

LAB 2 BLOCK 2 - Group A15

Zuxiang Li (*zuxli371*), Marcos F. Mourao (*marfr825*), Agustín Valencia (*aguva779*)

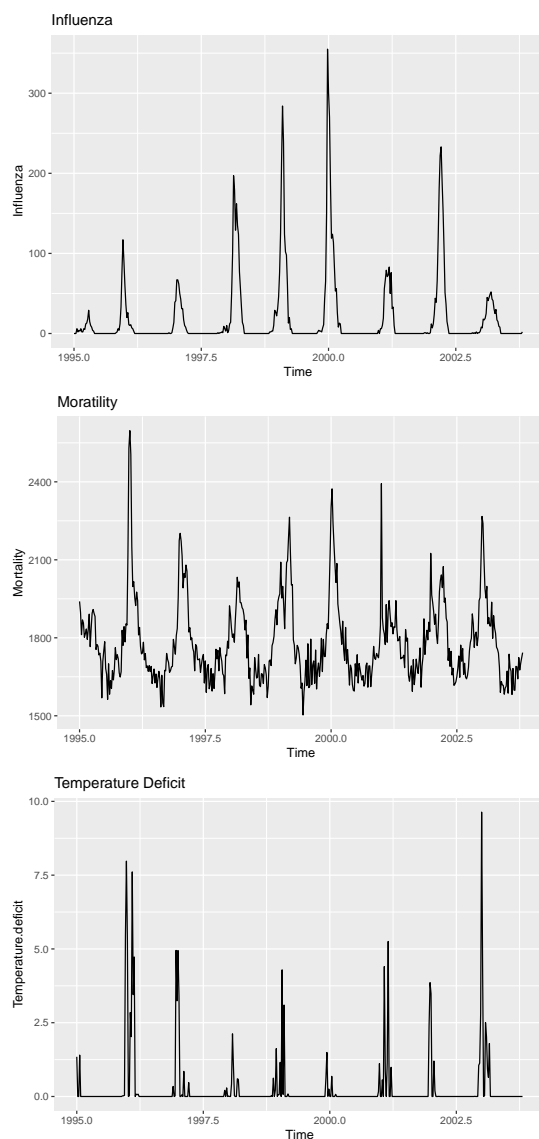
16 December 2019

Assignment 1. Using GAM and GLM to examine mortality rates

The Excel document *influenza.xlsx* contains weekly data on the mortality and the number of laboratory-confirmed cases of influenza in Sweden. In addition, there is information about population-weighted temperature anomalies (temperature deficits).

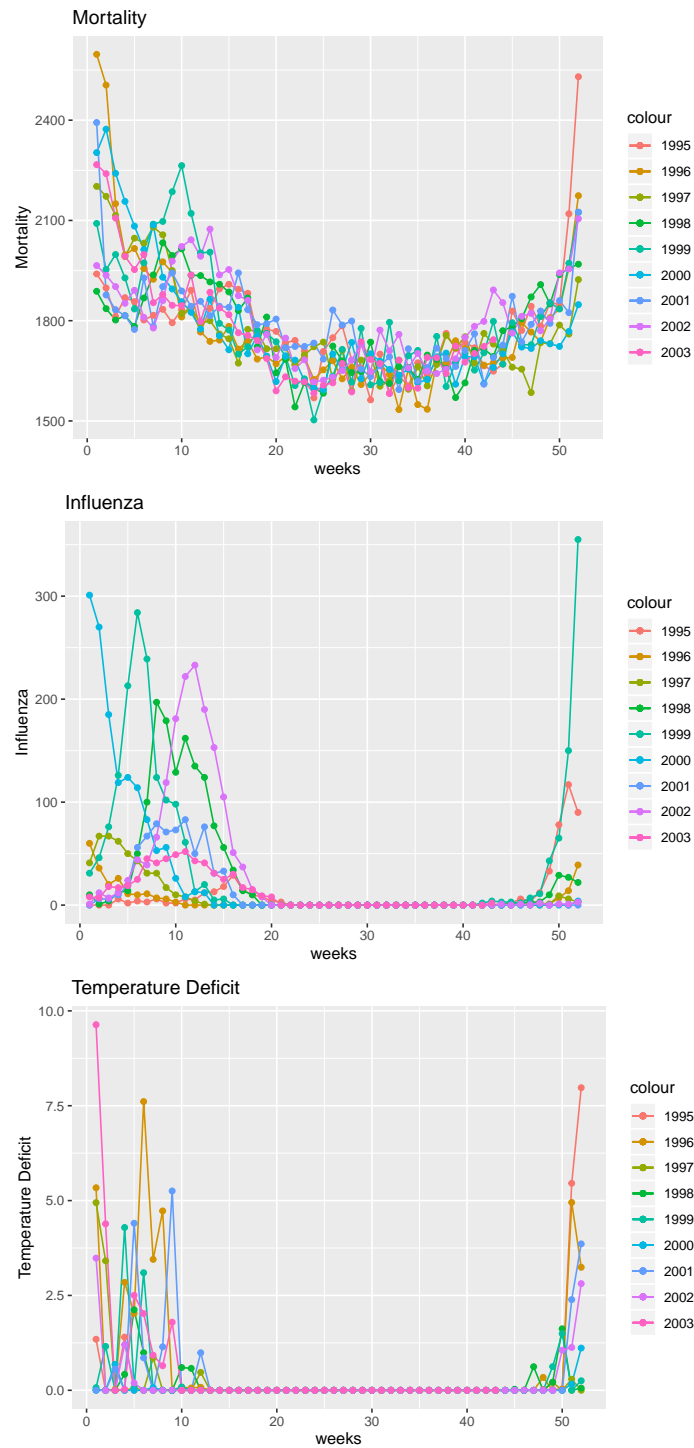
1. Use time series plots to visually inspect how the mortality and influenza number vary with time (use Time as X axis). By using this plot, comment how the amounts of influenza cases are related to mortality rates.

Plotting the data against the time:



It seems to be a modality in the data. It can be seen that Influenza peaks seems to have a correlation with peaks in mortality and temperature deficit

As the model is meant to learn from weeks variations along the year, it is better to separate the data into several time series per year and see how the data looks now.



Now it is quite clear that the data shows that during winter weeks Influenza cases and Mortality related to increases, also the temperature deficit seems to follow that trend.

2. Use `gam()` function from `mgcv` package to fit a GAM model in which Mortality is normally distributed and modelled as a linear function of Year and spline function of Week, and make sure that the model parameters are selected by the generalized cross-validation. Report the underlying probabilistic model.

