Jinchao Xu

# Deep Learning Algorithms and Analysis

Summer 2020

## Contributors:

This set of notes are based on contributions from many of graduate students, post-doctoral fellows and other collaborators. Here is a partial list:

Juncai He, Qingguo Hong, Li Jiang.....

# Contents

## 0.1 Binary LR and SVM and their relations

Given a binary linealy separable classification dataset $(x_i, y_i)_{i=1}^{N}$, where $x_i \in \mathbb{R}^d, y_i \in \{-1, +1\}$. We use $A_1, A_2$ to denote the data with label $+1, -1$ respectively. Our goal is to find a $\theta = (w, b)$ where $w \in \mathbb{R}^{1 \times d}, b \in \mathbb{R}$ such that the hyperplane $H_\theta = \{x : wx + b = 0\}$ can separate $A_1, A_2$.

### 0.1.1 Binary SVM

Binary SVM wants to find the classifiable hyperplane which has the biggest distance with $A_1$ and $A_2$.
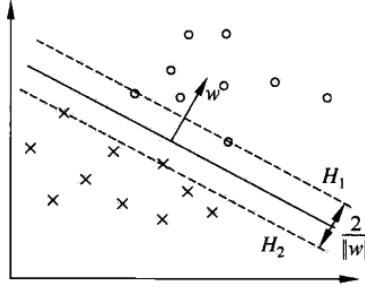
$$(0.1) \qquad \max_{w,b} \frac{\min_i y_i(wx_i + b)}{\|w\|_2}$$

Intuitively, the best separating hyperplane $H$ are only determined by those data points who are closest to $H$. Those data points are called support vector, and this method are called support vector machine.

Without loss of generality, we may restrict the norm of $\|w\|$ to be 1, which leads to a equivalent optimization problem

$$(0.2) \qquad \max_{\|w\|_2=1} \min_i y_i(wx_i + b)$$

Actually, we can prove $\mathrm{argmax}_{\|w\|_2=1} \min_i y_i(wx_i + b)$ is nonempty, but here we just admit this fact and only prove the uniqueness of the solution.

**Lemma 1.** *If $A_1, A_2$ are linearly separable, then*

$$\underset{\|w\|_2=1}{\text{argmax}} \min_i y_i(wx_i + b) \tag{0.3}$$

*is a singleton set.*

*Proof.* Denote $m(w, b) = \min_i y_i(wx_i + b)$. Notice that $m(w, b)$ is a concave homogeneous function w.r.t $w, b$ and $\|\cdot\|_2$ is a strictly convex norm. Suppose there are two solution $(w_1, b_1)$ and $(w_2, b_2)$ such that $w_1 \neq w_2$, take $\overline{w} = \frac{w_1+w_2}{2}, \overline{b} = \frac{b_1+b_2}{2}$, we must have

$$m(\overline{w}, \overline{b}) \geq \frac{m(w_1, b_1) + m(w_2, b_2)}{2} = \underset{\|w\|_2=1}{\max} m(w, b), \tag{0.4}$$

and

$$\|\overline{w}\|_2 < 1. \tag{0.5}$$

So

$$m(\frac{\overline{w}}{\|\overline{w}\|_2}, \frac{\overline{b}}{\|\overline{w}\|_2}) = \frac{m(\overline{w}, \overline{b})}{\|\overline{w}\|_2} > \underset{\|w\|_2=1}{\max} m(w, b), \tag{0.6}$$

which leads to a contradiction. So all the solution must have the same $w$, we denote it as $w^*$. Then if $(w^*, b^*)$ is a solution of problem (0.3), we must have

$$b^* \in \underset{b}{\text{argmax}}\, m(w^*, b) \tag{0.7}$$

Actually,

$$m(w^*, b) = \min\{b + \min_{x \in A_1} w^* x, -b + \min_{x \in A_2}(-w^* x)\}, \tag{0.8}$$

easy to observe that $\text{argmax}_b\, m(w^*, b)$ is a singleton set and

$$b^* = \frac{\min_{x \in A_2}(-w^* x) - \min_{x \in A_1} w^* x}{2}. \tag{0.9}$$

□

Denote

$$\theta^*_{SVM} = (w^*_{SVM}, b^*_{SVM}) = \underset{\|w\|=1}{\text{argmax}} \min_i y_i(wx_i + b). \tag{0.10}$$

4

**Theorem 1.** $w^*_{SVM}$ *must be a linear combination of* $x^T_i, i = 1, 2, \cdots, N$.

*Proof.* Denote

$$(0.11) \qquad\qquad S = \text{span}\{x^T_i\}^N_{i=1}$$

Then we have

$$(0.12) \qquad\qquad \mathbb{R}^{1\times d} = S \oplus^\perp S^\perp$$

So $w^*_{SVM}$ can be uniquely decomposed as $w^*_{SVM} = w^*_S + w^*_{S^\perp}$ where $w_S \in S$ and $w^*_{S^\perp} \in S^\perp$. We will prove that $w^*_{S^\perp} = 0$. Suppose not, we have

$$(0.13) \qquad\qquad \|w^*_S\|_2 < \|w^*\|_2 = 1.$$

Notice that

$$(0.14) \qquad\qquad w^*_{SVM} x_i = w^*_S x_i, \ \forall i = 1, 2, \cdots, N.$$

