

vector case: given a kernel, we know the mapping $h: \mathbb{R}^m \rightarrow \mathbb{R}^{m'}$

SMM Kernel Method

Chanwoo Lee, April 23, 2020

1 Kernel method

Claim: if we define h_matrix in this way, then the induced matrix-valued kernel satisfies the properties in Remark 2

I am suggesting new kernel method which makes optimization easier. Define feature mapping

$h: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m' \times n}$ where $m < m'$. Kernels for matrix case can be define as

$h_matrix: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m' \times n}$

$$K(X, X') = h(X)^T h(X') \in \mathbb{R}^{n \times n}.$$

$X=[x_1, \dots, x_n] \rightarrow h_matrix(X)=[h(x_1), \dots, h(x_n)]$ where $h(\cdot)$ is the classical mapping in vector space

Our objective primal problem is

$$\begin{aligned} \min_{U \in \mathbb{R}^{m' \times r}, V \in \mathbb{R}^{n \times r}, \xi} & \frac{1}{2} \|UV^T\|^2 + c \sum_{i=1}^N \xi_i \\ \text{subject to} & y_i (\langle UV^T, h(X_i) \rangle + b) \leq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, N. \end{aligned} \quad (1)$$

First, fix V and solve (1) with respect to U . We have the following dual problem.

$$\begin{aligned} \min_{\alpha} & - \sum_{i=1}^N \alpha_i + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle h(X_i) H_V, h(X_j) H_V \rangle \\ \text{subject to} & \sum_{i=1}^N y_i \alpha_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, N. \end{aligned}$$

Notice $\langle h(X_i) H_V, h(X_j) H_V \rangle = \text{tr}(H_V h(X_i)^T h(X_j)) = \text{tr}(H_V K(X_i, X_j))$. Therefore, we can update U as

$$U = \sum_{i=1}^N \alpha_i y_i h(X_i) V (V^T V)^{-1} \quad (2)$$

where h function is not known. We are going to borrow this formula to update V in the next step. Now, fix U and solve (1) with respect to V . We have the following dual problem.

$$\begin{aligned} \min_{\beta} & - \sum_{i=1}^N \beta_i + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \beta_i \beta_j y_i y_j \langle H_U h(X_i), H_U h(X_j) \rangle \\ \text{subject to} & \sum_{i=1}^N y_i \beta_i = 0 \\ & 0 \leq \beta_i \leq C, \quad i = 1, \dots, N. \end{aligned} \quad (3)$$

To get an optimal β in (3), we need the information of $\langle H_U h(X_i), H_U h(X_j) \rangle$. Notice

$$\begin{aligned} \langle H_U h(X_i), H_U h(X_j) \rangle &= \text{tr}(H_U h(X_j) h(X_i)^T) = \text{tr}(U (U^T U)^{-1} U^T h(X_j) h(X_i)^T) \\ &= \text{tr}((U^T U)^{-1} U^T h(X_j) h(X_i)^T U) \end{aligned} \quad (4)$$

Using the (2), we have the following expression of the component in (4).

$$\begin{aligned}
U^T U &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (V^T V)^{-1} V^T h(X_i)^T h(X_j) V (V^T V)^{-1} \\
&= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (V^T V)^{-1} V^T K(X_i, X_j) V (V^T V)^{-1} \\
&= (V^T V)^{-1} V^T \left(\sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(X_i, X_j) \right) V (V^T V)^{-1},
\end{aligned} \tag{5}$$

$$\begin{aligned}
U^T h(X_j) &= \sum_{l=1}^N \alpha_l y_l (V^T V)^{-1} V^T h(X_l)^T h(X_j) \\
&= \sum_{l=1}^N \alpha_l y_l (V^T V)^{-1} V^T K(X_l, X_j) \\
&= (V^T V)^{-1} V^T \sum_{l=1}^N \alpha_l y_l K(X_l, X_j).
\end{aligned}$$

Therefore, we can get an optimal β in (3) with (5). Finally, we update V with the help of Equation (5) as

$$V^T = \sum_{i=1}^N \beta_i y_i (U^T U)^{-1} U^T h(X_i).$$

Remark 1. We can get explicit update formula V while U cannot be expressed as explicit value. This does make sense because in feature mapping $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m' \times n}$, dimension m' can be increased to arbitrary dimension while n being fixed.