It is fitted one model with the following formula:

$$\text{Mortality} \approx \text{Year} + \text{spline}(\text{Week})$$

The summary of the obtained probabilistic model is:

```
##
## Family: gaussian
## Link function: identity
##
## Formula:
## Mortality ~ Year + s(Week, k = length(unique(data$Week)))
##
## Parametric coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept) -680.598    3367.760  -0.202    0.840
## Year         1.233      1.685    0.732    0.465
##
## Approximate significance of smooth terms:
##               edf Ref.df    F p-value
## s(Week) 14.32  17.87 53.86 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Rank: 52/53
## R-sq.(adj) =  0.677   Deviance explained = 68.8%
## GCV = 8708.6   Scale est. = 8398.9    n = 459
```

The model Mean and variance are :

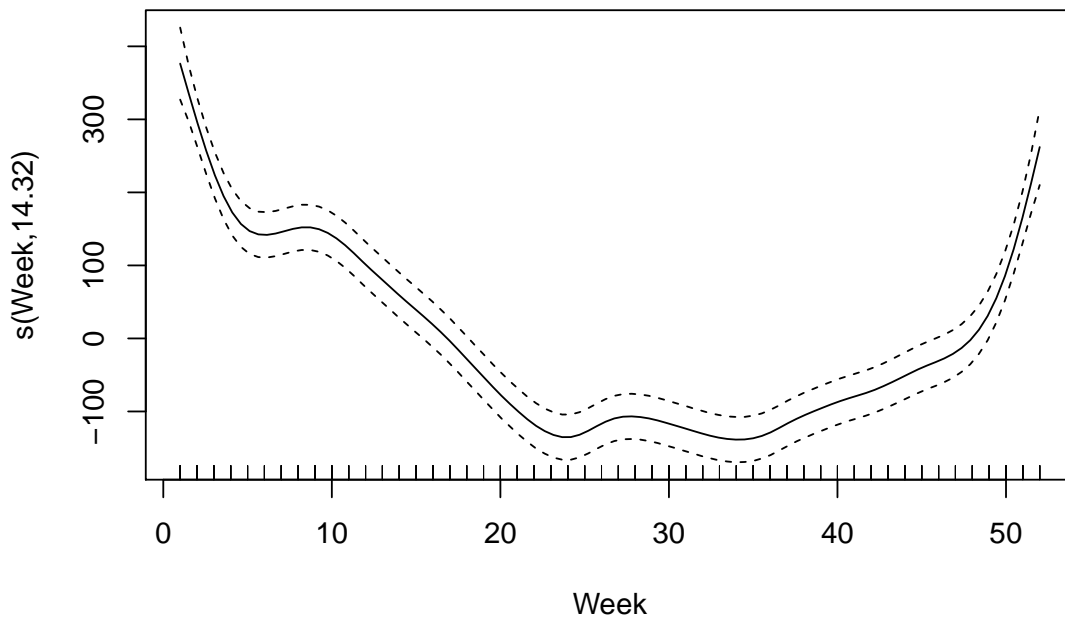
```
## [1] "Train predictions mean:"
## [1] 1783.765
## [1] "Train predictions variance:"
## [1] 25981.62
```

Thus as the summary states that the model family is Gaussian, it can be said that :

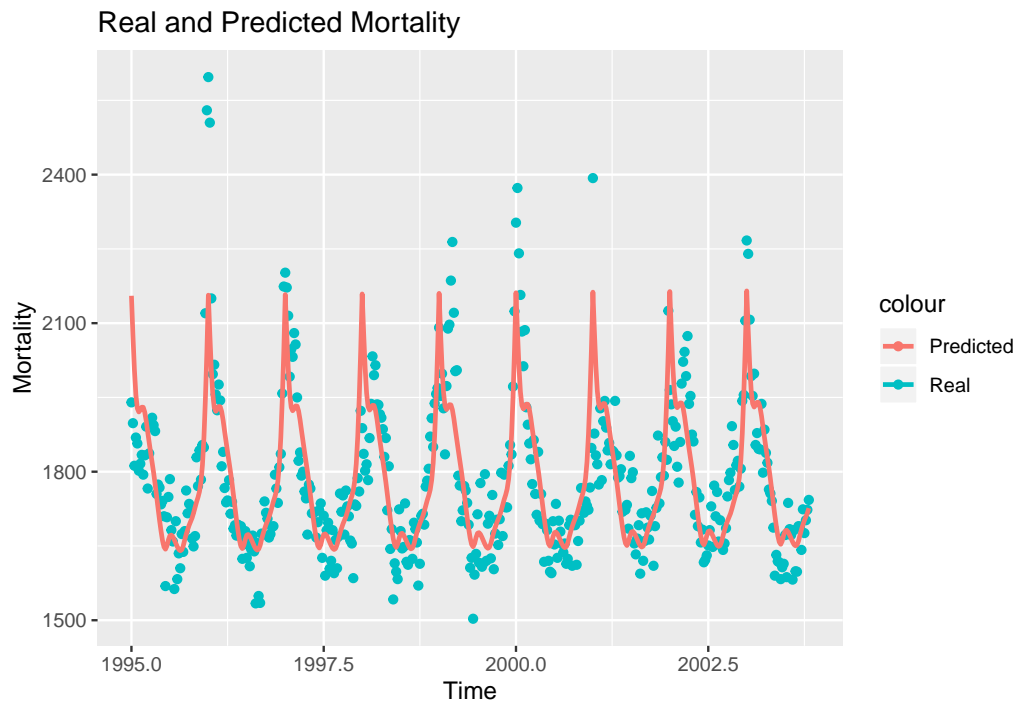
$$\begin{aligned} \text{Mortality} &\sim \text{Year} + \text{spline}(\text{Week}) \\ \text{Mortality} &\sim N(\mu = 1783.765, \sigma = 25981.62) \end{aligned}$$

The smoother curve approximated by the GAM model is:

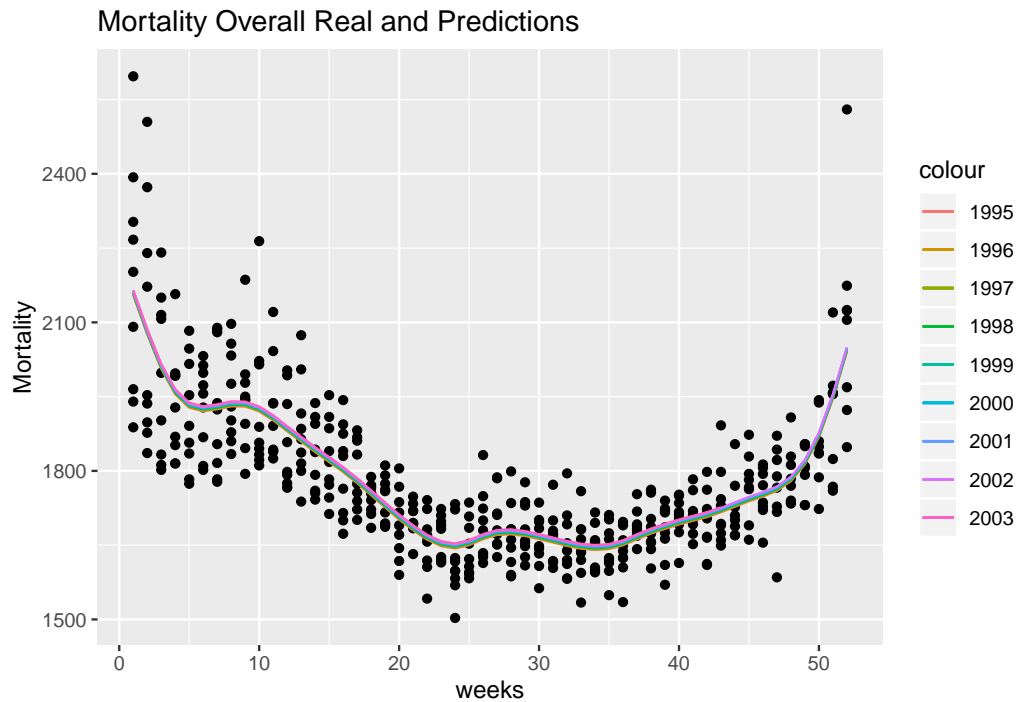
Smooth Approximation of Mortality



- Plot predicted and observed mortality against time for the fitted model and comment on the quality of the fit. Investigate the output of the GAM model and report which terms appear to be significant in the model. Is there a trend in mortality change from one year to another? Plot the spline component and interpret the plot.



