

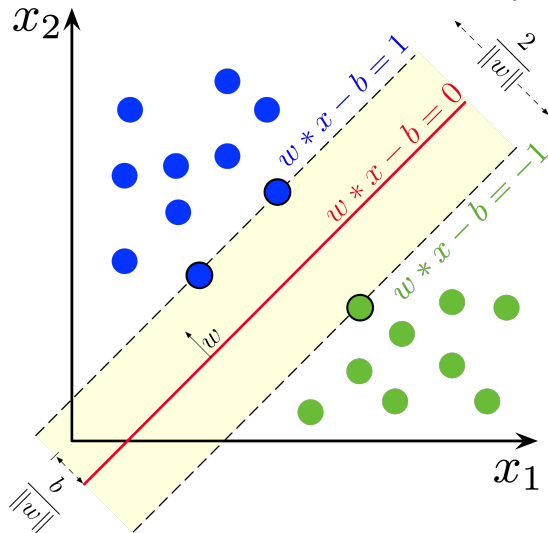
# Nonparametric learning with matrix-valued predictors in high dimensions

Chanwoo Lee, Miaoyan Wang

Summary of research work

May 26, 2020

## A successful story: Support vector machine (SVM)



Maximum-margin hyperplane for an SVM trained with 2-d vector predictors  
(Picture source: Wiki)

# SVM methods are powerful, however...

- How to **efficiently classify high-dimensional matrices** with limited sample size:

$$n \ll d_1 d_2 = \text{dimension of feature space} ?$$

- How to **robustly predict the label probability** when little is known about the function:

$$\mathbb{P}(\text{Label} = 1 | \text{Matrix}) \stackrel{\text{def}}{=} g(\text{Matrix}) ?$$

- How to **effectively reduce the sufficient dimensions** of the feature space without information lost:

$$\text{Label} \perp\!\!\!\perp \text{Matrix} | \phi(\text{Matrix}) ?$$

- For the purpose of presentation, we focus on **matrix** predictors and **binary** outcomes.

# Previous methods for nonparametric multivariate learning

Method	SDR*	Robust prediction	Binary outcome	Tensor predictor
<b>Weighted large-margin</b> (Wang et al, JCGS'19) (Zhang et al, JASA'10)	X	✓	✓	X
<b>Deep Neural Network</b> (...)	X	✓	✓	X
<b>Principal SVM</b> (Li et al, AOS'11) (Shin et al, Biometrika'17)	✓	X	X	X
<b>Nonparametric regression</b> (Imaizumi et al ICML'16)	X	✓	X	✓
<b>Our method</b>	✓	✓	✓	✓

\*SDR: sufficient dimension reduction. See later slides for definition.

# Overview

## 1 Motivation and Problem

## 2 Linear Learning with low-rank kernels

- Classification
- Probability function estimation
- Sufficient dimension reduction
- Simulations

## 3 Nonlinear learning with low-rank kernels

- Nonlinear matrix-valued kernels
- Nonlinear learning with matrix predictors
- Simulations

# Problem

Let  $\{(\mathbf{X}_i, y_i) \in \mathbb{R}^{d_1 \times d_2} \times \{-1, 1\} : i = 1, \dots, n\}$  denote an i.i.d. sample from distribution  $\mathcal{X} \times \mathcal{Y}$ .

- (Probability estimation.) Estimate the conditional probability function

$$\mathbb{P}(y^{\text{new}} = 1 | \mathbf{X}^{\text{new}}) = g(\mathbf{X}^{\text{new}}),$$

where  $g: \mathbb{R}^{d_1 \times d_2} \mapsto [0, 1]$  is the regression function of interest.

- (Sufficient dimension reduction, SDR.) Find a transformed predictor  $\phi(\mathbf{X})$  such that

$$y \perp\!\!\!\perp \mathbf{X} | \phi(\mathbf{X}),$$

where  $\phi: \mathbb{R}^{d_1 \times d_2} \mapsto \mathbb{R}^{r_1 \times r_2}$  is the dimension reduction of interest.

- Probability estimation concerns **conditional mean**, whereas SDR concerns **conditional independence** (stronger than mean).

