

Lab3

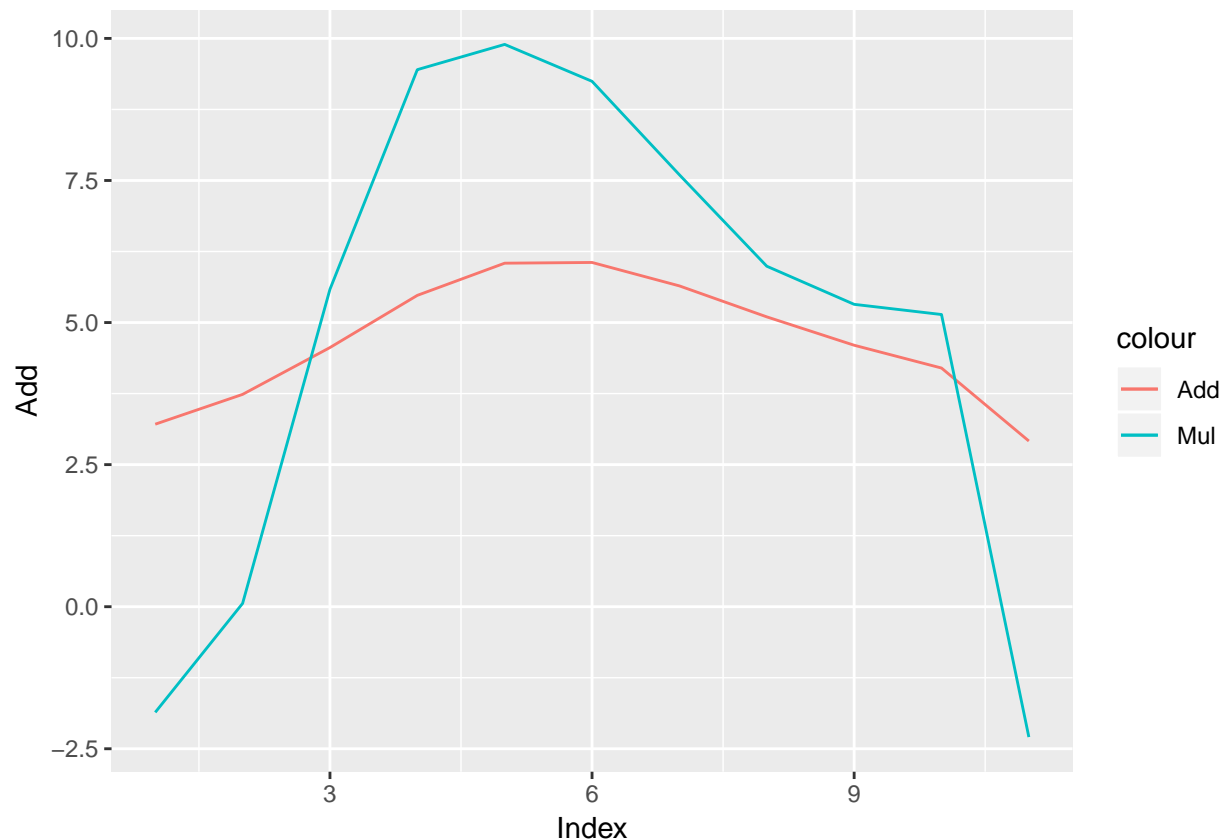
Obaid, Sridhar, Naveen

18 December 2018

Contents

Assignment 1 - Kernel Methods (Temperatue Forecast)	2
Assignment2 : Support Vector Machine	3
1. SVM with $C=0.5,1,5$	3
2. Generalization Error	3
3. SVM returned to user	3
4. Purpose of the parameter C	3
Appendix	4

Assignment 1 - Kernel Methods (Temperatue Forecast)



These are the two predictions we are getting on Adding all the kernels and Multiplying all the kernels. They are pretty close to each other for this given parameter values of the kernel. This prediction is been made for the 4th of May 2011, from 4 in the morning to 24 at night for that day. The temperature predicted by both the kernels is almost similar, and the temperatures looked to be correct for a day in month of May, these are some of the reasons we selected these values for the parameters.

The parameter value for distance we chose 30000 meters, that is 30 kms. We chose this value as the temperature in an area would be similar for an approximate 30kms radius. This is the reason we gave “h_distance=30000”. The function “distHaversine” gives the result in meters, so we had to specify the value in meters.