Although it seems to be a fair general approximation it seems to be missing some abnormal positive and negative peaks.



It can be seen that given the variability of mortality per week each year, the GAM fit acceptably well the data, though still too general, thus not as good as we would like to.

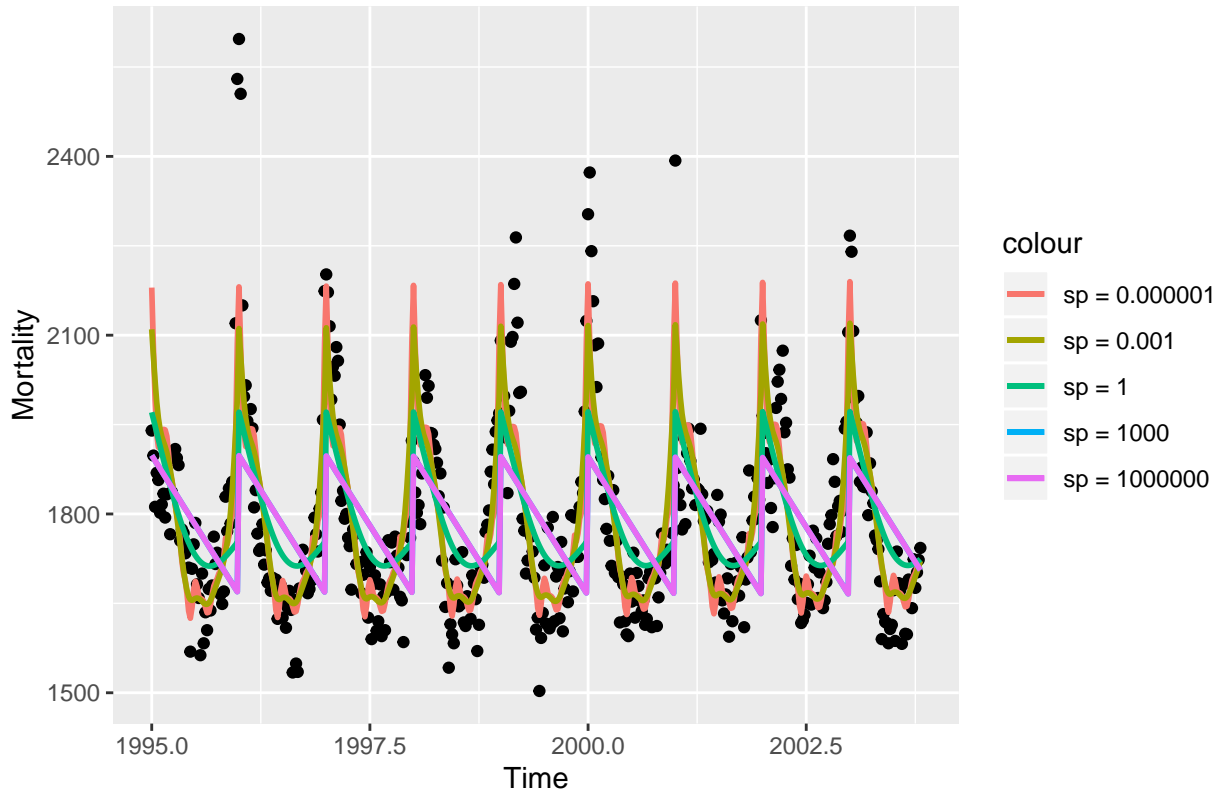
The 10 most significant terms in the model are:

```
## [1] "The 10 most important coefficients:"
## s(Week).26 s(Week).36 s(Week).25 s(Week).20 s(Week).22 s(Week).39 s(Week).18
## 6962.767 6367.790 6224.860 3619.116 2193.093 2038.671 1588.407
## s(Week).16 s(Week).40 s(Week).12
## 1478.685 1468.449 1406.516
```

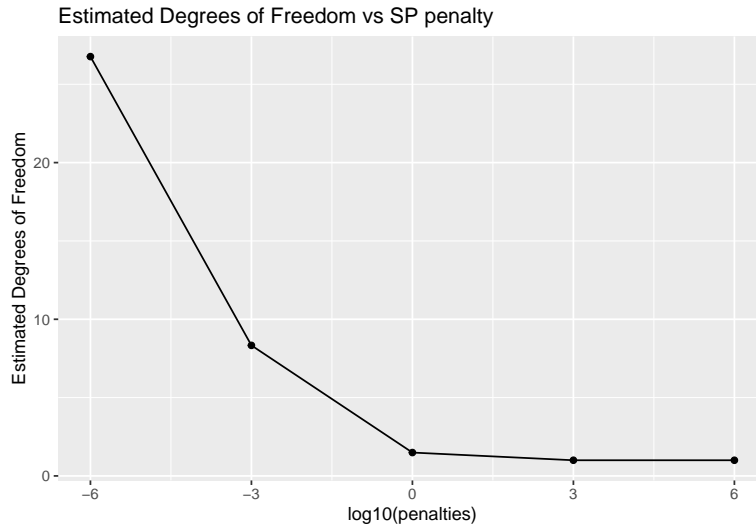
4. Examine how the penalty factor of the spline function in the GAM model from step 2 influences the estimated deviance of the model. Make plots of the predicted and observed mortality against time for cases of very high and very low penalty factors. What is the relation of the penalty factor to the degrees of freedom? Do your results confirm this relationship?

