

Support Vector Machine review

Chanwoo Lee, March 23, 2020

1 Linear SVM

1.1 Margin and Hard-SVM

Let $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ be a training set, where each $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{-1, 1\}$. The goal is to find a halfspace (\mathbf{w}, b) , such that $y_i = \text{sign}(\langle \mathbf{w}, \mathbf{x}_i \rangle + b)$ for all i if there exists. Alternatively this condition can be written as

$$\forall i \in [m], \quad y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0.$$

Hard-SVM is the learning rule in which we return an hyperplane that separates the training set with the largest possible margin.

Lemma 1. *The distance between a point \mathbf{x} and the hyperplane defined by (\mathbf{w}, b) where $\|\mathbf{w}\| = 1$ is $|\langle \mathbf{w}, \mathbf{x} \rangle + b|$*

From this claim, the Hard-SVM rule is

$$\arg \max_{(\mathbf{w}, b): \|\mathbf{w}\|=1} \min_{i \in [m]} |\langle \mathbf{w}, \mathbf{x}_i \rangle + b| \quad \text{s.t.} \quad \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0.$$

We can write an equivalent problem as follows

$$\arg \max_{(\mathbf{w}, b): \|\mathbf{w}\|=1} \min_{i \in [m]} y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b). \quad (1)$$

The algorithm is as follows

Algorithm 1: Hard-SVM

Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

Solve:

$$(\mathbf{w}_0, b_0) = \arg \min_{(\mathbf{w}, b)} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad (2)$$

Output: $\hat{\mathbf{w}} = \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}, \quad \hat{b} = \frac{b_0}{\|\mathbf{w}_0\|}$

Lemma 2. *Equation (1) and Equation (2) are equivalent.*

In homogenous case where the bias term b set to be zero, Equation (2) becomes

$$\min_{\mathbf{w}} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1. \quad (3)$$

1.1.1 The sample complexity of Hard-SVM

Notice the VC-dimension of halfspace in \mathbb{R}^d is $d + 1$. This VC-dimension implies that the sample complexity of learning halfspaces grows with the dimensionality of the problem. The fundamental theorem of learning tells us that if the number of examples is significantly smaller than d/ϵ , then

no algorithm can learn an ϵ -accurate halfspace. For Hard-SVM, we assume the data is separable with margin γ then the sample complexity is bounded from above by a function of $1/\gamma^2$. To specify this more accurately, we define,

Definition 1. Let \mathcal{D} be a distribution over $\mathbb{R}^d \times \{-1, 1\}$. We say \mathcal{D} is separable with a (γ, ρ) -margin if there exists (\mathbf{w}^*, b^*) such that $\|\mathbf{w}^*\| = 1$ and such that with probability 1 over the choice of $(\mathbf{x}, y) \sim \mathcal{D}$ we have that $y(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*) \geq \gamma$ and $\|\mathbf{x}\| \geq \rho$.

From Definition 1, we have the following sample complexity.

Theorem 1.1 (Sample complexity in Hard-SVM). *Let \mathcal{D} be a distribution over $\mathbb{R}^d \times \{-1, 1\}$ that satisfies the (γ, ρ) -separability with margin assumption. Then with probability of at least $1 - \delta$ over the choice of a training set of size m , the 0-1 error of the output of Hard-SVM is at most*

$$\sqrt{\frac{4(\rho/\gamma)^2}{m}} + \sqrt{\frac{2 \log(2/\delta)}{m}}.$$

1.2 Soft-SVM and norm regularization

Soft-SVM is a relaxation of the Hard-SVM rule which assumes that the training set is linearly separable. The algorithm is follows

Algorithm 2: Soft-SVM

Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

Parameter: $\lambda > 0$

Solve:

$$\begin{aligned} \min_{(\mathbf{w}, b, \xi)} & \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \\ \text{s.t. } & \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \text{and } \xi_i \geq 0 \end{aligned} \tag{4}$$

Output: \mathbf{w}, b

Define the hinge loss:

$$\ell^{\text{hinge}}((\mathbf{w}, b), (\mathbf{x}, y)) = \max\{0, 1 - y(\langle \mathbf{w}, \mathbf{x} \rangle + b)\}.$$

Given (\mathbf{w}, b) and a training set S , the averaged hinge loss on S is denoted by $L_S^{\text{hinge}}((\mathbf{w}, b))$.

