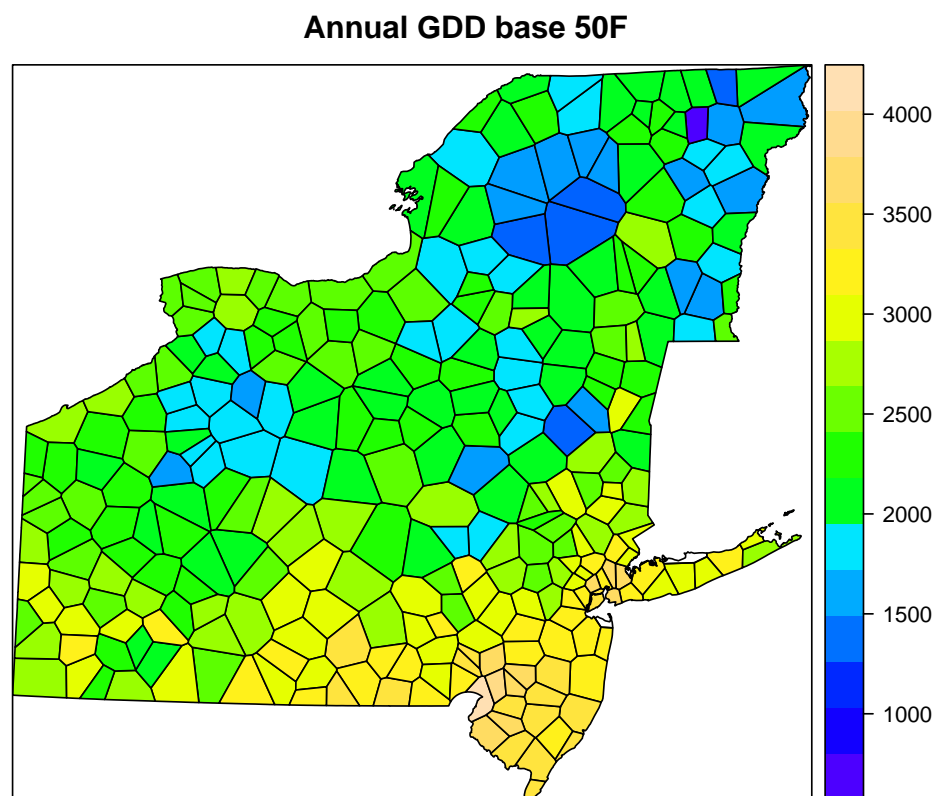
the attributes from the overlay as the data, and the geometry of the Thiessen polygons as the polygons.

```
v4 <- SpatialPolygonsDataFrame(v2, data=v3)
names(v4)

## [1] "STATION_ID"      "STATE"           "STATION_NAME"
## [4] "LATITUDE_D"      "LONGITUDE_D"     "ELEVATION_D"
## [7] "OID_"            "COOP_ID"         "STATE_1"
## [10] "STN_NAME"        "LAT_DD"          "LONG_DD"
## [13] "ELEV_FT"         "JAN_GDD50"       "FEB_GDD50"
## [16] "MAR_GDD50"       "APR_GDD50"       "MAY_GDD50"
## [19] "JUN_GDD50"       "JUL_GDD50"       "AUG_GDD50"
## [22] "SEP_GDD50"       "OCT_GDD50"       "NOV_GDD50"
## [25] "DEC_GDD50"       "ANN_GDD50"       "ols.resid"
## [28] "gls.resid"       "diff.gls.ols.resid" "resid.m.g.xy"
## [31] "rf.resid"        "rf.resid.oob"
```

TASK 164 : Plot the Thiessen polygons with their predicted values. •

```
sppplot(v4, zcol="ANN_GDD50", col.regions=topo.colors(64),
main="Annual GDD base 50F")
```



The prediction at any location can easily be found with a point overlay. Suppose a point is identified in Google Earth or using field GPS by its geographic coordinates. We make a `SpatialPoints` object, assign the correct CRS, and transform it to the CRS used in this project.

```
pt <- SpatialPoints(coords=cbind(-76.402175, 42.453271))
proj4string(pt) <- CRS("+init=EPSG:4326")
pt <- spTransform(pt, CRS(proj4string(ne.m)))
coordinates(pt)

##      coords.x1 coords.x2
## [1,] -33055.67 -5118.213
```

We then overlay with the predictive map:

```
pt.sp <- SpatialPointsDataFrame(pt, over(pt, v4))
pt.sp$ANN_GDD50

## [1] 2160
```

This shows the nearest station, its coordinates, and all the climate variables. So for the selected location, the predicted annual GDD50 is 2160.

Compare with the various other predictions:

```
dem.ne.m.spdf <- SpatialPixelsDataFrame(coordinates(dem.ne.m.sp),
                                          data=dem.ne.m.df)
proj4string(dem.ne.m.spdf) <- proj4string(ne.m)
pt.pred <- over(pt, dem.ne.m.spdf)
names(pt.pred)

## [1] "ELEVATION_"      "E"                "N"
## [4] "pred.ols"         "pred.gls"         "diff.gls.ols"
## [7] "ok.gls.resid"     "ok.gls.resid.var" "pred.rkgls"
## [10] "pred.ked"         "pred.ked.sd"      "diff.rkgls.ked"
## [13] "pred.ked.nn"      "pred.ked.nn.sd"   "diff.ked"
## [16] "diff.okrk.ked.pred" "pred.gam"         "pred.gam.se"
## [19] "diff.gls.gam"     "diff.rkgls.gam"   "pred.rt"
## [22] "pred.rf"          "diff.gls.rf"      "pred.ranger"
## [25] "diff.ranger.rf"   "pred.cubist"       "pred.ok"
## [28] "sd.ok"            "pred.idw"

pt.pred[,c("pred.ols", "pred.gls", "pred.rkgls", "pred.ked", "pred.rf")]

##   pred.ols pred.gls pred.rkgls pred.ked pred.rf
## 1 1998.056 1997.099  1978.911 1976.564 1958.331
```

These are all quite a bit lower than the nearest-neighbour prediction, because the elevation is not taken into account by the nearest-neighbour, just the planar distance. The selected point has a much higher elevation than the nearest station:

```
pt.sp$ELEVATION_

## [1] 960

pt.pred$ELEVATION_

## [1] 1459.784
```

14.5.1 Accuracy assessment

How accurate is a Thiessen polygon map? Ideally we would have a randomly-distributed sample of independent evaluation stations, but since we do not, we can apply the Thiessen polygon version of Leave-One-Out cross-validation (LOOCV). This is simply finding the nearest neighbour of each known station, and predicting at the known station from its neighbour.

TASK 165 : Build a matrix of inter-station distances. •

This can be done with the `spDists` “spatial distances” function of the `sp` package.

```
dm <- spDists(ne.m)
str(dm)

##  num [1:305, 1:305] 0 14343 64767 34083 159361 ...
```

TASK 166 : Find the nearest neighbour of each station, i.e., the one at closest distance. •

We use the `apply` function to apply the `which.min` “which minimum?” function across the rows of the distance matrix. However, all the diagonals are zero (distance between a station and itself), so we first have to put a large distance on diagonal so the stations themselves won’t come out as minima.

```
diag(dm) <- max(dm)*1.1
nn <- apply(dm, 1, which.min)
str(nn)

##  int [1:305] 2 1 248 24 267 26 11 4 38 31 ...
```

This gives the index of each station’s nearest neighbour. For example, station 1’s nearest neighbour is station 2, and vice-versa.

TASK 167 : Predict each station from its nearest neighbour, and plot the regional predictions. •

```
nn.gdd <- ne.m@data[nn, "ANN_GDD50"]
str(nn.gdd)

##  int [1:305] 3534 3310 3834 3553 2652 2692 2879 3450 2874 3250 ...
```

TASK 168 : Compare to the actual value and compute evaluation statistics. •

```
obs <- ne.m@data[, "ANN_GDD50"]
summary(obs)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      795   2100   2463   2518   2930   4021

summary(nn.gdd)

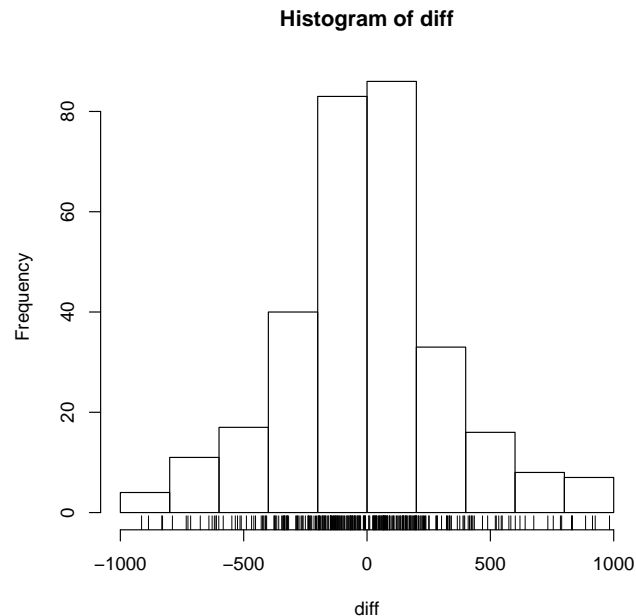
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      795   2113   2447   2522   2982   3834

summary(diff <- (obs - nn.gdd))

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -914.000 -193.000   -9.000   -4.393  175.000  983.000

hist(diff)
rug(diff)
(me <- mean(diff))
```

```
## [1] -4.393443
(rmse <- sqrt(sum(diff^2)/length(diff)))
## [1] 332.5643
```



We can compare these to the evaluation statistics at the same points from the Random Forest out-of-bag computed above (§12.2.1).

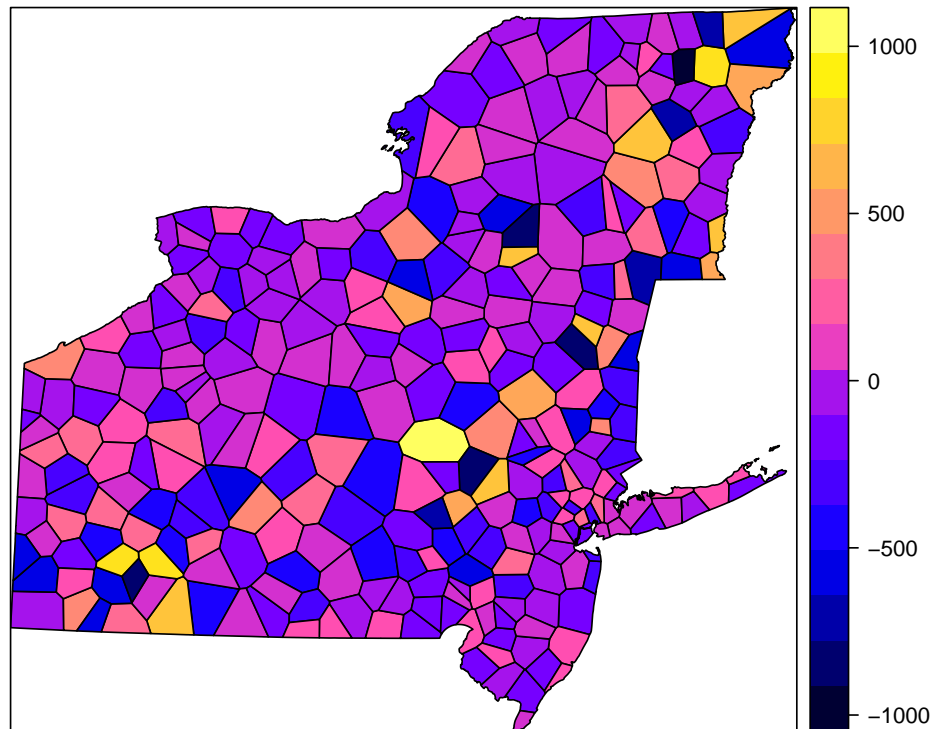
```
rf.oob.me; rf.oob.rmse
## [1] -2.06279
## [1] 11.58191
```

Clearly, the predictions from Thiessen polygons, with this density of climate stations, are quite poor, compared to the data-driven model using coördinates and elevation.

TASK 169: Plot the LOOCV errors. •

```
v4$resid <- diff
spplot(v4, zcol="resid", col.regions=bpy.colors(32),
       main="Cross-validation error, Thiessen polygon prediction")
```

Cross-validation error, Thiessen polygon prediction



The largest errors are where the topography changes rapidly between stations, for example at and next to Mt. Mansfield (VT) - these two polygons have the largest over- and under-predictions, respectively. The smallest errors are where elevations do not change much between neighbour stations, for example on the Lake Ontario plain.

