

Nonparametric learning with matrix-valued predictors in high dimensions

Abstract

We consider the problem of learning the relationship between a binary label response and a high-dimensional matrix-valued predictor. Such data problems arise commonly in brain imaging studies, sensor network localization, and personalized medicine. Existing regression analysis often takes a parametric procedure by imposing a pre-specified relationship between variables. However, parametric models are insufficient in capturing complex regression surfaces defined over high-dimensional matrix space. Here, we propose a flexible nonparametric framework for various learning tasks, including classification, level set estimation, and regression, that specifically accounts for the matrix structure in the predictors. Unlike classical approaches, our method adapts to the possibly non-smooth, non-linear pattern in the regression function of interest. The proposal achieves prediction and interpretability simultaneously via a joint optimization of prediction rules and dimension reduction in the matrix space. Generalization bounds, estimation consistency, and convergence rate are established. We demonstrate the advantage of our method over previous approaches through simulations and applications to **XXX** data analyses.

Keywords: Nonparametric learning, matrix-valued predictors, high dimension, classification, level-set estimation, regression.

1 Introduction

2 Methods

2.1 Models and Motivation

Consider a statistical learning problem where we would like to model the relationship between a feature $\mathbf{X} \in \mathcal{X}$ and a response $Y \in \mathcal{Y}$. Suppose that we observe a sample of n data points, $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$, identically and independently distributed (i.i.d.) according to a unknown distribution $\mathbb{P}(\mathbf{X}, Y)$ over $\mathcal{X} \times \mathcal{Y}$. We are interested in predicting a new response Y_{n+1} from a new feature value \mathbf{X}_{n+1} . The observations $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ are called the training data and $(\mathbf{X}_{n+1}, Y_{n+1})$ the test point. When no confusion arises, we often omit the subscript $(n+1)$ and simply write (\mathbf{X}, Y) for the prototypical test point. The test point is assumed independent of the training data and is drawn from the same unknown distribution \mathbb{P} . Our goal is to make accurate prediction under a wide range of distributions. In particular, we consider a non-parametric, distribution-free setting with no strong assumptions on the data generative distribution other than i.i.d.

We focus on the scenario with matrix-valued predictors and binary label response; that is, $\mathcal{X} = \mathbb{R}^{d_1 \times d_2}$ and $\mathcal{Y} = \{-1, 1\}$. Matrix-valued predictors ubiquitously arise in modern applications. One example is from electroencephalography studies of alcoholism. The data set records voltage value measured from 64 channels of electrodes on 256 subjects for 256 time points [3]. Each feature is a 256×64 matrix and the response is a binary indicator of subject being alcoholic or control. Another example is pedestrian detection from image data. Each image is divided into 9 regions where local orientation statistics are generated with a total of 22 numbers per region. This yields a 22×9 matrix-valued feature and a

binary label response indicating whether the image is pedestrian [1].

In the above two examples and many other studies, researchers are interested in *interpretable prediction*, where the goal is to not only make accurate prediction but also identify features that are informative to the prediction. While classical learning method has been ..., a lack. A key challenge with matrix data is the complex structure in the feature space. Matrix space is a tensor product space in which both rows and columns. The notions of features are not well defined for matrix space. For example, one may be interested in entrywise significance, or important vectors, or the spaces that corresponds to the matrix. Find the right representations. interpretable, adaptive, informative, ..., More relevant in non-linear learning. Bad representation leads to complicated data..., sensitive to noise. A model that is complex in the original space may projected into simpler surface in manifold. (complex function surface). Another challenge is ...Notably, the ambient dimension with matrix-valued feature is d_1d_2 , while the number of sample is n . Modern applications are often in the high dimensional regime where d_1d_2 is comparable to, or even larger, than the sample size n . (do we allow d to increase with n in the asymptotic result?). The challenge of high-dimensionality is notably hard especially for matrix-valued features. A second challenge is the notation of informative features. While vectors ... Matrices encodes more information. such as columns space, row space, local structure, global structure, entrywise, spacewise, vs. volumne, etc. One of our contribution is to have a general framework that ...

