

Large margin based regression for matrix feature

Abstract

The text of your abstract. 200 or fewer words.

Keywords: 3 to 6 keywords, that do not appear in the title

1 Introduction

2 Methods

2.0 Three main problems

Assume that we are given a set of training data $\{(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_n, y_n)\}$ with data matrix $\mathbf{X}_i \in \mathbb{R}^{d_1 \times d_2}$ and class labels $y_i \in \{+1, -1\}$ drawn according to some unknown probability density function $p_{\mathbf{X}, y}(\mathbf{X}, y)$. Define $p(\mathbf{X}) = \mathbb{E}(y|\mathbf{X})$. We consider three major problems on this setting: classification, level set estimation, and regression estimation.

2.1 *The problem of classification:* We aim to address learning a model with a low error on the training data. One successful approach is large-margin classifier. We consider a large margin classifier that 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 regularization 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 $L(z) = \log(1 + e^{-z})$.

2.2 *The problem of level set estimation:* The π -level set of p given a fixed $\pi \in [0, 1]$ is the set

$$S^* = S(p, \pi) = \{\mathbf{X} \in \mathbb{R}^{d_1 \times d_2} : p(\mathbf{X}) > \pi\}.$$

We consider the problem of estimating the π -level set of the regression function without assuming any certain models but utilizing the large-margin classifier.

2.3 The problem of regression estimation: Based on level set estimation problem, we aim to establish a model-free method for estimating $p(\mathbf{X}) = \mathbb{E}(y|\mathbf{X})$. In our setting, the regression $\mathbb{E}(y|\mathbf{X})$ is equivalent to the conditional probability $\mathbb{P}(y = 1|\mathbf{X})$ because the class label y is binary. This problem is also known as soft classification, since the estimated probability can be used to determine the classification boundary.

When utilizing classical methods based on feature vectors to solve the above problems, we have to transform the feature matrices to vectors. However, this vectorization would destroy the structural information of the data matrices. Moreover, the reshaping matrices to vectors results in very high dimensionality which leads to overfitting problem.

We propose a methodology for those problems. Our method exploits the structural information of the data matrix. We take advantage of low-rank assumption to describe the correlation within a matrix and overcome overfitting problem.

2.1 Classification

It is well known that the hinge loss enjoys sparseness and robustness, which are two desirable properties for a good classifiers. For demonstration, we focus on the hinge loss case in Equation (1). However, our estimation schemes and theorems are applicable to general large-margin classifiers. We propose a linear predictor with low-rank coefficient matrix 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 can extend the linear model to non-linear case in Section 4. We impose low-rank structure on

the coefficient matrix \mathbf{B} ,

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

We see that the condition provides a reasonable trade-off between model complexity and model flexibility. This low-rankness makes distinction from classical classification problem for feature vectors and considers structural information of feature matrices. Based on the considered function class, we present the following formulation from Equation (1):

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

Notice that when the rank of the coefficient matrix is full rank, the optimization problem (2) degenerates to the conventional linear SVM with vectorized feature matrices.

We show that the solution of (2) has the form (one can check Supplement for the details)

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

where $\{\alpha_i\}_{i=1}^n$ are solution (spare) of the dual problem in (2) 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:

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

The solution function $f(\cdot)$ can be written $f(\cdot) = \sum_{i=1}^n \alpha_i y_i \langle \mathbf{X}_i, \cdot \rangle_{\mathbf{P}_r}$, which can be viewed as an element in the reproducing kernel Hilbert space induced by rank- r linear kernels $\{\langle \cdot, \cdot \rangle_{\mathbf{P}} : \mathbf{P} \in \mathbb{R}^{r \times d_1} \text{ is a projection matrix}\}$. Therefore, we estimate the optimal hyperplane that separate the training data the best on projected feature matrices. The classification rule $\hat{G}(\mathbf{X})$ can be written as

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

where $\{\hat{\alpha}_i\}_{i=1}^n$ are estimated coefficients and $\hat{\mathbf{P}}_r$ is an estimated projection matrix. The detailed algorithm for the estimation will appear in Section 3.

2.2 Level set estimation

Since the large-margin classifier showed good performance in classification, one might expect to extract any information about level set $S(p, \pi)$ from the large-margin classifier. Lin et al. (2002) showed that the solution \hat{f} to (1) targets directly at $\text{sign}(p(\mathbf{X}) - \frac{1}{2})$. Moreover, Wang et al. (2008) proved that the minimizer \hat{f}_π to π -weighted loss function is a consistent estimate of $\text{sign}(p(\mathbf{X}) - \pi)$. Therefore, we solve the regularization problem with weighed loss function from (1),

$$\min_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n \omega_\pi(y_i) L(y_i f(\mathbf{X}_i)) + \lambda J(f), \quad (3)$$

where $\omega_\pi(y) = 1 - \pi$ if $y = 1$ and π if $y = -1$. From the solution \hat{f}_π to (3), we estimate the level set $\hat{S}(p, \pi) = \{\mathbf{X} : \mathbb{R}^{d_1 \times d_2} : \text{sign}(\hat{f}_\pi(\mathbf{X})) = 1\}$.

2.3 Regression function estimation

Our proposed method is designed to estimate the regression function $p(\mathbf{X}) \stackrel{\text{def}}{=} \mathbb{E}(y = 1 | \mathbf{X})$ at any \mathbf{X} which does not necessarily belong to the observed training data set. Consider the following two steps of approximation to the target function.

$$\begin{aligned} p(\mathbf{X}) &\stackrel{\text{step1}}{\approx} \sum_{h=1}^H \frac{1}{H} \mathbb{1} \left\{ \mathbf{X} : p(\mathbf{X}) \leq \frac{h}{H} \right\} \\ &\stackrel{\text{step2}}{\approx} \sum_{h=1}^H \frac{1}{H} \mathbb{1} \left\{ \mathbf{X} \in \hat{S}^c(p, h/H) \right\}, \end{aligned}$$

where $\hat{S}(p, \pi) = \{\mathbf{X} : \mathbb{R}^{d_1 \times d_2} : \text{sign}(\hat{f}_\pi(\mathbf{X})) = 1\}$. Step 1 approximates the target probability by linear combination of step functions and Step 2 uses the fact that estimated level set in Section 2.2 is consistent. The probability estimation scheme can be summarized as follows.

Shceme

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 (3) with $\omega_{\pi_h}(y)$

S.3 Denote the sequence of solutions and estimated level sets

$$\{\hat{f}_h\}_{h=1}^H \text{ and } \{\hat{S}(p, h/H)\}_{h=1}^H.$$

S.4 Estimate the target probability function by

$$\hat{p}(\mathbf{X}) = \sum_{h=1}^H \frac{1}{H} \mathbb{1} \left\{ \mathbf{X} \in \hat{S}^c(p, h/H) \right\}.$$

3 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

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

We optimize Equation (2) 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

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

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.

Algorithm 1: Linear classification algorithm

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

Parameter: U, V

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

Do until converges

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

4 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 $\{\mathbf{h}(\mathbf{X}_i), y_i\}_{i=1}^n$. 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 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.

5 Theory

6 Conclusion

SUPPLEMENTARY MATERIAL

References

- Lin, Y., Y. Lee, and G. Wahba (2002). Support vector machines for classification in nonstandard situations. *Machine learning* 46(1-3), 191–202.
- Wang, J., X. Shen, and Y. Liu (2008). Probability estimation for large-margin classifiers. *Biometrika* 95(1), 149–167.