The smoothing penalty in `gam()` is given by the `sp` parameter sent to the constructor to be used in training. For this testing purpose it have been used the following penalties $\{10^{-6}, 10^{-3}, 1, 10^3, 10^6\}$

Effects of variations on SP penalization on Spline training



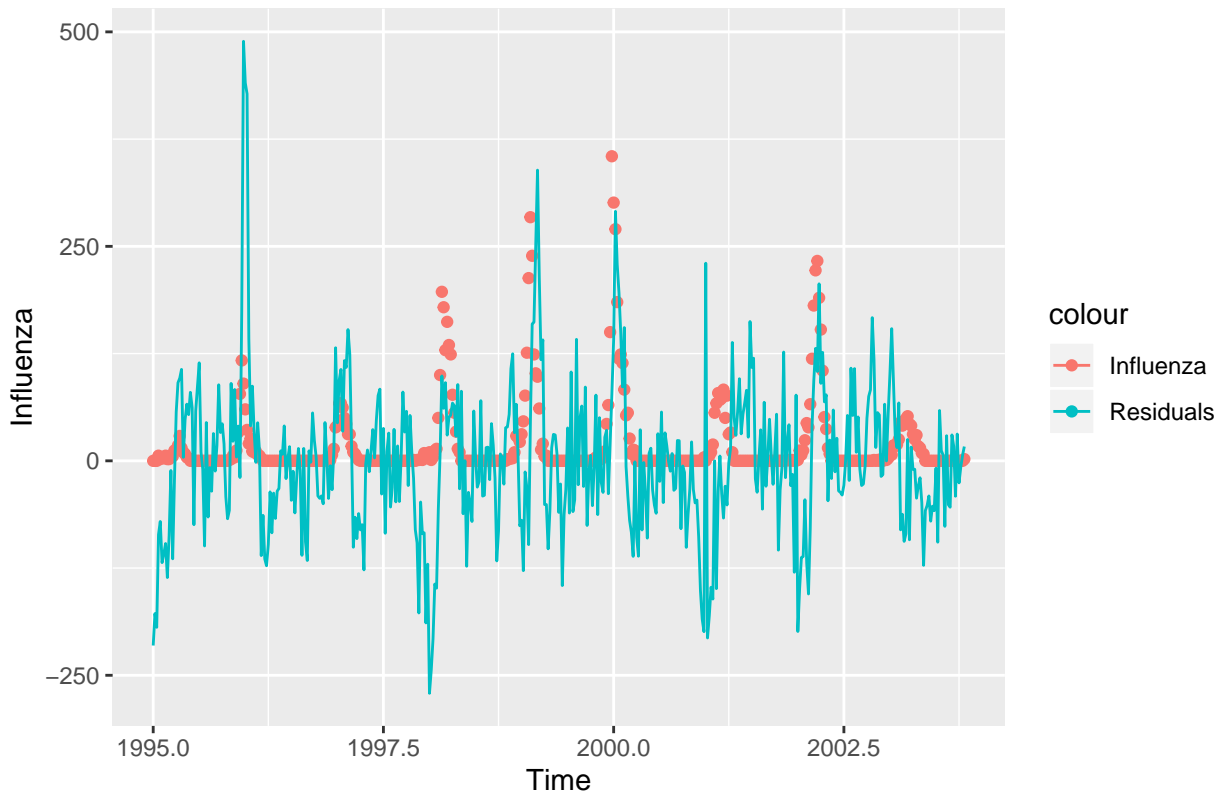
It can be seen that the lower the penalty the better the fit. In this case it looks good to have such a low penalty, though in other cases it might overfit the data.



Extracting the degrees of freedom from each penalization iteration it can be seen that the lower the penalization the more relaxed is the model, thus, the higher the degrees of freedom. The penalty factor is inversely proportional to the degrees of freedom. Indeed, by looking at the above plots, this relationship is confirmed. A high penalty essentially drops a higher number of parameters in the linear smoother matrix S_λ (by approximating them to zero).

5. Use the model obtained in step 2 and plot the residuals and the influenza values against time (in one plot). Is the temporal pattern in the residuals correlated to the outbreaks of influenza?

At a first glance Residuals seem to be noise compared to Influenza behavior.



Nonetheless visual assumptions might be incorrect. Testing numerically by getting their correlation statistic is low, thus there is no evidence to say that the residuals from the model and the influenza measures are correlated.

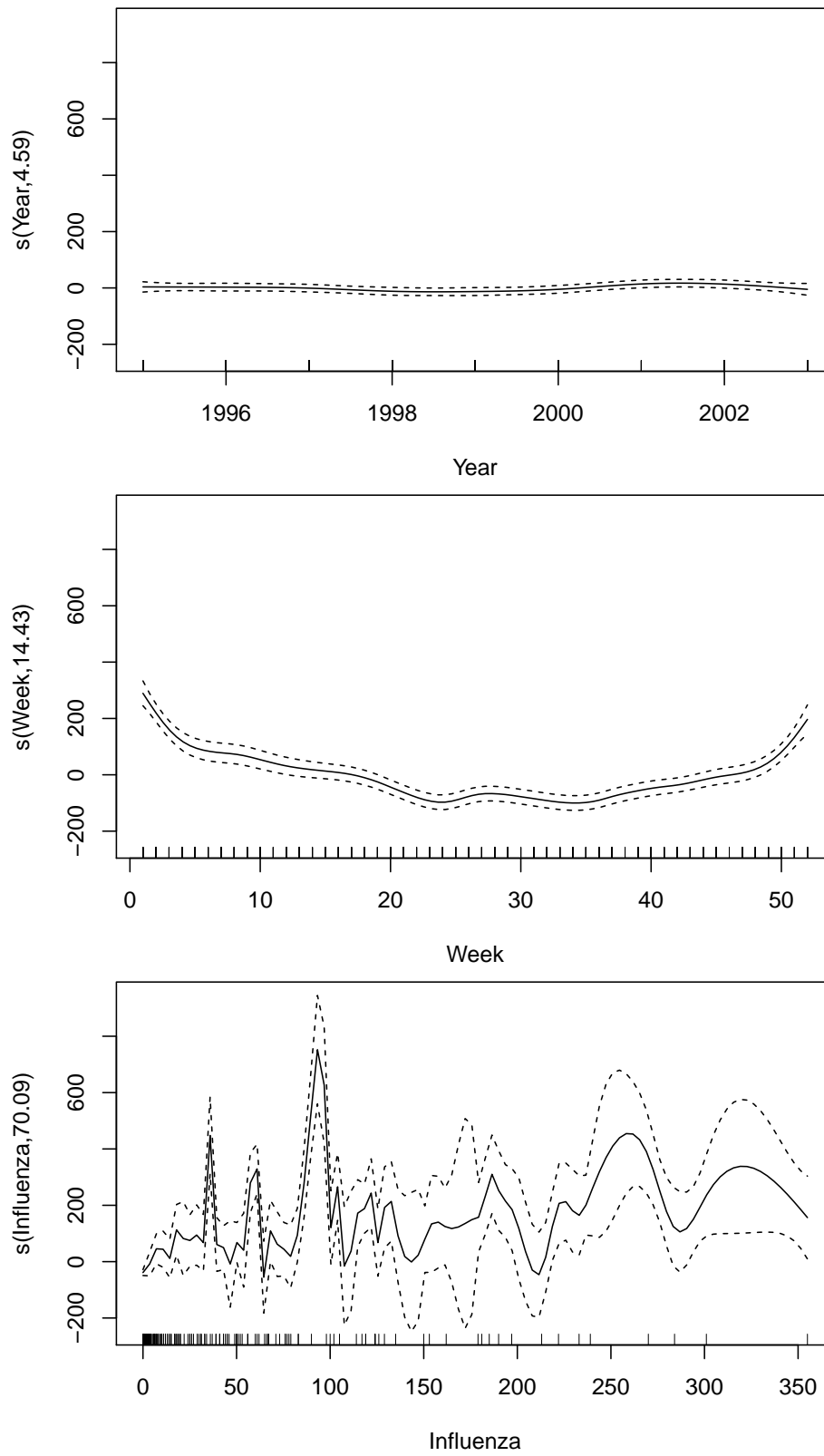
```
## [1] 0.3397395
```

6. Fit a GAM model in R in which mortality is be modelled as an additive function of the spline functions of year, week, and the number of confirmed cases of influenza. Use the output of this GAM function to conclude whether or not the mortality is influenced by the outbreaks of influenza. Provide the plot of the original and fitted Mortality against Time and comment whether the model seems to be better than the previous GAM models.