Before we propose methods to address the challanges. Our framework is general in that it . We focus on three concrete three learning problems and the connection between them. As we will see, combinations

2.2 Three Main Learning Problems

The distribution $y|\mathbf{X}$ is binary, so its distribution is completely determined by the class probability $\mathbb{P}(Y = 1|\mathbf{X})$. Based on the properties $\mathbb{E}(y|\mathbf{X}) = 2\mathbb{P}(Y = 1|\mathbf{X}) - 1$, it is equivalent to estimate $\mathbb{E}(Y|\mathbf{X})$. We denote the function $\eta(\mathbf{X}) = \mathbb{E}(y|\mathbf{X})$. Three questions in increasing difficulty.

Classification. Find the level set $S_0 = \{\mathbf{X} : \eta(\mathbf{X}) \geq 0\}$.

Level-set estimation. Given $\pi \in [-1, 1]$, find the level set $S_\pi = \{\mathbf{X} : \eta(\mathbf{X}) \geq \pi\}$.

Regression: Estimate the function $\eta(\mathbf{X}) : \mathbb{R}^{d_1 \times d_2} \mapsto [-1, 1]$.

Each of them does not depend on the unknown function η , and furthermore, the minimizer of decision rule is the desired properties. Problem 1 is a special case of Problem 2. Problem 3 can be addressed by.

Because log-likelihood for η is unknown, we consider the loss function that does not depend on the ground truth and distribution. The following loss can be treated as non-likelihood loss.

Example 1. Define classification error $\mathbb{P}[y \neq g(\mathbf{X})]$. Find a decision rule g such that minimizes classification error in the space $g \in \mathcal{F}$.

Example 2. Define weighted classification error: $\mathbb{E}[|y - \pi| \mathbb{1}\{y \neq g(\mathbf{X})\}]$. Find a decision rule g such that minimizes classification error in the space $g \in \mathcal{F}$.

Example 3. Define least-squared loss $\mathbb{E}[f(\mathbf{X}) - y]^2$. Find a decision rule g such that minimizes classification error in the space $f \in \mathcal{F}$.

3 Example: Trace regression model

3.1 Example: Non-linear matrix completion

$$y_{ij} = \text{sign}(f(x_{ij})), \text{ where } \mathbf{X} = \llbracket x_{ij} \rrbracket.$$

Estimate \mathbf{B} itself recovers the binary matrix completion problem.

3.2 Example: Graph classification

Both node level and edge level. 1. projection of nodes. 2. of edges (edges are generated from eigen model.. angle??? different homophily effect to the response label)

Difference in topology. absence of presence of small hubs. Output is the projection matrices.

3.2.1 Choice of loss function

Large-margin loss which is convex of the 0-1 loss. The sample version is

$$L(g) = \sum_{i=1}^n (1 - y_i g(\mathbf{X}_i))_+.$$

Other fisher-consistent loss is also possible. Such as Similarly, the weighted version is

$$L(g) = \sum_{y_i=1} (1 - \pi)(1 - g(\mathbf{X}_i))_+ + \sum_{y_i=-1} (1 + \pi)(1 - g(\mathbf{X}_i))_+$$

Consider the optimization problem

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n L(f(\mathbf{X}_i), y_i) + \lambda J(f)$$

where depending on the application, the cost function $L(\cdot, \cdot)$ can be selected to be, e.g., the least-squares (LS), the logistic or the hinge loss; $J(\cdot)$ is an increasing function; and, $\lambda > 0$ is a regularization parameter that controls overfitting.

4 Matrix kernels for subspace embedding

The aforementioned three problems are formulated as empirical risk minimization over functions defined over matrices. In this section, we will propose a family of matrix kernels as building blocks for constructing functions in matrix space. Informally, a matrix kernel is a distance measure between matrices based on their similarity. Kernels defined over non-vector objects have recently evolved into a rapidly developing area of structured learning. Unlike vectors, matrix-value feature encodes a two-way relationship across rows and columns. Taking into account the two-way relationship is essential in the matrix kernel formulation.

We now propose a matrix kernel inspired by low-rank latent factor models. The use of matrix kernel in learning will be described in Section XX.

Definition 1 (Matrix Kernel). *Given two matrices \mathbf{M} and \mathbf{M}' from the space of matrix \mathcal{M} . The problem of matrix comparison is to find a mapping s*

$$s: \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$$

such that $s(\mathbf{M}, \mathbf{M}')$ quantifies the similarity (or dissimilarity) of \mathbf{M} and \mathbf{M}' .