The parameter value for date we chose is 12 days. This means the temperature of the date we chose to predict for will be dependent on 12 days prior to that date. Those prior dates will be given higher weights for prediction of temperature. This seems reasonable for us as the tempreatre usually does not fluctuate much and is quiet similar along the 10 - 12 days margin.

The parameter value for time we chose is 4 hours. This means the temperature of the time of date we chose to predict for will be dependent on 4 hours prior to that time of the day. We thought the temperature value would not fluctuate much in 3-4 hours time period, it would be similar.

The multiplicative kernel peaks at the mid day and falls down again, but the additive kernel has a large standard deviation and is smooth in predictions. There are no sudden rise and fall in temperatures in a day according to the additive kernel, which seems right. Therefor we would say the additive kernel performs better than the multiplicative kernel.

But the additive kernel seems to give a lower bound on the temperatue and the multiplicative kernel gives

a upper bound on the temperature prediction. So if we somehow take an average on the prediction of the kernel that could be a better prediction of temperature.

The prediction of the additive kernel moves smoothly because it is not much affected by sudden changes in temperature, as we are adding the values and averaging over it. But in case of multiplicative kernel the kernel gets affected a lot with the sudden changes in temperature and this is the reason it has such sudden movements in the plot and is not smooth whereas the additive plot is smooth.

Assignment2 : Support Vector Machine

1. SVM with C=0.5,1,5

I have divided the data into train, validation and test in the ratio of 60%-20%-20%. The model is trained using a train dataset and then it is validated. By looking at the missclassification error on validation dataset, it seems C=1 gives the minimum error so I chose C=1 for SVM model.

Table 1: Missclassification Error on Validation Set

	C=0.5	C=1	C=5
Missclass_Error	0.0967391	0.0793478	0.0804348

2. Generalization Error

Choosing C=1, the model is trained on training and validation set. The prediction is done on test data. The generalization error reported is 6.9%.

Confusion Matrix for C=1 :

	Predicted	
Expected	nonspam	spam
nonspam	526	29
spam	35	331

Missclassification Error with C=1: 0.06948969

3. SVM returned to user

The model is trained on whole data. Below model would be returned to user.

```
svm_model <- ksvm(type=.,data=spamdata,kernel="rbfdot",kpar=list(sigma=0.05)
,C=1)
```

4. Purpose of the parameter C

Parameter C signifies the cost for residuals. Higher cost would result in high variance and less bias. This means that higher the cost C the margins will become more narrow to decrease the missclassification rate.