Remark 2. There are some kernel functions that might be used often.

$$\begin{aligned}
\text{Linear: } K(X, X') &= X^T X' \\
\text{Polynomial: } K(X, X') &= (X^T X' + I_n)^d \\
\text{Radial: } K(X, X') &= \exp((X - X')^T (X - X') / \sigma),
\end{aligned}$$

where $\exp(A) = \sum_{n=0}^{\infty} \frac{1}{n!} A^n$. Notice that when $X \in \mathbb{R}^{m \times 1}$ i.e. X is a vector, all those definitions are reduced to SVM case. From this way, we can generalize linear SMM method to Kernel SMM with tractable algorithm.

2 Algorithm construction and simulation

2.1 Algorithm construction: objective function

Since we do not have explicit U formula in SMMK algorithm, we do not evaluate the following objective function in the SMM algorithm.

$$L(U, V, b) = \frac{1}{2} \|UV^T\|^2 + C \sum_{i=1}^N (1 - y_i \langle UV^T, h(X_i) \rangle + b)_+ \quad (6)$$

Instead, we are using the same objective function with different parameter,

$$L(V, \alpha, b) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \text{tr}(H_V K(X_i, X_j)) + C \sum_{i=1}^N \left(1 - y_i \left(\sum_{j=1}^N \alpha_j y_j \text{tr}(H_V K(X_i, X_j)) \right) + b \right)_+ \quad (7)$$

This objective function in Equation (7) can be obtained by plugging $U = \sum_{i=1}^N \alpha_i y_i h(X_i) V (V^T V)^{-1}$ in Equation (6). In addition, I evaluate the objective value after updates of U because if we evaluate after updates of V , U expressed by V is changed resulting in wrong evaluation.

2.2 Simulation

I generate random matrix $X_1, \dots, X_{100} \in \mathbb{R}^{5 \times 4}$ whose entries are from i.i.d. $\text{Unif}(-1, 1)$. I set ground truth matrix $B = UV^T \in \mathbb{R}^{5 \times 4}$ such that $U \in \mathbb{R}^{5 \times 3}$, $V \in \mathbb{R}^{4 \times 3}$ whose entries are from i.i.d. $\text{Unif}(-1, 1)$. Our classification rule is

$$y = \text{sign}(\langle B, X \rangle + 0.1).$$

First, I perform five folded cross validation. Table 1 shows the miss classification rate (MCR) of SMM and SMMK (linear kernel) method at each test. Table 2 shows the loss function value at the last iteration at each test. This shows that two methods are pretty much the same. To be more accurate, only thing that makes two methods different is random initialization. I checked they have the same update of V from the same initialization.

	1st	2nd	3rd	4th	5th	average
SMMK(linear)	0.90	0.85	0.90	0.90	0.90	0.89
SMM	0.95	0.85	0.85	0.95	0.85	0.89

Table 1: Miss classification rate in Simulation

	1st	2nd	3rd	4th	5th	average
SMMK(linear)	32.41501	17.92673	18.18015	25.83209	24.44279	23.75935
SMM	22.31593	15.56181	17.99429	26.34058	32.12252	22.86703

Table 2: Loss function value in Simulation

2.3 Simulation with new kernel

I generate feature data matrix with the following rule.

$$\begin{aligned}\mathbb{P}((X_{\cdot 1}^T, X_{\cdot 2}^T)^T | y = 1) &\sim N((1, 1, -1, -1)^T, I_4), \\ \mathbb{P}((X_{\cdot 1}^T, X_{\cdot 2}^T)^T | y = -1) &\sim N((0, 0, 0, 0)^T, I_4).\end{aligned}$$

With this rule, we have a test data set $(X^1, 1), \dots, (X^{20}, 1), (X^{21}, -1), \dots, (X^{40}, -1)$.

$$Z_1 = X_{11} + X_{21} \text{ and } Z_2 = X_{12} + X_{22}.$$

Figure 1 shows how generated data looks like with $Z1$ and $Z2$ axes.