It has been trained a GAM using the following formula:

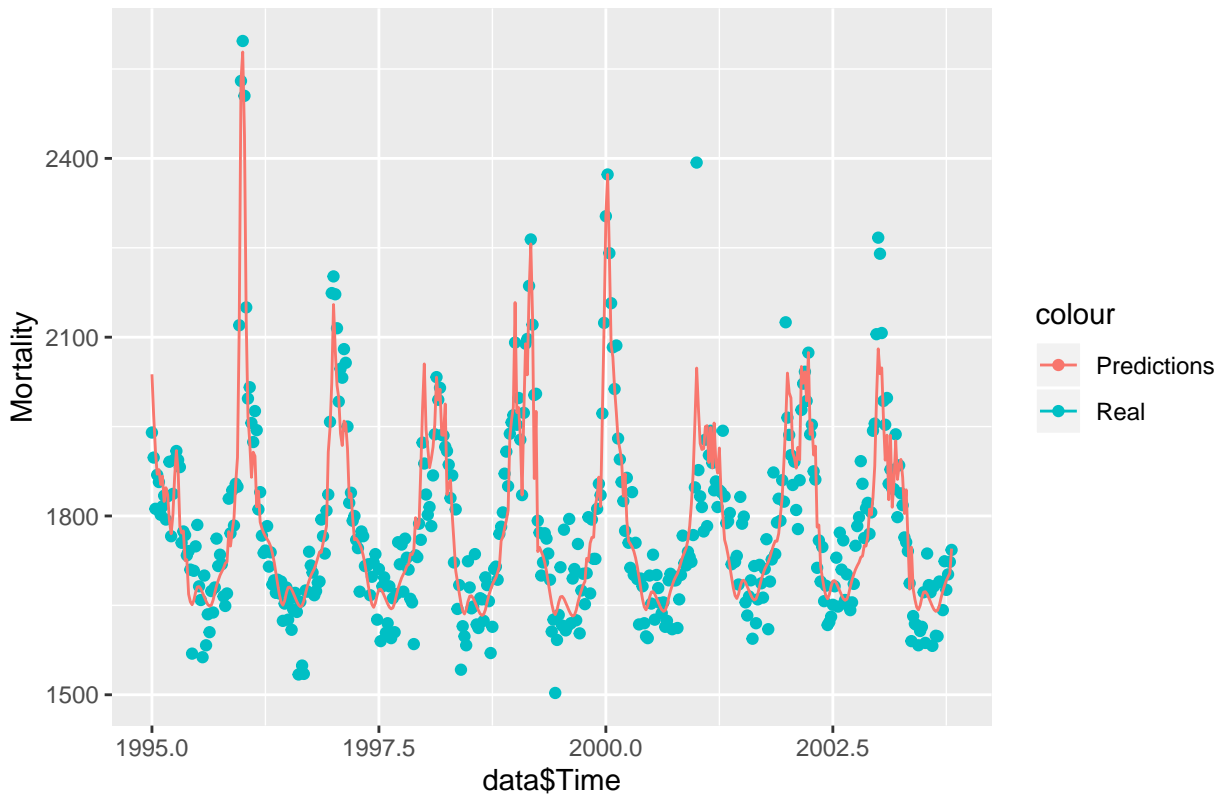
$$Mortality \approx \text{spline}(Year) + \text{spline}(Week) + \text{spline}(Influenza)$$

The obtained smoothers are the following.



Now, using the model to predict the Mortality :

Second GAM predictions



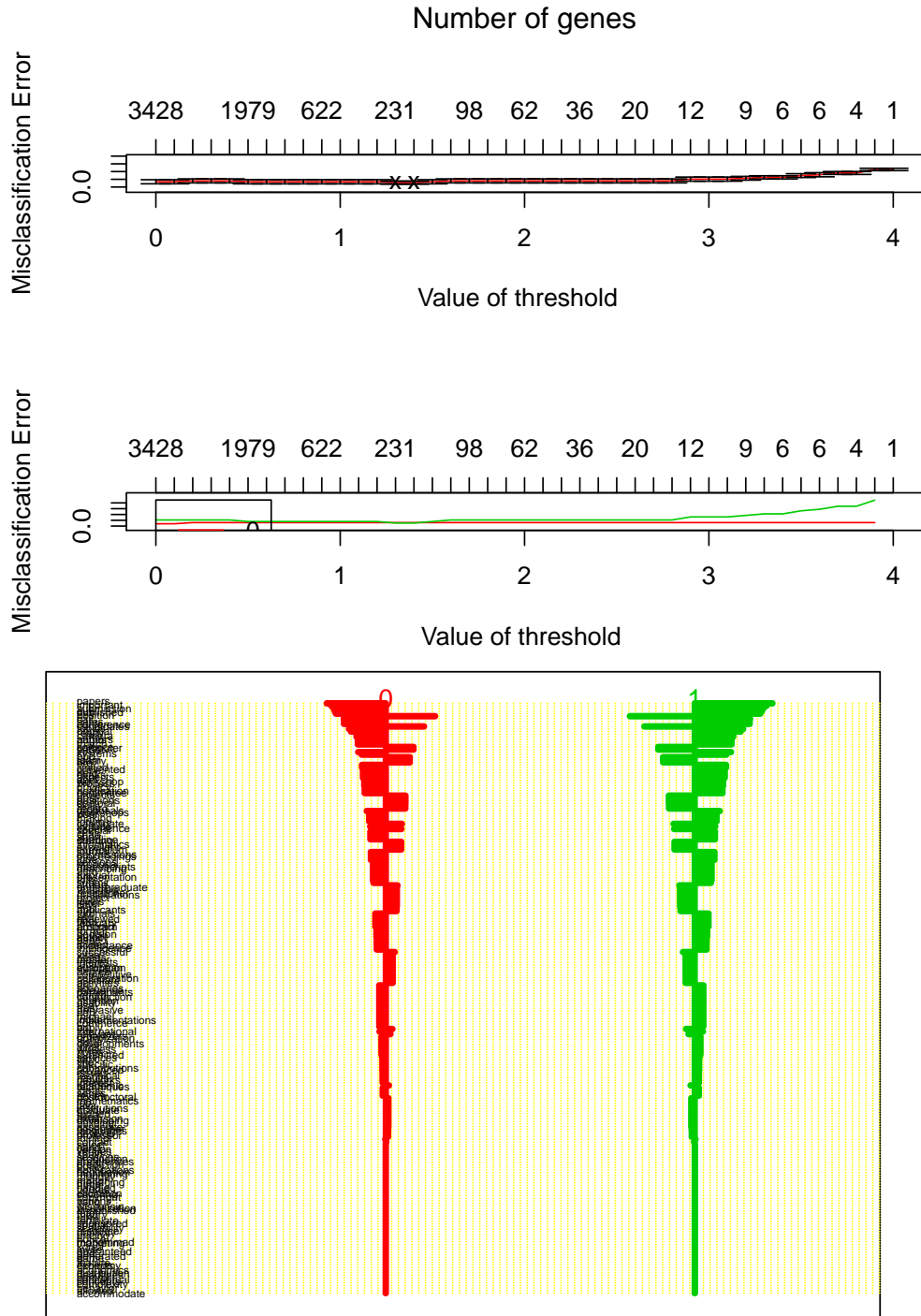
The model seems to fit better the data, though there are still some data points not being well approximated. This is sign that the predictor is not overfitted, so compared against the previous predictors, this is the best one.

