

---

# DEEP LEARNING

---

April 17, 2019

Zhihua Zhang (张志华)

Ruohua Shi (史若画)

Deep Learning Courses in Peking University

## Contents

1	Overview	3
1.1	References . . . . .	3
1.2	Machine Learning . . . . .	3
2	Chapter 2 Basics	4
2.1	Optimization Target . . . . .	4
2.2	Basic Models . . . . .	4
2.3	About Derivation . . . . .	7
2.4	Loss Functions Design . . . . .	11
3	Chapter 3 Kernel	15
3.1	Motivation . . . . .	15
3.2	Neural Network Methods . . . . .	16
3.3	Underfit VS Overfit & Bias VS Variance . . . . .	17
3.4	PCA . . . . .	18
4	Neural Network	19
4.1	Six Keys . . . . .	19

# 1 Overview

## 1.1 References

《Numerical Linear Algebra》

《Linear Algebra and Learning from Data》 <https://math.mit.edu/~gs/learningfromdata/>

《Numerical Optimization》 <http://www.bioinfo.org.cn/wangchao/maa/NumericalOptimization>  
<https://github.com/ShiqinHuo/Numerical-Optimization-Books>

## 1.2 Machine Learning

### 1) 频率派/贝叶斯派

贝叶斯派的 MCMC 是目前最常用的采样方法，但是其计算问题仍需解决

### 2) 参数化/非参数

参数化是指参数的数量和数值不依赖于数据，是固定的，如一个高斯分布就两个参数。

非参数指参数的数量和数值依赖于数据，如最近邻算法，数据的个数就是参数的个数。

kernel 方法：这种方法可以将参数与非参数方法转化，“主问题->对偶问题”就是一种“参数->非参数”的方法。

### 3) 生成模型/判别模型

$(x, y) \sim \mathbb{D}$ ,  $x$  和  $y$  服从联合分布。——生成模型

$(x, y), y \sim \mathbb{D}$ ,  $x$  固定，指关注  $y$  的分布。——判别模型

## 2 Chapter 2 Basics

### 2.1 Optimization Target

对于一个分类器 classification  $h$ , 它的泛化误差 The general error 定义为:

$$Pr(h(x) = y(x)) = \mathbb{E}_{x \in \mathbb{D}}(1_{(h(x) \neq y(x))}) \quad (1)$$

对分类问题 Classification: Input:  $x \in R^d$ . Output:  $y \in \{-1, 1\}$

对回归问题 Regression: Input:  $x \in R^d$ . Output:  $y \in R$

已知一个训练数据序列:  $(x_i, y_i), i = 1 \dots n, x_i \in R^d, y_i \in \{-1, 1\}$ , 它的经验误差 the experience error of  $h$  定义为:

$$\hat{R}(h) = \mathbb{E}_{x \in \mathbb{D}}(1_{(h(x) \neq y(x))}) \quad (2)$$

$$\hat{D} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad (3)$$

为了最小化这个经验误差, 自然的得到一个 0-1 loss, 它是高维不连续的, 在数据量较大的时候很难做离散的梯度下降来优化。所以需要找到一个损失函数来代替, 这个损失函数需要满足两点: 1) 凸函数 (Convex); 2) 是一个上界函数 (Upper bound of the original loss)。

### 2.2 Basic Models

SVM logistic regression

目标是  $\max\{0, 1 - y_i w^T x_i\}$ , 多层的感知机就是一个嵌套的过程  $f_0 = f_1(f_2(\dots))$ .

Linear Classification

首先做一个仿射函数 (Affine Function):

$$L_d = \{f_{w,b}, w \in R^d, b \in R\} \quad (4)$$

$$f_{w,b}(x) = \langle w, x \rangle + b = \frac{1}{n} \sum_{i=1}^n w_i x_i + b_i = \langle \bar{w}, \bar{x} \rangle \doteq \langle w, x \rangle \quad (5)$$

其中  $\bar{w} = (b, w) \in R^{d+1}$ ;  $\bar{x} = (1, x^T) \in R^{d+1}$ .  $h$  是在  $L_d$  上的映射  $R \rightarrow y$ . 对回归问题:  $h(x) = f_w(x)$ , 对分类问题:  $h(x) = \text{sign}(f_w(x))$ 。

**【生成模型】** Generation Model

对  $p$  维数据  $(x, y)$ ,  $P(x, y) = (P(y)P(x|y))$ , 设  $y$  的先验是伯努利分布,  $y \in \{0, 1\}$ ,  $x$  服从正态分布。

$$P(y|\pi) = \pi^y (1 - \pi)^{-y}; \pi \in (0, 1) \quad (6)$$

$$P(x|y=0) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right\} \quad (7)$$

$$P(x|y=1) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right\} \quad (8)$$

接下来做极大似然估计  $\min - \sum_{i=1}^n \log P(y_i) P(x_i|y_i)$ . 假设已知求得的所有参数极大似然估计  $\hat{\theta}_{MLE}$ . 那么去求  $y$  属于那一类其实就是求对  $y=0/1$  的后验概率, 如果大于  $1/2$  则属于这一类。

$$\begin{aligned} P(y=1|x, \theta) &= \frac{P(x|y=1)P(y=1|\pi)}{P(x|y=1, \theta)P(y=1|\pi) + P(x|y=0, \theta) \cdot P(y=0|\pi)} \\ &= \frac{\pi \exp\left\{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right\}}{\pi \exp\left\{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right\} + (1 - \pi) \exp\left\{-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right\}} \\ &= \frac{1}{1 + \frac{1-\pi}{\pi} \exp\left\{(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) - (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right\}} \\ &= \frac{1}{1 + \exp\left\{-(\mu_1 - \mu_0)^T \Sigma^{-1} x + \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 + \mu_0) - \log \frac{\pi}{1-\pi}\right\}} \\ &= \frac{1}{1 + \exp\{-w^T x + b\}} \\ & \quad w = \Sigma^{-1}(\mu_1 - \mu_0); b = \log \frac{\pi}{1-\pi} \end{aligned} \quad (9)$$

