

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

<b>1</b>	<b>Logistic Regression</b> .....	5
1.1	Introduction to logistic regression .....	5
1.1.1	Plain logistic regression .....	5
1.1.2	Regularized logistic regression .....	8



## Logistic Regression

### 1.1 Introduction to logistic regression

#### 1.1.1 Plain logistic regression

We first introduce the next definition of the set of linearly classifiable weights.

**Definition 1 (the set of linearly classifiable weights).** Assume that we are given  $k$  linearly separable sets  $A_1, A_2, \dots, A_k \in \mathbb{R}^d$ , we define the set of classifiable weights as

$$(1.1) \quad \Theta = \{\theta = (W, b) : w_i x + b_i > w_j x + b_j, \forall x \in A_i, j \neq i, i = 1, \dots, k\}$$

which means those  $(W, b)$  can separate  $A_1, A_2, \dots, A_k$  absolutely correctly.

Our linearly separable assumption implies that  $\Theta \neq \emptyset$ . Now we know the existence of linearly classifiable weights. But how can we find one element in  $\Theta$ ?

**Definition 2 (soft-max).** Given parameter  $\theta = (W, b)$ , a soft-max mapping  $p : \mathbb{R}^d \rightarrow \mathbb{R}^k$  is a mapping with the following formulation

$$(1.2) \quad p(x; \theta) = \frac{e^{Wx+b}}{e^{Wx+b} \cdot \mathbf{1}} = \frac{1}{\sum_{i=1}^k e^{w_i x + b_i}} \begin{pmatrix} e^{w_1 x + b_1} \\ e^{w_2 x + b_2} \\ \vdots \\ e^{w_k x + b_k} \end{pmatrix} = \begin{pmatrix} p_1(x; \theta) \\ p_2(x; \theta) \\ \vdots \\ p_k(x; \theta) \end{pmatrix}$$

$$\text{where } e^{Wx+b} = \begin{pmatrix} e^{w_1 x + b_1} \\ e^{w_2 x + b_2} \\ \vdots \\ e^{w_k x + b_k} \end{pmatrix}, \mathbf{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^k, \text{ and the } i\text{-th component}$$

$$(1.3) \quad p_i(x; \theta) = \frac{e^{w_i x + b_i}}{\sum_{i=1}^k e^{w_i x + b_i}}.$$

The soft-max mapping have several important properties.

1.  $0 < p_i(x; \theta) < 1, \sum_i p_i(x; \theta) = 1$ .
2.  $p_i(x; \theta) > p_j(x; \theta) \Leftrightarrow w_i x + b_i > w_j x + b_j$ . This implies that the linearly classifiable weights have an equivalent description as

$$(1.4) \quad \Theta = \left\{ \theta : p_i(x; \theta) > p_j(x; \theta), \forall x \in A_i, j \neq i, i = 1, \dots, k \right\}$$

3. We usually use the max-out method to do classification. For a given data point  $x$ , we first use a soft-max mapping to map it to  $p(x; \theta)$ , then we attached  $x$  to the class  $i = \arg \max_j p_i(x; \theta)$ .

*Remark 1.* The first property implies that  $p(x; \theta)$  can be regarded as a probability distribution of data points which means given  $x \in \mathbb{R}^d$ , we have  $x \in A_i$  with probability  $p_i(x; \theta)$ ,  $i = 1, \dots, k$ .

The last properties means we pick the label  $i$  as the class of  $x$  such that  $x \in A_i$  has the biggest probability  $p_i(x; \theta)$ .

More detailed discussion of logistic regression from the probability perspective will be presented in the nearly future.

From the above properties, we can define the next likelihood function to help find elements in  $\Theta$ :

$$(1.5) \quad P(\theta) = \prod_{i=1}^k \prod_{x \in A_i} p_i(x; \theta)$$

where this likelihood function comes from its probabilistic interpretation which we will discuss later. Based on the property that

$$(1.6) \quad p_i(x; \theta) = \max_{1 \leq j \leq k} p_j(x; \theta), \forall x \in A_i,$$

if  $\theta \in \Theta$ . This somehow means that if

$$(1.7) \quad P(\theta) = \prod_{i=1}^k \prod_{x \in A_i} p_i(x; \theta) = \max,$$

if  $\theta \in \Theta$ . Or we say that we may use the next optimization problem

$$(1.8) \quad \max_{\theta} P(\theta).$$

to find an element in  $\Theta$ .

More precisely, let us introduce the next lemmas (properties) of  $P(\theta)$ .