# Key ingredient: low-rank support matrix machine (SMM)

- Classification is an easier task than probability estimation.
- First, we develop a large-margin classifier for **matrix predictors in high dimensions**.
- Consider the classifier, **sign** $\{\hat{f}(\mathbf{X})\}$ , where  $\hat{f}(\cdot): \mathbb{R}^{d_1 \times d_2} \mapsto \mathbb{R}$  is the solution to the following optimization over a function class  $\mathcal{F}$ :

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \sum_i [1 - y_i f(\mathbf{X}_i)]_+ + \lambda \|f\|_{\mathcal{F}}^2 \right\}. \quad (1)$$

Here\*  $\mathcal{F} = \{f: f(\cdot) = \langle \cdot, \mathbf{B} \rangle \text{ where } \operatorname{rank}(\mathbf{B}) \leq r\}$  and  $\|f\|_{\mathcal{F}}^2 \stackrel{\text{def}}{=} \|\mathbf{B}\|_{\mathcal{F}}^2$ .

- Assume  $r$  is fixed and known (for now). The classifier (1) is referred to as **low-rank support matrix machine (SMM)**.

\*For presentation convenience, we omit the intercept  $b$  in the representation of  $f$ .

## Key ingredient: low-rank SMM

- The proposed low-rank SMM finds the **most separable hyperplane** among **all possible  $r$ -dimensional representations**.
- Note that the low-rank SMM

$$\min_{\text{rank}(\mathbf{B}) \leq r} \left\{ \sum_i \left[ 1 - y_i \left\langle \underbrace{\mathbf{X}_i}_{\text{ambient dimension } d_1 \times d_2}, \mathbf{B} \right\rangle \right]_+ + \lambda \|\mathbf{B}\|_F^2 \right\}.$$

is equivalent to

$$\min_{\substack{\mathbf{P}_r \mathbf{P}_r^T = \mathbf{P}_c \mathbf{P}_c^T = \mathbf{I}_r \\ \mathbf{C} \in \mathbb{R}^{r \times r}}} \left\{ \sum_i \left[ 1 - y_i \left\langle \underbrace{\mathbf{P}_r \mathbf{X}_i \mathbf{P}_c^T}_{\text{intrinsic dimension } r \times r}, \mathbf{C} \right\rangle \right]_+ + \lambda \|\mathbf{C}\|_F^2 \right\},$$

where  $\mathbf{P}_r \in \mathbb{R}^{r \times d_1}$ ,  $\mathbf{P}_c \in \mathbb{R}^{r \times d_2}$  are row- and column-wise projections, respectively.



## Key ingredient: low-rank SMM

- The solution to the low-rank SMM is represented as

$$\hat{f}(\mathbf{x}) = \sum_i \hat{\alpha}_i y_i \langle \hat{\mathbf{P}}_r \mathbf{x}_i, \hat{\mathbf{P}}_r \mathbf{x} \rangle,$$

where  $\{\hat{\alpha}_i\}$  are estimated (sparse) coefficients in the dual problem and  $\hat{\mathbf{P}}_r$  is the estimated projection matrix.

- We develop an alternating minimization algorithm to jointly estimate  $\mathbf{P}_r$  and  $\{\alpha_i\}$ . (The projection  $\mathbf{P}_c$  is absorbed into  $\{\alpha_i\}$ .)
- We are particularly interested in the high-dimensional regime when both  $n$  and  $\min\{d_1, d_2\} \rightarrow \infty$  while  $r = \mathcal{O}(1)$ . (*An illustration of failures for classical SVMs in high dimensions; JMLR 18(45), 1-21, 2017*).
- The **low-rankness** efficiently prevents overfitting in high dimensions.

# Theory: Low-rank SMM in high dimensions

- Let  $\langle \cdot, \cdot \rangle_{\mathbf{P}}$  denote the rank- $r$  projection kernel for a pair of matrices:

$$\langle \mathbf{X}, \mathbf{X}' \rangle_{\mathbf{P}} \stackrel{\text{def}}{=} \langle \mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{X}' \rangle, \quad \text{for all } \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{d_1 \times d_2},$$

where  $\mathbf{P} \in \mathbb{R}^{r \times d_1}$  is a (unknown) rank- $r$  projection matrix.