(Add a figure)

Goal: topological description. Map each matrix to a feature vector, and use distance and metrics on vectors for learning on graphs.

Let $\mathbf{X} \in \text{Sym}_d(\mathbb{R})$ be the feature matrix. We define a feature mapping Φ that sends a matrix in $\text{Sym}_d(\mathbb{R})$ to $\text{Sym}_d(\mathcal{H}^2)$. Here $\text{Sym}_d(\mathcal{H}^2)$ denote the collection of d -by- d symmetric matrices with each matrix entry takes value in \mathcal{H}^2 . Let $\mathbf{B} \in \mathcal{H}^{d \times d}$ be a matrix defined in \mathcal{H} . Replace product in \mathbb{R} to inner product in \mathcal{H} .

5 Mahalanobis metric

1. Sum. $\mathbf{B} + \mathbf{B}' = \llbracket b_{ij} + b'_{ij} \rrbracket \in \mathcal{H}^{d \times d}$
2. Liner combination. Let $\mathbf{P} \in \mathbb{R}^{d \times r}$ be a real-valued matrix. Define a $\mathbf{BP} = \llbracket c_{ij} \rrbracket \in \mathcal{H}^{d \times r}$, where $c_{ij} = \sum_s p_{sj} b_{is} \in \mathcal{H}^2$.
3. Inner product. $\langle \mathbf{B}, \mathbf{B}' \rangle = \sum_{ij} \langle b_{ij}, b'_{ij} \rangle$.
4. Matrix product $\mathbf{BB}' = \llbracket c_{ij} \rrbracket \in \mathbb{R}^{d \times d}$, where $c_{ij} = \sum_s \langle b_{is}, b'_{sj} \rangle_H$.

Here $\mathcal{H}^2 = \mathcal{H} \times \mathcal{H}$ denotes the Cartesian product of two Hilbert space. Specifically, the mapping $\mathbf{X} \mapsto \Phi(\mathbf{X}) = \llbracket f_{ij} \rrbracket$, where each element is a pair of possibly infinite dimensional features defined as

$$f_{ij} = [\Phi(\mathbf{X})]_{ij} \stackrel{\text{def}}{=} (\phi(\mathbf{X}_{i:}), \phi(\mathbf{X}_{:j})) \in \mathcal{H}^2, \quad \text{for all } i \geq j, (i, j) \in [d] \times [d],$$

and $f_{ji} = f_{ij}$. Note that the entry f_{ij} takes value in \mathcal{H}^2 for all $(i, j) \in [d] \times [d]$. Furthermore, $\Phi(\mathbf{X})$ is a symmetric matrix in the sense that $[\Phi(\mathbf{X})]_{ij} = [\Phi(\mathbf{X})]_{ji}$ for all (i, j) . Then the decision function is defined in the high dimensional features,

$$f(\mathbf{X}) = \langle \mathcal{B}, \Phi(\mathbf{X}) \rangle, \quad \text{where } \mathcal{B} \in (\mathcal{H}^2)^{d \times d}.$$

The parameter $\mathcal{B} = \llbracket \mathbf{b}_{ij} \rrbracket$ is a d -by- d matrix defined over \mathcal{H}^2 ; that is, each entry $\mathbf{b}_{ij} = (b_{ij1}, b_{ij2}) \in \mathcal{H}^2$ for all $(i, j) \in [d]^2$. We consider low-rank structure in $\mathcal{B} = \mathbf{P}^T \mathcal{C} \mathbf{P}$. Then low rank function space

$$\begin{aligned} \mathcal{F} &= \{f: \mathbf{X} \mapsto \langle \mathbf{C}, \mathbf{P}^T \Phi(\mathbf{X}) \mathbf{P} \rangle \mid \mathbf{P} \mathbf{P}^T = \mathbf{I}, \mathbf{P} \in \mathbb{R}^{d \times r}, \mathcal{C} \in (\mathcal{H}^2)^{r \times r}\} \\ &= \{f \in \text{RKHS generated by } \mathcal{K}(\mathbf{P}) \mid \mathbf{P} \mathbf{P}^T = \mathbf{I}, \mathbf{P} \in \mathbb{R}^{d \times r}\} \\ &= \{f \in \text{RKHS generated by } \mathcal{K}(\mathbf{W}) \mid \mathbf{W} \succeq 0, \text{rank}(\mathbf{W}) \leq r\} \end{aligned}$$