**Lemma 1.** Assume that the sets  $A_1, A_2, \dots, A_k$  are linearly separable. Then we have

$$(1.9) \quad \left\{ \theta : P(\theta) > \frac{1}{2} \right\} \subset \Theta.$$

*Proof.* It suffices to show that if  $\theta \notin \Theta$ , we must have  $P(\theta) \leq \frac{1}{2}$ . For any  $\theta \notin \Theta$ , there must exist an  $i_0$ , an  $x_0 \in A_{i_0}$  and a  $j_0 \neq i_0$  such that

$$(1.10) \quad w_{i_0}x_0 + b_{i_0} \leq w_{j_0}x_0 + b_{j_0}.$$

Then we have

$$(1.11) \quad p_{i_0}(x_0; \theta) \leq \frac{e^{w_{i_0}x_0 + b_{i_0}}}{e^{w_{i_0}x_0 + b_{i_0}} + e^{w_{j_0}x_0 + b_{j_0}}} \leq \frac{1}{2}.$$

Notice that  $p_i(x; \theta) < 1$  for all  $i = 1, \dots, k, x \in A$ . So

$$(1.12) \quad P(\theta) < p_{i_0}(x_0; \theta) \leq \frac{1}{2}.$$

□

**Lemma 2.** If  $A_1, A_2, \dots, A_k$  are linearly separable and  $\theta \in \Theta$ , we have

$$(1.13) \quad \lim_{\alpha \rightarrow +\infty} p_i(x; \alpha\theta) = 1 \Leftrightarrow x \in A_i.$$

*Proof.* We first note that if  $x \in A_i$ ,

$$(1.14) \quad p_i(\theta, x) = \frac{1}{1 + \sum_{j \neq i} e^{\alpha[(w_j x + b_j) - (w_i x + b_i)]}} \rightarrow 1, \quad \text{as } \alpha \rightarrow \infty.$$

On the other hand, if  $x \notin A_i$ ,

$$(1.15) \quad p_i(x; \alpha\theta) = \frac{1}{1 + \sum_{j \neq i} e^{\alpha[(w_j x + b_j) - (w_i x + b_i)]}} \leq \frac{1}{2}.$$

This implies that if  $x \notin A_i$ ,  $\lim_{\alpha \rightarrow \infty} p_i(x; \alpha\theta) \neq 1$  which is equivalent to the proposition that if  $\lim_{\alpha \rightarrow \infty} p_i(x; \alpha\theta) = 1$ , then  $x \in A_i$ . □

**Lemma 3.** If  $A_1, A_2, \dots, A_k$  are linearly separable,

$$(1.16) \quad \Theta = \left\{ \theta : \lim_{\alpha \rightarrow +\infty} P(\alpha\theta) = 1 \right\}.$$

*Proof.* We first note that if  $\theta \in \Theta$ , we have  $\lim_{\alpha \rightarrow +\infty} p_i(x; \alpha\theta) = 1$  for all  $x \in A_i$ . So

$$(1.17) \quad \lim_{\alpha \rightarrow +\infty} H(\alpha\theta) = \lim_{\alpha \rightarrow +\infty} \prod_{i=1}^k \prod_{x \in A_i} p_i(x; \alpha\theta) = \prod_{i=1}^k \prod_{x \in A_i} \lim_{\alpha \rightarrow +\infty} p_i(x; \alpha\theta) = 1.$$

On the other hand, if  $\lim_{\alpha \rightarrow +\infty} P(\alpha\theta) = 1$ , there must exist one  $\alpha_0 > 0$  such that  $P(\alpha_0\theta) > \frac{1}{2}$ . From Lemma 1, we have  $\alpha_0\theta \in \Theta$ , which means  $\theta \in \Theta$ .  $\square$

These properties above imply that if we can obtain a classifiable weight through maximizing  $P(\theta)$ , while lemma 3 tells us that  $P(\theta)$  will not have a minimum actually.

More specifically, we just need to find some  $\theta \in \Theta$  such that

$$(1.18) \quad P(\Theta) > \frac{1}{2} \Leftrightarrow L(\theta) := -\log P(\theta) < \log(2).$$

**Question:** how to find these element?

### 1.1.2 Regularized logistic regression

Here, we start from the regularization term  $e^{-\lambda R(\|\theta\|)}$  with these next properties:

1.  $\lambda > 0$ .
2.  $R(t)$  is a strictly increasing function on  $\mathbb{R}^+$  with  $R(0) = 0$ ,  $\lim_{t \rightarrow +\infty} R(t) = +\infty$ . For example,  $R(t) = t^2$ .
3.  $\|\cdot\|$  is a norm on  $\mathbb{R}^{k \times (d+1)}$ , a commonly used norm is the following Frobenius norm:

$$(1.19) \quad \|\theta\|_F = \sqrt{\sum_{i,j} w_{ij}^2 + \sum_i b_i^2}.$$

Based on this regularization term, we may consider the following regularized likelihood function  $P_\lambda(\theta)$  as

$$(1.20) \quad P_\lambda(\theta) = P(\theta)e^{-\lambda R(\|\theta\|)}.$$

Here, let us define

$$(1.21) \quad \Theta_\lambda = \arg \max_{\theta} P_\lambda(\theta),$$

where we have the next definition of  $\arg \max$

$$(1.22) \quad \arg \max_{\theta} P_\lambda(\theta) = \left\{ \theta : P_\lambda(\theta) = \max_{\theta} P_\lambda(\theta) \right\}.$$

The next lemma show that the maximal set of modified objective is not empty.

