

SMM conditional probability and kernel

Chanwoo Lee, April 09, 2020

1 Strong duality

Our primal problem for SMM is

$$\begin{aligned} \min_{U, V, \xi} & \frac{1}{2} \|UV^T\|^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to} & y_i (\langle UV^T, X_i \rangle + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, N. \end{aligned} \quad (1)$$

If we consider U and V at the same time, there is no guarantee that the strong duality holds. However, if we consider U fixing V and vice versa, we have the strong duality for each sub iteration. Suppose that we fix V to update U , then (5) becomes

$$\begin{aligned} \min_U & \frac{1}{2} \text{Vec}(V)^T \text{diag}(\underbrace{V^T V, \dots, V^T V}_{\text{total m}}) \text{Vec}(V) + C \xi^T \mathbf{1} \\ \text{subject to} & \xi + Z^T \text{Vec}(U) + b\mathbf{y} \geq \mathbf{1}, \\ & \xi \geq 0, \\ \text{where } Z &= (\text{Vec}(y_1 X_1 V), \dots, \text{Vec}(y_N X_N V))^T \end{aligned} \quad (2)$$

Since Equation (2) is a Quadratic programming problem with linear constraints, we have strong duality for the alternating update algorithm.

2 SVM conditional probability calculation

In simulation, I made inseparable data set generating $(x_1, y_1), \dots, (x_{30}, y_{30})$ where $x_i \in \mathbb{R}^2$ and $y \in \{-1, 1\}$. I randomly assigned y as 1 or -1 and add $(1.5, 1.5)^T$ to feature vectors x if $y = 1$. Figure 1 shows the generated data. First, I calculated conditional probability $\mathbb{P}(y = 1|x)$ with the linear kernel. Figure 3 shows that the weird phenomenon we talked before still happens in the right bottom part. In nonlinear SVM case with radial kernel, classification boundary changes as in Figure 3 according to π values. Based on the weighted hinge loss classification, I calculate the conditional probability and Figure 4 shows the result.

3 SMM conditional probability calculation

To calculate conditional probability, we solve the regularization problem based on weighted hinge loss.

$$\begin{aligned} \min_{U, V, \xi} & \frac{1}{2} \|UV^T\|^2 + C \left[(1 - \pi) \sum_{y_i=1} \xi_i + \pi \sum_{y_i=-1} \xi_i \right] \\ \text{subject to} & y_i (\langle UV^T, X_i \rangle + b) \geq 1 - \xi_i \end{aligned} \quad (3)$$

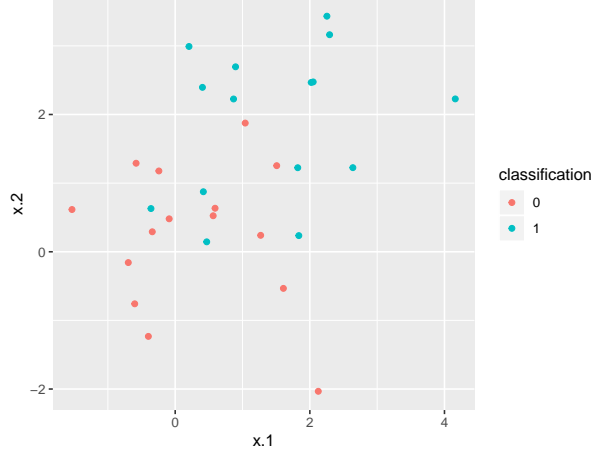


Figure 1: Simulation data points. Red colored points has $y = 1$ and blue colored points has $y = -1$.

$$\xi_i \geq 0, \quad i = 1, \dots, N.$$

To solve this problem, we take the alternating update approach. We update U fixing V with the following dual problem.