15 Comparing the spatial patterns of two climate variables

An interesting question for geographers is how the spatial patterns of two related climate variables are related. The previous part of this exercise have analyzed growing degree days, base 50°F (GDDG50), which is related to the amount of heat available for crop development. The provided dataset includes several other variables clearly related to heat, in the following files:

- gdd40_7100j : growing degree days, base 40°F (applies to C3 crops such as spring wheat and barley)
- maat7100 : mean air temperature °F
- frz32_7100j : length of frost-free period, consecutive days above 32 °F (for frost-sensitive crops)

frz28_7100j : length of frost-free period, consecutive days above 28 °F (for frost-tolerant crops)

extmin_7100j : extreme minimum temperature °F

These all have annual and monthly records, averages over 1971-2000. Since they are all related to heat, the overall geographic pattern should be similar, but there are physical differences between them. In this section we see how to compare regional maps to understand the spatial patterns of these differences.

We choose to compare mean annual temperature, °F, which is an overall measure of heat over the year, with GDD50.

TASK 170 : Import the temperature records for the entire USA into a temporary data frame. •

```
tmp <- readOGR(dsn=".", layer="maat7100",
               integer64="allow.loss")

## OGR data source with driver: ESRI Shapefile
## Source: "/Users/rossiter/data/edu/dgeostats/ex/ds/NEweather", layer: "maat7100"
## with 5556 features
## It has 19 fields
```

TASK 171 : Restrict this dataframe to the the four selected states. •

```
ix <- (tmp$STATE %in% c("NY", "NJ", "PA", "VT"))
tmp <- tmp[ix,]
names(tmp)

## [1] "STATION_ID" "STATE" "STATION_NA" "LATITUDE_D" "LONGITUDE_"
## [6] "ELEVATION_" "JAN_TEMP" "FEB_TEMP" "MARCH_TEMP" "APRIL_TEMP"
## [11] "MAY_TEMP" "JUNE_TEMP" "JULY_TEMP" "AUG_TEMP" "SEP_TEMP"
## [16] "OCT_TEMP" "NOV_TEMP" "DEC_TEMP" "ANN_TNORM_"
```

TASK 172 : Transform this frame to to the same metric CRS as the frame used for the GDD50. •

```
strwrap(proj4string(dem.ne.m))

## [1] "+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76"
## [2] "+ellps=WGS84 +units=m"

tmp.m <- spTransform(tmp, ne.crs)
proj4string(tmp.m)

## [1] "+proj=aea +lat_0=42.5 +lat_1=39 +lat_2=44 +lon_0=-76 +ellps=WGS84 +units=m"
```

TASK 173 : Transfer the temperature records to the same data frame with the GDD50, and then remove the temporary data frame. •

We use the cbind “bind columns” function for this.

```
names(ne.m)

## [1] "STATION_ID" "STATE" "STATION_NA"
## [4] "LATITUDE_D" "LONGITUDE_" "ELEVATION_"
```

```
## [7] "OID_" "COOP_ID" "STATE_1"
## [10] "STN_NAME" "LAT_DD" "LONG_DD"
## [13] "ELEV_FT" "JAN_GDD50" "FEB_GDD50"
## [16] "MAR_GDD50" "APR_GDD50" "MAY_GDD50"
## [19] "JUN_GDD50" "JUL_GDD50" "AUG_GDD50"
## [22] "SEP_GDD50" "OCT_GDD50" "NOV_GDD50"
## [25] "DEC_GDD50" "ANN_GDD50" "ols.resid"
## [28] "gls.resid" "diff.gls.ols.resid" "resid.m.g.xy"
## [31] "rf.resid" "rf.resid.oob"

ne.m <- cbind(ne.m, tmp.m[,7:19])
names(ne.m)

## [1] "STATION_ID" "STATE" "STATION_NA"
## [4] "LATITUDE_D" "LONGITUDE_" "ELEVATION_"
## [7] "OID_" "COOP_ID" "STATE_1"
## [10] "STN_NAME" "LAT_DD" "LONG_DD"
## [13] "ELEV_FT" "JAN_GDD50" "FEB_GDD50"
## [16] "MAR_GDD50" "APR_GDD50" "MAY_GDD50"
## [19] "JUN_GDD50" "JUL_GDD50" "AUG_GDD50"
## [22] "SEP_GDD50" "OCT_GDD50" "NOV_GDD50"
## [25] "DEC_GDD50" "ANN_GDD50" "ols.resid"
## [28] "gls.resid" "diff.gls.ols.resid" "resid.m.g.xy"
## [31] "rf.resid" "rf.resid.oob" "JAN_TEMP"
## [34] "FEB_TEMP" "MARCH_TEMP" "APRIL_TEMP"
## [37] "MAY_TEMP" "JUNE_TEMP" "JULY_TEMP"
## [40] "AUG_TEMP" "SEP_TEMP" "OCT_TEMP"
## [43] "NOV_TEMP" "DEC_TEMP" "ANN_TNORM_"

rm(tmp, tmp.m)
```

TASK 174 : Make a data frame version of the dataset.

```
ne.df <- as(ne.m, "data.frame")
```

To compare two variables on the same scale, it is necessary to **standardize** them, by subtracting the mean and dividing by the standard deviation.

TASK 175 : Compute standardized versions of the GDD50 and ANN_TNORM_ fields.

```
summary(ne.m$ANN_GDD50)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      795   2100   2463   2518   2930   4021

summary(ne.m$ANN_TNORM_)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      34.70  45.20  47.80  48.03  51.00  56.30

gdd50.std <- (ne.m$ANN_GDD50 - mean(ne.m$ANN_GDD50))/sd(ne.m$ANN_GDD50)
t.ann.std <- (ne.m$ANN_TNORM_ - mean(ne.m$ANN_TNORM_))/sd(ne.m$ANN_TNORM_)
summary(gdd50.std)

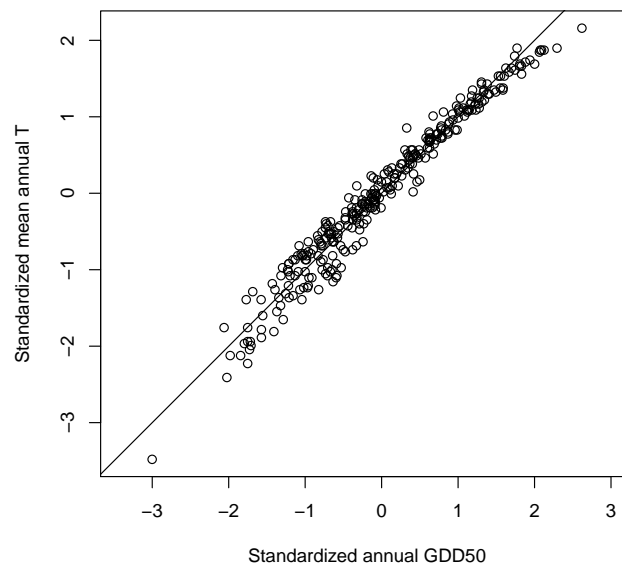
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##     -3.00125 -0.72747 -0.09499  0.00000  0.71869  2.61961

summary(t.ann.std)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##     -3.4800 -0.7388 -0.0600  0.0000  0.7754  2.1591
```

TASK 176 : Compare these to see how closely they are correlated. •

```
cor(ne.m$ANN_GDD50, ne.m$ANN_TNORM_)  
## [1] 0.9811989  
cor(gdd50.std, t.ann.std)  
## [1] 0.9811989  
plot(gdd50.std, t.ann.std, asp=1,  
      xlab="Standardized annual GDD50",  
      ylab="Standardized mean annual T")  
abline(0,1)
```



Note that the correlation is the same for both the standardized and un-standardized variables. It's interesting that the GDD50 are a bit higher than the mean annual T at both extremes of the 1:1 plot. So they are closely correlated, but not identical.

We choose to use RK-GLS (§8) and Random Forests (§12.2) as the prediction methods; the other methods could all be used.

15.1 Comparing variables with RK-GLS

TASK 177 : Build OLS models of both standardized variables from the two coördinates and square root of elevation. Compare their adjusted R^2 and coefficients. •

```
summary(m.ols.t <- lm(t.ann.std~ sqrt(ELEVATION_)+N+E, data=ne.df))  
##  
## Call:  
## lm(formula = t.ann.std ~ sqrt(ELEVATION_) + N + E, data = ne.df)  
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.6861 -0.1944 -0.0391  0.1842  0.8651
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.096e+00  5.505e-02  19.913 < 2e-16 ***
## sqrt(ELEVATION_) -5.321e-02  1.829e-03 -29.088 < 2e-16 ***
## N             -3.614e-06  1.199e-07 -30.148 < 2e-16 ***
## E             -9.150e-07  1.134e-07  -8.071 1.68e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2821 on 301 degrees of freedom
## Multiple R-squared:  0.9212, Adjusted R-squared:  0.9204
## F-statistic: 1173 on 3 and 301 DF, p-value: < 2.2e-16

summary(m.ols.gdd <- lm(gdd50.std ~ sqrt(ELEVATION_)+N+E, data=ne.df))

##
## Call:
## lm(formula = gdd50.std ~ sqrt(ELEVATION_) + N + E, data = ne.df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.9447 -0.2365 -0.0221  0.2303  1.0474
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.374e+00  6.709e-02  20.483 < 2e-16 ***
## sqrt(ELEVATION_) -6.184e-02  2.230e-03 -27.735 < 2e-16 ***
## N             -2.895e-06  1.461e-07 -19.814 < 2e-16 ***
## E             -9.386e-07  1.382e-07  -6.793 5.88e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3439 on 301 degrees of freedom
## Multiple R-squared:  0.8829, Adjusted R-squared:  0.8818
## F-statistic: 756.6 on 3 and 301 DF, p-value: < 2.2e-16

coef(m.ols.t)

##      (Intercept) sqrt(ELEVATION_)      N      E
##      1.096150e+00   -5.320974e-02   -3.613649e-06   -9.149859e-07

coef(m.ols.gdd)

##      (Intercept) sqrt(ELEVATION_)      N      E
##      1.374233e+00   -6.183727e-02   -2.894689e-06   -9.386282e-07

coef(m.ols.t)/coef(m.ols.gdd)

##      (Intercept) sqrt(ELEVATION_)      N      E
##      0.7976448    0.8604801    1.2483723    0.9748119
```

Q44: *Do the two models explain the same amount of spatial variability by the same predictors? If the two variables have the same spatial structure, what should be the ratio of the coefficients? Is that the case here?*

[Jump to A44](#) •

TASK 178: Display bubble plots of the residuals, and 1:1 plots of actual vs. fitted, for both variables. •

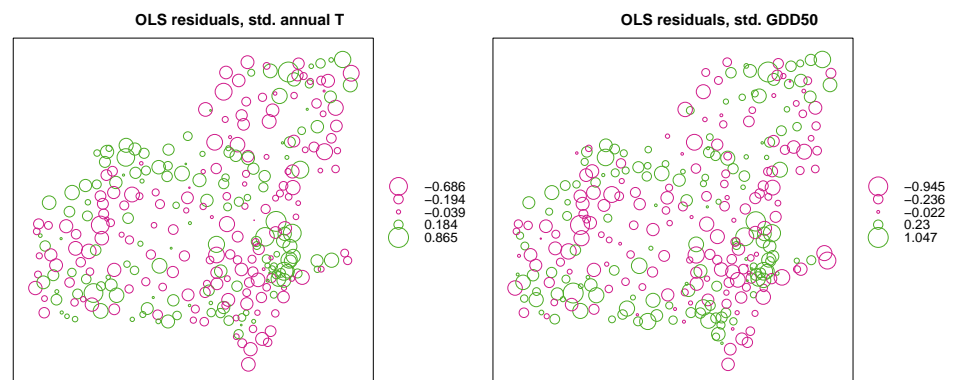
```
summary(ne.m$ols.resid.t <- residuals(m.ols.t))

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -0.6861 -0.1944 -0.0391  0.0000  0.1842  0.8651

summary(ne.m$ols.resid.gdd <- residuals(m.ols.gdd))

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -0.9447 -0.2365 -0.0221  0.0000  0.2303  1.0474

p1 <- bubble(ne.m, zcol="ols.resid.t", pch=1,
             main="OLS residuals, std. annual T")
p2 <- bubble(ne.m, zcol="ols.resid.gdd", pch=1,
             main="OLS residuals, std. GDD50")
print(p1, split=c(1,1,2,1), more=T)
print(p2, split=c(2,1,2,1), more=F)
```



TASK 179 : Compare the residuals with a bubble plot of their differences. •

```
summary(ne.m$ols.resid.diff <- ne.m$ols.resid.t - ne.m$ols.resid.gdd)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -0.583312 -0.103252  0.005628  0.000000  0.101299  0.610121

bubble(ne.m, zcol="ols.resid.diff", pch=1,
       main="Difference of residuals, T - GDD",
       col=c("#EF8A62", "#67A9CF"))
```




Now we begin to see the pattern, where the regional trend of annual mean temperature is higher or lower than that for GDD50.

