

STA6106 Statistical Computing Project 3

Robert Norberg, Jung-Han Wang

Sunday, November 23, 2014

Problem 1

```
##Robert's Path
mydat <- read.table('/home/robert/cloud/Classes/STA6106 Stat Computing/Project2/Project3/training data.csv')
mydat2 <- read.table('/home/robert/cloud/Classes/STA6106 Stat Computing/Project2/Project3/data_project3.csv')

##Jung-Han's Path
# mydat <- read.table("E:/Cloud Storage/Dropbox/Life long study/Ph.D/Lecture/2014 Fall/Statistical Computing/Project2/Project3/training data.csv")
# mydat2 <- read.table("E:/Cloud Storage/Dropbox/Life long study/Ph.D/Lecture/2014 Fall/Statistical Computing/Project2/Project3/data_project3.csv")
```

This problem is to get some codes to perform the support vector data description (SVDD)

a. Write an *R* function to perform the SVDD.

First, we want to compute the kernel matrix

$$\begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \dots & \dots & \dots & \dots \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix}$$

To do this, we must define a kernel function. This function essentially calculates the distance between each pair of data vectors. For simplicity, we begin by using the simplest distance, the euclidean distance. The euclidean distance between two data vectors is just their dot product. The `kernlab` package includes a function `vanilladot()` that when called, creates another function that will compute these dot products.

```
my_kernel <- vanilladot()
```

We have now created a function `my_kernel()` that will calculate the linear distance between two data vectors for us. We check that this is equivalent to the dot product.

```
my_matrix <- as.matrix(mydat)
my_kernel(my_matrix[1, ], my_matrix[2, ]) # dot prod using kernel function
```

```
      [,1]
[1,] 37.49
```

```
crossprod(my_matrix[1, ], my_matrix[2, ]) # dot prod using base R function
```

```
      [,1]
[1,] 37.49
```

```
my_matrix[1, ] %*% my_matrix[2, ] # old school matrix multiplication operator
```

```
      [,1]  
[1,] 37.49
```

It is useful to note here that our new function `my_kernel()` is a special kind of function - that is, it belongs to the class `vanillakernel`.

```
class(my_kernel)
```

```
[1] "vanillakernel"  
attr("package")  
[1] "kernlab"
```

Regular R users are familiar with classes like `matrix` and `function`, but what is the `vanillakernel` class?

```
getClass('vanillakernel')
```

```
Class "vanillakernel" [package "kernlab"]
```

```
Slots:
```

```
Name:      .Data      kpar  
Class: function      list
```

```
Extends:
```

```
Class "kernel", directly  
Class "function", by class "kernel", distance 2  
Class "OptionalFunction", by class "kernel", distance 3  
Class "PossibleMethod", by class "kernel", distance 3  
Class "kfunction", by class "kernel", distance 3
```

The `getClass()` command reveals that objects of the class `vanillakernel` should have two slots, one named `.Data` that contains a function and one named `kpar` that contains a list. It also reveals that `vanillakernel` is an extension of the class `kernel`, which is itself an extension of the class `function`. More on this later.

Now that we have defined an R function for applying our kernel function to a pair of data vectors, we can easily create a kernel matrix from our data matrix. There is a handy function in the `kernlab` package called `kernelMatrix()` that does exactly this. It requires as arguments `kernel`, the kernel function to be used, and `x`, the data matrix from which to compute the kernel matrix. We pass the function our kernel function `my_kernel()` and our data (`mydat`) as a matrix.

```
K <- kernelMatrix(kernel=my_kernel, x=as.matrix(mydat))
```

The function returns a $n \times n$ (66×66) matrix of class `kernelMatrix`.

```
dim(K)
```

```
[1] 66 66
```

```
class(K)
```

```
[1] "kernelMatrix"
attr("package")
[1] "kernlab"
```

The SVDD problem can be stated mathematically as

$$\max_{\alpha} \sum_i \alpha_i k(x_i, x_i) - \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)$$

subject to $\alpha_i \geq 0$ and $\sum \alpha_i = 1$.

The quadratic solver in the `kernlab` package (a command called `ipop()`) solves quadratic programming problems in the form

$$\min(c'x + \frac{1}{2}x'Hx)$$

subject to $b \leq Ax \leq b + r$ and $l \leq x \leq u$.