由上可知,  $P(y=1|x, \theta) > \frac{1}{2}$  与  $\text{sign}(w^T x + b)$  一致。但是参数量减少了很多。

### 【感知算法】Perception Algorithm

#### 1. Batch Perception 批处理 (Mini Batch)

Input:  $\{(x_i, y_i)\}$

Initialize:  $w^{(0)} = (0, \dots, 0)$

for  $t = 0, 1, \dots, T$

if  $(\exists i, s.t. y_i < w^{(t)}, x_i > 0)$   $w^{(t+1)} = w^{(t)} + y_i x_i$

*else*  $output w^{(t)}$

$$\begin{aligned}
 y_t < w^{(t+1)}, x_t > &= y_i < w^{(t)} + y_t x_t, x_t > \\
 &= y_i < w^{(t)}, x_t > + y_t^2 < x_t, x_t > \\
 &= y_t < w^{(t)}, x_t > + \|x_t\|^2
 \end{aligned} \tag{10}$$

## 2. Online Perception Algorithm

Initialize:  $w^{(0)} = (0, \dots, 0)$

*for*  $t = 0, 1, \dots, T$

*Receive*  $x_t$

Predict  $C_r = \text{sign}(< w^{(t)}, x_t >)$

*if* ( $y_t < w^{(t)}, x_t > \geq 0$ )  $w^{(t+1)} = w^{(t)} + \eta y_t x_t$   $\eta \in [0, 1]$

*else*  $w^{(t+1)} = w^{(t)}$

现在我们的损失函数  $1_{[y < w, x \leq 0]}$  (Fig.1) 替换成  $\max\{0, -y < w, x >\}$  ( $\max\{-z\}$ ) (Fig.2)。这时由一个不连续的函数变成连续但不可导的函数。这时的优化目标变成  $\min\{y(w) = \frac{1}{n} \sum_{i=1}^n \max\{0, -y_i < w_i, x_i >\}\}$ 。虽然零点处不可导，但是可以用此梯度来代替，参数更新过程是：

$$w^{(t+1)} = w^{(t)} + \eta \nabla_w y(w^{(t)})$$

$$w^{(t+1)} = w^{(t)} - \eta \nabla_w y(w^{(t)}), \quad \text{if } < w^{(t)}, x_t > \neq 0$$

$$w^{(t+1)} = w^{(t)} + \eta y_t x, \quad \text{if } < w^{(t)}, x_t > = 0 \quad \nabla_w y(w) = -y_t x_t, \quad \text{if } y_t(w^{(t)}) < 0$$

$$\nabla_w y(w) = 0, \quad \text{if } y_t(w^{(t)}) > 0$$

如果将拐点右移，如 Fig.3 蓝色曲线，优化过程转为  $\max\{0, 1 - y < w, x >\}$ ，是 SVM 的优化过程，也就是到 Hinton Loss。

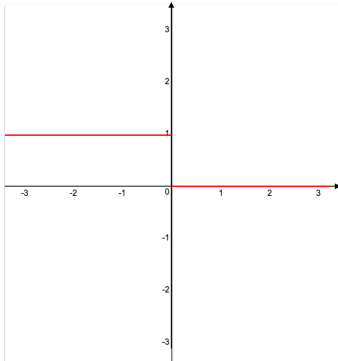


Figure 1: 0-1 Loss

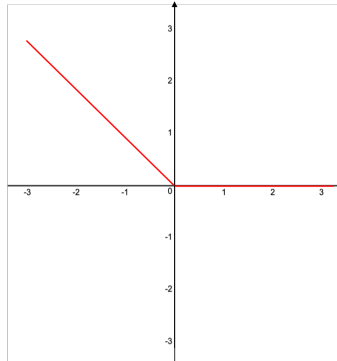


Figure 2: Replace Loss

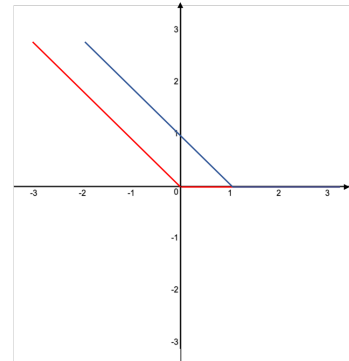


Figure 3: Hinton Loss

## 2.3 About Derivation

### 1. 梯度下降法

定义导数和次导数：

$$\begin{aligned}
 w^{(t+1)} &= w^{(t)} - \eta \nabla g(w^{(t)}) \\
 w^{(t+1)} &= \begin{cases} w^{(t)} - \eta \nabla g(w^{(t)}), & \text{if } \langle w^{(t)}, x_t \rangle \neq 0. \\ w^{(t)} + \eta y_t x_t, & \text{if } \langle w^{(t)}, x_t \rangle = 0. \end{cases} \\
 \eta \nabla g(w) &= \begin{cases} -y x_t, & \text{if } y_t \langle w, x_t \rangle < 0. \\ 0, & \text{if } y_t \langle w, x_t \rangle \geq 0. \end{cases}
 \end{aligned} \tag{11}$$

### 2. 判别模型

对  $y \in \{0, 1\}$ , 有：

$$\begin{aligned}
 P(x, y) &= P(y)P(x|y) \\
 P(y|\pi) &= \pi^y (1 - \pi)^{1-y} \quad \pi \in (0, 1) \\
 P(x|y=0) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right\} \\
 P(x|y=1) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right\}
 \end{aligned} \tag{12}$$