The orthogonality condition is based on the relationship $\mathbf{W} = \mathbf{P}^T \Lambda^2 \mathbf{P}$ and scale invariance between the function $\langle \mathcal{C}, \mathbf{P}^T \Phi(\mathbf{X}) \mathbf{P} \rangle = \langle \Lambda \mathbf{P} \mathcal{C}, \Phi(\mathbf{X}) \mathbf{P} \Lambda \rangle$. The kernel

$$\begin{aligned} \mathcal{K}(\mathbf{W}) &= \langle \mathbf{P}^T \Phi(\mathbf{X}) \mathbf{P}, \mathbf{P}^T \Phi(\mathbf{X}') \mathbf{P} \rangle \\ &= c_{\text{col}} \mathcal{K}_{\text{row}}(\mathbf{P} \mathbf{P}^T) + c_{\text{row}} \mathcal{K}_{\text{col}}(\mathbf{P} \mathbf{P}^T) \end{aligned}$$

Properties 1. Let $\mathbf{P} \in \mathbb{R}^{d \times r}$ be a projection matrix. The mapping $\mathbf{P} \circ \Phi = \mathbf{P} \Phi(\mathbf{X}) \mathbf{P} \in (\mathcal{H}^2)^{r \times r}$ defines a kernel $\mathcal{K}(\mathbf{X}, \mathbf{X}') =$

Linear function family: Represented by features. Let $\mathbf{X} \in \text{Sym}_d(\mathbb{R})$.

$$\begin{aligned} \mathcal{F}(r) &= \{f: \mathbf{X} \mapsto \langle \mathbf{W}, \mathbf{X} \rangle \mid \text{rank}(\mathbf{W}) \leq r, \mathbf{W} \in \text{Sym}_d(\mathbb{R})\} \\ &= \{f: \mathbf{X} \mapsto \langle \mathbf{C}, \mathbf{P}^T \mathbf{X} \mathbf{P} \rangle \mid \mathbf{P} \mathbf{P}^T = \mathbf{I}, \mathbf{P} \in \mathbb{R}^{d \times r}, \mathbf{C} \in \text{Sym}_r(\mathbb{R})\} \end{aligned}$$

Definition 2 (Kernel defined in matrix space). Let $d(\cdot, \cdot)$ be a kernel defined in vector space \mathbb{R}^d , and $\mathbf{W} = \llbracket w_{ij} \rrbracket \in \mathbb{R}^{d \times d}$ be a rank- r semi-positive definite matrix. Then \mathbf{K} and \mathbf{W} induce a low-rank projection kernel \mathcal{K} in matrix space:

$$\begin{aligned} \mathcal{K}: \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times d} &\mapsto \mathbb{R} \\ (\mathbf{X}, \mathbf{X}') &\mapsto \mathcal{K}(\mathbf{X}, \mathbf{X}') = \sum_{i,j \in [d]} w_{ij} d(\mathbf{x}_i, \mathbf{x}'_j), \end{aligned}$$

where $\mathbf{x}_i, \mathbf{x}'_j$ denote the i -th and j -th columns of \mathbf{X}, \mathbf{X}' , respectively.

Often, the kernel \mathbf{K} is specified by users, whereas the projection kernel \mathbf{W} is learned by our algorithm. In particular $\mathbf{W} \propto \mathbf{I}_{d \times d}$ corresponds to the classical vectorization case. We denote the low-rank projection kernel $\mathcal{K} = \mathcal{K}(\mathbf{W})$.

Properties 2 (Properties of projection kernel). The kernel family defines a set of kernels that indexed by \mathbf{W} . In particular,

- $\mathcal{K}(\mathbf{W}) + \mathcal{K}(\mathbf{W}') = \mathcal{K}(\mathbf{W} + \mathbf{W}')$ for all $\mathbf{W}, \mathbf{W}' \in \mathbb{R}^{d \times d}$.
- In the special rank-1 case $\mathbf{W} = \llbracket w_i w_j \rrbracket$, our kernel corresponds to the bandwidth selection.
- The kernel $\mathcal{K}(\mathbf{W})$ corresponds to a valid feature mapping $\Phi: \mathbb{R}^{d \times d} \rightarrow \mathcal{H}^r$.
- Rank- k kernel is a convex hull of rank-1 kernels.