Thus we have

$$(0.15) \qquad\qquad \min_i y_i(w^*_{SVM} x_i + b^*) = \min_i y_i(w^*_S x_i + b^*)$$

So

$$(0.16) \quad \min_i y_i(w^*_{SVM} x_i + b^*_{SVM}) < \frac{\min_i y_i(w^*_S x_i + b^*_{SVM})}{\|w^*_S\|} = \min_i y_i(\frac{w^*_S}{\|w^*_S\|_2} x_i + \frac{b^*_{SVM}}{w^*_S}),$$

which leads to a contradiction to the definition of $\theta^*_{SVM}$. $\square$

We may rewrite the SVM problem as

$$(0.17) \qquad\qquad \min_{w,b} \|w\|^2,$$

$$(0.18) \qquad\qquad s.t. \ y_i(wx_i + b) \geq 1, \ \forall i.$$

We can simply prove that the solution of (0.20) is $\theta^*_{SVM}$ multiplies a positive scalar. So it still satisfies the representer theorem. Thus we can restrict $w$ to be in the set $S$. Assume that

$$(0.19) \qquad\qquad w = \sum_{i=1}^N \alpha_i x^T_i,$$

Denote $\alpha = (\alpha_1, \cdots, \alpha_N)^T$, and $D \in \mathbb{R}^{N\times N}$ where $D_{ij} = <x_i, x_j>$. We can rewrite the problem (0.20) as

$$(0.20) \qquad\qquad \min_{w,b} \alpha^T D\alpha,$$

$$(0.21) \qquad\qquad s.t. \ y_i(\sum_{j=1}^N <x_j, x_i> \alpha_j + b) \geq 1, \ \forall i.$$

We can see that the whole problem is only determined by the inner product of data points but not the data itself. What we called kernel method is just use a symmetric positive definite kernel function to replace the inner product. Such kernel function can be regarded as a inner product of some feature space.

### 0.1.2 Binary Logistic Regression

For binary logistic regression, our score mapping can be written as $\left( \begin{array}{c} \frac{1}{1+e^{-(wx+b)}} \\ \frac{1}{1+e^{wx+b}} \end{array} \right)$. We can observe that, $(w, b)$ is classifiable if and only if

$$(0.22) \qquad \frac{1}{1 + e^{-y_i(wx+b)}} > \frac{1}{2}, \ \forall i = 1, 2 \cdots, N.$$

So we may consider to maximize following objetive

$$(0.23) \qquad P(\theta) = \prod_{i=1}^{N} \frac{1}{1 + e^{-y_i(wx+b)}},$$

which is equivalent to minimize

$$(0.24) \qquad L(\theta) = -\log P(\theta) = \sum_{i=1}^{N} -\log(1 + e^{-y_i(wx+b)}),$$

**Lemma 2.** *$L(\theta)$ is a strictly convex function without any global minima.*

To let the above problem have a global minima, we may add a $L_2$ regularization term as following

$$(0.25) \qquad \mathcal{L}(\theta, \lambda) = L(\theta) + \lambda \|w\|_2^2 = \sum_{i=1}^{N} -\log(1 + e^{-y_i(wx+b)}) + \lambda \|w\|_2^2,$$

Actually, we can prove $\operatorname{argmin}_{w,b} L(\theta, \lambda)$ is nonempty for $\lambda$ sufficiently small, but here we just admit this fact and only prove the uniqueness of the solution.

**Lemma 3.** *If $A_1, A_2$ are linearly separable, then*

$$(0.26) \qquad \operatorname*{argmin}_{w,b} L(\theta, \lambda)$$

*is a singleton set for $\lambda$ sufficiently small.*

*Proof.* Because $L(\theta)$ is strictly convex w.r.t. $\theta$ and $\|w\|^2$ is convex w.r.t. $\theta$, so $\mathcal{L}(\theta, \lambda) = L(\theta) + \lambda \|w\|_2^2$ is stricly convex w.r.t. $\theta$, which implies our result directly.
□

For $\lambda$ sufficiently small, denote

$$(0.27) \qquad \theta_{LR}(\lambda) = (w_{LR}(\lambda), b_{LR}(\lambda)) = \operatorname*{argmin}_{w,b} L(\theta, \lambda).$$

**Theorem 2.** *If $A_1, A_2$ are linearly separable, then $\frac{\theta_{LR}(\lambda)}{\|w_{LR}(\lambda)\|}$ converge to $\theta^*_{SVM}$ as $\lambda \to 0$, i.e.*

$$(0.28) \qquad \theta^*_{SVM} = \lim_{\lambda \to 0} \frac{\theta_{LR}(\lambda)}{\|w_{LR}(\lambda)\|}.$$