- The low-rank SMM considers the decision function of the form:

$$f(\cdot) = \sum_{i=1}^n \alpha_i y_i \langle \mathbf{X}_i, \cdot \rangle_{\mathbf{P}}. \quad (2)$$

- The function  $f \in \mathcal{F}_r$ , where  $\mathcal{F}_r$  is the reproducing kernel Hilbert Space (RKHS) induced by rank- $r$  kernels  $\{\langle \cdot, \cdot \rangle_{\mathbf{P}} : \text{projection } \mathbf{P} \in \mathbb{R}^{r \times d_1}\}$ .
- Column-wise projection kernel can be similarly introduced  $\Rightarrow$  same decision function. *A coincidence or something fundamental?*

# Theory: Low-rank SMM in high dimensions

## Generalization bound (Lee and Wang, 2020+)

With probability at least  $1 - \delta$ , the generalization error of low-rank SMM is

$$\mathbb{P}\{Y^{\text{new}} \neq \text{sign}[\hat{f}(\mathbf{X}^{\text{new}})]\} \leq \text{training error} + \mathbb{E}[\hat{R}_n(\mathcal{F}_r)] + \sqrt{\frac{\ln(\frac{1}{\delta})}{2n}},$$

where  $\hat{R}_n(\mathcal{F}_r)$  denotes the Rademacher complexity of rank- $r$  SMM classifiers. Roughly, in the case  $d_1 \asymp d_2 = \mathcal{O}(d)$ , we have

$$\mathbb{E}[\hat{R}_n(\mathcal{F}_r)] \leq 4\varepsilon + \frac{c_1 r}{\sqrt{n}} \int_{\varepsilon}^{\infty} \log^{1/2} \left( \frac{d/r}{c_2 \varepsilon'} \right) d\varepsilon',$$

where the bound increases with  **$d$**  and  **$r$** . Results highlight the role of  $r$  in preventing overfitting.

1. Check the proof more carefully; 2. consistency for relative risk.

See Varshney, K.R. and Willsky, A.S., *IEEE Trans. Signal Process.*, 59(6), 2496-2512, 2011.

## From classification to regression

Back to the probability estimation problem.

- Consider a piecewise-constant representation of the target probability function  $g(\mathbf{X}) \stackrel{\text{def}}{=} \mathbb{P}(Y = 1|\mathbf{X})$ :

$$g(\mathbf{X}) \approx \frac{1}{H} \sum_{h \in [H]} \mathbb{1} \left\{ \mathbf{X} : g(\mathbf{X}) \leq \frac{h}{H} \right\},$$

where  $H \in \mathbb{N}_+ \rightarrow \infty$  is the smoothing parameter.

- The classification problem has provided candidate decision regions:

$$\mathbb{1} \left\{ \mathbf{X} : \underbrace{\text{sign} [\hat{f}_h(\mathbf{X})] = -1}_{\text{decision region from classification}} \right\} \xrightarrow{P} \mathbb{1} \left\{ \mathbf{X} : \underbrace{\mathbb{P}(Y = 1|\mathbf{X}) \leq \frac{h}{H}}_{\text{target sublevel set}} \right\},$$

for any  $h = 1, \dots, H$ .

- This suggests a non-parametric approach to estimating  $g(\mathbf{X})$ .

# Algorithm

We develop the following algorithm to solve for the target function:

- Step 1. Choose a sequence of weights  $\pi_h = \frac{h}{M}$ , for  $h = 1, \dots, H$ .
- Step 2. For each weight  $\pi_h \in [0, 1]$ , solve the following **weighted low-rank support matrix machine (SMM)**:

$$\hat{\mathbf{B}}_h = \underset{\text{rank}(\mathbf{B}) \leq r}{\operatorname{argmin}} \left\{ \ell_h(\mathbf{B}) + \lambda \|\mathbf{B}\|_F^2 \right\}, \quad \text{where} \quad (3)$$

$$\ell_h(\mathbf{B}) = (1 - \pi_h) \sum_{y_i=1} [1 - \langle \mathbf{X}_i, \mathbf{B} \rangle]_+ + \pi_h \sum_{y_i=-1} [1 + \langle \mathbf{X}_i, \mathbf{B} \rangle]_+.$$

## Algorithm (cont.)

- Step 3. Denote the sequence of solutions and decision regions

$$\hat{f}_h(\cdot) \stackrel{\text{def}}{=} \langle \cdot, \hat{\mathbf{B}}_h \rangle \quad \text{and} \quad \hat{\mathcal{D}}_h = \{\mathbf{X} : \text{sign}[\hat{f}_h(\mathbf{X})] = -1\},$$

for all  $h = 1, \dots, H$ .