in total $d(d-1)/2$ terms involving interaction between row and column $K(\mathbf{x}_i, \mathbf{x}_j)$.
 $C = \mathcal{H}_1 + \dots + \mathcal{H}_{d(d-1)/2}$. Then rank- r projection over the left space.? \mathbf{X}
 Consider the family of kernels

$$\mathcal{F}(r) = \{f \in \text{RKHS generated by } \mathcal{K}(\mathbf{W}) \mid \mathbf{W} \succeq 0, \text{rank}(\mathbf{W}) = r\}.$$

We use $\mathcal{F}(r)$ to denote the set of kernels generated by rank- r projection.

Each kernel \mathcal{K} generates a RKHS. We consider the optimization over the union of RKHS.
 The union of RKHS is not RKHS.

$$\max_{f \in \mathcal{F}(r)} L(f) = \max_{\substack{\text{rank}(\mathbf{W})=r, \\ \mathbf{W} \succeq 0}} \max_{f \in \text{RKHS}(\mathcal{K}(\mathbf{W}))} L(f)$$

Connection to adaptive kernel learning The above can be viewed as an optimization of kernels

$$\begin{aligned} \max_{f \in \mathcal{F}(r)} L(f) &= \max_{f = \alpha_1 f_1 + \dots + \alpha_r f_r : f_i \in \mathcal{K}(1), \alpha \in \mathbb{R}^n} L(f) \\ &= \max_{f \in \mathcal{K}(1) \oplus \dots \oplus \mathcal{K}(1)} L(f) \end{aligned}$$

Connection to Metric learning

$$\max_{f \in \mathcal{F}(r)} L(f) = \max_{\alpha \in \Delta, f \in \text{RKHS}(\mathcal{K}(1))} L(f)$$

Ensemble methods.

Example 4 (Multiple kernel learning). Let $\mathbf{X} = \text{diag}(x_i)$, $\mathbf{X}' = \text{diag}(x'_i)$ are diagonal matrices and inner-product kernel $\mathbf{K}(\mathbf{x}, \mathbf{x}') = g(\langle \mathbf{x}, \mathbf{x}' \rangle)$. Then our rank-1 projection kernel reduces to

$$\mathcal{K}(\mathbf{X}, \mathbf{X}') = \sum_{i \in [d]} w_i^2 d(x_i, x'_i)$$

Example 5 (Feature selection with lasso penalty). Let \mathbf{X}, \mathbf{X}' are diagonal matrices.

$$\mathcal{K}(\mathbf{X}, \mathbf{X}') = \sum_{\alpha \in \Delta(r)} \alpha_i d(x_i, x'_i)$$

weighted features. Let $\mathcal{F}(1)$ denote the RKHS induced by rank-1 projection. Then rank $\mathcal{F}(r) = \mathcal{F}(1) \oplus \dots \oplus \mathcal{F}(1)$

Example 6. Group lasso on the features. Weight $\sim \text{Multinomial}(d, r)$. Then rank-kernel is r .

Example 7 (Convolution kernel).

Orthogonality between rank-1 kernels. If and only if $\mathcal{F}(1) \cap \mathcal{F}(1) = 0$

Each of the kernel induces the functions defined in matrix space.

$$\begin{aligned} \mathcal{F}(r) &= \{f \in \mathcal{F}(\mathbf{A}) \mid \mathbf{A} \succeq 0, \text{rank}(\mathbf{A}) \leq r\} \\ &= \left\{ \mathbf{X} \mapsto \sum_{i,j} a_{ij} \mathbf{K}(\mathbf{x}_i, \cdot) \mid \mathbf{A} = \llbracket a_{ij} \rrbracket \text{ is a rank-} r \text{ PSD matrix} \right\}. \end{aligned}$$

Example 8 (Rank-1 kernel).

