

Questions and my answers

Chanwoo Lee, June 3, 2020

algorithmic accuracy vs. statistical accuracy.

1. Why do you choose SVM instead of other non-parametric estimation?:

ans) SVM works the best when $\dim(\mathbf{X}) \gg \#$ of data points case. Especially in matrix feature (or higher dimensional) setting, feature dimension usually is greater than the number of data available. SVM works very well for this setting. One of the reason for this is that algorithm complexity of dual problem is fixed by the number of data points. Another benefit using SVM is that it is robust. Our obtained function from SVM has the form as

vs. K-NN, vs. neural network

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle,$$

where $\alpha \in \mathbb{R}^N$ is solution of dual problem. We can check this solution of dual problem is sparse having many 0 entries. This means that points around boundary only have impact on the classifier but others not, which gives SVM robustness.

independence among $(X, Y) \sim \text{unknown distribution}(\mathbf{x}, \mathbf{y})$; bernoulli distribution on $Y|X \sim \text{Bernoulli}(f(X))$?

2. Why is our method called non-parametric?

ans) First of all, we do not assume any distribution on the model i.e. it is distribution free model which is non-parametric. Secondly, if we check our classifier it has the form (including SMM) as

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle. \quad \text{agree with this.}$$

Therefore, as $N \rightarrow \infty$, we can think that our method has infinity parameter to estimate.

3. What is the practical benefit of using low rank structure in SMM?

ans) If we are using the full rank, we have too many parameters which can result in overfitting problem. Low rank structure prevent us from overfitting the data. Secondly, low rank consider matrix structure of feature matrix \mathbf{X} . Our classifier with this low rank has the form as

$$\hat{f}(\mathbf{X}) = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \hat{\mathbf{P}}_r \mathbf{X}, \hat{\mathbf{P}}_r \mathbf{X}_i \rangle = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \mathbf{X} \hat{\mathbf{P}}_c, \mathbf{X}_i \hat{\mathbf{P}}_c \rangle = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \hat{\mathbf{P}}_r \mathbf{X} \hat{\mathbf{P}}_c, \hat{\mathbf{P}}_r \mathbf{X}_i \hat{\mathbf{P}}_c \rangle. \quad (1)$$

Joint dimension reduction + supervised learning. can be viewed as a supervised PCA.

(1) shows that low rank structure consider column space and row space of the feature matrix.

4. What is intuition of (1)?, does it have something fundamental?

ans) I checked that this really happens from simulation i.e.

$\langle \mathbf{B}, \mathbf{X} \rangle = \langle \mathbf{Q}\mathbf{R}, \mathbf{X} \rangle$, then solve for Q and R

alpha are all the same.

$\langle \mathbf{B}, \mathbf{X} \rangle = \langle \mathbf{P}\mathbf{U}\mathbf{R}, \mathbf{X} \rangle$, then solve for P, U, R (matrix SVD parameterization)

$$\hat{\mathbf{B}} = \sum_{i=1}^N \hat{\alpha}_i y_i \hat{\mathbf{P}}_r \mathbf{X}_i = \sum_{i=1}^N \hat{\alpha}_i y_i \hat{\mathbf{X}}_i \hat{\mathbf{P}}_c = \sum_{i=1}^N \hat{\alpha}_i y_i \hat{\mathbf{P}}_r \mathbf{X}_i \hat{\mathbf{P}}_c. \quad (2)$$

Since this formula is derived from the first order condition in dual problem and linearity plays crucial role for this formula, (2) holds when the function is linear to feature matrix. Intuitively, we are looking for a linear coefficient \mathbf{B} that extract the information from \mathbf{X}_i the best. Since $\{\text{row}(\mathbf{X}_i)_j\}_{j=1}^{d_1}$ and $\{\text{column}(\mathbf{X}_i)_j\}_{j=1}^{d_2}$ contain the same information, the coefficient \mathbf{B} should be the same. But I am not clear what kinds of fundamental theorems are there mathematically.

reason: identifiability due to linearity in both d.r. and s.l

dimension reduction + supervised learning, 4 possible combinations:

linear d.r. + linear supervised learning

nonlinear d.r.+ linear s.l.

linear d.r + nonlinear s.l

nonlinear d.r + nonlinear s.l

1

$\langle \mathbf{P}\mathbf{X}\mathbf{P}, \mathbf{B} \rangle = \langle \mathbf{X}\mathbf{P}, \mathbf{P}\mathbf{B} \rangle = \langle \mathbf{X}, \mathbf{P}\mathbf{B}\mathbf{P} \rangle$

$\langle \mathbf{P} h(\mathbf{X}) \mathbf{P}, \mathbf{B} \rangle$

$\langle g(\mathbf{P}\mathbf{X}\mathbf{P}), \mathbf{B} \rangle$

$g(\langle \mathbf{P}h(\mathbf{X})\mathbf{P}, \mathbf{B} \rangle)$

complexity of SVM: number of sample size required for classification.