模型的参数为  $\pi, \mu_0, \mu_1, \Sigma$ , 已知数据对为  $\{(x_i, y_i)\}, i = 1, \dots, N$ ,

$$\begin{aligned}
 \prod_{n=1}^N P(x_i, y_i) &= \prod_{n=1}^N P(y_i)P(x_i|y_i) \\
 P(y=1|x, \theta) &= \frac{P(x|y=1)P(y=1|\pi)}{P(x|y=1, \theta)P(y=1|\pi) + P(x|y=0, \theta) \cdot P(y=0|\pi)} \\
 \min - \sum_{n=1}^N \log P(y_i)P(x_i|y_i) & \quad \hat{\theta}_{MLE}
 \end{aligned} \tag{13}$$

$$\begin{aligned}
&= \frac{\pi \exp\{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\}}{\pi \exp\{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\} + (1-\pi) \exp\{-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0)\}} \\
&= \frac{1}{1 + \frac{1-\pi}{\pi} \exp\{(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) - (x-\mu_0)^T \Sigma^{-1}(x-\mu_0)\}} \\
&= \frac{1}{1 + \exp\{-(\mu_1 - \mu_0)^T \Sigma^{-1} x + \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1}(\mu_1 + \mu_0) - \log \frac{\pi}{1-\pi}\}} \\
&= \frac{1}{1 + \exp\{-w^T x + b\}} \Leftrightarrow \text{sign}(w^T x + b)
\end{aligned} \tag{14}$$

## 2. 生成模型

已知数据  $(x, y), x \in R^p, y \in \{0, 1\}, \{-1, 1\}$

$$\begin{aligned}
P(x, y) &= P(y)P(x|y) \\
P(y|\tau) &= \tau^\theta (1-\tau)^{1-\theta} \quad \theta \in (0, 1) \\
P(x|y=0) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\{-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0)\} \\
P(x|y=1) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\}
\end{aligned} \tag{15}$$

模型的参数为  $\Theta = (\tau, \mu_0, \mu_1, \Sigma)$ ,

$$\begin{aligned}
P(y=1|x) &= \frac{P(x|y=1)P(y=1)}{P(x)} \\
P(x) &= \sum_y P(x, y) = P(y=0)P(x|y=0) + P(y=1)P(x|y=1)
\end{aligned} \tag{16}$$

### 1) 高斯混合模型, EM 算法

$$\begin{aligned}
P(y=1|x) &= \frac{1}{1 + \exp(-w^T x - b)} \\
w &\doteq \Sigma^{-1}(\mu_1 - \mu_0) \\
b &\doteq \frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1}(\mu_1 + \mu_0) - \log \frac{\tau}{1-\tau}
\end{aligned} \tag{17}$$

与用  $\text{sign}(f_{w,b} = w^T x + b)$  效果相同。

### 2) 贝叶斯

$$\begin{aligned}
P(x, y, \theta) &= P(\theta)P(x, y|\theta) \\
P(\theta|x, y) &= \frac{P(\theta)P(x, y|\theta)}{P(x, y)}
\end{aligned} \tag{18}$$



3) MAP

$$\begin{aligned} \underset{\theta}{\operatorname{argmax}} P(\theta|x, y) &\Leftrightarrow \underset{\theta}{\operatorname{argmax}} P(\theta)P(x, y|\theta) \\ \log P(\theta|x, y) &\Leftrightarrow \log(P(\theta)P(x, y|\theta)) \end{aligned} \quad (19)$$

4) Maximum Likelihood Estimation (MLE)

$$\begin{aligned} \hat{D} &= \{(x_n, y_n), n = 1, \dots, N\}. \quad x_n \in R^p, \quad y_n \in \{0, 1\} \\ l(\theta, 0) &= \sum_{n=1}^N \log[P(y_n|\tau)P(x_n, y_n, \mu, \Sigma)] \\ &= \sum_{n=1}^N \log P(y_n|\tau) \sum_{n=1}^N \log P(x_n|y_n, \mu, \Sigma) \\ \hat{\tau}_{MLE} &= \underset{\tau}{\operatorname{argmax}} \sum_{n=1}^N y_n \log \tau + (1 - y_n) \log(1 - \tau) = \frac{1}{N} \sum_{n=1}^N y_n \end{aligned} \quad (20)$$

$$\begin{aligned} &\sum_{n=1}^N \log P(x_n|y_n, \mu_0, \mu_1, \Sigma) \\ &= \sum_{n=1}^N \log(P(x_n|y_n = 1, \mu_1, \Sigma)^{y_n} P(x_n|y_n = 0, \mu_0, \Sigma)^{1-y_n}) \\ &= \sum_{n=1}^N [y_n \log(P(x_n|y_n = 1, \mu_1, \Sigma)) + (1 - y_n) \log(P(x_n|y_n = 0, \mu_0, \Sigma))] \\ &= \sum_{n=1}^N [y_n(1 - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1)) \\ &\quad + (1 - y_n)(-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_0)^T \Sigma^{-1} (x_n - \mu_0))] \\ &= -\frac{1}{2} \sum_{n=1}^N [\log |\Sigma| + y_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) + (1 - y_n) (x_n - \mu_0)^T \Sigma^{-1} (x_n - \mu_0)] (*) \end{aligned} \quad (21)$$

对与  $\mu_1$  有关的项  $-\frac{1}{2} \sum_{n=1}^N y_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1)$  求导, 得  $\sum_{n=1}^N y_n \Sigma^{-1} (x_n - \mu_1) = 0 \Rightarrow \sum_{n=1}^N y_n (x_n - \mu_1) = 0$ , 故  $\mu_1$  的估计为:  $\hat{\mu}_1 = \frac{\sum_{n=1}^N y_n x_n}{\sum_{n=1}^N y_n}$ , 同理  $\hat{\mu}_2 = \frac{\sum_{n=1}^N (1 - y_n) x_n}{\sum_{n=1}^N (1 - y_n)}$ .