$$\mathcal{K}(\mathbf{W}) = \sum_{ij} w_i w_j K(\mathbf{x}_i, \mathbf{x}'_j) \quad \text{for some } \mathbf{a} \in \mathbb{R}^d.$$

$$\mathcal{F} = \{ \mathbf{X} \mapsto \sum_{ij} a_i a_j \mathbf{K}(\mathbf{x}_i, \cdot) \mid \mathbf{a} \in \mathbb{R}^d \}$$

$$f \in \text{Span}(\mathcal{F}) \text{ and } \text{rank}(f) \leq r.$$

5.0.1 Non-Asymmetric case

$$\mathcal{K}(\mathbf{P}_1, \mathbf{P}_2) = \sum_{(i_1, i_2, j_1, j_2)} a_{i_1 j_1} b_{j_1 j_2} d(\mathbf{x}_{i_1 i_2}, \mathbf{x}'_{j_1 j_2}) = \langle D(\mathbf{X}, \mathbf{X}'), \mathbf{P}_1 \otimes \mathbf{P}_2 \rangle$$

5.1 Classification

We consider linear predictors on matrix feature of the form

$$f_{\mathbf{B}}(\mathbf{X}) = \langle \mathbf{B}, \mathbf{X} \rangle,$$

where $\mathbf{B}, \mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ and $\langle \mathbf{X}, \mathbf{X}' \rangle = \text{Tr}(\mathbf{X}^T \mathbf{X}')$. We extend the linear model to non-linear case in Section 7. This matrix representation has advantage of regularizing the coefficient matrix \mathbf{B} and restricting the number of parameters. Specifically, we assume that \mathbf{B} has low-rank structure such that

$$\mathbf{B} = \mathbf{U}\mathbf{V}^T \text{ where } \mathbf{U} \in \mathbb{R}^{d_1 \times r}, \mathbf{V} \in \mathbb{R}^{d_2 \times r} \text{ and } r \leq \min(d_1, d_2)$$

We show how low rankness let us consider column and row wise structure of feature matrices and prevents overfitting in Section 8. Assume we are given a set of training data and label pairs $\{\mathbf{X}_i, y_i\}_{i=1}^n$ where $\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}$ and $y_i \in \{-1, +1\}$. Our goal is to learn a model with a low error on the training data. One successful approach is large-margin classifiers. A large-margin classifier minimizes a cost function in f over a decision function class \mathcal{F} :

$$\min_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n L(y_i f(\mathbf{X}_i)) + \lambda J(f), \quad (1)$$

where $J(f)$ is a penalty term for model complexity and $L(z)$ is a margin loss that is a function of the functional margin $yf(\mathbf{X})$. Examples of such loss functions are the hinge loss function $L(z) = (1 - z)_+$ and the logistic loss function $L(z) = \log(1 + e^{-z})$. For

demonstration, we focus on the hinge loss functions. However, our estimation schemes and theorems are applicable to general large-margin classifiers. Consider linear decision function class with low rank coefficient $\mathcal{F} = \{f : f(\cdot) = \langle \mathbf{U}\mathbf{V}, \cdot \rangle \text{ where } \mathbf{U} \in \mathbb{R}^{d_1 \times r}, \mathbf{V} \in \mathbb{R}^{d_2 \times r}\}$. The solution $f(\mathbf{X})$ of Equation (1) is shown to have the form (check Supplement)

$$f(\mathbf{X}) = \sum_{i=1}^n \alpha_i y_i \langle \mathbf{P}_r \mathbf{X}_i, \mathbf{P}_r \mathbf{X} \rangle,$$

where $\{\alpha_i\}_{i=1}^n$ are solution (spare) in the dual problem of Equation (1) and $\mathbf{P}_r \in \mathbb{R}^{r \times d_1}$ is the projection matrix induced by low rank coefficient \mathbf{U}, \mathbf{V} . Let $\langle \cdot, \cdot \rangle_{\mathbf{P}_r}$ denote the low rank linear kernel for a pair of matrices:

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

Therefore, we consider the decision function of the form:

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

We can think of our considered classification function class as the reproducing kernel Hilbert space induced by rank- r linear kernels $\{\langle \cdot, \cdot \rangle_{\mathbf{P}} : \mathbf{P} \in \mathbb{R}^{r \times d_1} \text{ is a projection matrix}\}$. We estimate coefficient $\{\hat{\alpha}_i\}_{i=1}^n$ and the projection matrix $\hat{\mathbf{P}}_r$ from a given dataset. Detailed estimation algorithm appears in Section 6. We obtain our classification rule as

$$G(\mathbf{X}) = \text{sign} \left(\sum_{i=1}^n \hat{\alpha}_i y_i \langle \mathbf{X}_i, \mathbf{X} \rangle_{\hat{\mathbf{P}}_r} \right).$$