To re-state the SVDD problem in the form required by the quadratic solver we first note that the SVDD problem is a maximization, while the solver computes a minimization. Thus, we re-state the SVDD problem as

$$\min_{\alpha} - \sum_i \alpha_i k(x_i, x_i) + \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)$$

, subject to the same constraints $\alpha_i \geq 0$ and $\sum \alpha_i = 1$.

Thus, if we set

$$x' = [\alpha_1, \alpha_2, \dots, \alpha_n]$$

$$H = 2K = 2 * \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \dots & \dots & \dots & \dots \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix}$$

$$c' = (-1) [k(x_1, x_1), k(x_2, x_2), \dots, k(x_n, x_n)] = (-1) \text{diag}(K)$$

then

$$c'x + \frac{1}{2}x'Hx = (-1) [k(x_1, x_1), \dots, k(x_n, x_n)] \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_n \end{bmatrix} + \frac{1}{2}(2) [\alpha_1, \dots, \alpha_n] \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \dots & \dots & \dots & \dots \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_n \end{bmatrix}$$

To re-state the constraints of the SVDD problem in the form required by the quadratic solver, we set

$$b = 1, A = [1, 1, \dots, 1], x = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_n \end{bmatrix}, r = 0$$

$$l = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}, x = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_n \end{bmatrix}, u = \begin{bmatrix} \infty \\ \infty \\ \dots \\ \infty \end{bmatrix}$$

then $b \leq Ax \leq b + r$ is equivalent to

$$1 \leq [1, 1, \dots, 1] \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_n \end{bmatrix} \leq 1 + 0$$

or

$$1 \leq \sum \alpha_i \leq 1$$

and $l \leq x \leq u$ is equivalent to $0 \leq \alpha_i \leq \infty$.

Using the re-formulation of the problem, we pass the appropriate objects to the `ipop()` quadratic programming solver included in the `kernlab` package. this returns the vector of α 's that minimize the stated problem $\min(c'x + \frac{1}{2}x'Hx)$.

```
my_c <- (-1)*diag(K)
my_H <- (2)*K
my_A <- rep(1, nrow(my_matrix))
my_b <- 1
my_l <- rep(0, nrow(my_matrix))
my_u <- rep(1, nrow(my_matrix))
my_r <- 0
my_solution <- ipop(c=my_c, H=my_H, A=my_A, b=my_b, l=my_l, u=my_u, r=my_r, maxiter=300, margin=0.001)
my_alphas <- my_solution@primal # use @ symbol to access s4 slot
```

We check to make sure our α_i 's sum to one,

```
sum(my_alphas)
```

```
[1] 1
```

We believe we have found a solution to the SVDD problem, so we wrap the commands shown above into a function that takes as input a data matrix `x` and a kernel function `k`.

```
SVDD <- function(x, k, tol=1E-3){
  H <- kernelMatrix(kernel=k, x=x)
  n <- nrow(x)
  solution <- ipop(c=(-1)*diag(H),
                  H=2*H,
```

```

        A=rep(1, n),
        b=1,
        l=rep(0, n),
        u=rep(1, n),
        r=0,
        maxiter=300,
        margin=tol)

alphas <- solution@primal

# catch errors
if(signif(sum(alphas), 7) != 1){
  stop("An error has occurred! The solution values do not sum to one.")
}
## if(any(signif(alphas, 7) < 0)){
##   stop("An error has occurred! The solution values include at least one negative value, which is not allowed.")
## }

svs <- x[round(alphas, abs(log10(tol)))> 0, ]
return(list(alphas=alphas, support_vectors=svs, kernel=k, tolerance=tol))
}

test_solution <- SVDD(my_matrix, my_kernel)