Lemma 3. Equation (4) is equivalent to the following Equation (5)

$$\min_{\mathbf{w}, b} \left(\lambda \|\mathbf{w}\|^2 + L_S^{\text{hinge}}((\mathbf{w}, b)) \right). \tag{5}$$

Equation (5) is sometimes better to use when implementing algorithm.

1.2.1 The sample complexity of Soft-SVM

This section the sample complexity of Soft-SVM is for the case of homogenous halfspaces where $b = 0$.

Theorem 1.2. *Let \mathcal{D} be a distribution over $\mathcal{X} \times \{0, 1\}$, where $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\| \leq \rho\}$. Consider running Soft-SVM (Equation (5)) on a training set $S \sim \mathcal{D}^m$ and let $A(S)$ be the solution of Soft-SVM. Then, for every μ ,*

$$\mathbb{E}_{S \sim \mathcal{D}^m}[L_{\mathcal{D}}^{\text{hinge}}(A(S))] \leq L_{\mathcal{D}}^{\text{hinge}}(\mu) + \lambda \|\mu\|^2 + \frac{2\rho^2}{\lambda m}.$$

For every $B > 0$, if we set $\lambda = \sqrt{\frac{2\rho^2}{B^2 m}}$ then,

$$\mathbb{E}_{S \sim \mathcal{D}^m}[L_{\mathcal{D}}^{0.1}(A(S))] \leq \mathbb{E}_{S \sim \mathcal{D}^m}[L_{\mathcal{D}}^{\text{hinge}}(A(S))] \leq \min_{\mathbf{w}: \|\mathbf{w}\| \leq B} L_{\mathcal{D}}^{\text{hinge}}(\mathbf{w}) + \sqrt{\frac{8\rho^2 B^2}{m}}.$$

We therefore see that we can control the sample complexity of learning a half-space as a function of the norm of that halfspace, independently of the Euclidean dimension of the space over which the halfspace is defined.

1.2.2 Implementing Soft-SVM using SGD

The optimization problem of Soft-SVM is

$$\min_{\mathbf{w}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y\langle \mathbf{w}, \mathbf{x}_i \rangle\} \right),$$

We rely on the SGD framework for solving regularized loss minimization problems. We can write the update rule of SGD as

$$\mathbf{w}^{(t+1)} = -\frac{1}{\lambda t} \sum_{j=1}^t \mathbf{v}_j,$$

where \mathbf{v}_j is a subgradient of the loss function at $\mathbf{w}^{(j)}$ on the random example chosen at iteration j . For the hinge loss, given an example (\mathbf{x}, y) , we can choose \mathbf{v}_j to be 0 if $\langle \mathbf{w}^{(j)}, \mathbf{x} \rangle \geq 1$ and $\mathbf{v}_j = -y\mathbf{x}$ otherwise. Denoting $\boldsymbol{\theta} = -\sum_{j \leq t} \mathbf{v}_j$, we have Algorithm 4.

1.3 Optimal conditions and support vectors

The name ‘‘Support Vector Machine’’ stems from the fact that the solution of hard-SVM, \mathbf{w}_0 , is supported by the examples that are exactly at distance $1/\|\mathbf{w}_0\|$ from the separating hyperplane. These vectors are called support vectors.

Theorem 1.3. *Let \mathbf{w}_0 be as defined in equation (3) and let $I = \{i : |\langle \mathbf{w}_0, \mathbf{x}_i \rangle| = 1\}$. Then, there exist coefficients $\alpha_1, \dots, \alpha_m$ such that*

$$\mathbf{w}_0 = \sum_{i \in I} \alpha_i \mathbf{x}_i.$$

The examples $\{\mathbf{x}_i : i \in I\}$ are called support vectors.