5.2 Probability function estimation

Our proposed method is designed to estimate $p(\mathbf{X}) \stackrel{\text{def}}{=} \mathbb{P}(y = 1 | \mathbf{X})$ at any \mathbf{X} which does not necessarily belong to the observed training data set. We consider the weighed version

of (1),

$$\min_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n \omega_{\pi}(y_i) L(y_i f(\mathbf{X}_i)) + \lambda J(f), \quad (2)$$

where $\omega_{\pi}(y) = 1 - \pi$ if $y = 1$ and π if $y = -1$. **(author?)** [2] showed that The minimizer \hat{f}_{π} to Equation (2) is a consistent estimate of $\text{sign}(p(\mathbf{X}) - \pi)$. Therefore, for a given smoothing parameter $H \in \mathbb{N}_+$, We estimate the target probability through two main steps.

$$\begin{aligned} p(\mathbf{X}) &\stackrel{\text{step1}}{\approx} \sum_{h=1}^H \frac{h-1}{H} \mathbb{1} \left\{ \mathbf{X} : \frac{h-1}{H} \leq p(\mathbf{X}) < \frac{h}{H} \right\} \\ &\stackrel{\text{step2}}{\approx} \sum_{h=1}^H \frac{h-1}{H} \mathbb{1} \left\{ \mathbf{X} : \text{sign}(\hat{f}_{\frac{h-1}{H}}) = 1, \text{sign}(\hat{f}_{\frac{h}{H}}) = -1 \right\}, \end{aligned}$$

where Step 1 approximates the target probability by linear combination of step functions and Step 2 uses the fact that the solution of (2) is consistent to Bayes rule. The probability estimation scheme can be summarized as follows.

Scheme

S.1 Choose a sequence of weight $\pi_h = \frac{h}{H}$, for $h = 1, \dots, H$.

S.2 For each weight $\pi_h \in [0, 1]$, solve Equation (2) with $\omega_{\pi_h}(y)$

S.3 Denote the sequence of solutions and decision regions

$$\{\hat{f}_h\}_{h=1}^H \text{ and } \{\hat{\mathcal{D}}_h\}_{h=1}^H = \left\{ \left\{ \mathbf{X} : \text{sign}(\hat{f}_{\frac{h-1}{H}}) = 1, \text{sign}(\hat{f}_{\frac{h}{H}}) = -1 \right\} \right\}_{h=1}^H$$

S.4 Estimate the target probability function by

$$\hat{p}(mX) = \sum_{h=1}^H \frac{h-1}{H} \mathbb{1} \left\{ \mathbf{X} \in \hat{\mathcal{D}}_h \right\}.$$

Notice that this scheme can be applied to any large-margin classifiers though we focus on the hinge loss function in this paper. Solution of weighted hinge loss in Equation (2) is solved with simple modification from classification algorithm in Section 6.

6 Algorithm

In this Section, we describe the algorithm to seek the optimizer of Equation (1) in the case of hinge loss function $L(z) = (1 - z)_+$ and linear function class $\mathcal{F} = \{f : f(\cdot) = \langle \mathbf{U}\mathbf{V}^T, \cdot \rangle\}$, where $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$ $\mathbf{V} \in \mathbb{R}^{d_2 \times r}$. Equation (1) is written as

$$\min_{\{\mathbf{U}, \mathbf{V}\} \in \mathbb{R}^{d_1 \times r} \times \mathbb{R}^{d_2 \times r}} n^{-1} \sum_{i=1}^n (1 - y_i \langle \mathbf{U}\mathbf{V}^T, \mathbf{X}_i \rangle)_+ + \lambda \|\mathbf{U}\mathbf{V}^T\|_F^2 \quad (3)$$

We optimize Equation (3) with a coordinate descent algorithm that solves one block holding the other block fixed. Each step is a convex optimization and can be solved with quadratic programming. To be specific, when we fix \mathbf{V} and update \mathbf{U} we have the following equivalent dual problem

$$\begin{aligned} & \max_{\boldsymbol{\alpha} \in \mathbb{R}^n : \boldsymbol{\alpha} \geq 0} \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{X}_i, \mathbf{X}_j \mathbf{V} (\mathbf{V}^T \mathbf{V})^{-1} \mathbf{V}^T \rangle \right) \\ & \text{subject to} \quad \sum_{i=1}^N y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq \frac{1}{2\lambda n}, \quad i = 1, \dots, n, \end{aligned}$$

We use quadratic programming to solve this dual problem and update $\mathbf{U} = \sum_{i=1}^n \alpha_i y_i \mathbf{X}_i \mathbf{V} (\mathbf{V}^T \mathbf{V})^{-1}$. Similar approach is applied to update \mathbf{V} fixing \mathbf{U} . The Algorithm 1 gives the full description.