- Step 4. Estimate the target probability function by

$$\hat{g}(\mathbf{X}) = \frac{1}{M} \sum_{h \in [H]} \mathbb{1} \left\{ \mathbf{X} \in \hat{\mathcal{D}}_h \right\}. \quad (4)$$

*The estimator (4) is asymptotically equivalent to the original estimator in Wang et al [Biometrika'08].*

*I choose this form because of its good analytic properties.*

# Theory: Probability prediction via low-rank kernels

## High-dimensional consistency (Lee and Wang, 2020+)

Assume the true  $g \in \mathcal{F}$ , where  $\mathcal{F}$  is the RKHS induced by  $\langle \cdot, \cdot \rangle_{\text{rank-}r}$ .

- Given any  $\pi \in [0, 1]$ , the solution to (3) yields the Bayes rule:

$$\text{sign}[\hat{f}_\pi(\mathbf{X})] \xrightarrow{\text{in } P} \text{sign}[g(\mathbf{X}) - \pi], \text{ as } n, d \rightarrow \infty \text{ while } d/n \rightarrow 0.$$

- Our probability estimator (4) is consistent:

$$\hat{g}(\mathbf{X}) \xrightarrow{\text{in } P} g(\mathbf{X}), \text{ as } H, n, d \rightarrow \infty \text{ while } d/n \rightarrow 0.$$

- To the best of our knowledge, this is the first result for SVM-based prob. estimation in large dimension ( $d^2$ ), large sample size ( $n$ ) regime.
- Detailed proofs, assumptions? Convergence in terms of smoothness of  $g$ ,  $r$ ,  $d$ ,  $n$ , and  $H$ ? Do we really need the assumption,  $g \in \mathcal{F}$ ; i.e.,  $g(\cdot)$  linear in  $\mathbf{X}$ ? Perhaps not... composition of monotonic + linear functions is also fine.. indicator functions are dense in  $L[0, 1]^2$ ...*

# SDR via low-rank kernels

- A challenging problem is to identify the **sufficient features** with minimal modeling assumption in the prediction model.
- We develop a robust sufficient dimension reduction (SDR) method for matrix predictors in high dimensions.

Add two steps to Algorithm:

- 1 Pre-step 2 (Whitening):  $\mathbf{X}_i \leftarrow \hat{\Sigma}_r^{-1/2} [\mathbf{X}_i - \text{Mean}(\mathbf{X})] \hat{\Sigma}_c^{-1/2}$ , where  $\hat{\Sigma}_r, \hat{\Sigma}_c$  are empirical row- and column-wise covariance matrices.
- 2 Post-step 4 (Assembling): Arrange outputs  $\hat{\mathbf{B}}_h$  into an order-3 tensor:

$$\mathcal{B}(:, :, h) \leftarrow \hat{\Sigma}_r^{1/2} \hat{\mathbf{B}}_h \hat{\Sigma}_c^{1/2}, \quad \text{for } h = 1, \dots, (H-1).$$

Perform a rank- $(r_1, r_2, r_1 r_2)$  Tucker decomposition on  $\mathcal{B}$ . Let  $\hat{\mathbf{P}}_c, \hat{\mathbf{P}}_r$  denote the estimated factor matrices at the first two modes.



# Theory: SDR via low-rank kernels

## Bilinear sufficient dimension reduction (SDR)

Let  $(\mathbf{X}, y) \in \mathbb{R}^{d_1 \times d_2} \times \{0, 1\}$  be the pair of r.v.'s of interest. Bilinear SDR seeks a pair of matrices  $\mathbf{P}_r \in \mathbb{R}^{r_1 \times d_1}$ ,  $\mathbf{P}_c \in \mathbb{R}^{r_2 \times d_2}$  such that

$$y \perp\!\!\!\perp \mathbf{X} \mid \underbrace{\mathbf{X} \times_1 \mathbf{P}_r \times_2 \mathbf{P}_c}_{r_1\text{-by-}r_2}.$$

The minimum dimension of  $(r_1, r_2)$  is called the structure dimension.