Algorithm 3: SGD for solving Soft-SVM

Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ **Parameter:** $\mathbf{T}, \lambda > 0$ **For** $t = 1, \dots, T$ Let $\mathbf{w}^{(t)} = \frac{1}{\lambda t} \boldsymbol{\theta}^{(t)}$ Choose i uniformly at random from $[m]$ If $(y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle < 1)$ Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + y_i \mathbf{x}_i$

Else

 Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)}$ **Output:** $\bar{\mathbf{w}} = \frac{1}{T} \sum_{i=1}^T \mathbf{w}^{(t)}$

1.4 Duality

Many of the properties of SVM have been obtained by considering the dual of Equation (3). We derive the dual of Equation (3). Consider the function

$$g(\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \sum_{i=1}^m \alpha_i (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) = \begin{cases} 0 & \text{if } \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \\ \infty & \text{otherwise} \end{cases}.$$

Therefore, we can rewrite Equation (3) as

$$\begin{aligned} \min_{\mathbf{w}} (\|\mathbf{w}\|^2 + g(\mathbf{w})) &= \min_{\mathbf{w}} \max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \right) \\ &\geq \max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \min_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \right). \end{aligned}$$

By using weak duality (Lemma 4) we have the last inequality. It turns out that in our case, strong duality also holds; namely, the inequality holds with equality. Therefore the dual problem is

$$\max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \min_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) \right). \quad (6)$$

By first order necessary condition, we have

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i.$$

Therefore, Equation (6) can be written as

$$\max_{\alpha \in \mathbb{R}^m: \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right).$$

This property is important when implementing SVM with kernels.

Lemma 4 (Weak Duality). *For any function f , we have*

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) \geq \max_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y}).$$

2 Kernel SVM

2.1 The kernel trick

Given an embedding ψ for some domain space \mathcal{X} into some Hilbert space, define the kernel function $K(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle$. We can think of K as specifying similarity between instances and of the embedding ψ as mapping the domain set \mathcal{X} into a space where these similarities are realized as inner products. We aim to construct kernel SVM algorithm on the basis of the values of the kernel function over pairs of domain points.

Notice all versions of the SVM optimization problem we have derived in the previous chapter are instances of the following general problem:

$$\min_{\mathbf{w}} (f(\langle \mathbf{w}, \psi(\mathbf{x}_1) \rangle, \dots, \langle \mathbf{w}, \psi(\mathbf{x}_m) \rangle) + R(\|\mathbf{w}\|)), \quad (7)$$

where $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is an arbitrary function and $R : \mathbb{R}_+ \rightarrow \mathbb{R}$ is a monotonically nondecreasing function. The following theorem shows that there exists an optimal solution of Equation (7) that lies in the span of $\{\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_m)\}$.

Theorem 2.1 (Representer Theorem). *Assume that ψ is a mapping from \mathcal{X} to a Hilbert space. Then, there exists a vector $\boldsymbol{\alpha} \in \mathbb{R}^m$ such that $\mathbf{w} = \sum_{i=1}^m \alpha_i \psi(\mathbf{x}_i)$ is an optimal solution of Equation (7).*

Based on the representer theorem, we have the following equivalent problem.

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^m} f \left(\sum_{j=1}^m \alpha_j K(\mathbf{x}_j, \mathbf{x}_1), \dots, \sum_{j=1}^m \alpha_j K(\mathbf{x}_j, \mathbf{x}_m) \right) + R \left(\sqrt{\sum_{i,j=1}^m \alpha_i \alpha_j K(\mathbf{x}_j, \mathbf{x}_i)} \right). \quad (8)$$

To solve Equation (8) we solely need to know the value of the $m \times m$ matrix G s.t. $G_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$, which is often called the Gram matrix.

In particular, specifying the preceding to the Soft-SVM problem, we can write the problem given in Equation (5) as

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^m} \left(\lambda \boldsymbol{\alpha}^T G \boldsymbol{\alpha} + \frac{1}{2} \sum_{i=1}^m \max\{0, 1 - y_i(G\boldsymbol{\alpha}_i)\} \right).$$

2.2 Characterizing kernel function

Consider a given function of the form $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. The question that “Is it a valid kernel function?” is solved by the following lemma.