(\*) 式对  $\Sigma$  求导, 考虑  $\log |\Sigma|$  项,  $\log |\Sigma|: S^{p \times p} \rightarrow R$ ,  $\lim_{t \rightarrow 0} \frac{\log |\Sigma + tA| - \log |\Sigma|}{t} = \langle A, B \rangle = \operatorname{tr}(A^T B)$ . 则导数为  $B$ , 称为 G-导数。

几点性质：

$$\begin{aligned}
 X^T \Sigma^{-1} X &= \text{tr}(\Sigma^{-1} X X^T) & \Sigma \Sigma^{-1} &= I \\
 dX^T \Sigma^{-1} X &= \text{tr}(d\Sigma^{-1} X X^T) & d\Sigma \Sigma^{-1} + \Sigma d\Sigma^{-1} &= 0 \\
 &= \text{tr}(\Sigma^{-1} d\Sigma \Sigma^{-1} X X^T) & d\Sigma^{-1} &= -\Sigma^{-1} d\Sigma \Sigma^{-1} \\
 &= \text{tr}(\Sigma^{-1} X X^T \Sigma^{-1} d\Sigma) \\
 \frac{dX^T \Sigma^{-1} X}{d\Sigma} &= \Sigma^{-1} X X^T \Sigma^{-1}
 \end{aligned} \tag{22}$$

5) The Binomial and Bernoulli

对 Binomial 分布,  $C_n^k \theta^k (1-\theta)^{n-k}$

对 Bernoulli 分布,  $\theta^y (1-\theta)^{1-y}$ .

已知  $x \in R^p$ ,  $y_n \in \{1, \dots, k\}$ ,  $y \in R^k$ , eg.  $k=3$  时,  $y=(1,0,0), y=(0,1,0), y=(0,0,1)$  —One-Hot

6) The Multinomial

$y = (y_1, \dots, y_k)^T$ ,  $Mu(y) = \frac{n!}{y_i! - y_k!} \prod_{k=1}^K \theta_k^{y_k}$ ,  $\sum_{k=1}^K \theta_k = 1$

Multinoulli

$D(y) = \prod_{k=1}^K \tau_k^{y_k}$ ,  $\tau_k \in (0, 1)$

$$\begin{aligned}
 P(x|y_k = 1) &= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right) \\
 P(y_k = 1|x) &= \frac{\exp(w_k^T x)}{\sum_{k=1}^p \exp(w_k^T x)} \quad \text{--- (**) softmax}
 \end{aligned} \tag{23}$$

softmax 即找到  $\{< w_k, x >\}$  的最大值。

7) Gumbel-Max trick

$G - \text{Gumbel}(m)$   $P(G \leq y) = \exp(-\exp(y+m))$ , 若  $G$  是 iid 的  $\text{Gumbel}(0)$  分布, 则

$$(***) = \frac{\exp(\phi_m)}{\sum_{n=1}^N \exp(\phi_m)} = \underset{i}{\operatorname{argmax}} \{G_i + \phi_i\} \quad \text{--- softmax.}$$

8) The Naive Bayes Classifier

$x \in R^p$ ,  $x = (x_1, \dots, x_p)^T$ , 特征值是离散值, 不能用高斯分布。则  $P(x|y) = \sum_{j=1}^p P(y)P(x_j|y)$ ,  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p)$ 。

## 2.4 Loss Functions Design

对于一个 Discriminant Model, 我们看 Logistic Regression, 考虑输入是  $x \in R^p$ ,  $y \in \{0, 1\}$  是类标, 有如下分布:

$$P(y|x) = \mu(x)^y 1 - \mu(x)^{1-y}, \quad \mu \in (0, 1) \quad (24)$$

$$\mu(x) = \frac{1}{1 + \exp(-w^T x)} \quad (25)$$

现在想要训练  $w$ , 假设  $\{(x_n, y_n) : n = 1, \dots, N\}$  是训练数据, 可以通过最小化它的极大似然估计来优化:

$$\min_w L(w) = - \sum_{n=1}^N y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n), \quad \mu_n = \mu(x_n) \quad (26)$$

其中  $\mu_n$  对应真是分布, 这就是常见的交叉熵, 等价于极大似然, 也等价于算经验分布与真实分布的 KL。

对于一个数据矩阵  $X = [x_1, \dots, x_N]^T_{(N \times p)}$ ,  $Y = (y_1, \dots, y_N)^T_{(1 \times N)}$ ,  $\mu = \mu_1, \dots, \mu_N^T$ , 设:

对角阵  $D = \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_N(1 - \mu_N))$

L 的梯度  $\nabla_w L = X^T(\mu - Y)$

海森阵  $H = \frac{\partial L}{\partial w \partial w^T} = X^T D X \geq 0$ , 是一个半正定的, 也说明 L 关于  $w$  是凸的。

(这里插入一点矩阵求导的知识, 可以用极限近似的方法, 转化为求 log 或 trace, 如  $f(X), x \in R^{(p \times q)}, f: R^{(p \times q)} \rightarrow R$ , 可以转化为  $\lim_{t \rightarrow 0} \frac{f(X+tY) - f(X)}{t} = \langle Y, C \rangle$  这样的内积的形式, (G 导数) )