Assignment 2. High-dimensional methods

The data file data.csv contains information about 64 e-mails which were manually collected from DBWorld mailing list. They were classified as: 'announces of conferences' (1) and 'everything else' (0) (variable Conference)

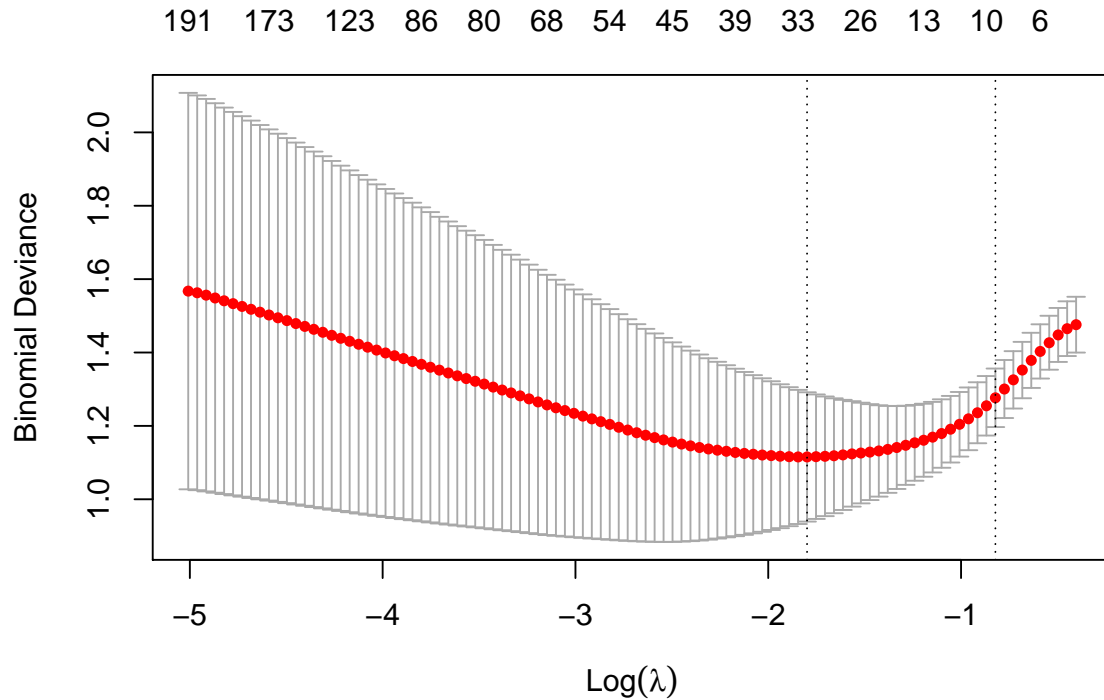
```
## Warning in RNGkind("Mersenne-Twister", "Inversion", "Rounding"): non-uniform
## 'Rounding' sampler used
```

1. Divide data into training and test sets (70/30) without scaling. Perform nearest shrunken centroid classification of training data in which the threshold is chosen by cross-validation. Provide a centroid plot and interpret it. How many features were selected by the method? List the names of the 10 most contributing features and comment whether it is reasonable that they have strong effect on the discrimination between the conference mails and other mails? Report the test error.



We get the threshold is 1.3 with using cross-validation. From the centroid plot we can see the contribution of each word made to the result(conference or not). There are 693 features selected in total. The 10 most contributing features are “papers”, “important”, “submission”, “due”, “published”, “position”, “call”, “conference”, “dates”, “candidates”. It’s clear that these word have a strong connection to conference. The test error is 5%.

2. Compute the test error and the number of the contributing features for the following methods fitted to the training data:
 - a. Elastic net with the binomial response and $\alpha = 0.5$ in which penalty is selected by the cross-validation



10 %

b. Support vector machine with “vanilladot” kernel.

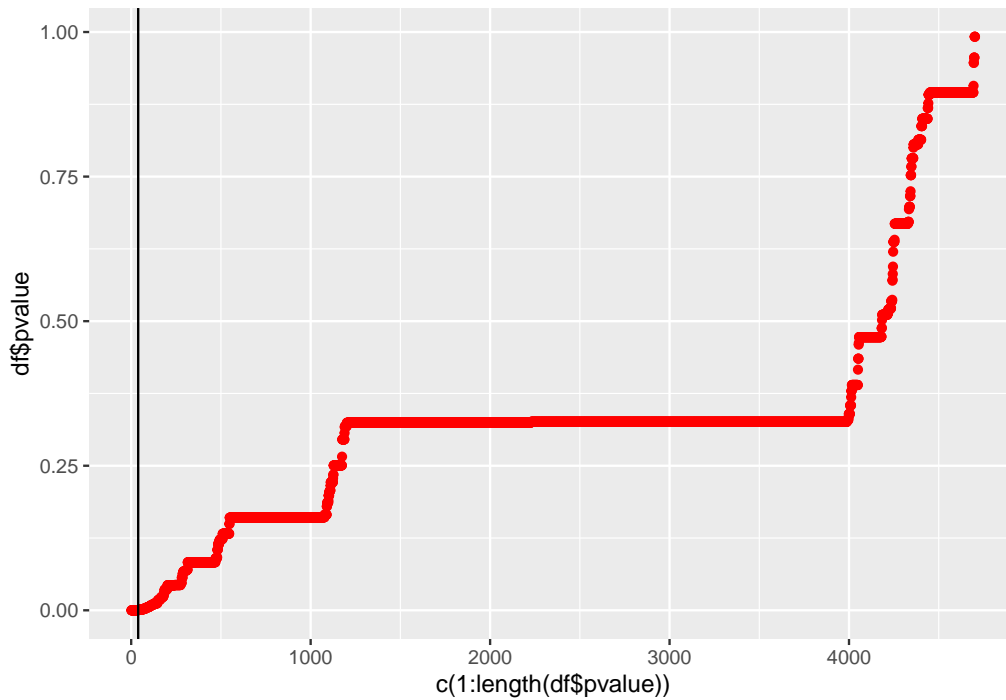
Setting default kernel parameters

5 %

Compare the results of these models with the results of the nearest shrunken centroids (make a comparative table). Which model would you prefer and why?

Error rate for Elastic net is 10% and for SVM is 5%. In this case we prefer to use SVM since it ignores the effect of high-dimensional data and it provides the lowest misclassification rate.