Q45 : *What is the pattern of differences between the residuals from the models for the two variables? Identify the most interesting areas. What could be an explanation?* *Jump to A45 •*

TASK 180 : Model the residual spatial dependence for both variables, i.e., fit variogram models to the OLS residuals, for both variables. •

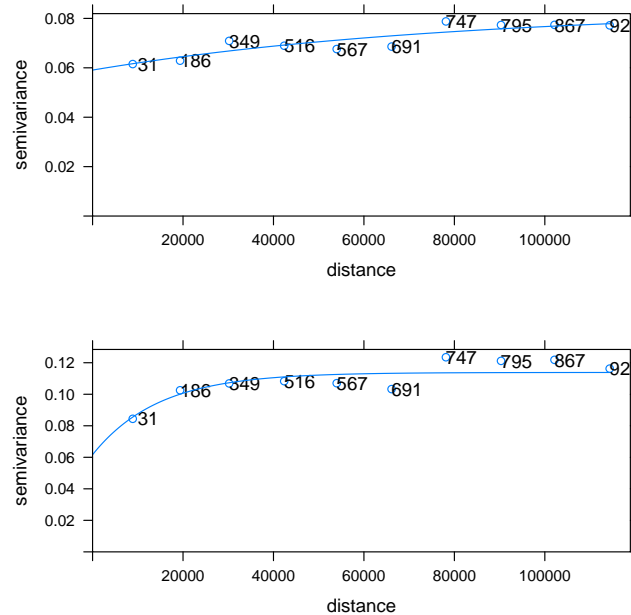
```
v.r.ols.t <- variogram(ols.resid.t ~ 1,
                      locations=ne.m, cutoff=120000, width=12000)
(vmf.r.ols.t <- fit.variogram(v.r.ols.t,
                             vgm(psill=0.1, model="Exp",
                                 range=20000, nugget=0.02)))

##    model      psill    range
## 1   Nug 0.05907886     0.0
## 2   Exp 0.02479735 80781.2

p1 <- plot(v.r.ols.t, pl=T, model=vmf.r.ols.t)
#
v.r.ols.gdd <- variogram(ols.resid.gdd ~ 1,
                       locations=ne.m, cutoff=120000, width=12000)
(vmf.r.ols.gdd <- fit.variogram(v.r.ols.gdd,
                               vgm(psill=0.12, model="Exp",
                                   range=20000, nugget=0.02)))

##    model      psill    range
## 1   Nug 0.06167406     0.00
## 2   Exp 0.05216480 14545.31

p2 <- plot(v.r.ols.gdd, pl=T, model=vmf.r.ols.gdd)
print(p1, split=c(1,1,1,2), more=T)
print(p2, split=c(1,2,1,2), more=F)
```



Q46 : *How strong is the local spatial dependence of the residuals from the two trend surfaces? Do the two trend surfaces have the same residual local spatial structure? If not, what is the difference? What does that imply about the trend surface model and spatial structure of the variables?*

[Jump to A46](#)

•

TASK 181 : Use this estimated spatial dependence among the residuals to re-fit the models with GLS.

•

We estimate starting values for the proportional nugget from the variogram fit.

```
# require(nlme)
(p.nugget <- vmf.r.ols.t[1,"psill"]/sum(vmf.r.ols.t[, "psill"]))

## [1] 0.7043578

m.gls.t <- gls(model=t.ann.std ~ sqrt(ELEVATION_) + N + E,
              data=ne.df,
              correlation=corExp(
                value=c(vmf.r.ols.t[2,"range"], p.nugget),
                form=~E + N,
                nugget=TRUE))

#
(p.nugget <- vmf.r.ols.gdd[1,"psill"]/sum(vmf.r.ols.gdd[, "psill"]))

## [1] 0.5417663

m.gls.gdd <- gls(model=gdd50.std ~ sqrt(ELEVATION_) + N + E,
                 data=ne.df,
                 correlation=corExp(
                   value=c(vmf.r.ols.gdd[2,"range"], p.nugget),
                   form=~E + N,
                   nugget=TRUE))
```

```

summary(m.gls.t)

## Generalized least squares fit by REML
## Model: t.ann.std ~ sqrt(ELEVATION_) + N + E
## Data: ne.df
##      AIC      BIC    logLik
## 159.9663 185.9161 -72.98317
##
## Correlation Structure: Exponential spatial correlation
## Formula: ~E + N
## Parameter estimate(s):
##      range      nugget
## 9.451201e+04 6.735544e-01
##
## Coefficients:
##              Value Std.Error   t-value p-value
## (Intercept)  1.0352866 0.07976232  12.97964  0e+00
## sqrt(ELEVATION_) -0.0510409 0.00217065 -23.51411  0e+00
## N            -0.0000035 0.00000027 -12.71021  0e+00
## E            -0.0000010 0.00000025  -4.02844  1e-04
##
## Correlation:
##              (Intr) s(ELEV N
## sqrt(ELEVATION_) -0.708
## N                0.323 -0.184
## E               -0.232  0.246 -0.335
##
## Standardized residuals:
##      Min      Q1      Med      Q3      Max
## -2.46015918 -0.62877289 -0.04851122  0.67873651  2.94128997
##
## Residual standard error: 0.2964007
## Degrees of freedom: 305 total; 301 residual

summary(m.gls.gdd)

## Generalized least squares fit by REML
## Model: gdd50.std ~ sqrt(ELEVATION_) + N + E
## Data: ne.df
##      AIC      BIC    logLik
## 274.5262 300.476 -130.2631
##
## Correlation Structure: Exponential spatial correlation
## Formula: ~E + N
## Parameter estimate(s):
##      range      nugget
## 7.741961e+04 6.175653e-01
##
## Coefficients:
##              Value Std.Error   t-value p-value
## (Intercept)  1.2724192 0.09520248  13.365400  0e+00
## sqrt(ELEVATION_) -0.0579702 0.00265657 -21.821480  0e+00
## N            -0.0000027 0.00000033  -8.161348  0e+00
## E            -0.0000011 0.00000030  -3.696353  3e-04
##
## Correlation:
##              (Intr) s(ELEV N
## sqrt(ELEVATION_) -0.734
## N                0.333 -0.188
## E               -0.246  0.254 -0.346
##
## Standardized residuals:
##      Min      Q1      Med      Q3      Max
## -2.847476938 -0.591834398  0.001698498  0.684143300  2.811395726
##
## Residual standard error: 0.3611266
## Degrees of freedom: 305 total; 301 residual

```

TASK 182 : Compare the model coefficients. •

```
coefficients(m.gls.t)

##      (Intercept) sqrt(ELEVATION_)      N      E
## 1.035287e+00 -5.104088e-02 -3.489755e-06 -1.018102e-06

coefficients(m.gls.gdd)

##      (Intercept) sqrt(ELEVATION_)      N      E
## 1.272419e+00 -5.797018e-02 -2.688525e-06 -1.117387e-06

coefficients(m.gls.t)/coefficients(m.gls.gdd)

##      (Intercept) sqrt(ELEVATION_)      N      E
## 0.8136364 0.8804679 1.2980184 0.9111453
```

This shows the ratio of the coefficients for the two variables. If they would have the same regional spatial structure (trend surface), these values would all be 1.

TASK 183 : Compare the correlation structures fit by REML. •

```
intervals(m.gls.t)$corStruct

##      lower      est.      upper
## range 2.350637e+04 9.451201e+04 3.800042e+05
## nugget 4.282451e-01 6.735544e-01 8.503852e-01
## attr("label")
## [1] "Correlation structure:"

intervals(m.gls.gdd)$corStruct

##      lower      est.      upper
## range 1.948023e+04 7.741961e+04 307686.16717
## nugget 4.076572e-01 6.175653e-01 0.79119
## attr("label")
## [1] "Correlation structure:"

intervals(m.gls.t)$corStruct/intervals(m.gls.gdd)$corStruct

##      lower      est.      upper
## range 1.206678 1.220776 1.235038
## nugget 1.050503 1.090661 1.074818
## attr("label")
## [1] "Correlation structure:"
```

Again, all these ratios would be 1 if the structures were identical. The mean annual temperature has proportionally longer range and higher nugget than GDD50. The range parameters are 77–95 km, for an effective range of about 240–300 km; this is quite a bit longer than what we estimated by eye from the residual variograms from the OLS fit.

TASK 184 : Add the GLS model residuals to the spatial data frame. •

```
ne.m$gls.resid.t <- residuals(m.gls.t)
ne.m$gls.resid.gdd <- residuals(m.gls.gdd)
```

TASK 185 : Display the fitted model of spatial correlation with the

empirical variogram of the GLS residuals.

We first convert the correlation structure found by `gls` to a variogram model; the partial sill is estimated as the variance of the residuals, adjusted for the proportional nugget. We estimate the nugget as a proportion of this total sill.

```
(p.nugget <- intervals(m.gls.t)$corStruct["nugget", "est."])
## [1] 0.6735544

(t.sill <- var(ne.m$gls.resid.t))
## [1] 0.08065724

(vmf.r.gls.t <- vgm(psill=t.sill*(1-p.nugget), model="Exp",
  range=intervals(m.gls.t)$corStruct["range", "est."],
  nugget=t.sill*p.nugget))

## model      psill      range
## 1   Nug 0.05432704      0.00
## 2   Exp 0.02633020 94512.01

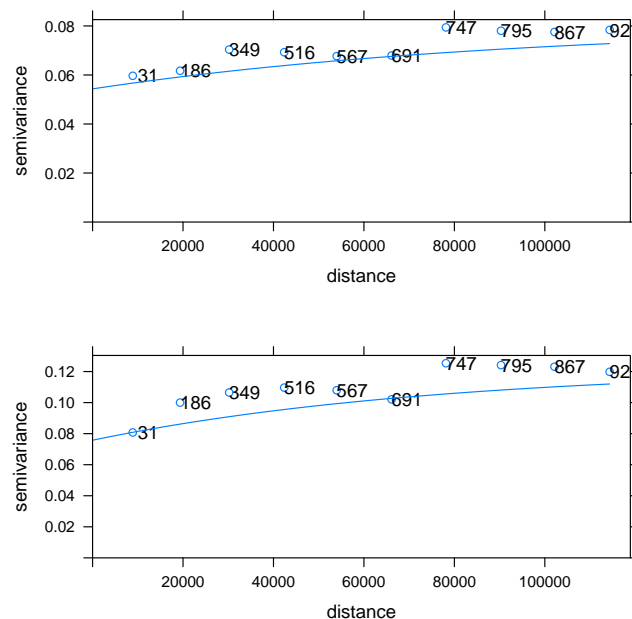
v.r.gls.t <- variogram(gls.resid.t ~ 1,
  locations=ne.m, cutoff=120000, width=12000)
p1 <- plot(v.r.gls.t, model=vmf.r.gls.t, pl=T)
#
(p.nugget <- intervals(m.gls.gdd)$corStruct["nugget", "est."])
## [1] 0.6175653

(t.sill <- var(ne.m$gls.resid.gdd))
## [1] 0.1226937

(vmf.r.gls.gdd <- vgm(psill=t.sill*(1-p.nugget), model="Exp",
  range=intervals(m.gls.gdd)$corStruct["range", "est."],
  nugget=t.sill*p.nugget))

## model      psill      range
## 1   Nug 0.07577136      0.00
## 2   Exp 0.04692231 77419.61

v.r.gls.gdd <- variogram(gls.resid.gdd ~ 1,
  locations=ne.m, cutoff=120000, width=12000)
p2 <- plot(v.r.gls.gdd, model=vmf.r.gls.gdd, pl=T)
print(p1, split=c(1,1,1,2), more=T)
print(p2, split=c(1,2,1,2), more=F)
```



TASK 186 : Predict over the regional grid with the GLS model and add the result to the dataframe. •

```
dem.ne.m.df$pred.gls.t <- predict(m.gls.t, newdata=dem.ne.m.df)
dem.ne.m.df$pred.gls.gdd <- predict(m.gls.gdd, newdata=dem.ne.m.df)
dem.ne.m.df$diff.gls.t.gdd <-
  dem.ne.m.df$pred.gls.t - dem.ne.m.df$pred.gls.gdd
```