## Unbiasedness (Lee and Wang, 2020+)

Under suitable assumptions\*, our proposed  $\hat{\mathbf{P}}_r$  and  $\hat{\mathbf{P}}_c$  are asymptotically unbiased estimators for  $\mathbf{P}_r$  and  $\mathbf{P}_c$ , respectively; i.e.

$$\Theta(\hat{\mathbf{P}}_r, \mathbf{P}_r) \rightarrow 0, \quad \text{and} \quad \Theta(\hat{\mathbf{P}}_c, \mathbf{P}_c) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Assume: (1)  $\mathbb{E}(\mathbf{X} | \mathbf{X} \times_1 \mathbf{P}_r \times_2 \mathbf{P}_c)$  is a linear function of  $\mathbf{X} \times_1 \mathbf{P}_r \times_2 \mathbf{P}_c$ ; (2)  $\text{Cov}(\text{vec}(\mathbf{X})) = \Sigma_r \otimes \Sigma_c$ .  
 $\text{cov}(X_{\{ip\}}, X_{\{jq\}}) = \text{sigma}_{\{ij\}} * \text{sigma}_{\{pq\}}$

# Simulation: Classification

- Linearly separable model:  $d_1 = 10, d_2 = 8, r = 5$ .

Error =  $\|\mathbf{B} - \hat{\mathbf{B}}\|_F$ , where both  $\mathbf{B}, \hat{\mathbf{B}}$  are normalized to have F-norm 1.

	N = 100	N = 200	N = 300	N = 400
Error (SMM)	0.4827625	0.1903077	0.1020411	0.0663940
Error (SVM)	0.5405979	0.3149391	0.1361625	0.0734705

Table 1: The columns are the number of sample size and rows are the Frobenius norm error according to the two methods.

- Linearly separable model:  $d_1 = 5, d_2 = 4, r = 3, N = 100$

	1st	2nd	3rd	4th	5th	average
SVM	0.90	0.95	0.9	0.75	0.9	0.88
SMM	0.95	0.95	1.0	0.75	0.9	0.91

Table 1: Miss classification rate in Simulation 1

- Linearly non-separable model:  $d_1 = 5, d_2 = 4, r = 2, N = 100$

	1st	2nd	3rd	4th	5th	average
SVM	0.55	0.65	0.75	0.6	0.75	0.66
SMM	0.50	0.80	0.75	0.7	0.75	0.70

Table 3: Miss classification rate in Simulation 2

Details in "041720.pdf". Table caption: Miss classification  $\Rightarrow$  1-MCR?

## Simulation: Probability function estimation

Ground truth  $g(\mathbf{X}) = \mathbb{P}(y = 1|\mathbf{X})$ : nonlinear in  $\mathbf{X}$ , but can be written as a composition of monotonic + linear functions.

$\mathbf{X}|y = -1 \sim_{i.i.d.} \text{MVN}(\mathbf{0}, \mathbf{I}_4)$ ,  $\mathbf{X}|y = 1 \sim_{i.i.d.} \text{MVN}(\mathbf{M}, \mathbf{I}_4)$ ,  $\text{rank}(\mathbf{M}) = 1$

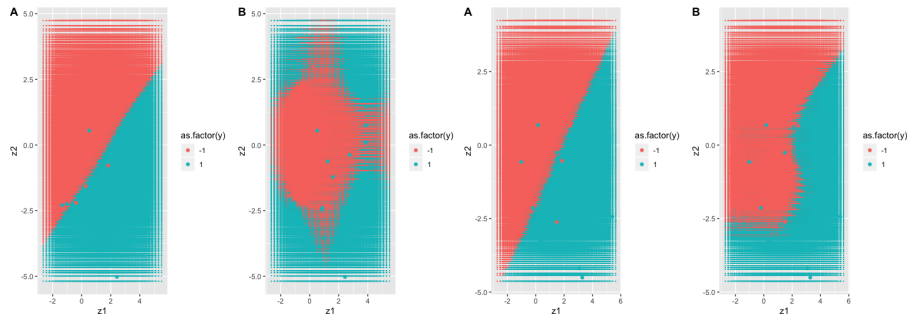


Figure A shows the classification rule of SMM with linear kernel and Figure B shows SMM with polynomial kernel.

Details in “042320.pdf”.

