

# SMMK conditional probability and SDR

Chanwoo Lee, May 21, 2020

## 1 SMMK conditional probability

Main changes of SMMK algorithm are as follow.

1. Weighted hinge loss based algorithm is available for conditional probability estimation.
2. Symmetric adjustment procedure is added in nonlinear kernel cases. But nothing is changed in linear case.

I check that the new SMMK algorithm gives us reasonable conditional probability output from simple data simulation. Training data  $\{(X_i, y_i)\}_{i=1}^{40}$  is generated with the following rule.  $X_i \in \mathbb{R}^2$  is from i.i.d. multivariate normal distribution  $N_2\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}\right)$  and  $y_i = \text{sign}\left(\phi(X_i) - \frac{1}{40} \sum_{i=1}^{40} \phi(X_i)\right)$  where  $\phi(\cdot)$  is density function of  $X_i$ . Then, true boundary is an ellipse in the plane. First, I fit the each kernel types (“linear”, “polynomial”, “exponential”) to training data. The following figure shows the boundary of classification.

good choice of simulation.

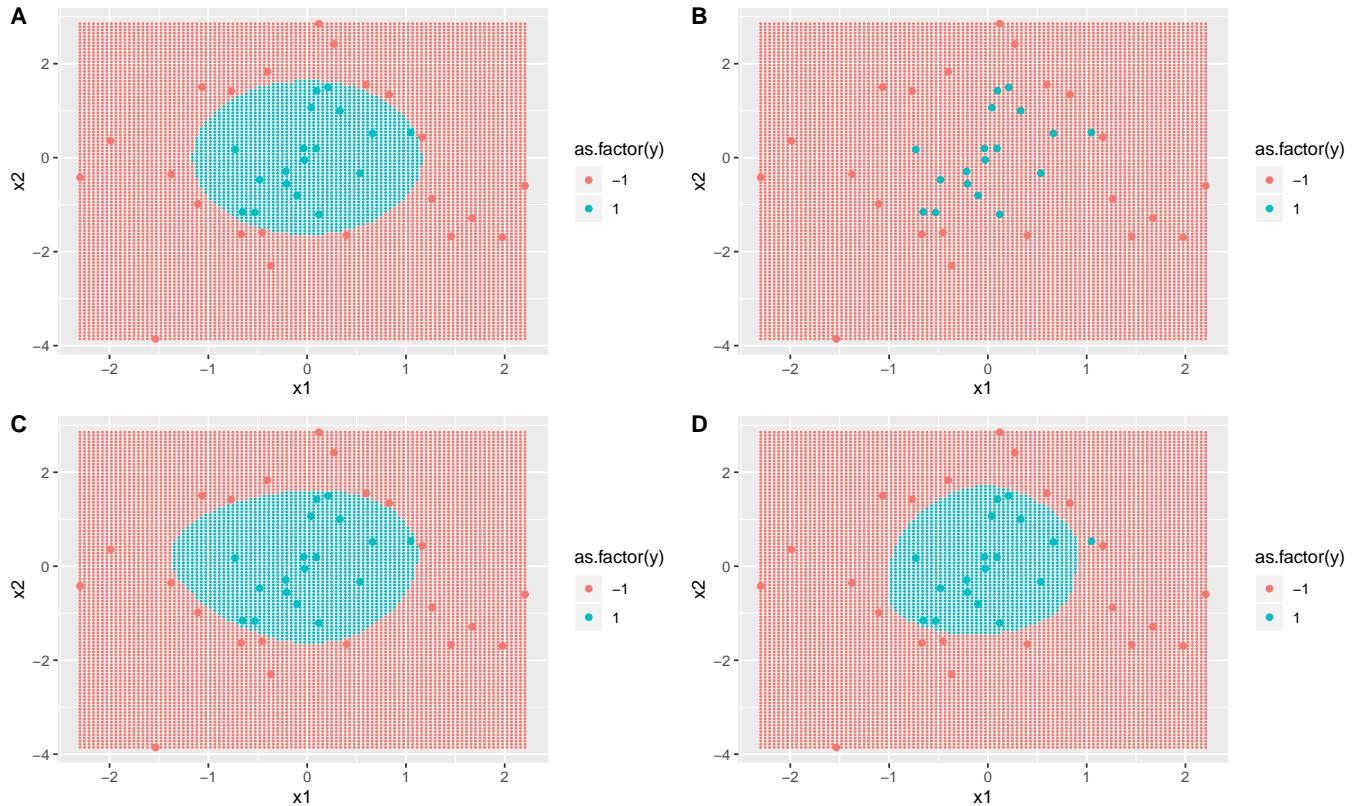


Figure 1: Subfigure A is true ellipsoid boundary. B is linear case boundary which assigns labels all 0. C and D show the boundary of polynomial and exponential kernel respectively.

Since there is no meaningful classification in linear case, I estimate conditional probability only in the case of polynomial and exponential kernel. Th2 shows the result of conditional probability

estimation which looks reasonable.

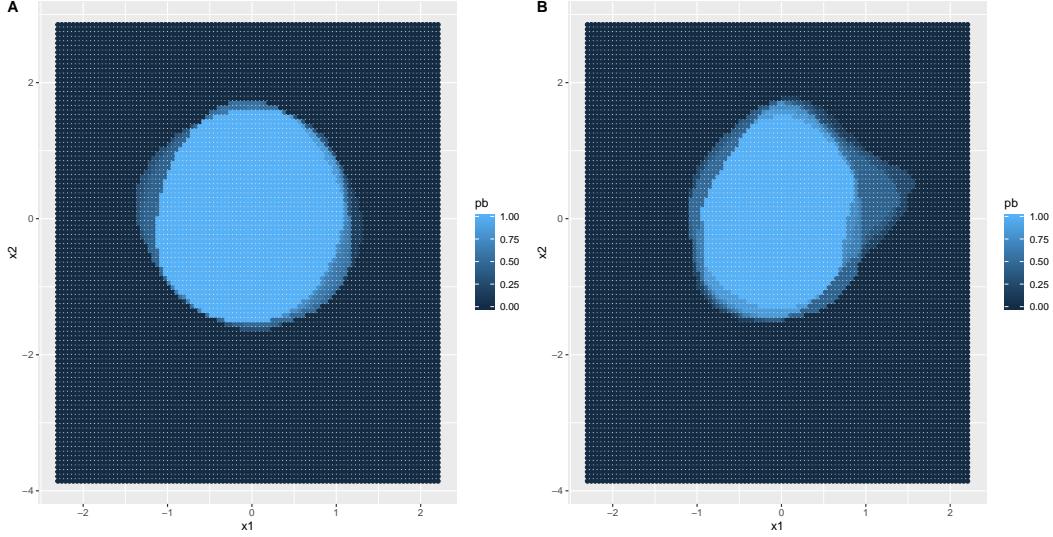


Figure 2: Figure A and B show the estimated probability with polynomial and exponential kernel respectively.

## 2 SDR summary and questions

### 2.1 SDR linear case

Sufficient dimension reduction here assumes that

$$Y \perp\!\!\!\perp X | B^T X.$$

We define the central subspace  $S_{Y|X}$  as

$$S_{Y|X} = \bigcap_{\{B: Y \perp\!\!\!\perp X | B^T X\}} \text{span}(B).$$

Our goal is to estimate basis of  $S_{Y|X}$  using  $S_{Y|X} = S_{P(X|Y=1)|X}$ .