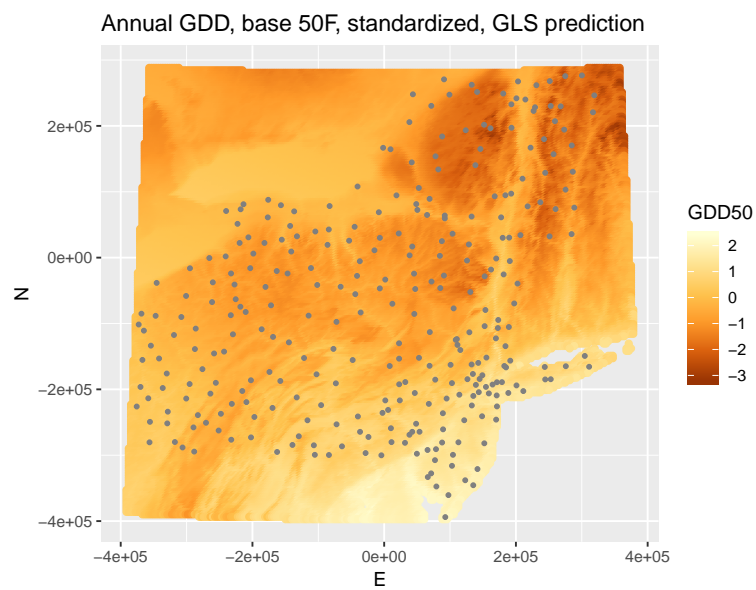
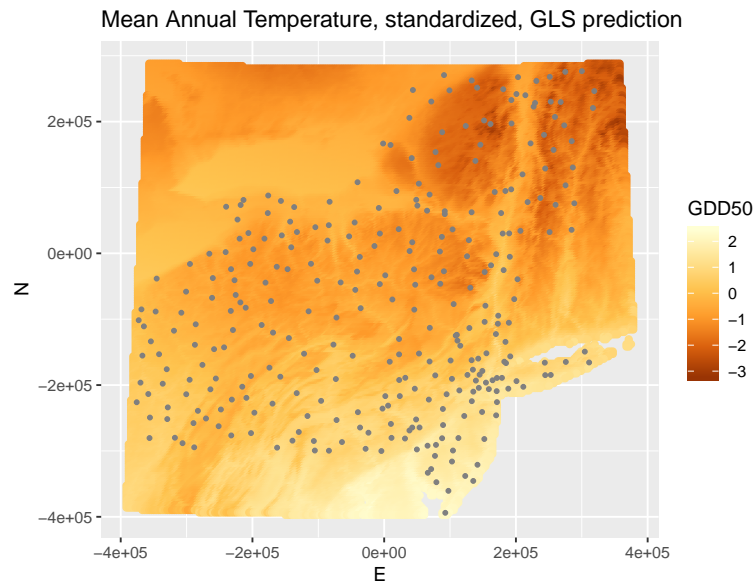
TASK 187 : Plot the two regional predictions on the same visual scale. •

We use a different palette from the non-standardized prediction maps, to emphasize that these are standardized values.

```
(std.pred.lim <- c(min(dem.ne.m.df[,c("pred.gls.t", "pred.gls.gdd")]),
  max(dem.ne.m.df[,c("pred.gls.t", "pred.gls.gdd")]))))

## [1] -3.203491 2.414317

display.prediction.map("pred.gls.t",
  "Mean Annual Temperature, standardized, GLS prediction",
  "GDD50", std.pred.lim, .palette="YlOrBr")
display.prediction.map("pred.gls.gdd",
  "Annual GDD, base 50F, standardized, GLS prediction",
  "GDD50", std.pred.lim, .palette="YlOrBr")
```

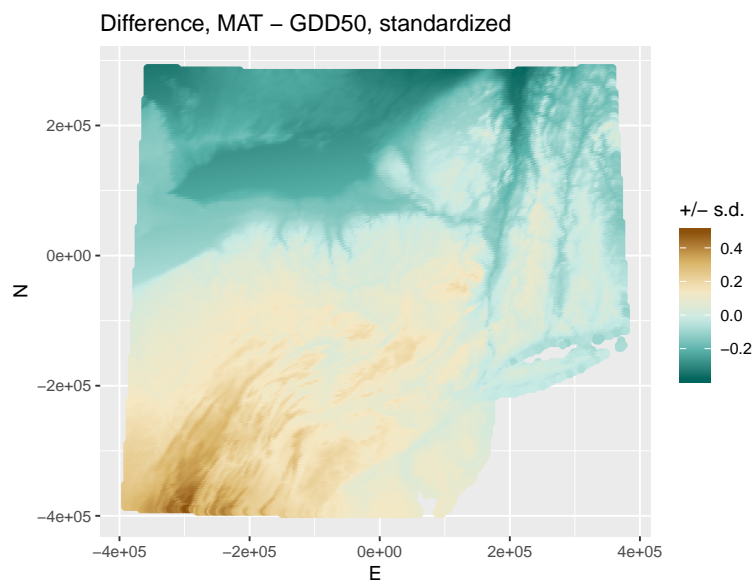


TASK 188 : Compute the differences between the two GLS predictions, add them to the data frame, summarize them, and display them as a map. •

```
summary(dem.ne.m.df$diff.gls.t.gdd <-
  dem.ne.m.df$pred.gls.t - dem.ne.m.df$pred.gls.gdd)

##      Min.   1st Qu.     Median       Mean   3rd Qu.      Max.
## -0.376528 -0.135178  0.014254 -0.009455  0.111029  0.492284

display.difference.map("diff.gls.t.gdd",
  "Difference, MAT - GDD50, standardized",
  "+/- s.d.",
  .palette="BrBG")
```



Q47 : Describe the spatial distribution of the differences between the MAT and GDD50 predictions. [Jump to A47](#) •

Although the GLS model residuals do not have strong spatial structure, still we can kriging them to improve the maps.

TASK 189 : Predict the deviations from the GLS trend surface at each location on the grid, using Ordinary Kriging (OK) of the GLS residuals for both standardized variables; display their summary, and display as maps. •

Recall we build variogram models of the residuals from the correlation structures found by `gls`.

```
ok.gls.resid.t <- krige(gls.resid.t ~ 1, loc=ne.m, newdata=dem.ne.m.sp,
  model=vmf.r.gls.t)

## [using ordinary kriging]

ok.gls.resid.gdd <- krige(gls.resid.gdd ~ 1, loc=ne.m, newdata=dem.ne.m.sp,
```



```

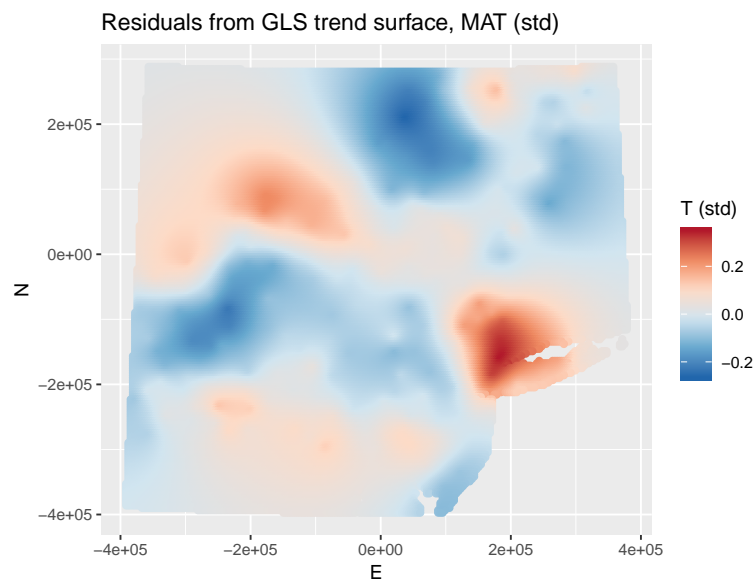
model=vmf.r.gls.gdd)

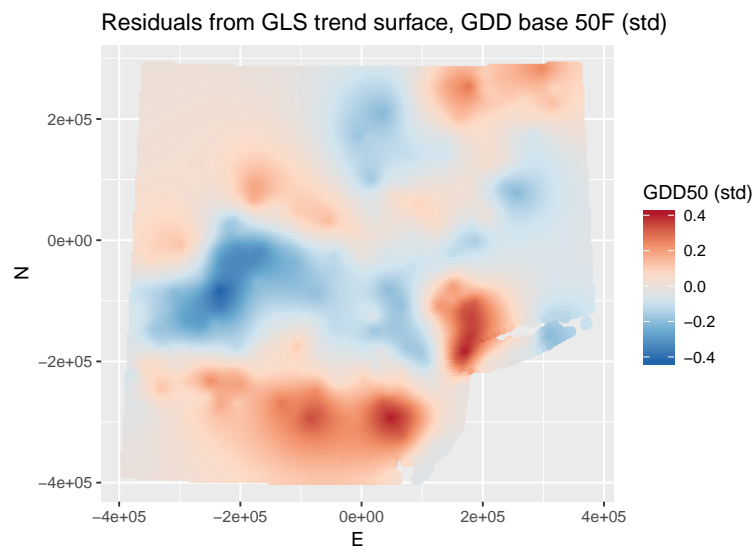
## [using ordinary kriging]

dem.ne.m.df$ok.gls.resid.t <- ok.gls.resid.t$var1.pred
dem.ne.m.df$ok.gls.resid.gdd <- ok.gls.resid.gdd$var1.pred

ggplot() +
  geom_point(aes(x=E, y=N, colour=ok.gls.resid.t), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Residuals from GLS trend surface, MAT (std)") +
  scale_colour_distiller(name="T (std)", space="Lab", palette="RdBu")
ggplot() +
  geom_point(aes(x=E, y=N, colour=ok.gls.resid.gdd), data=dem.ne.m.df) +
  xlab("E") + ylab("N") + coord_fixed() +
  ggtitle("Residuals from GLS trend surface, GDD base 50F (std)") +
  scale_colour_distiller(name="GDD50 (std)", space="Lab", palette="RdBu")

```





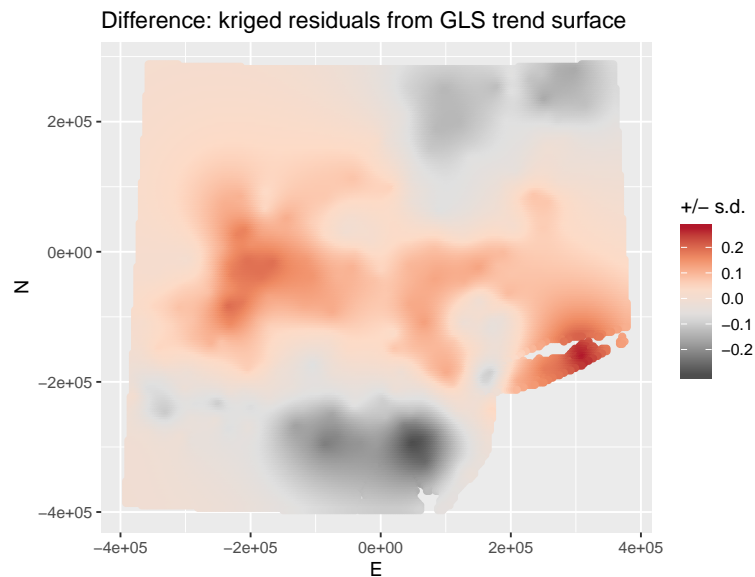
TASK 190 : Display the differences in the OK of the standardized residuals.