```

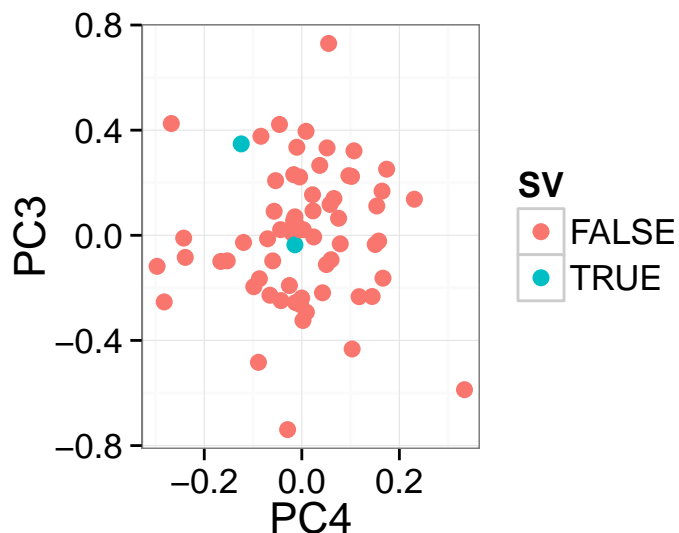
Lets check our solution to see if it makes sense.

```

my_pcs <- prcomp(my_matrix)
pc_data <- predict(my_pcs, my_matrix)
pc_data <- as.data.frame(pc_data)
pc_data$SV <- with(test_solution, round(alphas, abs(log10(tolerance))) > 0)

ggplot(pc_data, aes(x=PC4, y=PC3, color=SV))+
  geom_point(size=3)+
  theme_bw(base_size=16)

```



b. Write an R function to perform the prediction of a new observation using SVDD.

When x_s is Support Vector when $\alpha_s > 0$

Followed by, we want to select support vector x_s based on α_s .

Therefore, we combined α with original dataset, then make the selection.

```
test_solution$support_vectors
```

```
      V1 V2  V3  V4
[1,] 4.3  3 1.1 0.1
[2,] 7.6  3 6.6 2.1
```

The prediction of new observation can be calculated by solving R. The function of R can be written as:

$$R^2 = (x'_s \cdot x_s) - 2 \sum_{i=1}^N \alpha_i (x'_s \cdot x_i) + \sum_{i,j=1}^N \alpha_i \alpha_j (x'_i \cdot x_j)$$

```
find_term_3 <- function(SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance)))>0])
  k <- SVDD_obj$kernel
  term_3 <- 0
  for (i in 1:nrow(all_svs)){
    for (j in 1:nrow(all_svs)){
      tmp3 <- sv_alphas[i]*sv_alphas[j]*k(all_svs[i,], all_svs[j,])
      term_3 <- term_3+tmp3
    }
  }
  return(term_3)
}

find_r <- function(SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance)))>0])
  k <- SVDD_obj$kernel

  # pre-allocate a vector for each radius to be calculated
  r2s <- vector(length=nrow(all_svs), mode='numeric')

  # calculate term 3 just once for efficiency
  term_3 <- find_term_3(SVDD_obj)

  # term 1 and term 2 must be computed for each support vector separately
  for(i in 1:nrow(all_svs)){
    term_1 <- k(all_svs[i, ], all_svs[i, ])

    term_2 <- 0
    for (j in 1:nrow(all_svs)){
      tmp2 <- -2*sv_alphas[j]*k(all_svs[i, ], all_svs[j, ])
      term_2 <- term_2+tmp2
    }
  }
}
```

```

    r2s[i] <- term_1+term_2+term_3
  }

  # these r2s should all be similar
  if(var(r2s)>0.1){
    warning('The radii calculated have a variance larger than 0.1.')
  }
  return(mean(r2s))
}

r.sqr <- find_r(test_solution)

```

A test sample z is accepted when $\leq R^2$

```

z<-c(1,1,1,1)

predict_SVDD <- function(z, SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance)))>0])
  k <- SVDD_obj$kernel

  z1 <- k(z, z)
  z2 <- 0
  for (j in 1:length(sv_alphas)){
    tmp2 <- -2*sv_alphas[j]*k(z, z)
    z2 <- z2+tmp2
  }
  z3<-find_term_3(SVDD_obj)
  z_rad <- z1+z2+z3 # distance of point z from center
  return(z_rad)
}

predict_SVDD(z=z, SVDD_obj=test_solution)

```

c. Write an R function for detecting potential outliers for a new set of observations, along with the upper threshold.

```

detect_outliers <- function(SVDD_obj, newdata, plot=T){
  x_mat <- as.matrix(newdata)
  rad <- apply(x_mat, 1, predict_SVDD, SVDD_obj=SVDD_obj)
  SVDD_r2 <- find_r(SVDD_obj=SVDD_obj)
  outlier_ind <- rad>SVDD_r2
  num_outliers <- sum(outlier_ind)

  if(plot==T){
    plotdat <- as.data.frame(cbind(x_mat, rad))
    plot(rad, type='b')
    abline(h=SVDD_r2, col='red')
  }

  msg <- paste(num_outliers, 'outliers out of', nrow(x_mat), 'observations detected')