### 1. 一阶近似

对于一阶近似 (梯度迭代),  $w^{(t+1)} = w^{(t)} - \eta_t X^T(\mu^{(t)} - y)$ , 其中  $\eta_t$  是学习率, 满足  $m_{t=0}^\infty \eta_t = \infty$ ,  $\sum_{t=0}^\infty \eta_t^2 < \infty$ 。令  $L(w)$  在  $w^{(t)}$  点做一阶近似:

$$L(w) \approx L(w^{(t)}) + \langle \nabla_w L(w^{(t)}), w - w^{(t)} \rangle + \frac{1}{\eta} \|w - w^{(t)}\|^2 \quad (27)$$

这是一个用线性函数近似的过程, 最后一项是希望控制下一步迭代与这一步相差不太远 (近邻)。有时  $w$  有特殊的性质, 比如当  $w$  非负时, 用平方近邻就不太合适, 这时可以考虑  $+\eta KL(w^{(t)} \| w)$

### 2. 二阶近似

对于二阶近似，又称牛顿法（Newton-Raphson）

$$w^{(t+1)} = w^{(t)} - (X^T D^{(t)} X)^{-1} X^T (\mu^t - y) \quad (28)$$

$$L(w) \approx L(w^{(t)}) + \langle \nabla_w L(w^{(t)}), w - w^{(t)} \rangle + \frac{1}{2} (w - w^{(t)})^T H (w - w^{(t)}) \quad (29)$$

但二阶近似方法需对  $D_{(p \times p)}$  求逆，已知  $\text{rank}(X^T D X) \leq \min\{p, N\}, p \leq N$ ，对于少量数据或不满秩的情况，要满足可逆很难，所以需要加入一个正则化项，其中  $\lambda$  为超参数，交叉验证选取，即：

$$L(w) = \sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)) + \lambda \|w\|^2 \quad (30)$$

用贝叶斯方法来解释：

$$\begin{aligned} \underset{w}{\operatorname{argmin}} L(w) &\Leftrightarrow \underset{w}{\operatorname{argmax}} \sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)) - \lambda \|w\|^2 \\ &\Leftrightarrow \underset{w}{\operatorname{argmax}} \exp\left(\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))\right) \cdot \exp(-\lambda \|w\|^2) \quad (31) \\ &\Leftrightarrow \underset{w}{\operatorname{argmax}} \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{(1-y_n)} \frac{\lambda^{1/2}}{\sqrt{2\pi}^p} \exp(-\lambda \|w\|^2) \end{aligned}$$

也就是等价于最大后验， $p(y|w)p(w) \propto p(w|y)$ 。假设给  $w$  一个均匀分布  $[-m, m]$ ，相当于一个常数，也就是无信息先验（noninformation prior）。

下面考虑其中几个矩阵的运算，

$$\begin{aligned} H &= X^T D X + \lambda I_p \\ \nabla_w L &= X^T (\mu - y) + \lambda w \\ w^{(t+1)} &= w^{(t)} - (X^T D X + \lambda I_p)^{-1} [X^T (\mu^t - y) + \lambda w^{(t)}] \end{aligned} \quad (32)$$

这时  $H$  就是严格大于一个数， $I_p$  是一个  $p \times p$  的对称阵，现在就变成一个强凸的函数了。另外， $(X^T D X + \lambda I_p)$  是一个  $p \times p$  的矩阵，计算复杂度为  $O(p^3)$ 。假设  $p > N$ ，

有  $(Y^T Y + \lambda I_p)^{-1} Y^T = Y^T (Y Y^T + \lambda I_N)^{-1}$ ，通过变换把  $p \times p$  转换为  $N \times N$ 。

$$\begin{aligned}
& (X^T D X + \lambda I_p)^{-1} ((X^T D X + \lambda I_p) w - X^T (\mu - y) - \lambda w) \\
&= (X^T D X + \lambda I_p)^{-1} (X^T D X w - X^T (\mu - y)) \\
&= (X^T D X + \lambda I_p)^{-1} X^T D (X w - D^{-1} (\mu - y)) \\
&= ((D^{1/2} X)^T D^{1/2} + \lambda I_p)^{-1} (D^{1/2} X)^T D^{1/2} \\
&= X^T D^{1/2} (D^{1/2} X X^T D^{1/2} + \lambda I_N)^{-1} D^{1/2} \\
&= X^T (X X^T + \lambda D^{-1})^{-1}
\end{aligned} \tag{33}$$

虽然这是一个简化参数的方法，但参数量还是过大，导致神经网络很少用二阶方法。

### 3. 对于一阶方法的加速

现在考虑一般问题  $\min_w f(w)$ ，使用组合的思想。

① Player's heavy ball

$$\begin{aligned}
w^{(t+1)} &= w^{(t)} - \eta \nabla f(w^{(t)}) + \beta (w^{(t)} - w^{(t-1)}) \\
&= (1 + \beta) w^{(t)} - \beta w^{(t-1)} - \eta \nabla f(w^{(t)})
\end{aligned} \tag{34}$$

其中第一行的  $w^{(t)} - w^{(t-1)}$  就是常说的动量（momentum），第二行的  $(1 + \beta) w^{(t)} - \beta w^{(t-1)}$  就是外插的过程。

② Nesterous “蛙跳”

$$\begin{aligned}
y^{(t+1)} &= w^{(t)} - \eta \nabla f(w^{(t+1)}) \\
w^{(t+1)} &= y^{(t+1)} + \beta (y^{(t+1)} - y^{(t)}) \\
w^{(t+1)} &= (1 + \beta) w^{(t)} - \beta w^{(t-1)} - \eta ((1 + \beta) \nabla f(w^{(t)}) - \beta \nabla f(w^{(t+1)}))
\end{aligned} \tag{35}$$

