

Detecting patterns in space and time

Lab session

1. The dataset *SubstationRPD.RData* contains real power delivered (KW) for each 10-minute period, of every day during August, for 410 substations in the southwest of Wales, UK.
 - (a) Produce summaries of the dataset *SubstationRPD.RData* and produce histograms showing the distributions of real power delivered for the 410 substations.
 - (b) For each substation calculate the daily average demand and then plot these on the same plot, using a different colour for each substation.
 - (c) Add to your plot in 1(b) a thick, black line showing the overall median for the demand of all of the substations, with two grey, dashed lines to show the first and third quartiles of demand over all the substations.
 - (d) Split your plot in 1(c) into four separate plots representing; 1) All days, 2) Weekdays, 3) Saturdays and 4) Sundays
 - (e) Write a function that would create the plot from 1(d) for a given set of substations.

The dataset *Characteristics.csv* contains information on all of the aforementioned substations, including number of customers, number of feeder ends and transformer rating.

- (e) Summarise the data and find the distributions for the percentage of industrial and commercial customers (*Percentage_IC*), transformer ratings (*Transformer_RATING*) and pole or ground monitored substations (*TRANSFORMER_TYPE*).
 - (f) Reproduce 1(e) for the following:
 - Substations which have $> 80\%$ industrial and commercial customers;
 - Substations with $< 80\%$ industrial and commercial customers AND a transformer rating < 250
2. Using a unit observation of individual substations, perform hierarchical clustering for the daily average demand dataset you created in 1(b)
 - (a) Using your preferred choice of a dissimilarity function, create a distance matrix for these data.
 - (b) Use your distance matrix from 2(a) to produce a dendrogram.
 - (c) Choose an appropriate number of clusters and label each substation according to its cluster membership.
 - (d) For each cluster, as in question Q1(d), on the same plot, plot the daily average demand for 1) All days, 2) Weekdays, 3) Saturdays and 4) Sundays.
 3. Dataset *NewSubstations.csv* contains information for five new substations.

- (a) For each substation, on the same plot, plot the daily average demand for 1) All days, 2) Weekdays, 3) Saturdays and 4) Sundays.
 - (b) Using the *k-means* function (or if you fancy a challenge, by writing your own algorithm), which cluster would you allocated to each of the new substations?
 - (c) Is the cluster allocation as you expected?
4. In this question we explore how fitting a generalised additive model (*GAM*) to data allows us to forecast future data.
- (a) Reformat the *SubstationRPD.RData* dataset so that each row is the average of all demand data for each substation.
 - (b) Fit and plot a *GAM* which accounts for the underlying seasonal pattern. What are the degrees of freedom?
 - (c) Choose an appropriate model with which to predict the demand for the 21st to the 28th of July. Produce a plot showing these predictions against time.
5. In the *geoR* library there are data *ca20* which you should explore/analyze using geostatistical techniques. For example, you may:
- (a) Look at empirical semi-variograms (clouds and binned).
 - (b) Examine Monte Carlo intervals of no spatial dependence.
 - (c) Fit variogram models to the data.
 - (d) Carry out kriging and examine the resultant surfaces.
6. In this question we will fit several theoretical variogram to a variable of your choice in the *meuse* data set from *gstat* package. We will find the best fitted model based on the SSE criteria and by using cross validation.
- (a) Use the `fit.variogram()` function from *gstat* package. Set the option `print.SSE` of this function to `TRUE`. Read the help page for this function carefully. Concentrate on one of the metal variables in the *meuse* data set and fit at least four different families of variogram models to the empirical variogram computed by the `variog()` function. You may do the analysis on the original or make a transformation if you like.
- ```
library(gstat) data(meuse)
vgm1 <- variogram(log(zinc)~1, ~x+y, meuse,
print.SSE=TRUE) plot(vgm1)
meuse.vfit <- fit.variogram(vgm1, vgm(1,"Sph",300,1))
plot(vgm1,model=meuse.fit)
```
- Based on the SSE criteria choose the best fitted model.
- (b) Now we will use cross validation to choose between a set of models. We will use the `krige.cv()` function from the *gstat* package. Read the help page carefully. When doing cross validation choose to use the method of one-leave-out by specifying `nfold=1`. For example you can do is like this,

```

data(meuse)
m <- vgm(.59, "Sph", 874, .04)
x <- krige.cv(log(zinc)~1, ~x+y,
model = m, data = meuse, nmax = 40, nfold=1)

```

Use the following functions to calculate the mean error (ME), the mean squared error (MSE), and the mean squared deviation ratio (MSDR) diagnostics.

```

ME <- function(xv.obj){ tmp <- xv.obj$error
return(sum(tmp)/length(tmp))
}
MSE <- function(xv.obj){ tmp <- xv.obj$error
return(sum(tmp^2)/length(tmp))
}
MSDR <- function(xv.obj){ e2 <- xv.obj$error^2
s2 <- xv.obj$krige.var
msdr <- sum(e2/s2)/length(e2) return(msdr)
}

```

To get the diagnostics do the following on the cross-validation object x computed above  
ME(x), MSE(x), MSDR(x)