```

```

cat(msg, '\n')
return(outlier_ind)
}

detect_outliers(SVDD_obj=test_solution, newdata=mydat2)

```

Problem 2

The goal of problem 2 is to perform the support vector data description (SVDD) using the Mahalanobis kernel function. We will simplify the problem by using the identity function for g .

a. Write an R function to compute the Mahalanobis kernel distance $d_g(x)$

```

malhalanobis <- function(covmat){
  rval <- function(x, y = NULL) {
    if (!is(x, "vector"))
      stop("x must be a vector")
    if (!is(y, "vector") && !is.null(y))
      stop("y must be a vector")
    if (is(x, "vector") && is.null(y)) {
      t(x)%*%solve(covmat)%*%x
    }
    if (is(x, "vector") && is(y, "vector")) {
      if (!length(x) == length(y))
        stop("number of dimension must be the same on both data points")
      t(x)%*%solve(covmat)%*%y
    }
  }
  return(new('kernel', .Data=rval, kpar=list(covmat=covmat)))
}

my_mal_kernel <- malhalanobis(covmat=cov(mydat))

my_mal_kernel(c(1, 1, 1, 1), c(1, 1, 1, 1)) # test case

```

```

      [,1]
[1,] 17.01475

```

b. Write an R function to perform the Mahalanobis kernel SVDD.

Already did that.

```

mal_solution <- SVDD(my_matrix, k=my_mal_kernel)

```


- c. Write an R function to perform the prediction of a new observation using the Mahalanobis kernel SVDD.
- d. Write an R function for detecting potential outliers for a new set of observations, along with the upper threshold.

Appendix with R code

```
# Clear working environment
rm(list=ls())
library(ggplot2) # for plots
library(kernlab)

# Options for document compilation
knitr::opts_chunk$set(warning=FALSE, message=FALSE, comment=NA, fig.width=4, fig.height=3)
##Robert's Path
mydat <- read.table('/home/robert/cloud/Classes/STA6106 Stat Computing/Project2/Project3/training dataset')
mydat2 <- read.table('/home/robert/cloud/Classes/STA6106 Stat Computing/Project2/Project3/data_project3')

##Jung-Han's Path
# mydat <- read.table("E:/Cloud Storage/Dropbox/Life long study/Ph.D/Lecture/2014 Fall/Statistical Computing")
# mydat2 <- read.table("E:/Cloud Storage/Dropbox/Life long study/Ph.D/Lecture/2014 Fall/Statistical Computing")
my_kernel <- vanilladot()
my_matrix <- as.matrix(mydat)
my_kernel(my_matrix[1, ], my_matrix[2, ]) # dot prod using kernel function
crossprod(my_matrix[1, ], my_matrix[2, ]) # dot prod using base R function
my_matrix[1, ] %*% my_matrix[2, ] # old school matrix multiplication operator
class(my_kernel)
getClass('vanillakernel')
K <- kernelMatrix(kernel=my_kernel, x=as.matrix(mydat))
dim(K)
class(K)
my_c <- (-1)*diag(K)
my_H <- (2)*K
my_A <- rep(1, nrow(my_matrix))
my_b <- 1
my_l <- rep(0, nrow(my_matrix))
my_u <- rep(1, nrow(my_matrix))
my_r <- 0
my_solution <- ipop(c=my_c, H=my_H, A=my_A, b=my_b, l=my_l, u=my_u, r=my_r, maxiter=300, margin=0.001)
my_alphas <- my_solution@primal # use @ symbol to access slot
sum(my_alphas)
SVDD <- function(x, k, tol=1E-3){
  H <- kernelMatrix(kernel=k, x=x)
  n <- nrow(x)
  solution <- ipop(c=(-1)*diag(H),
                  H=2*H,
                  A=rep(1, n),
                  b=1,
                  l=rep(0, n),
                  u=rep(1, n),
                  r=0,
                  maxiter=300,
```

```

        margin=tol)

alphas <- solution@primal

# catch errors
if(signif(sum(alphas), 7) != 1){
  stop("An error has occurred! The solution values do not sum to one.")
}
## if(any(signif(alphas, 7) < 0)){
##   stop("An error has occurred! The solution values include at least one negative value, which is not allowed")
## }

svs <- x[round(alphas, abs(log10(tol)))> 0, ]
return(list(alphas=alphas, support_vectors=svs, kernel=k, tolerance=tol))
}

test_solution <- SVDD(my_matrix, my_kernel)
my_pcs <- prcomp(my_matrix)
pc_data <- predict(my_pcs, my_matrix)
pc_data <- as.data.frame(pc_data)
pc_data$SV <- with(test_solution, round(alphas, abs(log10(tolerance)))) > 0)

ggplot(pc_data, aes(x=PC4, y=PC3, color=SV))+
  geom_point(size=3)+
  theme_bw(base_size=16)
test_solution$support_vectors

find_term_3 <- function(SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance))))>0])
  k <- SVDD_obj$kernel
  term_3 <- 0
  for (i in 1:nrow(all_svs)){
    for (j in 1:nrow(all_svs)){
      tmp3 <- sv_alphas[i]*sv_alphas[j]*k(all_svs[i,], all_svs[j,])
      term_3 <- term_3+tmp3
    }
  }
  return(term_3)
}

find_r <- function(SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance))))>0])
  k <- SVDD_obj$kernel

  # pre-allocate a vector for each radius to be calculated
  r2s <- vector(length=nrow(all_svs), mode='numeric')

  # calculate term 3 just once for efficiency
  term_3 <- find_term_3(SVDD_obj)

  # term 1 and term 2 must be computed for each support vector separately