这里的  $\beta$  可以不  $>0$ ，若  $\beta < 0$ ，则变成内插形式（非凸）。在强凸的情况下可以将时间复杂度从  $O(\frac{1}{\sqrt{T}})$  变为  $O(\frac{1}{T})$ 。

### 4. 隐变量处理

还是考虑开头讲的  $\mu(x) = \frac{1}{1 + \exp(-w^T x)}$ ，这时一个可求导的 Logistic link。它的 Probit link:

$$P(y = 1|x, w) = \mu(x) = \Phi(w^T x) = \int_{-\infty}^{w^T x} \frac{\exp(1 - \frac{1}{2} t^2)}{(2\pi)^{1/2}} dt \tag{36}$$

一般用隐变量处理。设 latent variable  $z = w^T x + \epsilon, \epsilon \in N(0, 1)$

$$P(y=1|z, w) = \begin{cases} 1, & z > 0. \\ 0, & z \leq 0. \end{cases}, \quad \text{and} \quad P(y=1|z, w) = 1 - P(y=0|z, w) \quad (37)$$

这时把  $z$  边界掉, 有:

$$\begin{aligned} P(y=1|z, w) &= P(y=1|z > 0, w) \cdot P(z > 0|w) + P(y=1|z < 0, w) \cdot P(z < 0|w) \\ &= P(z > 0|w) \\ &= Pr\{\epsilon > -w^T X\} \\ &= \int_{-w^T X}^{+\infty} \frac{1}{(2\pi)^{1/2}} \exp(-\frac{1}{2}t^2) dt \\ &= \int_{-\infty}^{w^T X} \frac{1}{(2\pi)^{1/2}} \exp(-\frac{1}{2}t^2) dt \\ &= \Phi(w^T X) \end{aligned} \quad (38)$$

现在回顾一下这个问题, 本质上是解决 0-1 问题, 想找一个凸的上界函数, 已知一个变换:

$$\mu(x)^y (1 - \mu(x))^{1-y} = \left( \frac{1}{1 + \exp(-w^T X)} \right)^y \left( \frac{\exp(-w^T X)}{1 + \exp(-w^T X)} \right)^{1-y} \quad (39)$$

那么观察式子:

$$\begin{aligned} -\log P(y|x) &= y \log(1 + \exp(-w^T X)) + (1-y) w^T X + (1-y) \log(1 + \exp(-w^T X)) \\ &= \log(1 + \exp(-w^T X)) + (1-y) w^T X \\ &= \begin{cases} \log(1 + \exp(w^T X)), & y = 0. \\ \log(1 + \exp(-w^T X)), & y = 1. \end{cases} \end{aligned} \quad (40)$$

若  $y \in \{-1, 1\}$ , 则上式  $= \log(1 + \exp(-yw^T X))$ 。由于它是大于 0 的, 是 0-1 问题的一个上界函数, 已知  $\beta > 0$ , 若:

$$\lim_{\beta \rightarrow \infty} \frac{1}{\beta} \log(1 + \exp(-\beta y w^T X)) = \begin{cases} 0, & y w^T X = 0. \\ 0, & y w^T X > 0. \\ -y w^T X, & y w^T X < 0. \end{cases} \quad (41)$$

这时候量化为 SVM。

若

$$\lim_{\beta \rightarrow \infty} \frac{1}{\beta} \log(1 + \exp(\beta(-yw^T X))) = \begin{cases} 0, & yw^T X = 0. \\ 0, & yw^T X > 0. \\ yw^T X, & yw^T X < 0. \end{cases} \quad (42)$$

这时候退化为 ReLU。

若记  $z = yw^T X$ ，回到隐变量  $z$ ：

$$\max_{w \in (0,1)} wz + \lambda(w \log w + (1-w) \log(1-w)) \quad (43)$$

这就是熵 (entropy)，作用是把一个东西变光滑。

——— 这里推荐两本书：Leo Breiman 《Statistical Modeling》 David A. Freedman 《The Two Cultures Statistical Models Theory and Practice》 ———

## 3 Chapter 3 Kernel

### 3.1 Motivation

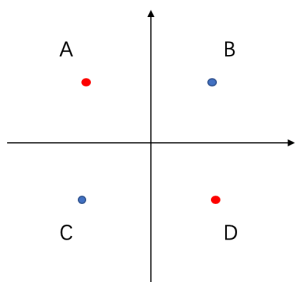


Figure 4: XOR

前面讨论了生成模型（对数据有假设）和判别模型（计算）。都是在平面内线性可分的，考虑 Fig.4 中的情况：若想要把 AD 分为一类，BC 分为一类，显然在这个空间线性是不可分的，这时候就需要引入一个核方法 (Kernel)。这种思想是用最简单的方法达到线性结果。

也就是从输入空间  $\phi$  特征空间做一个非线性映射， $\phi: R^p \rightarrow R^r$ ， $r \gg p$ 。让数据在跟高维的空间线性可分。 $\phi(x)$  作为特征，做一个线性模型。神经网络的做法是直接近似  $\phi$ ，而核函数的方法是直接定义在这个特征空间里的内积形式，也就是只关心空间的内积性质，具体是什么空间，有什么其他性质并不关心。

定义  $\langle \phi(x_1), \phi(x_2) \rangle = k(x_1, x_2)$  是两个数据点  $x_1, x_2$  的核函数。这个核函数的映射过程  $k: R^p \times R^p \rightarrow R$  必须满足两点性质：

1)  $k(x, y) = k(y, x)$  对称性 2)  $k(x, y) \geq 0$  半正定性

常用的两种核函数是：

1) 多项式核：  $k(x, y) = 1 + \langle x, y \rangle^d$

2) 高斯核:  $k(x, y) = \exp(-\frac{\|x-y\|^2}{2\sigma^2})$  ( $\phi(x) \in R^\infty$  , 因为可以拓展到无限维所以常用)