Lemma 5. *A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ implements an inner product in some Hilbert space if and only if it is positive semi-definite; namely, for all $\mathbf{x}_1, \dots, \mathbf{x}_m$, the Gram matrix, $G_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$, is a positive semi-definite matrix.*

2.3 Implementing Soft-SVM with kernels

While we could have designed an algorithm for solving Equation (8), there is an even simpler approach that directly tackles the Soft-SVM optimization problem in the feature space,

$$\min_{\mathbf{w}} \left(\frac{\lambda}{2} \|\omega\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y \langle \mathbf{w}, \psi(\mathbf{x}_i) \rangle\} \right),$$

while only using kernel evaluations. The basic observation is that the vector $\mathbf{w}^{(t)}$ maintained by the SGD procedure

Algorithm 4: SGD for solving Soft-SVM with Kernels

Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

Parameter: $T, \lambda > 0$

Initilize: $\beta^{(0)} = 0$

For $t = 1, \dots, T$

 Let $\alpha^{(t)} = \frac{1}{\lambda t} \beta^{(t)}$

 Choose i uniformly at random from $[m]$

 For all $j \neq i$ set $\beta_j^{(t+1)} = \beta_j^{(t)}$

 If $(y_i \sum_{j=1}^m \alpha_j^{(t)} K(\mathbf{x}_j, \mathbf{x}_i) < 1)$

 Set $\beta^{(t+1)} = \beta^{(t)} + y_i$

 Else

 Set $\beta^{(t+1)} = \beta^{(t)}$

Output: $\bar{\mathbf{w}} = \sum_{i=1}^m \alpha_i \psi(\mathbf{x}_i)$ where $\bar{\alpha} = \frac{1}{T} \sum_{t=1}^T \alpha^{(t)}$

3 Generalization

3.1 Binary class classification

SVM objective function can be expressed more generally as

$$\min_{f \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m [1 - y_i f(\mathbf{x}_i)]_+ + \lambda J(f), \quad (9)$$

where \mathcal{H} is the structured space of functions, and $J(f)$ an appropriate regularizer on that space. For example of \mathcal{H} , we can think $f(x) = b + \sum_{i=1}^m \alpha_i K(\mathbf{x}, \mathbf{x}_i)$. Lin (2002) showed that the SVM solution \hat{f} to Equation (9) targets directly at $\text{sign}[p_1(\mathbf{x}) - \frac{1}{2}]$, where $p_1(\mathbf{x}) = P(Y = 1 | \mathbf{X} = \mathbf{x})$. One can solve the regularization problem based on the weighted hinge loss

$$\min_{f \in \mathcal{H}} \frac{1}{m} \left[(1 - \pi) \sum_{y_i=1} [1 - y_i f(\mathbf{x}_i)]_+ + \pi \sum_{y_i=-1} [1 - y_i f(\mathbf{x}_i)]_+ \right] + \lambda J(f), \quad (10)$$

for $\pi \in (0, 1)$. Wang, Shen, and Liu (2008) showed that the minimizer to Equation (10) is a consistent estimate of $\text{sign}[p_1(\mathbf{x}) - \pi]$. Therefore, one can repeatedly solve Equation (10) using different $0 = \pi_1 < \dots < \pi_{k+1} = 1$ and search \hat{j} such that $\text{sign}[p_1(\mathbf{x}) - \pi_{\hat{j}}] \leq \text{sign}[p_1(\mathbf{x}) - \pi_{\hat{j}+1}]$. The probability can be estimated as $\hat{p}_1(\mathbf{x}) = \frac{1}{2}(\pi_{\hat{j}} + \pi_{\hat{j}+1})$.

3.2 Multiclass classification

Here observation space of y is $\{1, 2, \dots, K\}$. A classifier seeks the function vector $\mathbf{f} = (f_1, \dots, f_K)$, where f_k is a map from the input domain S to \mathbb{R} representing the class k . To ensure uniqueness of the solution, a sum to zero constraint $\sum_{k=1}^K f_k = 0$ is usually employed. For any new input vector \mathbf{x} , its label is estimated via a decision rule $\hat{y} = \arg \max_{k=1, \dots, K} f_k(\mathbf{x})$. Let us define

$$\mathbf{g}(\mathbf{f}(\mathbf{x}), y) = \{f_y(\mathbf{x}) - f_k(\mathbf{x}), k \neq y.\}$$

The quantity $\min \mathbf{g}(\mathbf{f}(\mathbf{x}), y)$ is known as the generalized functional margin and can be reduced to $y f(\mathbf{x})$ in the binary case. For multiclass classification, we have the following regularization problem.

$$\begin{aligned} \min_{\mathbf{f}} \frac{1}{m} \sum_{i=1}^m \ell(\min(\mathbf{g}(\mathbf{f}(\mathbf{x}_i), y_i)) + \lambda \sum_{k=1}^K J(f_k) \\ \text{subject to } \sum_{k=1}^K f_k(\mathbf{x}) = 0. \end{aligned}$$