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

From this (3), we have updated $U = \sum_{i=1}^N \alpha_i y_i X_i V (V^T V)^{-1}$. We can update V fixing U by the same way. If we have theoretical guarantee that

$$\text{sign}(\langle X, UV^T \rangle + b) \rightarrow \text{sign}(\mathbb{P}(y = 1|X) - \pi),$$

we can obtain the conditional probability. My guess is that if the targeted rank is large enough, we have the result. To build this guess, consider the regularization problem of minimizing

$$\frac{1}{N} \sum_{i=1}^N L(y_i)[1 - y_i f(X_i)]_+ + \frac{1}{C} \|h\|_{H_k}^2, \quad (4)$$

over all the functions of the form $f(X) = h(X) + b$, with $h \in H_k$. In our situation, $L(1) = 1 - \pi$ and $L(-1) = \pi$. There are two components: a data fit functional component and a regularization penalty component. The data fit component usually approaches a limiting functional as $N \rightarrow \infty$. In general, the target function is the minimizer of the limiting functional. Under the assumption that the target function can be approximated by the elements in the Reproducing Kernel Hilbert Space (RKHS) with other general regularity conditions, the solution of (4) approaches the target function as $N \rightarrow \infty$ [Lin, 2002]. The following lemma suggests that the solution of (4) approaches $\text{sign}[\mathbb{P}(Y = 1|X) - \frac{L(Y=-1)}{L(-1)+L(1)}]$ [Lin et al., 2002]. In this sense, as long as considered reproducing

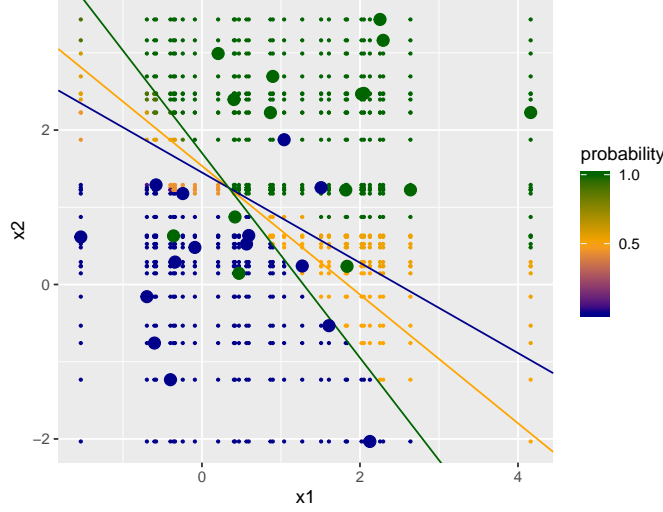


Figure 2: Green large points are from given data with $y = 1$. Blue large points are from $y = -1$. The yellow line is classification boundary when $\pi = 0.5$. The blue line and green line are when $\pi = 0.01$ and $\pi = 0.99$ respectively. Small dots are calculated conditional probability from the grid.

kernel Hilbert space is rich enough to approximate the target function, we have consistent estimator. This explains better performance in estimating the conditional probability in nonlinear case (you can compare Figure and Figure).

Lemma 1. *The minimizer of $E[L(Y)(1 - Yf(X))_+]$ is $\text{sign}[\mathbb{P}(Y = 1|X) - \frac{L(Y=-1)}{L(-1)+L(1)}]$.*

4 Kernel SMM method

The suggested kernel method is to use $K(X, X') = g(\langle H_U X H_V, H_U X' H_V \rangle)$. In SVM case, we are looking for good hyper plane that separate the feature well including the kernel method (in kernel, hyper plane is in the extended dimension space). To make it clear, let us consider kernel method in SVM. For nonlinear SVM, we set the basis function $h : \mathbb{R}^{\dim(\mathbf{x})} \rightarrow \mathbb{R}^d$ where $d > \dim(\mathbf{x})$ is the usual case. Once the basis function are selected, the procedure is the same as before. We fit the SV classifier using input feature $h(\mathbf{x}_i) = (h_1(\mathbf{x}_i), \dots, h_d(\mathbf{x}_i))$, $i = 1, \dots, N$ and the produce the nonlinear function $\hat{f}(\mathbf{x}) = \langle h(\mathbf{x}), \hat{\mathbf{w}} \rangle + \hat{b} = \sum_{i=1}^N \hat{\alpha}_i y_i K(\mathbf{x}, \mathbf{x}_i) + \hat{b}$, where $K(\mathbf{x}, \mathbf{x}') = \langle h(\mathbf{x}), h(\mathbf{x}') \rangle$. Therefore, even in nonlinear case, we are finding the hyper plane in extended dimension. However, if we use $K(X, X') = g(\langle H_U X H_V, H_U X' H_V \rangle)$, we are obtaining classification function as $f(X) = \sum_{i=1}^N \alpha_i y_i g(\langle H_U X_i H_V, H_U X H_V \rangle) + b$. From the SVM case, we can interpret

$$g(\langle H_U X H_V, H_U X' H_V \rangle) = \langle H_{U'} h(X) H_{V'}, H_{U'} h(X') H_{V'} \rangle.$$