#### 2.1.1 Population level

For a pair of random variables  $(X, Y) \in \mathbb{R}^p \times \{-1, +1\}$ , the linear principal weighted support vector machine minimizes

$$\Lambda_\pi(\boldsymbol{\beta}, \alpha) = \boldsymbol{\beta}^T \Sigma \boldsymbol{\beta} + \lambda \mathbb{E} [w_\pi(Y)(1 - Y f(X; \theta))_+], \quad (1)$$

where  $\Sigma = \text{cov}(X)$ , and  $f(X; \boldsymbol{\beta}, \alpha) = \alpha + \langle \boldsymbol{\beta}, X - \mathbb{E}(X) \rangle$ . Note (1) is reduced to SVM with weighted hinge loss with the population version when  $\mathbb{E}(X) = 0$  and  $\text{cov}(X) = I_p$ . Fisher consistency of the weighted support vector machine ensures that a hyperplane  $\{X : f(X; \boldsymbol{\beta}_{0,\pi}, \alpha_{0,\pi}) = 0\}$  optimally separates  $S_\pi^+ = \{X : p(X) \geq \pi\}$  and  $S_\pi^- = \{X : p(X) < \pi\}$  where  $(\boldsymbol{\beta}_{0,\pi}, \alpha_{0,\pi}) = \arg \min \Lambda_\pi(\boldsymbol{\beta}, \alpha)$ .

**Theorem 2.1.** Assume that  $\mathbb{E}(X|B^T X)$  is a linear function of  $B^T X$ . Then for any given weighted  $\pi \in (0, 1)$ ,  $\beta_{0,\pi} \in S_{Y|X}$ .

The assumption in Theorem 2.1 is known as the linearity condition. It implies that  $\mathbb{E}(\beta^T X|B^T X) = \beta^T P_B(\Sigma)X$  where  $P_B(\Sigma) = B(B^T \Sigma B)^{-1}B^T \Sigma$ . The condition holds when  $X$  is elliptically symmetric and approximately holds when  $p$  is large. We can assume that  $\text{span}(\beta_{0,1}, \dots, \beta_{0,H}) = S_{Y|X}$  whenever  $H$  is large enough.

### 2.1.2 Finite sample estimation

For the finite sample case,  $\Lambda_\pi$  in (1) changes to,

$$\hat{\Lambda}_{n,\pi}(\beta, \alpha) = \beta^T \hat{\Sigma}_n \beta + \frac{\lambda}{n} \sum_{i=1}^n w_\pi(y_i)(1 - y_i \hat{f}_n(X_i; \beta, \alpha))_+,$$

where  $\hat{f}_n(X_i; \beta, \alpha) = \alpha + \langle \beta, X_i - \bar{X}_n \rangle$ ,  $\bar{X}_n$  is the sample mean, and  $\hat{\Sigma}_n$  denotes the sample covariance matrix. Let  $(\hat{\beta}_{n,\pi}, \hat{\alpha}_{n,\pi}) = \arg \min_{\beta, \alpha} \Lambda_\pi(\beta, \alpha)$ . The candidate matrix of the linear principal weighted support vector machine is

$$\hat{M}_n = \sum_{i=1}^H \hat{\beta}_{n,h} \hat{\beta}_{n,h}^T.$$

The first  $k$  eigenvectors of  $\hat{M}_n$ , denoted by  $\hat{V}_n = (\hat{v}_1, \dots, \hat{v}_k)$ , estimate a basis of  $S_{Y|X}$ . To determine  $k$ , we consider

$$G_n(m; \rho, \hat{M}_n) = \sum_{j=1}^m \ell_j - \rho \frac{m \log n}{\sqrt{n}} \ell_1,$$

where  $\ell_j$  is the  $j$ -th leading eigenvalues of  $\hat{M}_n$  and  $\rho$  is a tunning parameter. It is known that  $\hat{k} = \arg \max_m G_n(m; \rho, \hat{M}_n)$  is a consistent estimator of  $k$ .

## 2.2 SDR nonlinear case

Sufficient dimension reduction here assumes that

$$Y \perp\!\!\!\perp X | \phi(X), \quad (2)$$

where  $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^k$  is an unknown vector valued function of  $X$ . We define unbiasedness notion in nonlinear case for the later use.