# Simulation: Probability function estimation

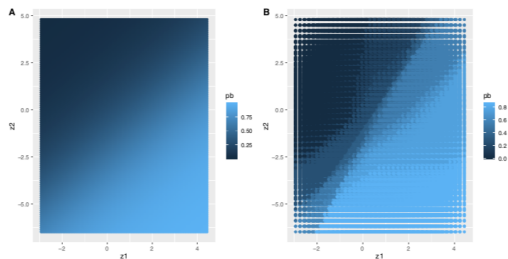


Figure 2: Figure A shows the true conditional probability of  $y$  given  $(Z_1, Z_2)$ . Figure B shows the estimated probability using SMM weighted hinge loss approach.

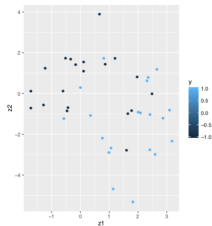


Figure 1: Visualization of the feature matrices.  $Z_1$  is the sum of the first column and  $Z_2$  is of the second one.

Details in “042020.pdf”.

## Extension to nonlinear kernels

- Let  $\mathbf{X}, \mathbf{X}' \in \mathbb{R}^{d_1 \times d_2}$  be a pair of matrices. Given a projection matrix  $\mathbf{P} \in \mathbb{R}^{r \times d_1}$ , a **linear low-rank kernel** is defined as

$$\begin{aligned}\langle \mathbf{X}, \mathbf{X}' \rangle_{\mathbf{P}} &\stackrel{\text{def}}{=} \langle \mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{X}' \rangle \\ &= \text{trace} \left[ \underbrace{\mathbf{X}\mathbf{X}'^T}_{\text{bilinear map}} \underbrace{\mathbf{P}^T \mathbf{P}}_{\text{low-rank projection}} \right].\end{aligned}$$

- Similarly, we propose a **nonlinear low-rank kernel**

$$\begin{aligned}\langle \mathbf{X}, \mathbf{X}' \rangle_{\mathbf{P}, h} &\stackrel{\text{def}}{=} \langle \mathbf{P}h(\mathbf{X}), \mathbf{P}h(\mathbf{X}') \rangle \\ &= \text{trace} \left[ \mathbf{K}(\mathbf{X}, \mathbf{X}') \mathbf{P}^T \mathbf{P} \right],\end{aligned}$$

where  $h: \mathbb{R}^{d_1 \times d_2} \mapsto \mathbb{R}^{d_1' \times d_2'}$  denotes a nonlinear map between two matrix spaces, and  $\mathbf{K}(\mathbf{X}, \mathbf{X}') \stackrel{\text{def}}{=} h(\mathbf{X})h^T(\mathbf{X}') \in \mathbb{R}^{d_1' \times d_1'}$  denotes the matrix product of mapped features.

# Nonlinear kernels for matrix predictors

- Each entry of  $\mathbf{K}(\mathbf{X}, \mathbf{X}')$  is the inner product between two vectors.
- We refer to  $\mathbf{K}(\cdot, \cdot) \in \mathbb{R}^{d_1 \times d_1}$  as the *lifted (matrix-valued) kernel* induced by the nonlinear map  $h$ .
- Examples:
  - ▶ Linear kernel:  $\mathbf{K}(\mathbf{X}, \mathbf{X}') = \mathbf{X}\mathbf{X}'^T$ .
  - ▶ Polynomial kernel with degree  $m$ :  $\mathbf{K}(\mathbf{X}, \mathbf{X}') = (\mathbf{X}\mathbf{X}'^T + \lambda \mathbf{I})^{\circ m}$ .
  - ▶ Gaussian kernel: the  $(i, j)$ -th entry of  $\mathbf{K}(\mathbf{X}, \mathbf{X}')$  is

$$[\mathbf{K}(\mathbf{X}, \mathbf{X}')]_{(i,j)} = \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{X}[i, :] - \mathbf{X}'[j, :]\|_2^2 \right\}$$

for all  $(i, j) \in [d_1] \times [d_1]$ .

- Implementation: (1) symmetrization trick; (2) A single projection matrix is enough to guide the optimization. Update at every other step.
- *Perhaps use a different name for  $\mathbf{K}$ ? Not really a kernel...*  
 $\mathbf{K}(\mathbf{X}, \mathbf{X}') \neq \mathbf{K}(\mathbf{X}', \mathbf{X})$ ;  $\mathbf{K}(\mathbf{X}, \mathbf{X}') \neq \mathbf{K}(\mathbf{X}^T, \mathbf{X}'^T)$ .