This interpretation is possible because as long as gram matrix K is positive semi definite and symmetric, we can find induced U' , V' , and $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m' \times n'}$ from K (need to verify). From this h , our primal SMM problem becomes,

$$\begin{aligned} \min_{U', V', \xi} \quad & \frac{1}{2} \|U'(V')^T\|^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i (\langle h(X_i), U'(V')^T \rangle + b) \geq 1 - \xi_i \end{aligned} \tag{5}$$

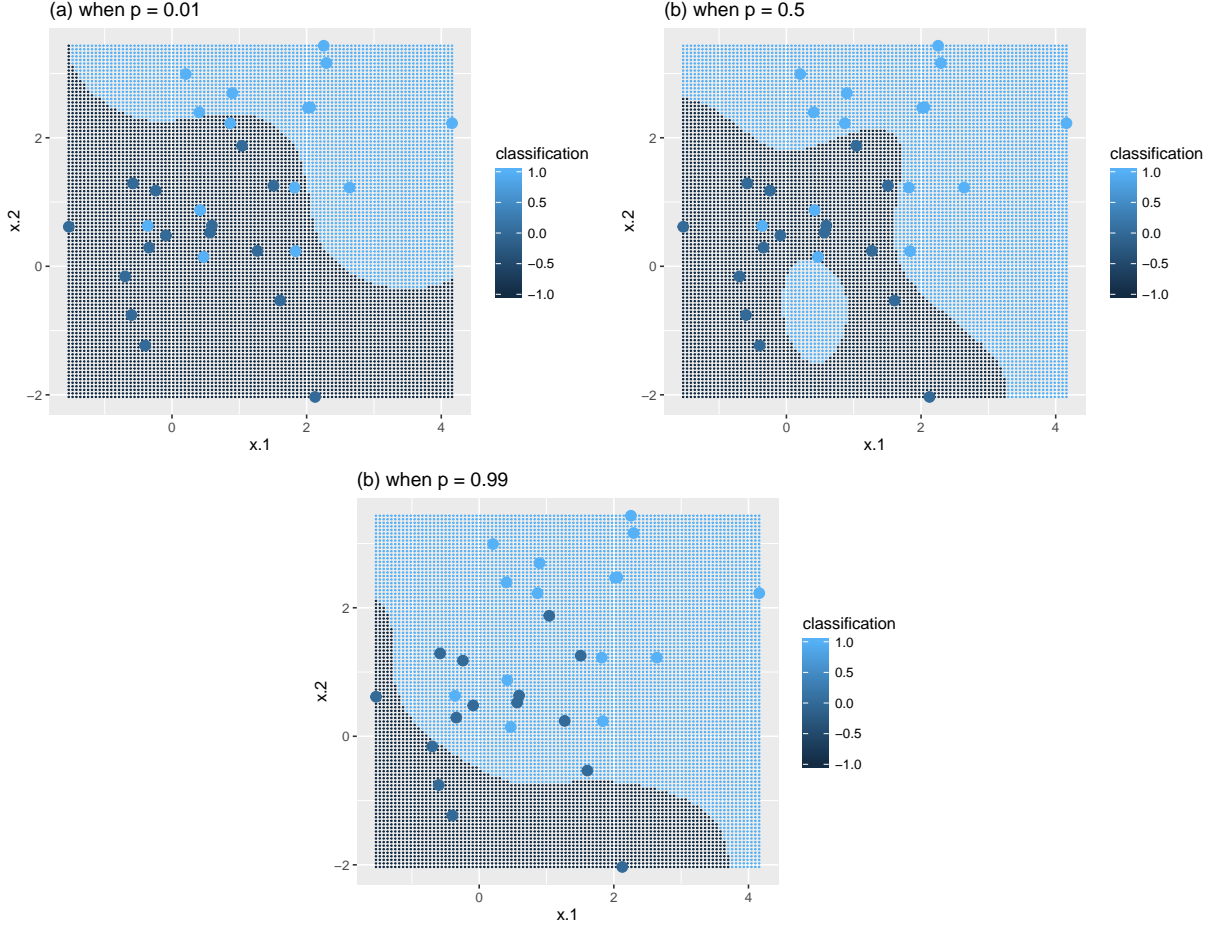


Figure 3: The figures show weighted hinge loss SVM classifier with radial kernel when $\pi \in \{0.01, 0.5, 0.99\}$. Black dots are classified as -1 and blue dots as 1. Sub figure (b) is the regular SVM with radial kernel.