**Definition 1.** A function  $\psi \in \mathcal{H}$  is unbiased for nonlinear sufficient dimension reduction (2) if it has a version that is measurable  $\sigma\{\phi(X)\}$

### 2.2.1 Population level

The principal weighted support vector machine minimizes

$$\Lambda_\pi(\psi, \alpha) = \text{var}(\psi(X)) + \lambda \mathbb{E} [\omega_\pi(Y)(1 - Y f(X; \psi, \alpha))_+], \quad (3)$$

where  $f(X; \psi, \alpha) = \alpha + (\psi(X) - \mathbb{E}(\psi(X)))$ . We can equivalently express (3) as

$$\Lambda_\pi(\psi, \alpha) = \langle \psi, \Sigma \psi \rangle_{\mathcal{H}} + \lambda \mathbb{E} [\omega_\pi(Y)(1 - Y f(X; \psi, \alpha))_+], \quad (4)$$

where  $\Sigma : \mathcal{H} \rightarrow \mathcal{H}$  is an operator such that  $\langle f_1, \Sigma f_2 \rangle = \text{cov}[f_1(X), f_2(X)]$

**Theorem 2.2.** Suppose the mapping

$$\mathcal{H} \rightarrow L_2(P_X), \quad f \mapsto f$$

is continuous and :

1.  $\mathcal{H}$  is a dense subset of  $L_2(P_X)$ ,
2.  $Y \perp\!\!\!\perp X | \phi(X)$ .

If  $(\psi^*, \alpha^*)$  minimizes (4) among all  $(\psi, \alpha) \in \mathcal{H} \times \mathbb{R}$ , then  $\psi^*(X)$  is unbiased.

### 2.2.2 Finite sample estimation

We define  $K_n(\cdot) = (K(\cdot, X_1), \dots, K(\cdot, X_n))^T$  from a given kernel  $K$  and feature predictors  $\{X_i\}_{i=1}^n$ . Considering Hilbert space  $\mathcal{H} = \{\beta^T K_n(\cdot) = \sum_{i=1}^n \beta_i K(\cdot, X_i) : \beta \in \mathbb{R}^n\}$  is too rich, so that the solution often overfits the data. So we consider new Hilbert space. Define basis of new  $\mathcal{H}$  as

High-level explanation:

$K(X)$ : function space

$w_j$ : canonical space

$$\psi_j(X) = \tilde{K}_n(X)^T \omega_j / \lambda_j$$

formula in blue  $\rightarrow$  operated in matrix space  $\rightarrow$  canonical basis in  $\mathbb{R}^n$  is of interest

formula in yellow  $\rightarrow$  operated in function space  $\rightarrow$  function basis in  $\mathcal{H}$  is of interest

green  $\rightarrow$  convert canonical basis to function basis for the purpose of interpretation.

$$\text{change of basis: from } \mathbb{R}^n \text{ (canonical basis) to } \mathcal{H} \text{ (function basis)} \\ \text{where } \tilde{K}_n(X) = K_n(X) - \frac{1}{n} \sum_{i=1}^n K_n(X_i)$$

$$\omega_j, \lambda_j \text{ are j-th components s.t. } (\mathbf{I}_n - \mathbf{J}_n/n) K_n (\mathbf{I}_n - \mathbf{J}_n/n) \omega_j = \lambda_j \omega_j.$$

The algorithm relies on only this matrix and its top k eigen-vectors.

Finally, we have new Hilbert space as  $\mathcal{H} = \{\sum_{i=1}^k \beta_i \psi_i(\cdot) : \beta \in \mathbb{R}^k\}$ . Consider finite sample case of (4).

$$\Lambda_\pi(\beta, \alpha) = \beta^T \Psi^T \Psi \beta + \frac{\lambda}{n} \sum_{i=1}^n \omega_\pi(y_i) (1 - y_i(\Psi_i^T \beta + \alpha))_+, \quad (6)$$

$$\text{where } \Psi = \begin{pmatrix} \psi_1(X_1) & \cdots & \psi_k(X_1) \\ \vdots & & \vdots \\ \psi_1(X_n) & \cdots & \psi_k(X_n) \end{pmatrix} \text{ and } \Psi_i^T = (\psi_1(X_i), \dots, \psi_k(X_i)).$$