3. Implement Benjamini-Hochberg method for the original data, and use `t.test()` for computing p-values. Which features correspond to the rejected hypotheses? Interpret the result.



```
## [1] 39
```

##		name	pvalue
##	3036	papers	1.116910e-10
##	4060	submission	7.949969e-10
##	3187	position	8.219362e-09
##	3364	published	1.835157e-07
##	2049	important	3.040833e-07
##	596	call	3.983540e-07
##	869	conference	5.091970e-07
##	607	candidates	8.612259e-07
##	1045	dates	1.398619e-06
##	3035	paper	1.398619e-06
##	4282	topics	5.068373e-06
##	2463	limited	7.907976e-06
##	606	candidate	1.190607e-05
##	599	camera	2.099119e-05
##	3433	ready	2.099119e-05
##	389	authors	2.154461e-05
##	3125	phd	3.382671e-05
##	3312	projects	3.499123e-05
##	2974	org	3.742010e-05
##	681	chairs	5.860175e-05
##	1262	due	6.488781e-05
##	2990	original	6.488781e-05
##	2889	notification	6.882210e-05
##	3671	salary	7.971981e-05
##	3458	record	9.090038e-05
##	3891	skills	9.090038e-05
##	1891	held	1.529174e-04
##	4177	team	1.757570e-04
##	3022	pages	2.007353e-04

```

## 4628      workshop 2.007353e-04
## 810       committee 2.117020e-04
## 3285    proceedings 2.117020e-04
## 272        apply 2.166414e-04
## 4039      strong 2.246309e-04
## 2175 international 2.295684e-04
## 1088        degree 3.762328e-04
## 1477    excellent 3.762328e-04
## 3191         post 3.762328e-04
## 3243    presented 3.765147e-04

```

There are 39 features correspond to the rejected hypotheses with $\alpha = 0.05$.

We are testing:

H_0, j : word j has no effect on classification.

H_1, j : word j has effect on classification.

The variables correspondent to the rejected NULL hypothesis, therefore, are the ones that are significant to the classification as conference mail.

The 39 words selected are significant and are ranked in order of importance, similarly to what can be seen in the centroid plot from item 1.

Appendix A : Code for Assignment 1

```
#####  
##                               Assignment 1  
#####  
  
# Import data  
dataPath <- "data/influenza.xlsx"  
data <- read.xlsx(dataPath)  
  
# Time Series plotting  
mortPlot <- ggplot(data) +  
  geom_line(aes(x=Time, y=Mortality), color="black") + ggtitle("Moratality")  
infPlot <- ggplot(data) +  
  geom_line(aes(x=Time, y=Influenza), color="black") + ggtitle("Influenza")  
tempPlot <- ggplot(data) +  
  geom_line(aes(x=Time, y=Temperature.deficit), color="black") +  
  ggtitle("Temperature Deficit")  
infPlot  
mortPlot  
tempPlot  
  
# Time series per year  
years <- unique(data$Year)  
weeks <- unique(data$Week)  
mortData <- list()  
infData <- list()  
tempData <- list()  
for(i in 1:length(years)) {  
  year <- years[i]  
  mortData[[i]] <- data$Mortality[which(data$Year == year)]  
  infData[[i]] <- data$Influenza[which(data$Year == year)]  
  tempData[[i]] <- data$Temperature.deficit[which(data$Year == year)]  
}  
names(mortData) <- years  
names(infData) <- years  
names(tempData) <- years  
  
# create data.frames for ggplot  
plotData <- function(d, title) {  
  shortWeeks <- 1:length(d$'2003')  
  p <- ggplot() +  
    geom_line(aes(x=weeks, y=d$'1995', color="1995")) +  
    geom_line(aes(x=weeks, y=d$'1996', color="1996")) +  
    geom_line(aes(x=weeks, y=d$'1997', color="1997")) +  
    geom_line(aes(x=weeks, y=d$'1998', color="1998")) +  
    geom_line(aes(x=weeks, y=d$'1999', color="1999")) +  
    geom_line(aes(x=weeks, y=d$'2000', color="2000")) +  
    geom_line(aes(x=weeks, y=d$'2001', color="2001")) +  
    geom_line(aes(x=weeks, y=d$'2002', color="2002")) +  
    geom_line(aes(x=shortWeeks, y=d$'2003', color="2003")) +  
    geom_point(aes(x=weeks, y=d$'1995', color="1995")) +  
    geom_point(aes(x=weeks, y=d$'1996', color="1996")) +  
    geom_point(aes(x=weeks, y=d$'1997', color="1997")) +
```

```

    geom_point(aes(x=weeks, y=d$'1998', color="1998")) +
    geom_point(aes(x=weeks, y=d$'1999', color="1999")) +
    geom_point(aes(x=weeks, y=d$'2000', color="2000")) +
    geom_point(aes(x=weeks, y=d$'2001', color="2001")) +
    geom_point(aes(x=weeks, y=d$'2002', color="2002")) +
    geom_point(aes(x=shortWeeks, y=d$'2003', color="2003")) +
    ggtitle(title) + ylab(title)
  return(p)
}

mortPlot <- plotData(mortData, "Mortality")
infPlot <- plotData(infData, "Influenza")
tempPlot <- plotData(tempData, "Temperature Deficit")
mortPlot
infPlot
tempPlot

# Training GAM
model <- gam (
  Mortality ~ Year +
    s(Week, k = length(unique(data$Week))),
  data=data,
  method = "GCV.Cp"
)
summary(model)
modelMean <- mean(model$y)
modelVar <- var(model$y)
print("Train predictions mean:")
modelMean
print("Train predictions variance:")
modelVar
plot(model, main="Smooth Approximation of Mortality")

# Predictions
predictions <- predict(model, data)
ggplot(data) +
  geom_point(aes(x=Time, y=Mortality, color="Real")) +
  geom_line(aes(x=Time, y=predictions, color="Predicted"), size=1) +
  ggtitle("Real and Predicted Mortality")

# Plot GAM predictions per year
plotAll <- function(d, p, title) {
  shortWeeks <- 1:length(d$'2003')
  a <- ggplot() +
    geom_point(aes(x=weeks, y=d$'1995')) +
    geom_point(aes(x=weeks, y=d$'1996')) +
    geom_point(aes(x=weeks, y=d$'1997')) +
    geom_point(aes(x=weeks, y=d$'1998')) +
    geom_point(aes(x=weeks, y=d$'1999')) +
    geom_point(aes(x=weeks, y=d$'2000')) +
    geom_point(aes(x=weeks, y=d$'2001')) +
    geom_point(aes(x=weeks, y=d$'2002')) +
    geom_point(aes(x=shortWeeks, y=d$'2003')) +
    geom_line(aes(x=weeks, y=p$'1995', color="1995")) +