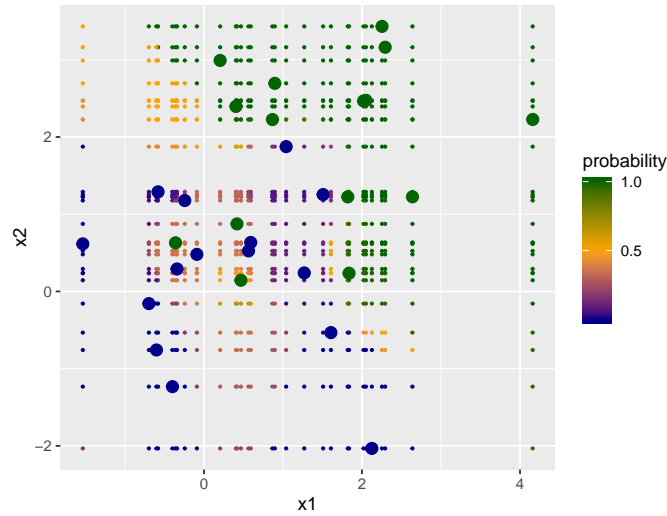


Figure 4: Green large points are from given data with $y = 1$. Blue large points are from $y = -1$. Small dots are calculated conditional probability from the grid.

$$\xi_i \geq 0, \quad i = 1, \dots, N.$$

If we have the relationship

$$g(\langle H_U X, H_U X \rangle) = \langle H_U h(X), H_U h(X') \rangle \quad \text{and} \quad (\langle X H_V, X H_V \rangle) = \langle h(X) H_{V'}, h(X') H_{V'} \rangle.$$

we can make use of alternating algorithm we do on linear SMM case.

5 Updated rcodes

I updated svm function to be available for arbitrary kernel methods.

```

1 # SVM with kernel functions and weighted cost function
2 svm = function(X,y,cost = 10, kernels = function(x1,x2) sum(x1*x2), p = .5){
3   if (p==.5) {
4     cost = 2*cost
5   }
6   result = list()
7   error = 10
8   iter = 0
9   # SVM
10  m= nrow(X[[1]]); n = ncol(X[[1]]); N = length(X)
11
12  x = matrix(unlist(X),nrow = N,byrow = T)
13  dvec = rep(1,length(X))
14  Dmat = kernelmat(x,y,kernels)
15  Amat = cbind(y,diag(1,N),-diag(1,N))
16  bvec = c(rep(0,1+N),ifelse(y==1,-cost*(1-p),-cost*p))
17  alpha = solve.QP(Dmat,dvec,Amat,bvec,meq =1)
18  coef = y*alpha$solution
19
20  Bhat=matrix(t(coef)%*%x,nrow = m)
21  # b0hat = -(min(unlist(lapply(X,function(x) sum(Bhat*x)))[which(y==1)]))+
22  #          max(unlist(lapply(X,function(x) sum(Bhat*x)))[which(y==1)]))/2
23  b0hat = -(min((Dmat%*%alpha$solution)[which(y==1)]*y[which(y==1)])+
24            max((Dmat%*%alpha$solution)[which(y==1)]*y[which(y==1)]))/2
25  obj = objv(Bhat,b0hat,X,y,cost,prob = p)
26
27  predictor = function(y){
28    return(sign(sum(coef*apply(x,1, function(x) kernels(x,y)))+b0hat))
29  }
30  # predictor = function(x) sign(sum(Bhat*x)+b0hat)
31  result$B = Bhat; result$b0 = b0hat; result$obj = obj;
32  result$predict = predictor
33  return(result)
34 }
```

References

- Yi Lin. Support vector machines and the bayes rule in classification. *Data Mining and Knowledge Discovery*, 6:259–275, 2002.
- Yi Lin, Yoonkyung Lee, and Grace Wahba. Support vector machines for classification in nonstandard situations. *Machine learning*, 46(1-3):191–202, 2002.