Define the unit K -cube hyperplane as

$$A_K = \{(\pi_1, \dots, \pi_K) : \sum_{k=1}^K \pi_k = 1, \pi_k \geq 0, k = 1, 2, \dots, K\}.$$

For any given $\boldsymbol{\pi} \in A_K$, one can solve the regularization problem based on the weighted loss

$$\begin{aligned} \min_{\mathbf{f}} \frac{1}{m} \sum_{i=1}^m \pi_{y_i} \ell(\min \mathbf{g}(\mathbf{f}(\mathbf{x}_i), y_i)) + \lambda \sum_{k=1}^K J(f_k) \\ \text{subject to } \sum_{k=1}^K f_k(\mathbf{x}) = 0. \end{aligned} \tag{11}$$

To construct a good probability estimate from the classification rules, we require that the loss function ℓ in Equation (11) is consistent in the following sense.

Definition 2. A functional margin base loss ℓ is called weighted Fisher-consistent for weighted classification problem if the minimizer \mathbf{f}^* of $E[\pi_Y \ell(\min(\mathbf{g}(\mathbf{f}(\mathbf{X}), Y)) | \mathbf{X} = \mathbf{x})]$ satisfies

$$\arg \max_{k=1, \dots, K} f_k^*(\mathbf{x}) = \arg \max_{k=1, \dots, K} \pi_k p_k(\mathbf{x}), \quad \forall \mathbf{x}, \forall \boldsymbol{\pi} \in A_K.$$

For any $\ell(\cdot)$, we define its truncated loss at a location $s \leq 0$ by

$$l_{T_s}(\cdot) = \min(\ell(\cdot), \ell(s)).$$

The following theorem shows that the truncated loss ℓ_{T_s} is weighted Fisher-consistent.

Theorem 3.1. Let $\ell(\cdot)$ be a nonincreasing loss function satisfying $\ell'(0^+) < \ell'(0^-) \leq 0$. Then a sufficient condition for the weighted truncated loss $\pi_y \ell_{T_s}(\min \mathbf{g}(\mathbf{f}(\mathbf{x}), y))$ with $K > 2$ and $s \leq 0$ to be weighted Fisher-consistent for estimating $\arg \max_j \pi_j p_j$ is that the truncation location s satisfies $\sup_{\{u: u \geq -s \geq 0\}} \frac{\dot{\ell}(0) - \dot{\ell}(u)}{\dot{\ell}(s) - \dot{\ell}(0)} \geq K - 1$. This condition is also necessary if $\ell(\cdot)$ is convex.

Example 1. The hinge loss, $\ell(u) = [1 - u]_+ : \text{the condition becomes } s \in [-\frac{1}{K-1}, 0]$.

Example 2. The logistic loss, $\ell(u) = \log(1 + e^u) : \text{the condition becomes } s \in [-\log(2^{K/(k-1)} - 1), 0]$.

4 Review of Robust Model-Free Multiclass Probability Estimation

In this paper, linear SVM is used where $f(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle + \beta_0$.

4.1 Direct scheme for probability recovery

Define the truncated hinge loss as

$$H_{T_s}(u) = \min(H_1(s), H_1(u)),$$

where H_1 is hinge loss and $s = -\frac{1}{K-1}$. Denote $\hat{\mathbf{f}}^\pi$ as the solution of the $\boldsymbol{\pi}$ -weighted truncated hinge loss SVM, obtained from,

$$\begin{aligned} \min_{\mathbf{f}} \frac{1}{n} \sum_{i=1}^n \pi_{y_i} H_{T_s}(\min \mathbf{g}(\mathbf{f}(\mathbf{x}_i), y_i)) + \lambda \sum_{k=1}^K J(f_k) \\ \text{subject to } \sum_{k=1}^K f_k(\mathbf{x}) = 0. \end{aligned}$$

where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \in A_K$. By Theorem 3.1, we have that $\arg \max_k \hat{f}_k^\pi$ converges to $\max_k p_i p_k$ as $n \rightarrow \infty$ and $\lambda \rightarrow 0$. The following proposition gives a key result for estimating the probabilities for each $\mathbf{x} \in S$.