**Lemma 4.** Suppose that  $A_1, A_2, \dots, A_k$  are linearly separable, then

1. if  $\lambda = 0$ ,  $\Theta_0 = \emptyset$ ,
2.  $\Theta_\lambda$  must be nonempty for all  $\lambda > 0$ .



*Proof.* Lemma 3 shows the first proposition. For the second proposition, we notice that

1.  $P_\lambda(\mathbf{0}) = \frac{1}{k^N}$ .
2.  $\exists M_\lambda > 0$  such that  $e^{-\lambda R(\|\boldsymbol{\theta}\|)} < \frac{1}{k^N}$  whenever  $\|\boldsymbol{\theta}\| > M_\lambda$  because of the properties of  $R(\|\boldsymbol{\theta}\|)$ .

So a maxima on  $\{\boldsymbol{\theta} : \|\boldsymbol{\theta}\| < M_\lambda\}$  must be a global maxima. Then we can easily obtain the result in the lemma from the boundedness and closeness of  $\{\boldsymbol{\theta} : \|\boldsymbol{\theta}\| < M_\lambda\}$ .  $\square$

Furthermore, we have the next theorem which shows that we can indeed get  $\boldsymbol{\theta}$  by maximizing  $P_\lambda(\boldsymbol{\theta})$ .

**Theorem 1.** *If  $A_1, A_2, \dots, A_k$  are linearly separable,*

$$(1.23) \quad \boldsymbol{\Theta}_\lambda \subset \boldsymbol{\Theta},$$

*when  $\lambda > 0$  and sufficiently small.*

*Proof.* By Lemma 1, we can take  $\boldsymbol{\theta}_0 \in \boldsymbol{\Theta}$  such that  $P(\boldsymbol{\theta}_0) > \frac{3}{4}$ . Then, for any  $\lambda < \frac{\log \frac{3}{2}}{R(\|\boldsymbol{\theta}_0\|)}$ ,  $\boldsymbol{\theta}_\lambda \in \boldsymbol{\Theta}_\lambda$ , we have

$$P(\boldsymbol{\theta}_\lambda) \geq P_\lambda(\boldsymbol{\theta}_\lambda) \geq P_\lambda(\boldsymbol{\theta}_0) = P(\boldsymbol{\theta}_0)e^{-\lambda R(\|\boldsymbol{\theta}_0\|)} > \frac{3}{4} \cdot \frac{2}{3} = \frac{1}{2},$$

which implies that  $\boldsymbol{\theta}_\lambda \in \boldsymbol{\Theta}$ . Thus, for any  $0 < \lambda < \frac{\log \frac{3}{2}}{R(\|\boldsymbol{\theta}_0\|)}$ ,  $\boldsymbol{\Theta}_\lambda \subset \boldsymbol{\Theta}$ .  $\square$

The design of logistic regression is that maximize  $P_\lambda(\boldsymbol{\theta})$  is equivalent to minimize  $-\log P_\lambda(\boldsymbol{\theta})$ , i.e.,

$$(1.24) \quad \max_{\boldsymbol{\theta}} \{P_\lambda(\boldsymbol{\theta})\} \Leftrightarrow \min_{\boldsymbol{\theta}} \{-\log P_\lambda(\boldsymbol{\theta})\},$$

while the second one is more convenient to evaluate the gradient. Meanwhile, we add a regularization term  $R(\boldsymbol{\theta})$  to the objective function which makes the optimization problem has a unique solution.

Mathematically, we can formulate Logistic regression as

$$(1.25) \quad \min_{\boldsymbol{\theta}} L_\lambda(\boldsymbol{\theta}),$$

where

$$(1.26) \quad L_\lambda(\boldsymbol{\theta}) := -\log P_\lambda(\boldsymbol{\theta}) = -\log P(\boldsymbol{\theta}) + \lambda R(\|\boldsymbol{\theta}\|) = L(\boldsymbol{\theta}) + \lambda R(\|\boldsymbol{\theta}\|),$$

with

$$(1.27) \quad L(\boldsymbol{\theta}) = -\sum_{i=1}^k \sum_{x \in A_i} \log p_i(x; \boldsymbol{\theta}).$$

Then we have the next logistic regression algorithm.

## 1.1. INTRODUCTION TO LOGISTIC REGRESSION

---

---

**Algorithm 1** Logistic Regression

---

Given data  $A_1, A_2, \dots, A_k$ , find

$$(1.28) \quad \theta^* = \arg \min_{\theta} L_{\lambda}(\theta),$$

for some sufficient small  $\lambda > 0$ .

---

*Remark 2.* Here

$$(1.29) \quad L(\theta) = -\log P(\theta),$$

is known as the loss function of logistic regression model. The next reasons may show that why  $L(\theta)$  is popular.

1. It is more convenient to take gradient for  $L(\theta)$  than  $P(\theta)$ .
2.  $L(\theta)$  is related the so-called cross-entropy loss function which will be discussed in the next section.
3.  $L(\theta)$  is a convex function which will be discussed later.