This shows where each variable is more or less adjusted.

```
summary(dem.ne.m.df$ok.gls.resid.diff.t.gdd <-
  dem.ne.m.df$ok.gls.resid.t - dem.ne.m.df$ok.gls.resid.gdd)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -0.297577 -0.053065  0.016131  0.003648  0.061518  0.273665

display.difference.map("ok.gls.resid.diff.t.gdd",
  "Difference: kriged residuals from GLS trend surface",
  "+/- s.d.",
  .palette="RdGy")
```



TASK 191 : Add the kriged GLS residuals to the trend surfaces for a final GLS-RK prediction. Display the two maps. •

```
summary(dem.ne.m.df$pred.rkgls.std.t <-
  dem.ne.m.df$pred.gls.t + dem.ne.m.df$ok.gls.resid.t)

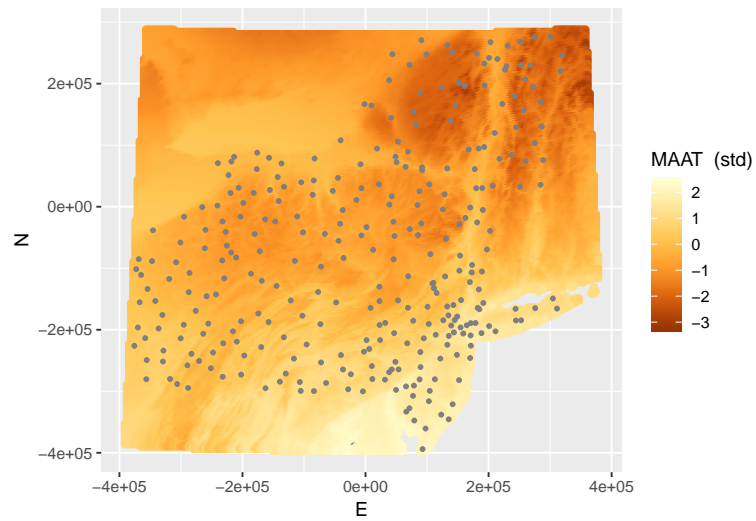
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -3.1610 -0.9539 -0.3710 -0.3092  0.2745   2.4335

summary(dem.ne.m.df$pred.rkgls.std.gdd <-
  dem.ne.m.df$pred.gls.gdd + dem.ne.m.df$ok.gls.resid.gdd)

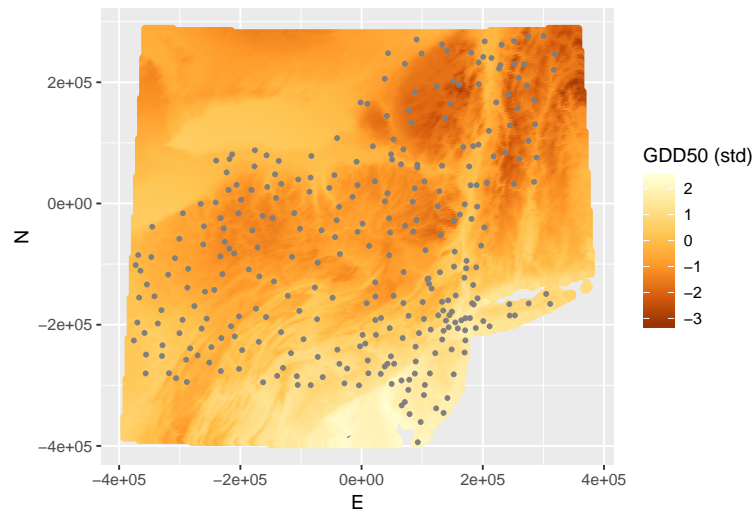
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## -3.1660 -0.9465 -0.4005 -0.3033  0.2757   2.4660

display.prediction.map("pred.rkgls.std.t",
  "GLS-RK prediction, Mean Annual Temperature, standardized",
  "MAAT (std)", std.pred.lim, .palette="YlOrBr")
display.prediction.map("pred.rkgls.std.gdd",
  "GLS-RK prediction, Annual GDD, base 50F, standardized",
  "GDD50 (std)", std.pred.lim, .palette="YlOrBr")
```

GLS-RK prediction, Mean Annual Temperature, standardized



GLS-RK prediction, Annual GDD, base 50F, standardized



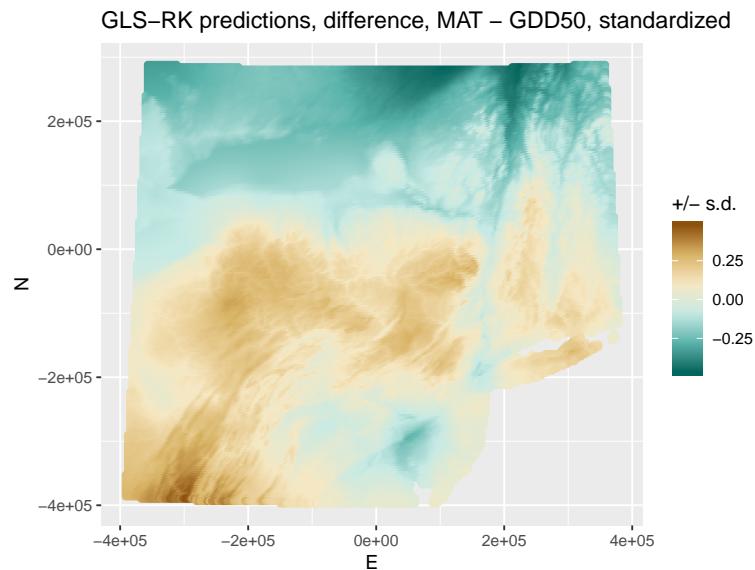
TASK 192 : Compute and display the difference between the two GLS-RK predictions. •

```
#
summary(dem.ne.m.df$diff.rkgls.std.t.gdd <-
```

```
(dem.ne.m.df$pred.rkgls.std.t - dem.ne.m.df$pred.rkgls.std.gdd))

##      Min.   1st Qu.    Median      Mean   3rd Qu.      Max.
## -0.462388 -0.142287  0.013942 -0.005807  0.134904  0.477189

display.difference.map("diff.rkgls.std.t.gdd",
  "GLS-RK predictions, difference, MAT - GDD50, standardized",
  "+/- s.d.",
  .palette="BrBG")
```



There is a considerable difference between the two normalized variables. Normalized MAT is higher than normalized GDD50 in the Appalachian, Taconic and Allegheny mountains, and on eastern Long Island; the reverse is the case in Adirondacks and Green Mountains, and especially in the Lake Champlain valley. To the north, the mean annual temperature is lowered by the very cold winter months. If the comparison were with growing season mean temperature, perhaps the results would be more similar.

15.2 Comparing variables with Random Forests

Another way to compare the variables is with their RF models. This does not show linear model coefficients or local spatial correlation structure, but does show variable importance, and also produces two maps.

TASK 193 : Build RF models of the two normalized variables. •

```
m.rf.std.t <- randomForest(t.ann.std ~ ELEVATION_ + N + E,
  data=ne.df, ntree=1200,
  importance=TRUE)
```

```
m.rf.std.gdd <- randomForest(gdd50.std ~ ELEVATION_ + N + E,
                             data=ne.df, ntree=1200,
                             importance=TRUE)
```

TASK 194 : Compare the variable importance of the two models. •

We use the `importance` function of the `randomForest` package:

```
randomForest::importance(m.rf.std.t)

##              %IncMSE IncNodePurity
## ELEVATION_  66.87351      104.73476
## N           91.23734      143.96688
## E           55.61519       48.56437

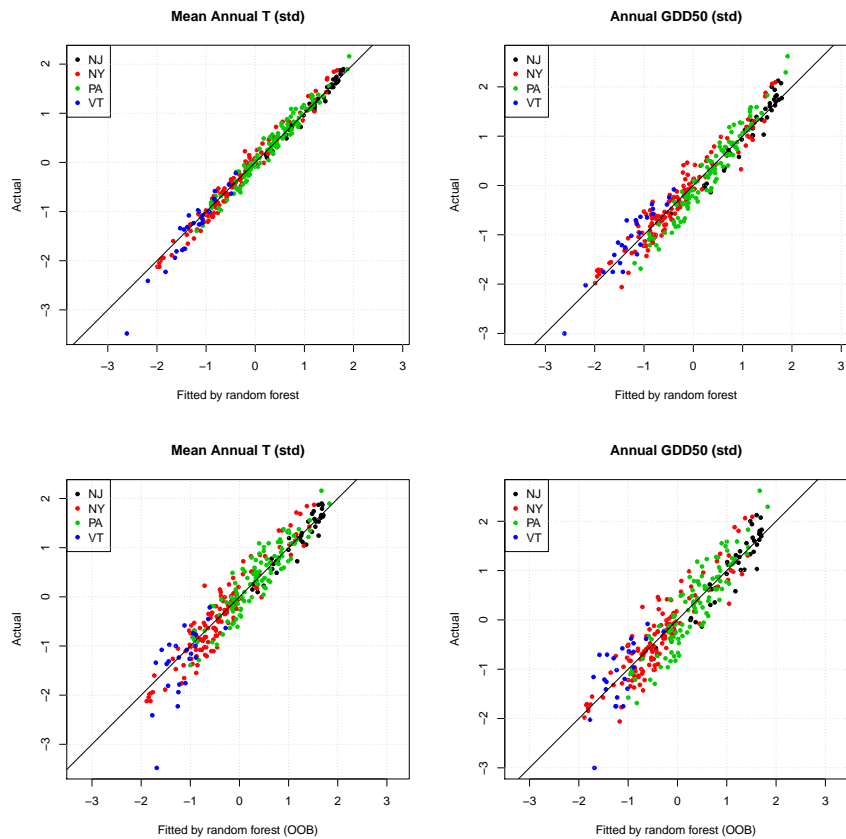
randomForest::importance(m.rf.std.gdd)

##              %IncMSE IncNodePurity
## ELEVATION_  69.83814      120.73660
## N           80.76900      129.97032
## E           58.74444       46.66238
```

The Northing is more influential in the MAT model than in the GDD50 model, whereas the elevation is slightly more influential in the GDD50 model. This agrees with the results of the GLS model (§15.1), where the absolute Northing coefficient was larger for the MAT model, and the absolute elevation coefficient larger for the GDD50 model.

TASK 195 : Plot the two fits, and the two OOB fits, side-by-side. •

```
plot(t.ann.std ~ predict(m.rf.std.t, newdata=ne.m), col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest", ylab="Actual",
     main="Mean Annual T (std)")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
plot(gdd50.std ~ predict(m.rf.std.t, newdata=ne.m), col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest", ylab="Actual",
     main="Annual GDD50 (std)")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
plot(t.ann.std ~ predict(m.rf.std.t), col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest (OOB)", ylab="Actual",
     main="Mean Annual T (std)")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
plot(gdd50.std ~ predict(m.rf.std.t), col=ne.m$STATE, pch=20, asp=1,
     xlab="Fitted by random forest (OOB)", ylab="Actual",
     main="Annual GDD50 (std)")
legend("topleft", levels(ne.m$STATE), pch=20, col=1:4)
grid(); abline(0,1)
```



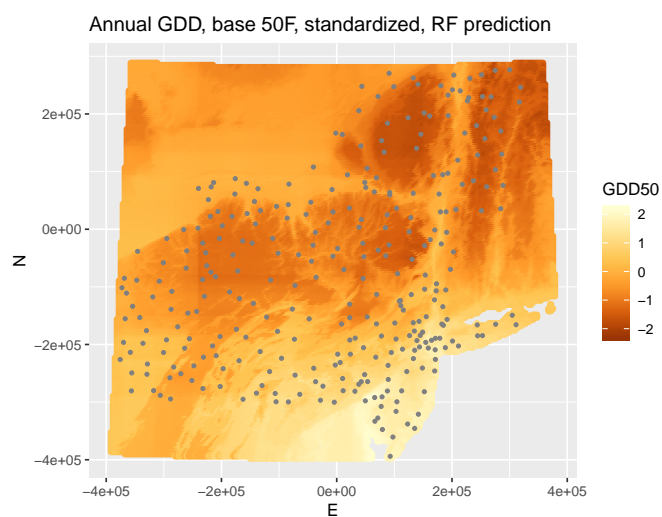
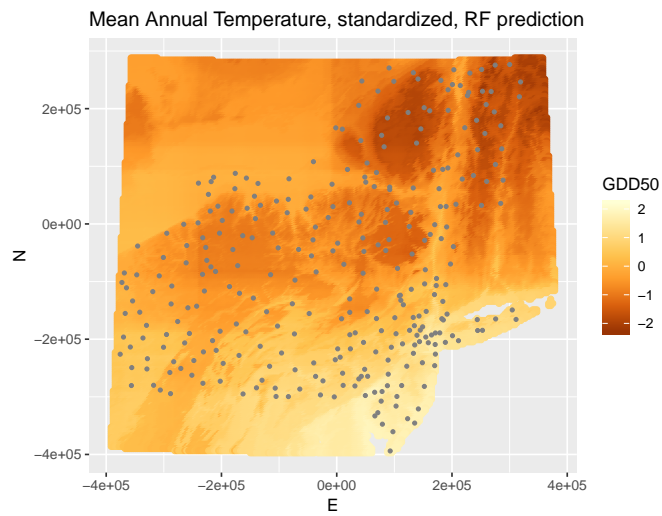
There is not much difference between these. The GDD50 for Mount Mansfield (VT) is better fit than the MAT.