**例 1:** 设  $x, y \in R^2, d = 2, x = (x_1, x_2)^T, y = (y_1, y_2)^T$ 。它的多项式核将 2 维扩展到 6 维:

$$\begin{aligned} k(x, y) &= (1 + x_1 y_1 + x_2 y_2)^2 \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 1 + 2x_1 x_2 y_1 y_2 + 2x_1 y_1 + 2x_2 y_2 \\ &= (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2)(1, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2, y_1^2, y_2^2)^T \\ &= \langle \phi(x), \phi(y) \rangle \end{aligned} \quad (1)$$

对于一个线性模型:  $\sum_{n=1}^N \log(1 + \exp(-y_n w_n^T x_n)) + \lambda \|w\|^2$ ,  
 设有  $\phi(x) \Rightarrow \sum_{n=1}^N \log(1 + \exp(-y_n \langle \phi(x), \phi(x_n) \rangle)) + \lambda \|w\|^2$ , 但是  $\phi(x)$  不知道, 只知道  $k(x_1, x_2)$ 。

**例 2:** 设  $x_1, x_2 \in R$ , 对于一个高斯核

$$\begin{aligned} k(x_1, x_2) &= \exp(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}) \\ &= \exp(-\frac{\|x_1\|^2 + \|x_2\|^2 - 2\langle x_1, x_2 \rangle}{2\sigma^2}) \\ &= \exp(-\frac{\|x_1\|^2}{2\sigma^2}) \exp(-\frac{\|x_2\|^2}{2\sigma^2}) \exp(\frac{\langle x_1, x_2 \rangle}{\sigma^2}) \end{aligned} \quad (2)$$

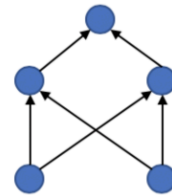
$$\exp(\frac{\langle x_1, x_2 \rangle}{\sigma^2}) = \sum_{k=0}^{\infty} \frac{1}{k!} (\frac{\langle x_1, x_2 \rangle}{\sigma^2})^k = (x_1, x_1^2, \dots)(x_2, x_2^2, \dots)^T \quad (3)$$

这样高斯核就把二维数据扩展到无限维。

### 3.2 Neural Network Methods

现在看神经网络的解决方法。神经网络希望通过最简单的方法让  $X$  真实地表现出来。任意给定一个两层的神经网络, 用 ReLU 做激活函数, 假设输入四个数据分别为  $(1,1), (-1,-1), (-1,1), (1,-1)$ , 设  $h = (X^T w + b)$ , 权重为  $B = [[1, -1], [-1, 1]]$ , 则通过一个 ReLU 的输出为

$$ReLU(X \cdot B) = ReLU \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ -2 & 2 \\ 2 & -2 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 2 \\ 2 & 0 \end{pmatrix}$$



**Figure 5:** Neural Network



当再经过一个  $(1,1)^T$  的滤波，最后输出的结果是  $(0,0,2,2)^T$ ，这样就把四个数据分开了。

### 3.3 Underfit VS Overfit & Bias VS Variance

#### 1. Underfit VS Overfit

给定  $N$  个数据  $\{x_1, \dots, x_N\} \in R^N$ ,  $b_n \in R$ , 定义  $y_n = f(x_n) + \epsilon_n$  是一个预测器，其中  $\epsilon_n$  是噪音，且  $E(\epsilon_n) = 0, D(\epsilon_n) = \sigma^2$ 。预测器未知，通过数据找到一个估计  $F(x)$  尽可能与  $f$  接近，且在测试数据上表现良好。

a).  $F$  does poorly on training sample with large error  $\rightarrow$  Underfitting

b).  $F$  does well on training sample but not well on testing sample  $\rightarrow$  Overfitting

考虑  $\|Xb - y\|^2$ ,  $X$  是一个  $N \times P$  的矩阵，若  $N \gg P$ ，即数据个数远大于参数个数，容易造成 underfit，反之若  $N \ll P$ ，即参数个数远大于数据个数，容易造成 overfit。

#### 2. Bias VS Variance

$Bias = E(F(x) - f(x))$ ，若  $Bias = 0$  则  $F(x)$  是无偏的；若  $Bias \neq 0$ ，但是当  $N$  趋近于无穷大时等于 0，则是渐近无偏。

$$Variance = E(E(x) - E^2(F(x))) = E(F(x))^2 - (E(F(X)))^2$$

【欠拟合】：有偏

【过拟合】：方差很大

$$\begin{aligned} & E((y - F(x))^2) \\ &= \underbrace{E(f(x) - F(x))^2}_{bias^2} + \underbrace{E((F(x)^2 - E(F(x)))^2)}_{variance} + \underbrace{E(y - f(x))^2}_{noise^2} \end{aligned} \quad (4)$$

对于  $F_\theta(x)$ ,  $\theta$  是未知的参数，利用  $\{x_n\}$  估计  $\hat{\theta}_N = g(x_1, \dots, x_N)$  统计量，

$$\begin{aligned} bias(\hat{\theta}_N) &= E(\hat{\theta}_N) - \theta \\ Var(\hat{\theta}_N) &= E(\hat{\theta}_N - E(\hat{\theta}_N))^2 \end{aligned} \quad (5)$$

若  $x_1, \dots, x_N \stackrel{iid}{\sim} bernulii(\theta) \sim \theta^x(1-\theta)^{1-x}$ :

$$\hat{\theta}_N = \frac{1}{N} \sum_{n=1}^N x_n \quad E(\hat{\theta}_N) = \frac{1}{N} \sum_{k=1}^N E(x_n) = \theta \text{ 无偏估计} \quad (6)$$