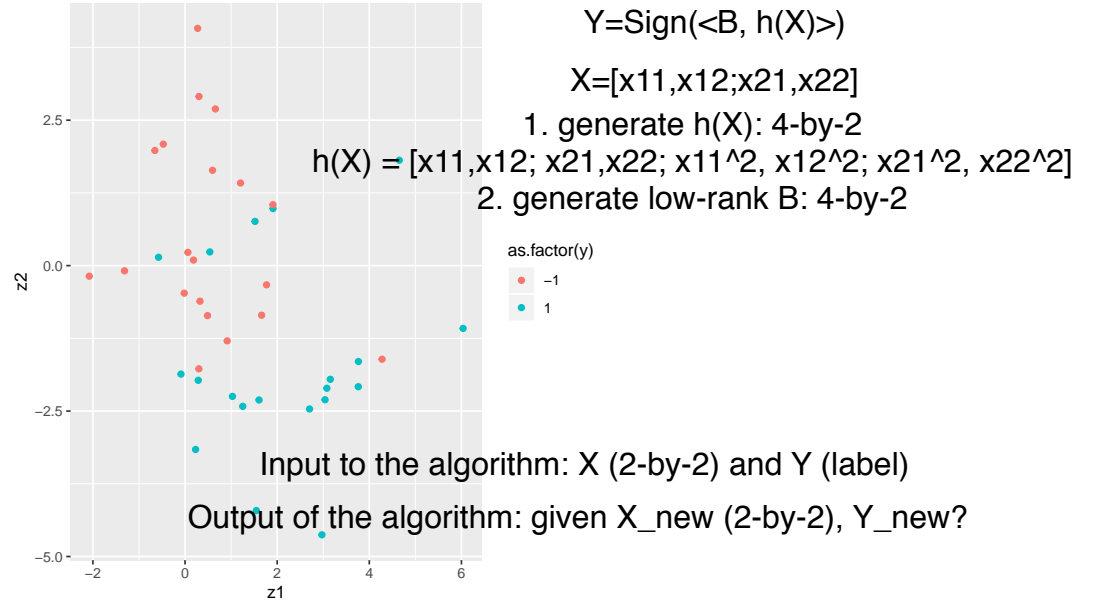


Figure 1: Visualization of the feature matrices. $Z1$ is the sum of the first column and $Z1$ is of the second one.

I implement SMMK method with linear kernel and polynomial kernel with degree three. I do not use exponential kernel I suggested above. The reason is that algorithm does not work well getting exploding objective values. I think that there is no feature h that corresponding to the exponential kernel. Figure 2 shows the classification results using the linear kernel and the polynomial kernel. One thing to note is that there are blurry part around the boundary in each figure. This is because for each point $(Z1, Z2)$, there are infinitely many possible combinations of $(X1, X2, X3, X4)$ which might give different value of y for each classifier.

3 Things to do

1. Construct algorithm that provides weighted loss function for conditional probability estimation.
2. Find good condition to make us easily check if given kernel is available i.e. has feature mapping such that $K(X, X') = h(X)^T h(X')$.