7 Extension to nonlinear case

We extend linear function class to non-linear class with kernel trick. We enlarge feature space through feature mapping $\mathbf{h} : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^{d_1 \times d'_2}$. Once this mapping fixed, the procedure is the same as before. We fit the linear classifier using pair of input feature and label

Algorithm 1: Linear classification algorithm

Input: $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_n, y_n)$, rank r

Parameter: U, V

Initilize: $U^{(0)}, V^{(0)}$

Do until converges

Update U fixing V :

 Solve $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{X}_i, \mathbf{X}_j \mathbf{V} (\mathbf{V}^T \mathbf{V})^{-1} \mathbf{V}^T \rangle$.

$\mathbf{U} = \sum_{i=1}^n \alpha_i y_i \mathbf{X}_i \mathbf{V} (\mathbf{V}^T \mathbf{V})^{-1}$.

Update V fixing U :

 Solve $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{X}_i, \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{X}_j \rangle$.

$\mathbf{V} = \sum_{i=1}^n \alpha_i y_i \mathbf{X}_i^T \mathbf{U} (\mathbf{U}^T \mathbf{U})^{-1}$.

Output: $\mathbf{B} = \mathbf{U} \mathbf{V}^T$

$\{\mathbf{h}(\mathbf{X}_i), y_i\}_{i=1}^n$. Define a nonlinear low rank kernel in similar way to linear case.

$$\langle \mathbf{X}, \mathbf{X}' \rangle_{\mathbf{P}_r, h} \stackrel{\text{def}}{=} \langle \mathbf{P}_r h(\mathbf{X}), \mathbf{P}_r h(\mathbf{X}') \rangle = \text{trace} [\mathbf{K}(\mathbf{X}, \mathbf{X}') \mathbf{P}_r^T \mathbf{P}_r] \quad \text{for all } \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{d_1 \times d_2},$$

where $\mathbf{K}(\mathbf{X}, \mathbf{X}') \stackrel{\text{def}}{=} h(\mathbf{X}) h^T(\mathbf{X}') \in \mathbb{R}^{d_1 \times d_1}$ denotes the matrix product of mapped features.

The solution function $f(\cdot)$ of (1) on enlarged feature can be written

$$f(\cdot) = \sum_{i=1}^n \alpha_i y_i \langle \mathbf{P}_r h(\mathbf{X}_i), \mathbf{P}_r h(\cdot) \rangle = \sum_{i=1}^n \alpha_i y_i \langle \mathbf{X}_i, \cdot \rangle_{\mathbf{P}_r, h} = \sum_{i=1}^n \alpha_i y_i \text{trace} [\mathbf{K}(\mathbf{X}_i, \cdot) \mathbf{P}_r^T \mathbf{P}_r],$$

which involves feature mapping $h(\mathbf{X})$ only thorough inner products. In fact, we need not specify the the transformation $h(\mathbf{X})$ at all but only requires knowledge of the $\mathbf{K}(\mathbf{X}, \mathbf{X}')$. A sufficient condition and a necessary condition for \mathbf{K} being reasonable appear in Supplement. Three popular choices for \mathbf{K} are

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

One can check detailed description for non-linear case algorithm in Supplement.

8 Theory

9 Conclusion

SUPPLEMENTARY MATERIAL

References

- [1] Amnon Shashua, Yoram Gdalyahu, and Gaby Hayun. Pedestrian detection for driving assistance systems: single-frame classification and system level performance. *IEEE Intelligent Vehicles Symposium, 2004*, pages 1–6, 2004.
- [2] Junhui Wang, Xiaotong Shen, and Yufeng Liu. Probability estimation for large-margin classifiers. *Biometrika*, 95(1):149–167, 2008.
- [3] Hua Zhou and Lexin Li. Regularized matrix regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(2):463–483, 2014.