Appendix

```
knitr::opts_chunk$set(
  echo = FALSE,
  message = FALSE,
  warning = FALSE
)
library(kernlab)
library(tidyverse)
library(geosphere)
#Assignment 1
kernal_methods <- function(st, date, lat, lon, time_seq, h_date, h_time, h_distance){
  st$h_dist = abs(distHaversine(p1 = c(lon, lat), p2 = st[,c("longitude", "latitude"])))
  st$h_dist = exp(-(st$h_dist/h_distance)^2)
  st$h_date = as.numeric(difftime(date, st$date, units = c("days")))
  st$h_date = ifelse(st$h_date>0, st$h_date, 0)
  st = subset(st, st$h_date!=0)
  st$h_date = exp(-(st$h_date/h_date)^2)
  times = c()
  for(t in 1:length(time_seq)){
    d = as.Date(time_seq[t])
    c_time = format(time_seq[t], "%H:%M:%S")
    times = append(times, c(c_time))
    st[c_time] = as.numeric(abs(difftime(strptime(paste(d, c_time),
                                                    "%Y-%m-%d%H:%M:%S"),
                                          strptime(paste(d, st$time),
                                                    "%Y-%m-%d%H:%M:%S"),
                                          units = c("hour"))))
    st[c_time] = exp(-(st[c_time]/h_time)^2)
  }
  temp = st[times]
  temp_add_d = temp + (st$h_date + st$h_dist)
  temp_mul_d = temp * (st$h_date * st$h_dist)
  temp_add_n = temp_add_d*st$air_temperature
  temp_mul_n = temp_mul_d*(st$air_temperature)

  temp_add = colSums(temp_add_n)/colSums(temp_add_d)
  temp_mul = colSums(temp_mul_n)/colSums(temp_mul_d)

  d = data.frame(Index = 1:length(times), Time = times, Add = temp_add, Mul = temp_mul)
  return(d)
}

set.seed(1234567890)
stations = read.csv("stations.csv")
temps = read.csv("temps50k.csv")
st <- merge(stations,temps,by="station_number")
#colnames(st)
h_distance <- 30000
h_date <- 12
h_time <- 4
lat <- 58.4274 # The point to predict (up to the students)
lon <- 14.826
```

```

date <- "2014-05-04" # The date to predict (up to the students)
start <- as.POSIXct(date)
interval <- 60
end <- start + as.difftime(1, units="days")
time_seq <- seq(from=start, by=interval*120, to=end)
time_seq = time_seq[3:length(time_seq)]

pred = kernal_methods(st, date, lat, lon, time_seq, h_date, h_time, h_distance)
ggplot(pred) + geom_line(aes(Index, Add, col="Add")) + geom_line(aes(Index, Mul, col="Mul"))
data("spam", package="kernlab")
spamdata <- spam

nr <- nrow(spamdata)
set.seed(12345)
id <- sample(1:nr, floor(0.6*nr))

spamdata_train <- spamdata[id,]

id2 <- setdiff(1:nr, id)
id3 <- sample(id2, floor(0.5*length(id2)))
id4 <- setdiff(id2, id3)
spamdata_valid <- spamdata[id3,]
spamdata_test <- spamdata[id4,]

#rbfdot is the radial basis function for the kernel type.
#C signifies the cost for residuals. Higher cost would result in
#high variance and less bias.

#With C=0.5
svm_model <- ksvm(type~., data=spamdata_train, kernel="rbfdot", kpar=list(sigma=0.05),
                  prob.model=TRUE, C=0.5)
predict1 <- predict(svm_model, newdata=spamdata_valid, type="response")

#With C=1
svm_model2 <- ksvm(type~., data=spamdata_train, kernel="rbfdot", kpar=list(sigma=0.05),
                  prob.model=TRUE, C=1)
predict2 <- predict(svm_model2, newdata=spamdata_valid, type="response")

#With C=5
svm_model3 <- ksvm(type~., data=spamdata_train, kernel="rbfdot", kpar=list(sigma=0.05),
                  prob.model=TRUE, C=5)
predict3 <- predict(svm_model3, newdata=spamdata_valid, type="response")

conf_table <- matrix(nrow=0, ncol=3)
colnames(conf_table) <- c("C=0.5", "C=1", "C=5")

conf_matrix1 <- table("Expected" = spamdata_valid$type, "Predicted"=predict1)

# cat("\n\nConfusion Matrix for C=0.5 :\n\n ")
# conf_matrix1

missclas1 <- sum(conf_matrix1[1,2]+conf_matrix1[2,1])/nrow(spamdata_valid)

```

```

conf_matrix2 <- table("Expected" = spamdata_valid$type, "Predicted"=predict2)

# cat("\n\nConfusion Matrix for C=0.5 :\n\n ")
# conf_matrix2

missclas2 <- sum(conf_matrix2[1,2]+conf_matrix2[2,1])/nrow(spamdata_valid)

conf_matrix3 <- table("Expected" = spamdata_valid$type, "Predicted"=predict3)

# cat("\n\nConfusion Matrix for C=0.5 :\n\n ")
# conf_matrix3

missclas3 <- sum(conf_matrix3[1,2]+conf_matrix3[2,1])/nrow(spamdata_valid)

conf_table <- rbind(conf_table, c(missclas1, missclas2, missclas3))
rownames(conf_table) <- "Missclass_Error"

knitr::kable(conf_table, caption="Missclassification Error on Validation Set")

spamdata_train2 <- spamdata[c(id, id3),]

svm_model <- ksvm(type~., data=spamdata_train2, kernel="rbfdot", kpar=list(sigma=0.05),
                  C=1)
predict <- predict(svm_model, newdata=spamdata_test)

conf_matrix <- table("Expected" = spamdata_test$type, "Predicted"=predict)

cat("Confusion Matrix for C=1 :\n")
conf_matrix

missclas <- sum(conf_matrix[1,2]+conf_matrix[2,1])/nrow(spamdata_test)

cat("\n\nMissclassification Error with C=1:", missclas)

svm_model <- ksvm(type~., data=spamdata, kernel="rbfdot", kpar=list(sigma=0.05),
                  C=1)

```