```

```

    geom_line(aes(x=weeks, y=p$'1996', color="1996")) +
    geom_line(aes(x=weeks, y=p$'1997', color="1997")) +
    geom_line(aes(x=weeks, y=p$'1998', color="1998")) +
    geom_line(aes(x=weeks, y=p$'1999', color="1999")) +
    geom_line(aes(x=weeks, y=p$'2000', color="2000")) +
    geom_line(aes(x=weeks, y=p$'2001', color="2001")) +
    geom_line(aes(x=weeks, y=p$'2002', color="2002")) +
    geom_line(aes(x=shortWeeks, y=p$'2003', color="2003")) +
    ggtitle(title)
  return(a)
}

preds <- list()
for (i in 1:length(years)) {
  year <- years[i]
  d <- data[which(data$Year == year),]
  preds[[i]] <- predict(model, d)
}
names(preds) <- years
p <- plotAll(mortData, preds, "Mortality Overall Real and Predictions")
p <- p + ylab("Mortality")
p

# significant terms.
orderedCoefficients <- model$coefficients[order(model$coefficients, decreasing = TRUE)]
print("The 10 most important coefficients:")
orderedCoefficients[1:10]

# Penalty factor analysis
penalties <- c(1e-6, 1e-3, 1, 1e3, 1e6)
predictions <- list()
estDegFreedom <- vector(length = length(penalties))
for (i in 1:length(penalties)) {
  model <- gam ( Mortality ~ Year +
                 s(Week, k = length(unique(data$Week))),
                 data=data,
                 method = "GCV.Cp",
                 sp = penalties[i]
               )
  predictions[[i]] <- predict(model, data)
  estDegFreedom[i] <- summary(model)$edf
}
df <- data.frame(
  time = data$Time,
  real = data$Mortality,
  sp_1u = predictions[[1]],
  sp_1m = predictions[[2]],
  sp_1 = predictions[[3]],
  sp_1k = predictions[[4]],
  sp_1M = predictions[[5]]
)

```



```

ggplot(df) +
  geom_point(aes(x=time, y=real), color="black") +
  geom_line(aes(x=time, y=sp_1u, color="sp = 0.000001"), size=1) +
  geom_line(aes(x=time, y=sp_1m, color="sp = 0.001"), size=1) +
  geom_line(aes(x=time, y=sp_1, color="sp = 1" ), size=1) +
  geom_line(aes(x=time, y=sp_1k, color="sp = 1000"), size=1) +
  geom_line(aes(x=time, y=sp_1M, color="sp = 1000000"), size=1) +
  ggtitle("Effects of variations on SP penalization on Spline training") +
  xlab("Time") + ylab("Mortality")

# Relation between penalty factors and degrees of freedom
ggplot() +
  geom_point(aes(x=log10(penalties), y=estDegFreedom)) +
  geom_line(aes(x=log10(penalties), y=estDegFreedom)) +
  ggtitle("Estimated Degrees of Freedom vs SP penalty") +
  ylab("Estimated Degrees of Freedom")

# Comparing model residuals against influenza cases
model <- gam (
  Mortality ~ Year +
  s(Week, k = length(unique(data$Week))),
  data=data,
  method = "GCV.Cp"
)
ggplot(data) +
  geom_point(aes(x=Time, y=Influenza, color="Influenza")) +
  geom_line(aes(x=Time, y=model$residuals, color="Residuals")) +
  ggtitle("")

# Correlation measurement
corr <- cor(model$residuals, data$Influenza)
print(corr)

# GAM 2 training
model <- gam( formula = Mortality ~ s(Year, k = length(unique(data$Year))) +
  s(Week, k = length(unique(data$Week))) +
  s(Influenza, k=length(unique(data$Influenza))),
  data = data,
  method = "GCV.Cp"
)
plot(model)

# GAM 2 predictions
predictions <- predict(model, data)
ggplot() +
  geom_point(aes(x=data$Time, y=data$Mortality, color="Real")) +
  geom_line(aes(x=data$Time, y=predictions, color="Predictions")) +
  ggtitle("Second GAM predictions") + ylab("Mortality")

```

Appendix B : Code for Assignment 2

```
#####  
##                               Assignment 2  
#####  
data<-read.csv2("data/data.csv",check.names = FALSE)  
names(data)<-iconv(names(data),to="ASCII")  
RNGversion("3.5.1")  
  
n=dim(data)[1]  
set.seed(12345)  
id=sample(1:n, floor(n*0.7))  
train=data[id,]  
test=data[-id,]  
  
x<-t(train[,-4703])  
  
y<-train[[4703]]  
  
x_test<-t(test[,-4703])  
  
y_test<-test[[4703]]  
  
my_data<-list(x=x,y=as.factor(y),geneid=as.character(1:nrow(x)),genenames=rownames(x))  
my_data_test<-list(x=x_test,y=as.factor(y_test),geneid=as.character(1:nrow(x_test)),genenames=rownames(x_test))  
  
mod<-pamr.train(my_data,threshold = seq(0,4,0.1))  
cvmodel<-pamr.cv(mod,my_data)  
  
thr<-cvmodel$threshold[which.min(cvmodel$error)]  
pamr.plotcv(cvmodel)  
pred<-pamr.predict(mod,my_data_test$x,threshold = thr,type="class")  
pamr.plotcen(mod,my_data,thr)  
  
#res<-as.data.frame(pamr.listgenes(mod,my_data,thr,genenames = TRUE))  
#listgene  
#my_data$genenames[as.numeric(listgene)]  
  
library(glmnet)  
  
x<-train[,-4703]  
y<-train[[4703]]  
  
x_test<-test[,-4703]  
y_test<-test[[4703]]  
  
mod<-cv.glmnet(as.matrix(x),y,alpha=0.5,family="binomial")  
plot(mod)  
penalty_min<-mod$lambda.min  
real_mod<-glmnet(as.matrix(x),y,alpha=0.5,lambda = penalty_min,family="binomial")  
pred<-predict(real_mod,as.matrix(x_test),type="class")  
cft<-table(pred,y_test)  
mis_rate<-1-(cft[1,1]+cft[2,2])/sum(cft)
```

```

cat(mis_rate*100,"%")
fit<-ksvm(as.matrix(x),y,data=train,kernel="vanilladot",type="C-svc",scale=FALSE)
pred<-predict(fit,x_test,type="response")

cft<-table(pred,y_test)

mis_rate<-1-(cft[1,1]+cft[2,2])/sum(cft)
cat(mis_rate*100,"%")
x=as.matrix(data[, -4703])
y=as.factor(data[[4703]])

df<-data.frame(name=c(),pvalue=c())

for(i in 1:ncol(x)){
  tmpv<-t.test(x[,i]~y,alternative="two.sided",conf.level=0.95)$p.value
  tdf<-data.frame(name=colnames(x)[i],pvalue=tmpv)
  df<-rbind(df,tdf)
}
df<-df[order(df$pvalue),]

a=0.05
max_i=1
for(i in 1:length(df$pvalue)){
  if(df$pvalue[i]<=a*i/length(df$pvalue)){
    max_i=i
  }
}
ggplot()+geom_point(aes(x=c(1:length(df$pvalue)),y=df$pvalue),col="red")+geom_vline(xintercept = 39)
print(max_i)

df[1:39,]

```

Appendix C : Environment setup Code

```
knitr::opts_chunk$set(echo = FALSE)
library(openxlsx)
library(ggplot2)
library(mgcv)
library(pamr)
library(glmnet)
library(kernlab)
library(readr)
RNGversion('3.5.1')
set.seed(12345)
```