= top k eigen-vectors of the above matrix in blue

Let  $(\hat{\beta}_{n,\pi}, \hat{\alpha}_{n,\pi}) = \arg \min_{\beta, \alpha} \Lambda_\pi(\beta, \alpha)$ , the minimizer of (6). It is shown that  $\hat{\beta}_{n,\pi} = \lambda \sum_{i=1}^n \hat{\gamma}_{i,\pi} y_i (\Psi^T \Psi)^{-1} \Psi^T / 2$  where  $\hat{\gamma}_\pi = (\hat{\gamma}_{1,\pi}, \dots, \hat{\gamma}_{n,\pi})^T$  solves

$$\max_{\gamma} \sum_{i=1}^n \gamma_i - \frac{1}{4} \sum_{i=1}^n \sum_{j=1}^n \gamma_i \gamma_j y_i y_j [P_\Psi]_{i,j}$$

subject to  $0 \leq \gamma_i \leq \lambda \omega_\pi(y_i)$ ,  $\sum_{i=1}^n \gamma_i y_i = 0$ ,

with  $[P_\Psi]_{i,j}$  the  $(i,j)$ th element of  $P_\Psi = \Psi (\Psi^T \Psi)^{-1} \Psi^T$ . By the similar way in linear SDR case, The candidate matrix of the non linear principal weighted support vector machine is

$$\hat{M}_n = \sum_{i=1}^H \hat{\beta}_{n,h} \hat{\beta}_{n,h}^T.$$

Denote first  $k$  eigenvectors of  $\hat{M}_n$ , by  $\hat{V}_n = (\hat{v}_1, \dots, \hat{v}_k)$ . Then, the  $s$ th sufficient predictor evaluated at  $X$  is  $\sum_{i=1}^n v_{si} \psi_i(X)$ .

## 2.3 Question

There are some questions I had while summarizing the papers. I will think about those questions until tomorrow meeting.

1. In linear case, we define the central subspace as

Good question.  
In general, not necessarily true. Under some conditions, yes.  
See Section 3.1 <http://users.stat.umn.edu/~rdcook/SDR/ASA94.pdf>

$$S_{Y|X} = \bigcap_{\{B: Y \perp\!\!\!\perp X | B^T X\}} \text{span}(B).$$

If we find the basis  $\{v_1, \dots, v_k\}$  such that  $\text{span}(\{v_1, \dots, v_k\}) = S_{Y|X}$ , is it true that  $Y \perp\!\!\!\perp X | (v_1, \dots, v_k)^T X$ ?

2. In linear case, obtained  $\hat{\beta}_{n,\pi}$  is normal vector of hyperplane that optimally separate  $S_1 = \{X_i : P(X_i|y_i = 1) > \pi\}$  and  $S_2 = \{X_i : P(X_i|y_i = 1) < \pi\}$ . Intuitively, the normals of these hyperplanes are roughly aligned with the directions that forms the central subspace. We use the principal components of these  $\hat{\beta}_{n,\pi}$ s to estimate the central subspace. However, I cannot get the intuition that how eigen-vectors of  $\hat{M}_n = \sum_{i=1}^H \hat{\beta}_{n,h} \hat{\beta}_{n,h}^T$  give us estimation of the central subspace.
  3. I understand overall procedures of nonlinear SDR. However, I actually, do not understand how we define new Hilbert space as in Equation (5) and under what geometric intuition.
2. In principle,  $\text{Span}\{\beta_1, \dots, \beta_H\}$  is already a good candidate for central subspace.  
The eigen-vector amounts to normalize the collection of beta's.

Suppose we sample two vectors,  $b_1, b_2$ , from a linear space  $S$  of interest.

Q: how to estimate  $S$  from  $\{b_1, b_2\}$ ?

A1: Estimate  $S$  using  $\text{Span}\{b_1, b_2\}$ . The pitfall is that  $b_1$  and  $b_2$  may not be orthonormal.

A2: Estimate  $S$  using normalized column space of  $B = [b_1, b_2]$

= left singular space of  $B$   
= eigen-space of  $B^* B^T$   
= eigen-space of  $(b_1 \ b_1^T + b_2 \ b_2^T)$