```

```

for(i in 1:nrow(all_svs)){
  term_1 <- k(all_svs[i, ], all_svs[i, ])

  term_2 <- 0
  for (j in 1:nrow(all_svs)){
    tmp2 <- -2*sv_alphas[j]*k(all_svs[i, ], all_svs[j, ])
    term_2 <- term_2+tmp2
  }
  r2s[i] <- term_1+term_2+term_3
}

# these r2s should all be similar
if(var(r2s)>0.1){
  warning('The radii calculated have a variance larger than 0.1.')
}
return(mean(r2s))
}

r.sqr <- find_r(test_solution)

z<-c(1,1,1,1)

predict_SVDD <- function(z, SVDD_obj){
  all_svs <- SVDD_obj$support_vectors
  sv_alphas <- with(SVDD_obj, alphas[round(alphas, abs(log10(tolerance)))>0])
  k <- SVDD_obj$kernel

  z1 <- k(z, z)
  z2 <- 0
  for (j in 1:length(sv_alphas)){
    tmp2 <- -2*sv_alphas[j]*k(z, z)
    z2 <- z2+tmp2
  }
  z3<-find_term_3(SVDD_obj)
  z_rad <- z1+z2+z3 # distance of point z from center
  return(z_rad)
}

predict_SVDD(z=z, SVDD_obj=test_solution)
detect_outliers <- function(SVDD_obj, newdata, plot=T){
  x_mat <- as.matrix(newdata)
  rad <- apply(x_mat, 1, predict_SVDD, SVDD_obj=SVDD_obj)
  SVDD_r2 <- find_r(SVDD_obj=SVDD_obj)
  outlier_ind <- rad>SVDD_r2
  num_outliers <- sum(outlier_ind)

  if(plot==T){
    plotdat <- as.data.frame(cbind(x_mat, rad))
    plot(rad, type='b')
    abline(h=SVDD_r2, col='red')
  }

  msg <- paste(num_outliers, 'outliers out of', nrow(x_mat), 'observations detected')

```

```

    cat(msg, '\n')
    return(outlier_ind)
}

detect_outliers(SVDD_obj=test_solution, newdata=mydat2)
malhalanobis <- function(covmat){
  rval <- function(x, y = NULL) {
    if (!is(x, "vector"))
      stop("x must be a vector")
    if (!is(y, "vector") && !is.null(y))
      stop("y must be a vector")
    if (is(x, "vector") && is.null(y)) {
      t(x)%*%solve(covmat)%*%x
    }
    if (is(x, "vector") && is(y, "vector")) {
      if (!length(x) == length(y))
        stop("number of dimension must be the same on both data points")
      t(x)%*%solve(covmat)%*%y
    }
  }
  return(new('kernel', .Data=rval, kpar=list(covmat=covmat)))
}

my_mal_kernel <- malhalanobis(covmat=cov(mydat))

my_mal_kernel(c(1, 1, 1, 1), c(1, 1, 1, 1)) # test case
mal_solution <- SVDD(my_matrix, k=my_mal_kernel)

```