TASK 196 : Predict over the regional grid with both RF models, and display the maps. •

```
dem.ne.m.df$pred.rf.std.t <- predict(m.rf.std.t, newdata=dem.ne.m.df)
dem.ne.m.df$pred.rf.std.gdd <- predict(m.rf.std.gdd, newdata=dem.ne.m.df)
(std.pred.lim <- c(min(dem.ne.m.df[,c("pred.rf.std.t", "pred.rf.std.gdd")],
max(dem.ne.m.df[,c("pred.rf.std.t", "pred.rf.std.gdd")]))))

## [1] -2.288294 2.174055

display.prediction.map("pred.rf.std.t",
  "Mean Annual Temperature, standardized, RF prediction",
  "GDD50", std.pred.lim, .palette="YlOrBr")
display.prediction.map("pred.rf.std.gdd",
  "Annual GDD, base 50F, standardized, RF prediction",
  "GDD50", std.pred.lim, .palette="YlOrBr")
```

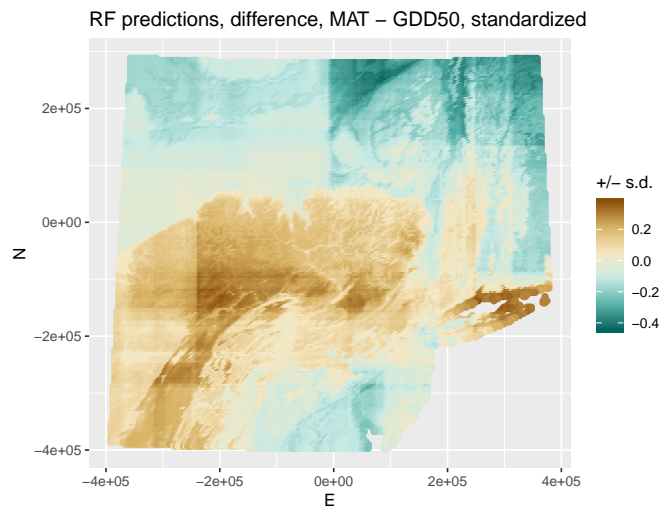


TASK 197 : Compute the differences between the two maps and display the difference map. •

```
summary(dem.ne.m.df$diff.rf.std.t.gdd <-
  dem.ne.m.df$pred.rf.std.t - dem.ne.m.df$pred.rf.std.gdd)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## -0.443410 -0.123015 -0.023286 -0.005924  0.114057  0.375542

display.difference.map("diff.rf.std.t.gdd",
  "RF predictions, difference, MAT - GDD50, standardized",
  "+/- s.d.",
  .palette="BrBG")
```

The summary show that half the differences between the standardized variables are quite small, about ± 0.115 standard deviations. Even the extremes are not too different. This agrees with the close correlation between the variables.

This map shows similar differences between the variables as the GLS model difference map (§15.1): normalized MAT is higher than normalized GDD50 in the Appalachian, Taconic and Allegheny mountains, and on eastern Long Island; the reverse is the case in Adirondacks and Green Mountains, and especially in the Lake Champlain valley.

16 Answers

A1 : *There is an obvious regional trend in the N-S dimension: degree-days tend to decrease going North, as shown by the smaller symbols in that direction. There does not seem to be a consistent trend in the E-W dimension, There is also a trend with elevation: degree-days tend to decrease as elevation increases, at the same N latitude.* [Return to Q1](#) •

A2 :

- Northing :** *This relation looks linear at the more southerly portions of the map (NJ and PA except the northern tier of counties) but quite spread out, with a less obvious relation, in the northerly portion (NY and VT, northern PA).*
- Easting :** *Note the much wider spread of GDD in the East, which ranges from southern NJ to northern VT, than in the continental climate of the West. Easting does not appear to be predictive of the expected value of GDD – it appears that a linear regression would have a (near) zero slope.*

ELEVATION_ : *There is a clear relation (higher elevations have fewer GDD) but it appears inverse-parabolic rather than linear.*

[Return to Q2](#) •

A3 : *Elevation has the strongest correlation; Northing is not much weaker. The correlation coefficients are both negative: GDD50 decreases with increasing Northing and elevation. There is essentially no relation of GDD50 with Easting.*

[Return to Q3](#) •

A4 : *Square root of elevation explains a bit more than half (55.6%) of the total variation in GDD50.*

[Return to Q4](#) •

A5 : *In general there is a good fit: most points are near the 1:1 actual:fit line. However two NY points (red) are poorly fit, one underfit at $\approx (1300, 1950)$ (fit, actual) and one overfit at $\approx (2700, 1700)$.*

[Return to Q5](#) •

A6 :

The model diagnostics look pretty good: (1) no relation between fits and residuals; (2) approximately equal variance at all fitted values; (3) residuals mostly normally-distributed; (4) the high-leverage point is consistent with the others. However two points are severely under-predicted (positive residuals) and one severely over-predicted (negative residual).

[Return to Q6](#) •

A7 :

There is obvious correlation, with some reasonable geographic interpretation. The model under-predicts GDD along the Lake Erie and Ontario plains (this due to the moderating effect of the lake) and over-predicts along the Atlantic coast, southeastern VT and the northern tier of PA. The Atlantic coast may stay cooler in spring than indicated by its very low elevation and southerly position. The model also under-predicts in SE PA and SW NJ (the Philadelphia area) where there may be influence from warm southerly winds.

•

A8 : *There appears to be some structure up to ≈ 40 km.*

[Return to Q8](#) •

A9 : *The confidence interval of the range parameter seems quite wide, almost double from the lower limit 1.285×10^4 to the upper limit 2.3715×10^4 .*

[Return to Q9](#) •

A10 :

The range of spatial dependence has been adjusted; from the variogram fit we estimated in §4.4, i.e., 1.2107×10^4 ; the REML estimate is somewhat longer, 1.7457×10^4 m. These are 1/3 of the effective range, since we fit an exponential

model. The effective range found by *gls* is thus 5.2371×10^4 m. [Return to Q10](#) •

A11 : The GLS fit does not remove spatial correlation, it just takes it into account when computing the regression parameters. [Return to Q11](#) •

A12 : This fits reasonably well; we conclude that we've detected the true correlation structure of the residuals. [Return to Q12](#) •

A13 :

The coefficients have changed a small amount; that for Northing was reduced by over 2%. This shows that some clustered far N and far S elevation points had higher leverage in the OLS model than the cluster warranted. [Return to Q13](#) •

A14 : We see the effect of the reduced coefficient for Northing: the GLS predicts somewhat lower in the N, so the residuals are lower than the OLS residuals (red circles in the bubble plot); the reverse is true in the S. [Return to Q14](#) •

A15 : There is quite some adjustment along the Great Lakes plain (NW) and near some cities (additional GDD) and the Atlantic coast (lower GDD). There are several “hotspots” of large adjustments near single climate stations. [Return to Q15](#) •

A16 : These areas are beyond the range of spatial correlation, here about 50 km, and so are not affected by “nearby” observations. [Return to Q16](#) •

A17 : Compared to the GLS surface this has more detail and is adjusted locally, for example along the Lake Ontario plain. The discrepancies (points that can be seen on the map) are much smaller. [Return to Q17](#) •

A18 : (1) The global trend on the covariates is adequate to compute residuals and their spatial structure; we do not expect it to change much if re-fit locally. (2) It is not practical because the trend surface would have to be re-computed, the residuals extracted, and the variogram re-fit at 1000's of prediction points; this would have to be done automatically without the possibility of checking for artefacts. In local KED we can use the single fitted residual variogram model, while the trend is adjusted locally and the predictions are dependent only on observations in the local neighbourhood. The model is not changed, the results of applying it are.

Note that with very dense point networks (e.g., precision agriculture) this procedure is applicable, and implemented by the VESPER computer program⁴³ for

⁴³ <https://sydney.edu.au/agriculture/pal/software/vesper.shtml>

Ordinary Kriging, but not KED.

[Return to Q18](#) •

A19 : The largest over-predictions by local KED are extrapolations, outside the area with points, e.g., west side of Lake Ontario, and the south-central Appalachians. Within the interpolation area, global KED is generally a bit higher, especially in north-central PA. This may be due to a stronger Northing effect in the global model, vs. the effect as found in local neighbourhoods. [Return to Q19](#) •

A20 : Local KED with the 61 neighbours is considerably better: Cross-validation RMSE is about 4% lower and the extreme errors are smaller.

This shows that both the effect of the covariates and any local effects represented by the stations are not uniform across the area, and some form of localized prediction is indicated. [Return to Q20](#) •

A21 : The differences are certainly geographically-consistent. The global model gives more positive residuals (excessively high predictions) in the south-west, Lake Ontario plain, and far North Country of NY and VT. It gives more negative residuals (excessively low predictions) in the east. This shows that the local fit substantially changed the predictions. [Return to Q21](#) •

A22 : The marginal relations are not well-fit by linear relations, although the square root of elevation is nearly linear, as we saw in the OLS/GLS modelling. However, in GAM modelling we do not need to select a single transformation over the whole range of a predictor to linearize the relation. One transformation may not be applicable to the whole range. Instead, we allow the smooth fit to determine the local adjustment. [Return to Q22](#) •

A23 : The model fits well; the adjusted R^2 is 0.908, and the residuals are less spread than those from the OLS and GLS models. The effective degrees of freedom, i.e., accounting for the many local regressions in the splines, was 23.53 for the 2D surface, and 8.52 for the 1D relation with elevation.

[Return to Q23](#) •

A24 : Both of these plots support the conclusion that there is no spatial dependence of the model residuals. (1) The bubbles appear to be randomly distributed, both in colour and size. (2) The variogram shows the same spatial correlation at all separations. [Return to Q24](#) •

A25 : The GAM 2D trend clearly differs from the linear trend surface, especially in the Lake Ontario plain (towards the right centre in the figure). It is also higher than a linear trend in the S Hudson valley, but lower along the Atlantic shore (upper left in the figure). These are areas we identified with large OLS and GLS residuals in the linear trend surfaces (§4, §5). [Return to](#)

Q25 •

A26 : This is almost linear now, as suggested by physical theory, once the smooth geographic trend is removed. However at the low elevations there is a wide spread of GDD50, which is well-fit by an almost vertical portion of the marginal smooth function. [Return to Q26 •](#)

A27 : The largest differences are to the E and W, out of the calibration area - this illustrates that GAM should not be extrapolated. The much lower GDD50 in New England is because the GAM trend surface was considerably below the GLS surface along the Atlantic coast, and this was extrapolated eastward. The GAM predicts higher values along the Great Lakes plains and lower Hudson valley, as we saw in the GLS residuals. [Return to Q27 •](#)

A28 : The largest differences are quite local, around certain weather stations where the local deviation from the trend could be accounted for in RK-GLS, but was somewhat averaged out in GAM. [Return to Q28 •](#)

A29 : The first (root) splitting variable is N. The split is at -1.55967×10^5 N. The mean value of GDD50 of the whole dataset is 2517.52; the mean value of the observations in the left branch (less than) is 2182.54 and of the right branch (greater than) is 3019.99. These branches have 183 and 122 observations, respectively. [Return to Q29 •](#)

A30 : Northing is most important; it explains almost 50% of the variance. Then elevation explains about another 40%, and Easting very little. This agrees with the linear correlations computed as preparation for fitting the OLS model. [Return to Q30 •](#)

A31 : Each run of the *rpart* function will give the same tree (if the same parameters are specified) but slightly different cross-validation statistics. The cross-validation error reaches an effective minimum around 15 splits, CP about 0.0045. So building the tree with CP=0.003 was overfitting. [Return to Q31 •](#)

A32 : The pruned tree has the same root and higher levels, but fewer splits and leaves. [Return to Q32 •](#)

A33 : There are 16 unique values predicted by the pruned regression tree. The fit to the actual values is better the OLS or GLS models; the RMSE from the regression tree is 11.16 GDD50; for the GLS model 12.08 In this case the linear model predicts more values but on average they deviate more from the true values. [Return to Q33 •](#)

A34 : The regression tree divides the area into "blocks" mostly with the Nor-

thing but in one place with the Easting. It also slices most of the “blocks” according to elevation zones. These give the maximum between-group variance at the leaves of the regression tree, without overfitting. [Return to Q34](#)

•

A35 : Permuting elevation, i.e., assigning elevations randomly to stations, would increase the MSE by about 69%, thus it is a very important predictor. Northing is also important as expected; surprisingly the Easting is important in the trees in which it was used. [Return to Q35](#) •

A36 : The out-of-bag mean errors are from 2 to 3 times that of the fits. This is a typical result for random forests. The OOB errors are indicative of the errors at unknown points, i.e., prediction accuracy. [Return to Q36](#) •

A37 : For most runs we see no or weak residual spatial structure, because of the averaging effect of the repeated bootstrap sampling; in many trees close point-pairs lose one of the points. [Return to Q37](#) •

A38 : The RF surface shows some irregular patches and abrupt transitions, whereas the RK-GLS surface is by construction smooth. [Return to Q38](#) •

A39 : The RF prediction is from the “box” containing the elevation (all the same in the lake), Northing and Easting, which was fit with the most similar points. These are presumably along the Lake shore. There is no extrapolation via a trend surface to modify the effect of the coördinates in the model. [Return to Q39](#) •

A40 : The GLS and GLS-RK models have coefficients for the coördinates, which here are far East, so the predictions change. The RF does not have any information in this area and so puts it all in the “boxes”. [Return to Q40](#)

•

A41 : There are no points in OH so the prediction by RF is made from the fitted model of the nearby PA stations. These apparently use the Easting. [Return to Q41](#) •

A42 : The RF model has no way to extrapolate to higher or lower elevations than in the calibration set. In the Alleghenies and Catskills there are a lot of higher-elevation area beyond the elevation of weather stations. In the Adirondacks the stations are all at low elevations, but apparently the Northing is here used to predict the GDD. This is why the RF over-predicts in the lowlands around Plattsburgh and Lake Champlain. [Return to Q42](#) •

A43 : The random forest model is clearly not suitable to show a regional trend, especially outside of the model calibration area. It does allow for non-

linearities and local combinations of factors, for example on the Lake Ontario and Erie plains. [Return to Q43](#) •

A44 : Slightly more variation of annual temperature is explained by the model. Since these are both standardized, if their regional relation with the predictors is the same, so should the coefficients. They are close but not identical. [Return to Q44](#) •

A45 : There is a clear spatial pattern and local spatial dependence. The residuals from the MAT model are lower than those from the GDD50 model in the N and S edges, and at the higher elevations (Adirondacks, Green Mountains, Allegheny Plateau). The reverse is the case in the centre and especially on Long Island (NY) and northern NJ.