Theorem 4.1. *For any given $\mathbf{x} \in S$ satisfying $\min_k p_k(\mathbf{x}) > 0$, there exists a unique weight vector $\tilde{\boldsymbol{\pi}}(\mathbf{x}) = (\tilde{\pi}_1(\mathbf{x}), \dots, \tilde{\pi}_K(\mathbf{x})) \in A_K$ such that*

$$\tilde{\pi}_1(\mathbf{x}) p_1(\mathbf{x}) = \tilde{\pi}_2(\mathbf{x}) p_2(\mathbf{x}) = \dots = \tilde{\pi}_K(\mathbf{x}) p_K(\mathbf{x}).$$

We can estimate $p_j(\mathbf{x})$ from the estimation of $\tilde{\boldsymbol{\pi}}(\mathbf{x})$.

Theorem 4.2. *For any given $\mathbf{x} \in S$, we assume that its associated border weight is estimated as $\hat{\boldsymbol{\pi}}(\mathbf{x})$. Then its class probability estimated as*

$$\hat{p}_k(\mathbf{x}) = \frac{\hat{\pi}_k(\mathbf{x})^{-1}}{\sum_{i=1}^K \pi_i(\mathbf{x})^{-1}}, \quad k = 1, \dots, K$$

This estimation of probability is consistent.

Theorem 4.3. *For any non-increasing loss function $\ell(\cdot)$ with $\ell'(0) < 0$, if the truncation location s is chosen such that $\sup_{\{u: u \geq -s \geq 0\}} \frac{\ell(0) - \ell(u)}{\ell(s) - \ell(0)} \geq K - 1$. When $\lambda \rightarrow 0$ and the grid size $d_\pi \rightarrow 0$ as $n \rightarrow \infty$, estimate $\hat{p}_k(\mathbf{x})$ based on the truncated loss ℓ_{T_s} is asymptotically consistent, i.e.,*

$$\hat{p}_k(\mathbf{x}) \rightarrow p_k(\mathbf{x}) \text{ for } k = 1, \dots, K \text{ as } n \rightarrow \infty.$$

5 Matrix predictor SVM

Suppose we have train data set such as $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{-1, 1\}$ for all i . We have generalized optimization problem with $f(x) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle + \beta_0$,

$$\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i)) + \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2,$$

where $L(y, f(x))$ is a loss function. We can characterize several loss functions in terms of what they are estimating at the population level. Table 1 shows several loss functions and corresponding minimization functions on population level.

Loss Function	$L(y, f(x))$	Minimization Function
Binomial Deviance	$\log(1 + e^{-yf(x)})$	$f(x) = \log \frac{P(Y=+1 x)}{P(Y=-1 x)}$
SVM Hinge Loss	$[1 - yf(x)]_+$	$f(x) = \text{sign}[P(Y = +1 x) - \frac{1}{2}]$
Squared Error	$[1 - yf(x)]^2$	$f(x) = 2P(Y = +1 x) - 1$

Table 1: The population minimizers for the different loss functions.

Consider we have a test data set such as $(X_1, y_1), \dots, (X_N, y_N)$ where $X_i \in \mathbb{R}^{m \times n}$ and $y_i \in \{-1, 1\}$ for all i . My proposal to train this kind of matrix predictors dataset is to define

$$f(x) = \langle X, B \rangle + \beta_0 \quad \text{where } B \in \mathcal{B} = \{B : \text{rank}(B) \leq r\}.$$

and solve the optimization problem,

$$\min_{\beta_0, B} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i)) + \frac{\lambda}{2} \|B\|_F^2.$$

The rank r can be chosen as $\max_{i=1, \dots, N} \text{rank}(X_i)$. This rank constrained space implies the belief that predictor space has low rank structure. This constraint also functions as reducing the number of parameters to estimate.

6 Discussion and to do list

- Comparison between kernel SVM and new matrix preserving linear SVM (benefit of using linear SVM with matrix predictor?)
- Algorithm derivation in matrix case.
- Extension to multiclass case (similar to Section 4).

References

- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.

Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

Yichao Wu, Hao Helen Zhang, and Yufeng Liu. Robust model-free multiclass probability estimation. *Journal of the American Statistical Association*, 105(489):424–436, 2010.