若  $x_1, \dots, x_N \stackrel{iid}{\sim} \text{Gaussian}(\mu, \sigma^2)$ :

$$\begin{aligned}\hat{\mu}_N &= \frac{1}{N} \sum_{n=1}^N x_n \quad E(\hat{\mu}_N) = \mu \\ \hat{\sigma}_N &= \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2 \quad E(\hat{\sigma}_N) = \frac{N-1}{N} \sigma^2 \text{渐近无偏} \\ \hat{\sigma}_N &= \frac{1}{N-1} \sum_{n=1}^N (x_n - \hat{\mu}_N)^2 \text{无偏}\end{aligned} \quad (7)$$

### 3.4 PCA

设  $X$  是一个  $N \times p$  维的数据，定义  $\phi(A, Z) = \frac{1}{2} \text{tr}((X - ZA^T)(X - ZA^T)^T)$ ， $A$  是  $p \times q$  维矩阵，其中  $q < p$ ， $A$  是列正交的，即  $A^T A = I_q$ 。这时我们称  $A$  是投影矩阵， $Z$  是一个低维表示。我们的目标是找到  $A$  和  $Z$ ，然后  $\min_{Z, A} \phi(A, Z)$ ，也就是最小化  $\|X - ZA^T\|_F^2$ ，若找到  $A_{p \times q}$ ，有  $\min(X - XAA^T)$  最小即可。

定义 Lagrange 函数  $(A, Z) = \frac{1}{2} \text{tr}((X - ZA^T)(X - ZA^T)^T) - \frac{1}{2} \text{tr}(L(A^T A - I_q))$ ， $A = [a_1, \dots, a_q]^T$ ， $L$  是拉格朗日乘子矩阵，是个对称阵，则要找  $\frac{\partial L}{\partial A} = 0, \frac{\partial L}{\partial Z} = 0$  的解。

$$\begin{aligned}dL &= \frac{1}{2} \text{tr}(d(XA^T)(X - ZA^T)^T) + \frac{1}{2} \text{tr}(X - ZA^T)d(X - ZA^T)^T - \frac{1}{2} \text{tr}(L(dA^T A + A^T dA)) \\ &= \dots + \dots - \frac{1}{2} (\text{tr}(LdA^T A) + \text{tr}(LA^T dA)) \\ &= \dots + \dots - \frac{1}{2} (\text{tr}(ALdA^T) + \text{tr}(dA^T AL)) \\ &= \dots + \dots - \frac{1}{2} (\text{tr}(ALdA^T) + \text{tr}(ALdA^T)) \quad (\text{tr}(AB) = \text{tr}(BA)) \\ &= \text{tr}((X - ZA^T)d(X - ZA^T)^T) - \text{tr}(ALdA^T) \\ &= \text{tr}(X - ZA^T)(AdZ^T + dAZ^T) - \text{tr}(ALdA^T) \\ &= \text{tr}((X - ZA^T)AdZ^T) + \text{tr}(X - ZA^T)dAZ^T - \text{tr}(ALdA^T) \\ &= \dots + \text{tr}(ZdA^T(X - ZA^T)^T) - \dots \\ &= \dots + \text{tr}((X - ZA^T)ZdA^T) - \dots \\ &= -\text{tr}((X - ZA^T)A(dZ^T)) - \text{tr}((X^T Z - AZ^T Z + AL)dA^T)\end{aligned} \quad (8)$$

$$\begin{cases} \frac{dL}{dZ} = -(X - 2A^T)A = 0 \\ \frac{dL}{dA} = -(X^T Z - AZ^T Z + AL) = 0 \\ A^T A - I_q = 0 \end{cases} \quad (9)$$

解方程：

$$\begin{aligned}
 (X - ZA^T)A &= 0 \Rightarrow Z = XA \\
 X^Z - AZ^T Z + AL &= 0 \Rightarrow A^T X^Z - A^T AZ^T Z + A^T AL = 0 \\
 L &= Z^Z - A^T X^Z \\
 &\Rightarrow X^T XA = AA^T X^T XA
 \end{aligned} \tag{10}$$

对  $A^T X^T XA$  可以做特征值分解写成  $S\Lambda S^T$ ,  $X^T XA = AS^T$ ,  $X^T XAS = AS\Lambda$ ,  $AS$  是特征向量,  $\Lambda$  是特征值。

回到开始的问题  $\min_{Z,A} \phi(A, Z)$ , 设  $Z = XA$ ,

$$\begin{aligned}
 \min \phi &= \min \text{tr}((X - XA^T)(X - XA^T)^T) \\
 &= \min \text{tr}(X(I - AA^T)X^T) \\
 &= \min \text{tr}(XX^T) - \text{tr}(XAA^T X^T) \\
 &= \dots - \text{tr}(A^T X^T XA) \\
 &= \dots - \text{tr}(\Lambda) \\
 &= \text{tr}(XX^T) - \text{tr}(\Lambda)
 \end{aligned} \tag{11}$$

则  $\Lambda$  为  $XX^T$  的最大前  $q$  个特征值才能让  $\phi$  最小,  $AS$  是对应的特征向量。

## 4 Neural Network

### 4.1 Six Keys

神经网络的基本思想是想要找一个  $F(x)$  来逼近函数  $f(x)$ , 用复合形式来找:  $F = F_L(F_{L-1} \dots 0)$ 。其中有六个关键:

1. Key Compound Operation
2. Key Chain Rule
3. Key Architecture (CNN/RNN)
4. Key Algorithm (SGD)
5. Key Subroutine (Back-Propagation to execute the channel)
6. Key Nonlinear (ReLU)