# Application of nonlinear kernels to matrix-based learning

- Nonparametric classifier:  $\text{sign}[\hat{f}(\cdot)]: \mathbb{R}^{d_1 \times d_2} \mapsto \{-1, 1\}$ , where

$$\hat{f}(\cdot) = \sum_i \hat{\alpha}_i y_i \text{tr} \left[ \mathbf{K}(\cdot, \mathbf{x}_i) \hat{\mathbf{P}}^T \hat{\mathbf{P}} \right].$$

- Nonparametric regression:

$$\hat{\mathbb{P}}(Y = 1 | \mathbf{X}^{\text{new}}) = \frac{1}{H} \sum_{h \in [H]} \mathbb{1} \left\{ \mathbf{X}^{\text{new}}: \text{sign}[\hat{f}_h(\mathbf{X}^{\text{new}})] = 1 \right\}$$

- Nonlinear SDR,  $Y \perp\!\!\!\perp \mathbf{X} | \phi(\mathbf{X})$ , where

$$\phi(\mathbf{X}) = h(\mathbf{X}) \times_1 \hat{\mathbf{P}}_c \times_2 \hat{\mathbf{P}}_r(??)$$

# Simulation: Linear vs. nonlinear classification

- linearly separable, homoscedastic case:  $d_1 = 10, d_2 = 8, r = 3, N = ?$

	1st	2nd	3rd	4th	5th	average
SVM	0.85	0.70	0.75	0.70	0.50	0.70
SMM	0.90	0.90	0.75	0.80	0.85	0.84
SMM(polynomial)	0.75	0.55	0.85	0.55	0.75	0.69
SMM(gaussian)	0.65	0.50	0.85	0.70	0.80	0.70

Table 1: Miss Classification Rate (MCR) on 5 folded Cross validation(CV)

- Non-separable, heteroscedastic case:  $d_1 = d_2 = 2, \text{Cov}(\mathbf{X}|y = -1) = 4\text{Cov}(\mathbf{X}|y = 1) = 4\mathbf{I}_2$ .

	sim 2.1	sim 2.2	sim 2.3
SVM	0.76	0.735	0.860
SMM	0.75	0.740	0.865
SMM(polynomial)	0.80	0.750	0.785
SMM(gaussian)	0.71	0.745	0.830

Table 3: This table shows the averaged MCR of 5 folded CV according to different methods.

Details in "042720.pdf"



# Simulation: Nonlinear learning

- Nonlinear classification

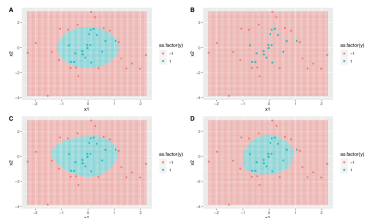


Figure 1: Subfigure A is true ellipsoid boundary. B is linear case boundary which assigns labels all 0. C and D show the boundary of polynomial and exponential kernel respectively.

- Nonlinear probability estimation

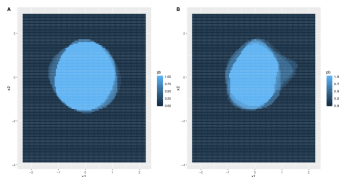


Figure 2: Figure A and B show the estimated probability with polynomial and exponential kernel respectively.

# Summary

- We have developed a robust nonparametric framework for classification, prediction and dimension reduction with **matrix-valued predictors**.
- The proposed **low-rank projection kernel** achieves statistical efficiency and adaptability in high dimension regimes.
- Our method can be thought of as **adaptive kernel learning** by jointly performing (nonlinear or linear) dimension reduction and classification/prediction.
- **Computational efficiency** can be further boosted by employing solution path algorithm and stochastic gradient descent.

# Questions

- [Regression] How does SVM-based learning compare to other nonparametric regression, such as local smoothing, K-NN, Neural Network?  
  
⇒ SVM-based learning is essentially a kernel smoothing method. Neighborhood is defined by support vectors (c.f. (2)). Intuition?
- [SDR] How does SVM-based SDR compare to other SDR methods?
- [Kernel] Connection between low-rank kernel vs. adaptive kernel, kernel learning, sum of rank-1 kernels?
- Extend to higher-order tensors. Data-driven choice of  $r$ .
- Real data analysis