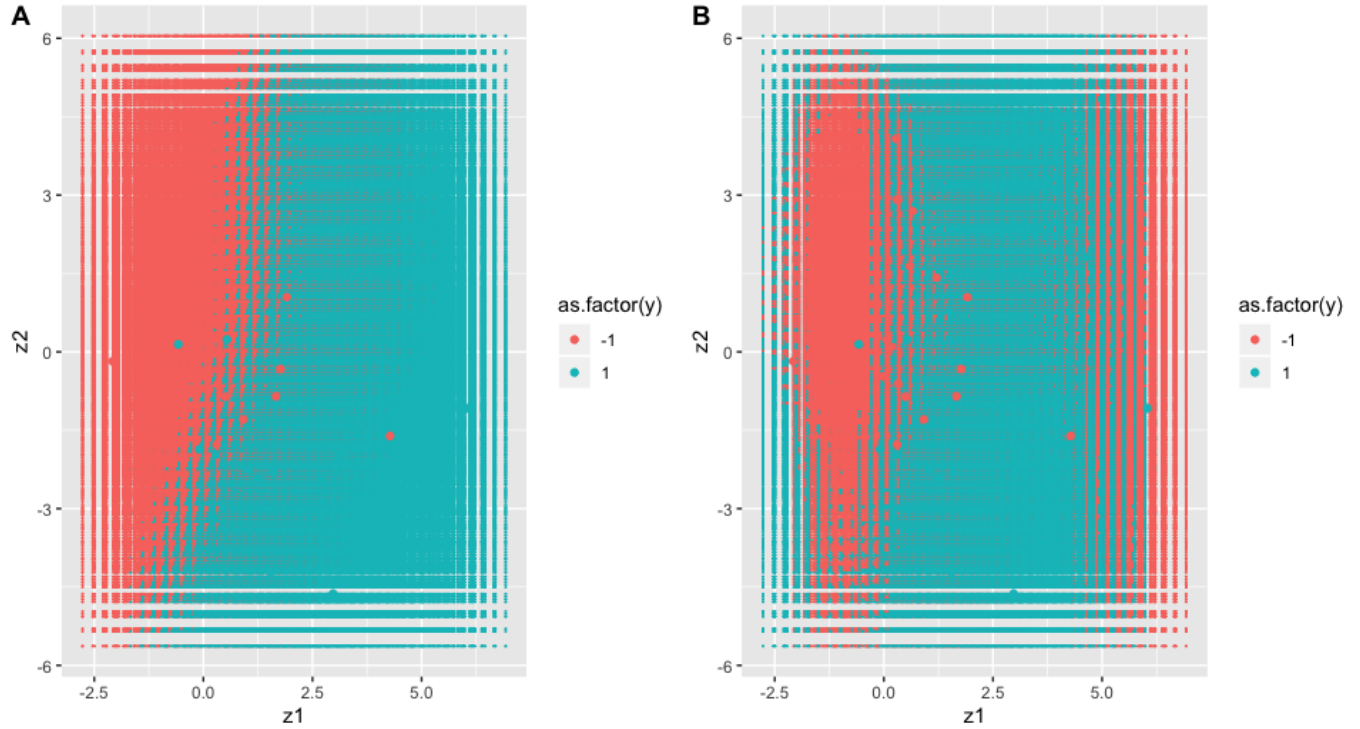


Figure 2: Figure A shows the classification rule of SMM with linear kernel and Figure B shows SMM with polynomial kernel.

4 New Algorithms

```

1 #SMMK
2
3 # Make sure Q matrix is positive definite
4 Makepositive = function(mat){
5   h = eigen(mat,symmetric = T)
6   nmat = (h$vectors)%*%diag(pmax(h$values,10^-4))%*%t(h$vectors)
7   return(nmat)
8 }
9
10 # save kernel values
11 Karray = function(X,kernel = function(X1,X2) t(X1)%*%X2){
12   m= nrow(X[[1]]); n = ncol(X[[1]]); N = length(X)
13   K = array(dim = c(N,N,n,n))
14   for(i in 1:N){
15     for(j in 1:N){
16       K[i,j, , ] = kernel(X[[i]],X[[j]])
17     }
18   }
19   return(K)
20 }
21
22
23 # objective value function
24 objm = function(X,y,alpha,V,b,K,cost = 10){
25   m= nrow(X[[1]]); n = ncol(X[[1]]); N = length(X)
26   Hv = V%*%solve(t(V)%*%V)%*%t(V)

```

```

27 Kv = matrix(nrow =N,ncol = N)
28 for(i in 1:N){
29   for(j in 1:N){
30     Kv[i,j] = sum(Hv*K[i,j,,])
31   }
32 }
33 coef = as.matrix(alpha*y)
34 obj = t(coef)%%Kv%%coef/2 + cost*sum(pmax(1-y*(Kv%%coef+b),0))
35 return(obj)
36 }
37
38
39
40 # Main function (not available for weighed loss yet)
41 SMM = function(X,y,r,kernel = function(X1,X2) t(X1)%%X2, cost = 10,rep = 1,p =
    .5){
42   result = list()
43   m= nrow(X[[1]]); n = ncol(X[[1]]); N = length(X)
44   K = Karray(X,kernel)
45
46   compareobj = 10^10
47   for(i in 1:rep){
48     error = 10; iter = 0; V = randortho(n)[,1:r,drop = F]
49     obj = compareobj
50
51
52     ### update U fixing V ###
53     vtv = solve(t(V)%%V)
54     Hv = V%%vtv%%t(V)
55     Kv = matrix(nrow =N,ncol = N)
56     for(i in 1:N){
57       for(j in 1:N){
58         Kv[i,j] = sum(Hv*K[i,j,,])
59       }
60     }
61     dvec = rep(1,length(X))
62     Dmat = Makepositive((y%%t(y))*Kv)
63     Amat = cbind(y,diag(1,N),-diag(1,N))
64     bvec = c(rep(0,1+N),ifelse(y==1,-cost*(1-p),-cost*p))
65     alpha = solve.QP(Dmat,dvec,Amat,bvec,meq =1)$solution
66
67
68     while((iter < 20)&(error >10^-3)){
69
70
71       ### update V fixing U ###
72
73       # sum_{i,j} alpha_i alpha_j y_i y_j K(X_i,X_j)
74       aayyK = 0
75       for(i in 1:N){
76         for(j in 1:N){
77           aayyK = aayyK + alpha[i]*alpha[j]*y[i]*y[j]*K[i,j,,]
78         }
79       }
80
81       # sum_j alpha_i y_i K(X_i,X_j)
82       ayK = array(0,dim = c(N,n,n))
83       for(i in 1:N){
84         for(j in 1:N){