5. $K(\mathbf{X}, \mathbf{X}') = h(\mathbf{X})h^T(\mathbf{X}')$ is not really kernel. How do we call this?

ans) asymmetric **kernel?** lifted matrix?

6. **Intuitive way to explain why SVM based probability estimation is better than histogram method.**

convergence rate: $n^{-(2\alpha/(2\alpha+d))}$: d is dimension of the feature, and n is the number of sample points, α = smooth parameter, e.g. for Holder function or Lipschitz function

ans) **The histogram probability estimation method is not working when the dimension of feature is large.** Notice that histogram is making grid on feature map space for estimation. Histogram method cannot estimate probability well in the area of which no data point is observed especially in high dimensional setting. SVM based probability estimation have advantage when dimension is large but have small data size.

reasonable.

7. **how $p(\mathbf{X}) = \mathbb{P}(Y = 1|\mathbf{X})$ affects the accuracy in probability estimation?**

The condition $p(\mathbf{X})$ is monotonic function of $\langle \mathbf{X}, \cdot \rangle$ is enough for consistent estimation. Consider the following probability estimation process

$$\frac{1}{H} \sum_{i=1}^N \mathbb{1} \left\{ \mathbf{X} : \text{sign}(\hat{f}_h(\mathbf{X})) = -1 \right\} \rightarrow \frac{1}{H} \sum_{i=1}^N \mathbb{1} \left\{ \mathbf{X} : p(\mathbf{X}) \leq \frac{h}{H} \right\} \approx p(\mathbf{X})$$

Lebesgue/Riemann

In the last approximation we can get a insight how geometric figure of $p(\mathbf{X})$ impact on our estimation. If $p(\mathbf{X})$ is smooth function that the last approximation works well has accurate estimator. On the other hand $p(\mathbf{X})$ is a wiggly function, the last approximation perform poorly resulting in bad estimation. **As long as $p(\mathbf{X})$ is Lebesgue measurable function, the last approximation is consistent as $H \rightarrow \infty$.** But in the finite approximation case, the geometric features of $p(\mathbf{X})$ do matter in estimation accuracy.

8. **Why whitening step is needed for SDR but not for prob.est?**

ans) We obtain prob.est from solving

Claim: does the prob. estimation theory go through with either formulation (3) or formulation (4)?

$$\text{Vec}(\mathbf{B})^T \text{Vec}(\mathbf{B}) + \frac{\lambda}{n} \sum_{i=1}^n \omega_{\pi}(y_i) \left| 1 - Y_i \hat{f}_n(\mathbf{X}_i; \mathbf{B}) \right|_+, \quad (3)$$

$p(Y=1|\mathbf{X})$

where $\hat{f}_n(\mathbf{X}_i; \mathbf{B}) = \langle \mathbf{B}, \mathbf{X}_i \rangle$.

However, we obtain SDR estimator from solving

we have to use (4). (3) is invalid for SDR

$$\text{Vec}(\mathbf{B}) \hat{\Sigma}_n \text{Vec}(\mathbf{B}) + \frac{\lambda}{n} \sum_{i=1}^n \omega_{\pi}(y_i) \left| 1 - Y_i \hat{f}_n(\mathbf{X}_i; \mathbf{B}) \right|_+, \quad (4)$$

where $\hat{f}_n(\mathbf{X}_i; \mathbf{B}) = \langle \mathbf{B}, \mathbf{X}_i - \bar{\mathbf{X}}_n \rangle$, $\bar{\mathbf{X}}_n = n^{-1} \sum_{i=1}^n \mathbf{X}_i$, and $\hat{\Sigma}_n$ denotes the sample covariance matrix of $\{\text{Vec}(\mathbf{X}_i)\}_{i=1}^n$. Therefore, by whitening step, we simplify (4) to (3) where we can use SMM algorithm.

Probability estimation: (\mathbf{X}, y) . 1. whiten $\mathbf{X} \leftarrow \mathbf{X} \Sigma^{-1/2}$

2. weighted SVM/SMM \rightarrow function estimation.

$p(Y=1|\mathbf{X}) \Sigma^{-1/2}$