[Return to Q45](#) •

A46 :

The variogram structures are both weak (short range, very high nugget proportion). That is, the regional trend explains most of the variation in both cases.

There is weaker residual spatial structure for mean annual temperature than for GDD50 (lower total sill, higher nugget proportion). This implies that the MAT trend (the predictors N, E and elevation) explains more of the spatial variability. Note however that the MAT trend surface has a lower R^2 than the GDD50 trend surface. This implies that MAT is more explained regionally and less by local variations than annual GDD50. [Return to Q46](#) •

A47 : The MAT is proportionally higher than GDD50 in the higher elevations, especially towards the SW. The opposite effect is seen in the lowlands, especially the Lake Champlain valley (NE). [Return to Q47](#) •

17 Challenge

Do a similar analysis *either* for:

- Over the same study area as the example:
 - the growing degree days in one of the growing-season months: May through September; or
 - the annual growing degree days at base 40°F
 - one of the other climate variables in one of the other shape-files.
- Over one of the States within the study area.
- Over some other region of the USA:
 - the growing degree days base 50°F, i.e., the same as used in this example. Note that for this you will have to define a suitable coordinate reference system (CRS) for that area.

If you have to define a suitable CRS:

- For E-W oriented regions, you can use the same Albers Equal Area projection as was used in §2.4, but the parameters will be different.
- For N-S oriented regions, you will need to select a different projection. A good choice is Transverse Mercator, PROJ.4 name `tmerc`. The parameters for this are⁴⁴:

```
+proj=tmerc +lat_0=Latitude of natural origin
            +lon_0=Longitude of natural origin
            +k=Scale factor at natural origin
            +x_0=False Easting
            +y_0=False Northing
```

1. Determining which covariables (elevation, Northing, Easting, their transformations or interaction) best model the target variable;
2. Estimating the coefficients of the linear model with OLS;
3. Examining the (non-spatial) OLS model diagnostics;
4. Modelling the spatial structure of the OLS model residuals;
5. Estimating the coefficients of the linear model and the spatial correlation structure (nuisance parameters) with GLS;
6. Fitting a random forest with all three predictors.
7. Mapping the various predictions and their differences.

Then answer these questions:

1. Is the linear model justified?

⁴⁴ http://geotiff.maptools.org/proj_list/transverse_mercator.html, see the 'PROJ.4 organization' subheading for the names and meanings of the parameters

2. Is there spatial structure in the OLS model residuals? So, is a GLS model needed?
3. How much do the OLS and GLS model coefficients differ?
4. What is the spatial correlation structure fitted by the GLS model? Does this agree with that estimated from the OLS variogram?
5. How does your model compare to that developed in this exercise for annual GDD50?
6. What could be the reason(s) for differences between the model in the exercise and your model?

References

- [1] D Bates. Fitting linear mixed models in R. *R News*, 5(1):27–30, 2005. 38
- [2] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and regression trees*. Wadsworth, 1983. 112
- [3] I. C. Briggs. Machine contouring using minimum curvature. *Geophysics*, 39(1):39–48, January 1974. ISSN 0016-8033, 1942-2156. doi: 10.1190/1.1440410. 148
- [4] Lev M Bugayevskiy and John Parr Snyder. *Map projections: a reference manual*. Taylor and Francis, 1995. 11
- [5] Bradley Efron and Gail Gong. A leisurely look at the bootstrap, the jackknife & cross-validation. *American Statistician*, 37:36–48, 1983. 123
- [6] A. Stewart Fotheringham, Chris Brunsdon, and Martin Charlton. *Geographically weighted regression: the analysis of spatially varying relationships*. Wiley, 2002. ISBN 0-471-49616-2. doi: 10.4135/9781849209755. 91
- [7] A. Stewart Fotheringham, Wenbai Yang, and Wei Kang. Multiscale geographically weighted regression (mgwr). *Annals of the American Association of Geographers*, 107(6):1247–1265, Nov 2017. ISSN 2469-4452. doi: 10.1080/24694452.2017.1352480. 91
- [8] P Goovaerts. *Geostatistics for natural resources evaluation*. Applied Geostatistics. Oxford University Press, New York; Oxford, 1997. 67, 153
- [9] P. Harris, A. S. Fotheringham, R. Crespo, and M. Charlton. The use of geographically weighted regression for spatial prediction: an evaluation of models using simulated data sets. *Mathematical Geosciences*, 42:657–680, Jun 2010. doi: 10.1007/s11004-010-9284-7. 91
- [10] Paul Harris, Chris Brunsdon, and A. Fotheringham. Links, comparisons and extensions of the geographically weighted regression model when used as a spatial predictor. *Stochastic Environmental Research and Risk Assessment*, 25(2):123–138, 2011. doi: 10.1007/s00477-010-0444-6. 91
- [11] T Hastie, R Tibshirani, and J H Friedman. *The elements of statistical learning data mining, inference, and prediction*. Springer series in statistics. Springer, New York, 2nd ed edition, 2009. ISBN 9780387848587. 79, 112, 122, 123, 147
- [12] Tomislav Hengl, Gerard B. M. Heuvelink, and David G. Rossiter. About regression-kriging: From equations to case studies. *Computers & Geosciences*, 33(10):1301–1315, 2007. doi: 10.1016/j.cageo.2007.05.001. 59

- [13] M. F. Hutchinson. Interpolating mean rainfall using thin plate smoothing splines. *International Journal of Geographical Information Science*, 9(4):385 – 403, 1995. 148
- [14] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning: with applications in R*. Number 103 in Springer texts in statistics. Springer, 2013. ISBN 9781461471370. 79, 112, 122
- [15] Max Kuhn. Building predictive models in R using the caret package. *Journal of Statistical Software*, 28(5):1–26, 2008. 133
- [16] Max Kuhn and Kjell Johnson. *Applied Predictive Modeling*. Springer, 2013 edition edition, Sep 2013. ISBN 978-1-4614-6848-6. 139
- [17] R. M. Lark and B. R. Cullis. Model based analysis using REML for inference from systematically sampled data on soil. *European Journal of Soil Science*, 55(4):799–813, 2004. 36, 38
- [18] Binbin Lu, Paul Harris, Martin Charlton, and Chris Brunsdon. The GWmodel R package: further topics for exploring spatial heterogeneity using geographically weighted models. *Geo-spatial Information Science*, 17(2):85–101, Apr 2014. ISSN 1009-5020. doi: 10.1080/10095020.2014.917453. 95
- [19] G. S. McMaster and W. W. Wilhelm. Growing degree-days: one equation, two interpretations. *Agricultural and Forest Meteorology*, 87(4):291–300, 1997. doi: 10.1016/S0168-1923(97)00027-0. 2
- [20] H. Mitasova and J. Hofierka. Interpolation by regularized spline with tension: II. Application to terrain modeling and surface geometry analysis. *Mathematical Geology*, 25(6):657–669, 1993. doi: 10.1007/BF00893172. 148
- [21] H. Mitasova and L. Mitas. Interpolation by regularized spline with tension: I. Theory and implementation. *Mathematical Geology*, 25(6):641–655, 1993. doi: 10.1007/BF00893171. 147
- [22] J C Pinheiro and D M Bates. *Mixed-effects models in S and S-PLUS*. Springer, 2000. ISBN 0387989579. 37, 38
- [23] J. R. Quinlan. *C4.5: programs for machine learning*. The Morgan Kaufmann series in machine learning. Morgan Kaufmann Publishers, 1993. ISBN 1-55860-238-0. 139
- [24] Cosma Shalizi. The bootstrap. *American Scientist*, 98(3):186–190, 2010. doi: DOI:10.1511/2010.84.186. URL <http://www.americanscientist.org/issues/pub/2010/3/the-bootstrap/3>. 123
- [25] John P. Snyder. *Map projections: a working manual*. USGS Professional Paper 1395. US Government Printing Office, 1987. URL <https://pubs.er.usgs.gov/publication/pp1395>. 11

- [26] John P. Snyder and Philip M. Voxland. *An album of map projections*. USGS Professional Paper 1453. U.S. Geological Survey, 1989. URL <https://pubs.er.usgs.gov/publication/pp1453>. 11
- [27] W. R. Tobler. A computer movie simulating urban growth in the Detroit region. *Economic Geography*, 46:234–240, 1970. ISSN 0013-0095. doi: 10.2307/143141. 152
- [28] Ku Wang, Chuanrong Zhang, and Weidong Li. Predictive mapping of soil total nitrogen at a regional scale: A comparison between geographically weighted regression and cokriging. *Applied Geography*, 42:73–85, Aug 2013. ISSN 0143-6228. doi: 10.1016/j.apgeog.2013.04.002. 91
- [29] R. Webster and M. A. Oliver. *Geostatistics for environmental scientists*. John Wiley & Sons Ltd., 2nd edition, 2008. 67, 153
- [30] Hadley Wickham. ggplot2. <http://ggplot2.org/>. URL <http://ggplot2.org/>. 19
- [31] Hadley Wickham. *ggplot2: Elegant graphics for data analysis*. Use R! Springer, August 2009. ISBN 0387981403. 19
- [32] Leland Wilkinson. *The grammar of graphics*. Statistics and computing. Springer, New York, 2nd ed edition, 2005. ISBN 9780387286952. 17
- [33] S. N. Wood. Thin plate regression splines. *Journal of the Royal Statistical Society Series B-Statistical Methodology*, 65:95–114, 2003. doi: 10.1111/1467-9868.00374. 147
- [34] Marvin N. Wright and Andreas Ziegler. ranger: a fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, 77(1):1–17, Mar 2017. doi: 10.18637/jss.v077.i01. 133
- [35] Yihui Xie. knitr: Elegant, flexible and fast dynamic report generation with R, 2011. URL <http://yihui.name/knitr/>. Accessed 04-Mar-2016. 2
- [36] Canying Zeng, Lin Yang, A-Xing Zhu, David G. Rossiter, Jing Liu, Junzhi Liu, Chengzhi Qin, and Desheng Wang. Mapping soil organic matter concentration at different scales using a mixed geographically weighted regression method. *Geoderma*, 281:69–82, Nov 2016. ISSN 0016-7061. doi: 10.1016/j.geoderma.2016.06.033. 91

A * Colour ramps with ggplot2

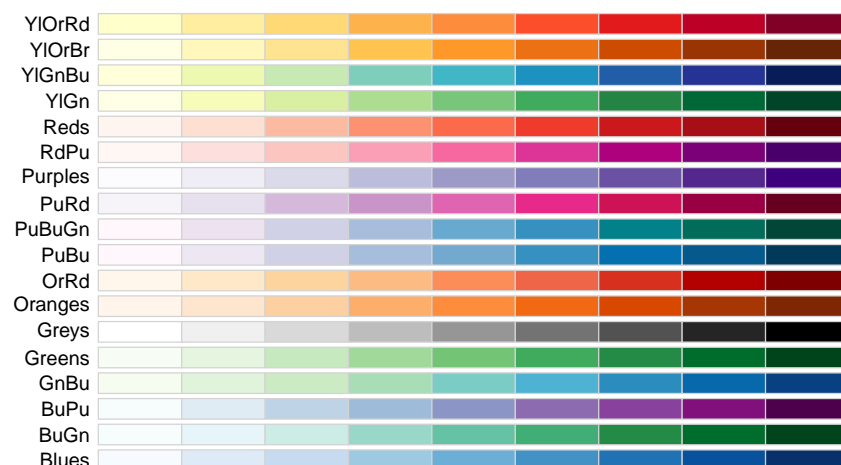
The `colour` argument to the `aes` “aesthetics” function has a default colour ramp.