```

```

85     ayK[i,,] = ayK[i,,] + alpha[j]*y[j]*K[j,i,,]
86   }
87 }
88
89 # (U^TU)^{-1}
90 utui = t(V)%%V%%solve(t(V)%%aayK%%V)%%t(V)%%V
91
92 # U^th(X_i)
93 uth = array(dim = c(N,r,n))
94 for(i in 1:N){
95   uth[i,,] = vtv_i%%t(V)%%ayK[i,,]
96 }
97
98 # Ku[i,j] = tr(H_uh(X_i),H_uh(X_j))
99 Ku = matrix(nrow = N,ncol = N)
100 for(i in 1:N){
101   for(j in 1:N){
102     Ku[i,j] = sum(uth[i,,]*(utui%%uth[j,,]))
103   }
104 }
105 dvec = rep(1,length(X))
106 Dmat = Makepositive((y%%t(y))*Ku)
107 Amat = cbind(y,diag(1,N),-diag(1,N))
108 bvec = c(rep(0,1+N),ifelse(y==1,-cost*(1-p),-cost*p))
109 beta = solve.QP(Dmat,dvec,Amat,bvec,meq =1)$solution
110
111 # update V
112 V = 0
113 for(i in 1:N){
114   V = V+as.matrix(uth[i,,],nrow = r)*beta[i]*y[i]
115 }
116 V = V%%utui
117
118
119 ### update U fixing V ###
120 vtv_i = solve(t(V)%%V)
121 Hv = V%%vtv_i%%t(V)
122 Kv = matrix(nrow =N,ncol = N)
123 for(i in 1:N){
124   for(j in 1:N){
125     Kv[i,j] = sum(Hv*K[i,j,,])
126   }
127 }
128 dvec = rep(1,length(X))
129 Dmat = Makepositive((y%%t(y))*Kv)
130 Amat = cbind(y,diag(1,N),-diag(1,N))
131 bvec = c(rep(0,1+N),ifelse(y==1,-cost*(1-p),-cost*p))
132 alpha = solve.QP(Dmat,dvec,Amat,bvec,meq =1)$solution
133
134
135
136 # slope part estimation
137
138 slope = function(nX){
139   coef <- as.matrix(alpha*y)
140   sp <- t(coef)%%unlist(lapply(X,function(x) sum(Hv*kernel(nX,x))))
141   return(sp)
142 }
143

```

```

144     # intercept part estimation (update b)
145     positiv = min(unlist(lapply(X,slope))[which(y==1)])
146     negativ = max(unlist(lapply(X,slope))[which(y==-1)])
147     if ((1-positiv)<(-1-negativ)) {
148         b0hat = -(positiv+negativ)/2
149     }else{
150         gridb0 = seq(from = -1-negativ,to = 1-positiv,length = 100)
151         b0hat = gridb0[which.min(sapply(gridb0,function(b) objm(X,y,alpha,V,b,K,
cost)))]
152     }
153     obj = c(obj,objm(X,y,alpha,V,b0hat,K,cost));obj
154     iter = iter+1
155     error = abs(-obj[iter+1]+obj[iter])/obj[iter];error
156 }
157 if (compareobj>obj[iter+1]) {
158     compareobj = obj[iter+1]
159     predictor = function(nX) sign(slope(nX)+b0hat)
160     result$slope = slope; result$b0 = b0hat; result$obj = obj[-1]; result$iter =
iter
161     result$error = error; result$predict = predictor; result$V = V
162 }
163
164
165 }
166 return(result)
167 }
168
169
170 ## Some kernels (Expkernel does not work)
171 Expkernel = function(Y,Z){
172     return(exp(-t(Y-Z)%*(Y-Z)))
173 }
174
175 polykernel = function(Y,Z,deg = 3){
176     n = ncol(Y)
177     return((t(Y)%*Z+diag(1,n))^deg)
178 }

```