In `ggplot2` terminology this is called a **scale**. It maps numbers (such as GDD50) to a scale in colour space, in the same way there is a mapping

from numbers to a position on an axis of a scatterplot. It is a **function** from a **domain** (data values, which can be classes or continuous) to a **range** (the colour for that value). This is a quite tricky – the user’s visual perception must match the change in colour. One package that has dealt with this is RColorBrewer⁴⁵, which palettes (colour choices) can be accessed with the `scale_colour_brewer` function for categorical variables and `scale_colour_distiller` for continuous variables. This package provides ready-made sequential, diverging, and qualitative palettes.

TASK 198 : Load the RColorBrewer package and display the ready-made palettes for continuous **sequences**. •

```
library(RColorBrewer)
display.brewer.all(type="seq")
```



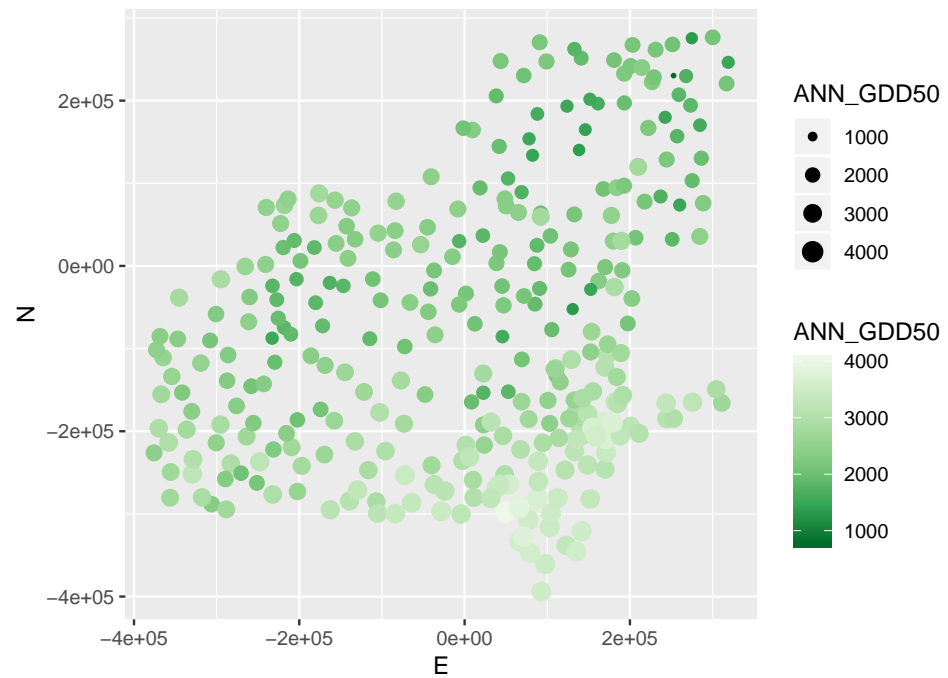
These should be used when the variable to display is a continuous value from some minimum to some maximum, where increasing values indicate increasing intensity of the variable. For example, the annual GDD50, from low to high.

TASK 199 : Select one of the palettes and re-display the GDD50 map with it, •

```
ggplot(data=ne.df) +
  aes(x=E, y=N) +
  geom_point(aes(size=ANN_GDD50, colour=ANN_GDD50),
             shape=20) +
  scale_colour_distiller(space="Lab",
```

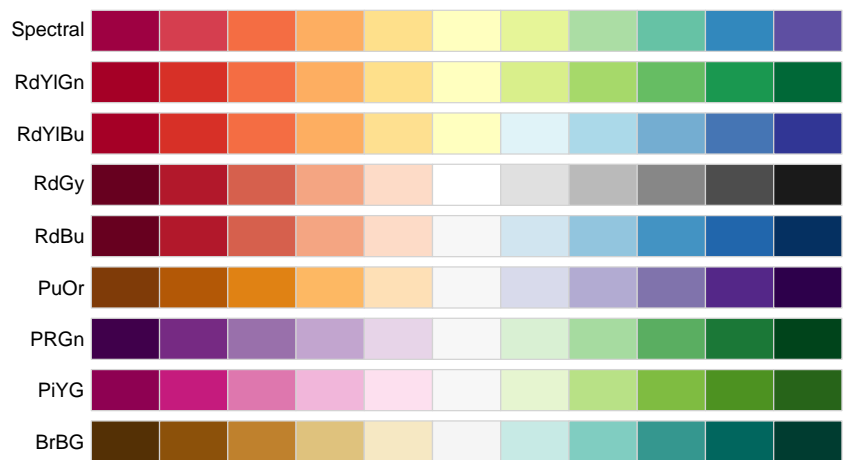
⁴⁵ <http://colorbrewer2.org/>

```
palette="Greens") +  
xlab("E") + ylab("N") + coord_fixed()
```



TASK 200 : Display the ready-made palettes for **diverging** palettes. •

```
display.brewer.all(type="div")
```



These should be used when the central value is a natural zero and we want to emphasize divergences from it in two directions, e.g., for residuals.

Index of R Concepts

+ formula operator, 114
+ operator, 19
- operator, 10
[] operator, 7
%in% operator, 7, 15
~ formula operator, 114

adapt argument (gwr function), 97
add argument (plot function), 17
aes (ggplot2 package), 19, 56, 204
aes_ (ggplot2 package), 56
aggregate (raster package), 53
apply, 171
as, 13
as.character, 5
as.data.frame, 34

bandwidth argument (gwr function), 97
block argument (krige function), 157
border argument (plot function), 17
bubble (sp package), 30
byid argument (gIntersection function), 167

caret package, 134, 141
cbind, 174
ceiling, 12
class, 50
cloud argument (variogram function), 34
colour argument (aes function), 19, 204
committees argument (cubist function), 145
control argument (rpart function), 118
coord_fixed (ggplot2 package), 19
coordnames (sp package), 13
coords slot (SpatialPointsDataFrame class), 10
cor, 23
corExp (nlme package), 39, 43
correlation argument (glS function), 38
corSpher (nlme package), 45
corStruct class, 39
cp argument (prune function), 119
cp argument (rpart function), 114, 133
crop (raster package), 51
CRS (rgdal package), 12
CRS (sp package), 13
Cubist package, 1, 140, 145

cubist (Cubist package), 145

data argument (ggplot function), 19
data slot (SpatialPointsDataFrame class), 10, 13
data.frame class, 13, 55, 159
deldir (deldir package), 165
deldir package, 165
dim, 8
dismo package, 165
dsn argument (readOGR function), 4, 5

expand.grid, 49, 134
ext argument (voronoi function), 165
extent (raster package), 50

fact argument (aggregate function), 53
factor, 8
fields package, 1, 148, 150, 152
fit.variogram (gstat package), 33, 67, 156
floor, 12
form argument (corExp function), 39
function, 56

gam (mgcv package), 82
geom_point (ggplot2 package), 19
geom_smooth (ggplot2 package), 80
getModelInfo (train package), 134
getwd, 4
ggplot (ggplot2 package), 19, 56, 159
ggplot2 package, 1, 19, 80, 204
gIntersection (rgeos package), 167
glS, 183, 186
glS (nlme package), 38–40, 43, 45, 67, 197
grid.arrange (gridExtra package), 80
gridExtra package, 80
gstat package, 1, 32, 33, 62, 67, 72, 154, 161
gUnaryUnion (rgeos package), 166
GWmodel package, 96
gwr (gwr package), 97, 102, 107
gwr class, 97, 98
gwr.sel (spgwr package), 96

hatmatrix argument (gwr function), 102

idp argument (idw function), 163
idw (gstat package), 163

image (raster package), 52, 54
 importance (randomForest package), 192
 importance argument (randomForest function), 124
 integer64 argument (readOGR function), 4
 intervals (nlme package), 40
 kml_close (plotKML package), 21
 kml_layer (plotKML package), 21
 kml_open (plotKML package), 21
 knitr package, 2
 krige (gstat package), 62, 67, 68, 73, 74, 77–79, 163
 krige.cv (gstat package), 72, 161
 labels argument (text function), 17
 layer argument (readOGR function), 4
 list.files, 3, 5
 lm, 24, 38, 78, 82
 load, 16, 54
 loess, 80
 log, 82
 lwd argument (plot function), 17
 matrix, 150
 max, 55
 maxdist argument (krige function), 73, 74
 mean, 53
 method argument (gwr.sel function), 96
 method argument (train function), 134, 142, 144
 mgcv package, 82, 84, 85
 min.node.size argument (ranger function), 135, 137
 minsplitt argument (rpart function), 114, 133
 model argument (krige function), 62
 mtry argument (randomForest function), 124, 133
 mtry argument (ranger function), 135, 137, 140
 ncol (raster package), 50
 neighbors argument (predict.cubist function), 145
 newdata argument (predict.rpart function), 121
 newdata argument (predict function), 125, 126
 nlme package, 1, 38, 40
 nmax argument (krige function), 73, 74, 157
 nmin argument (krige function), 74
 nodesize argument (randomForest function), 133
 nrow (raster package), 50
 ntree argument (randomForest function), 124, 133
 nugget argument (corExp function), 39, 43
 nugget argument (corSpher function), 45
 order, 34
 over (sp package), 153, 168
 palette, 17
 pattern argument (list.files function), 3
 plot, 17
 plot.gam (mgcv package), 84
 plot.SpatialPolygons (sp package), 17
 plot.xy, 17
 plotKML package, 1, 4, 21
 pmax, 55
 predict, 120
 predict (cubist package), 145
 predict (randomForest package), 125, 126
 predict argument (gwr function), 107
 predict.gam (mgcv package), 88
 predict.Krig (fields package), 152
 predict.rpart (rpart package), 120
 printcp (rpart package), 117
 proj4string (rgdal package), 5, 11, 95
 proj4string (sp package), 12, 15
 proj4string slot (SpatialPointsDataFrame class), 11, 12
 projection (raster package), 50
 projectRaster (raster package), 51
 projInfo (rgdal package), 11
 prune (rpart package), 119
 quote, 56
 randomForest (randomForest package), 124, 133, 137, 139
 randomForest package, 1, 124, 134, 140, 192

ranger (ranger package), 137
 ranger package, 1, 134, 138, 140
 raster (raster package), 50
 raster package, 1, 49–54
 RasterLayer class, 50
 RColorBrewer package, 204
 read.metadata (plotKML package), 4
 readOGR (rgdal package), 4, 14
 res (raster package), 50
 rgdal package, 1, 4, 5, 11
 rgeos package, 166
 round, 12
 rpart (rpart package), 113, 114, 118, 199
 rpart class, 120
 rpart package, 1, 113
 rpart.control (rpart package), 118
 rpart.plot (rpart.plot package), 115
 rpart.plot package, 115

 s (mgcv package), 82
 save, 16, 53, 54
 scale_colour_brewer (ggplot2 package), 204
 scale_colour_distiller (ggplot2 package), 204
 scheme argument (plot.gam function), 84
 se.fit argument (gwr function), 102
 se.fit argument (predict.gam function), 88
 select argument (plot.gam function), 84
 setwd, 3
 shape argument (geom_point function), 19
 size argument (aes function), 19
 sp class, 4, 13, 17
 sp package, 1, 4, 5, 9, 11, 13, 21, 30, 96, 151, 153, 158, 168, 171
 span argument (loess function), 80
 SpatialGridDataFrame (sp class), 158
 SpatialPixels (sp class), 153
 SpatialPixelsDataFrame (sp class), 153
 SpatialPixelsDataFrame class, 55
 SpatialPoints (sp class), 169
 SpatialPointsDataFrame (sp class), 11, 13, 15, 30, 168
 SpatialPointsDataFrame class, 98, 106
 SpatialPolygonsDataFrame (sp class), 14, 168
 spDists (' package), 171
 spgwr package, 96
 spplot (sp package), 158
 spsample (sp package), 49, 151
 spTransform (sp package), 11, 13, 15, 21, 95
 sqrt, 82
 strwrap, 5
 sum, 8

 text, 17
 theta argument (plot.gam function), 84
 topo.colors, 52
 Tps (fields package), 150
 train (caret package), 134, 137, 142
 trainControl (caret package), 134
 trControl argument (train function), 134
 tuneGrid argument (train function), 134
 type argument (projInfo function), 11

 unique, 121

 value argument (corExp function), 39
 value argument (corSpher function), 45
 variogram (gstat package), 32, 34, 67, 154
 verbose argument (load function), 16, 54
 vgm (gstat package), 42, 67, 156
 vis.gam (mgcv package), 85
 voronoi (dismo package), 165

 which, 8
 which.max, 26
 which.min, 26, 171

 x argument (train function), 134
 xlab argument (ggplot2 function), 19

 y argument (train function), 134
 ylab argument (ggplot2 function), 19

 zerodist (sp package), 9