

# Learning Theory

## bias/variance tradeoff

设  $x$  为预测样本,  $y_0$  为  $x$  在数据集中的标记,  $y$  是  $x$  的真实标记,

$f(x; D)$  为由数据集  $D$  学得模型  $f$  对  $x$  的预测输出,  $\bar{f}(x)$  为模型对不同数据集

预测值的期望, 则有:

1° generalization error.

$$\text{Err}(x) = E[(y - f(x))^2]$$

2° ~~bias~~ variance

$$\left\{ \begin{array}{l} \bar{f}(x) = E_D[f(x; D)] \quad \text{对不同数据集预测值的期望} \\ \text{var}(x) = E_D[(f(x; D) - \bar{f}(x))^2] \end{array} \right.$$

3° noise

$$\varepsilon^2 = E_D[(y_0 - y)^2]$$

4° bias.

$$\text{bias}^2(x) = (f(x) - y)^2$$

$$E[f; D] = E_D[(f(x; D) - y_D)^2]$$

$$= E_D[(f(x; D) - \bar{f}(x))^2 + (\bar{f}(x) - y_D)^2 + 2(f(x; D) - \bar{f}(x))(\bar{f}(x) - y_D)]$$

$$= E_D[(f(x; D) - \bar{f}(x))^2] + E_D[(\bar{f}(x) - y_D)^2] + \underbrace{2 E_D[(f(x; D) - \bar{f}(x))(\bar{f}(x) - y_D)]}_{=0}$$

$$= E_D[(f(x; D) - \bar{f}(x))^2] + E_D[(\bar{f}(x) - y)^2] + E_D[(y - y_D)^2] + \underbrace{2 E_D[(\bar{f}(x) - y)(y - y_D)]}_{=0}$$

$$= \underbrace{E_D[(f(x; D) - \bar{f}(x))^2]}_{\text{Variance}} + \underbrace{(\bar{f}(x) - y)^2}_{\text{bias}^2} + \underbrace{E_D[(y - y_D)^2]}_{\text{noise}}$$

↓  
算法受数据扰动的结果。

↓  
算法的拟合能力

↓  
学习问题本身的难度。

## learning theory

union bound

$$P(A_1 \cup A_2 \dots \cup A_K) \leq P(A_1) + \dots + P(A_K)$$

Hoeffding inequality

假设  $z_1, \dots, z_N$  iid.  $\sim \text{Bernoulli}(p)$ ,  $\hat{p} = \frac{1}{N} \sum_{i=1}^N z_i$

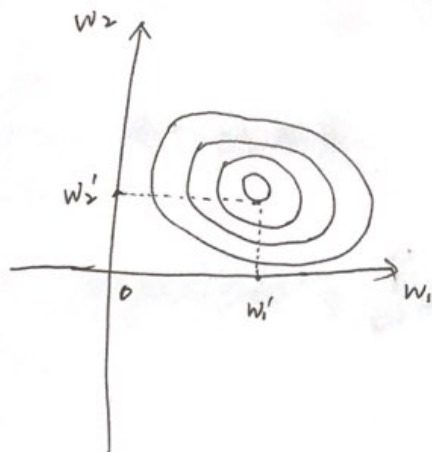
$$P(|p - \hat{p}| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

# Regularization

目的: 缓解过拟合问题, 限制参数过大或过多.

$l_1$  正则化:  $E = E_{in} + \lambda \sum_j |w_j|, \sum_j |w_j| \leq C.$

$l_2$  正则化:  $E = E_{in} + \lambda \sum_j \|w_j\|^2, \sum_j \|w_j\|^2 \leq C.$



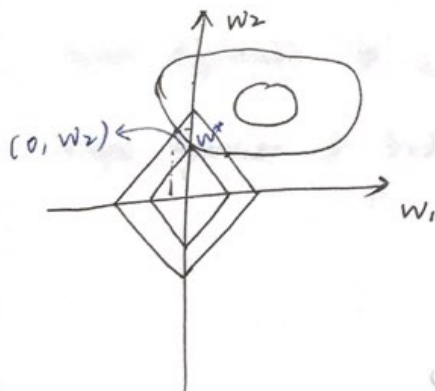
1) 误差函数等高线

此时我们希望找到  $(w_1, w_2)$  使得误差最小, 即最里面的等高线, 但此时得到的是  $(w_1, w_2)$  比较大.

2)  $l_1$  正则化.

此时我们希望最小化误差的同时,  $C$  也越小.

对于同一条等高线来说, 当  $\sum_j |w_j| \leq C$  与其只有一个交点, 即相切时  $(w_1, w_2)$  最小. 且此时很容易点在坐标轴上, 从而使得某些维度的坐标为 0, 从而得到稀疏解.

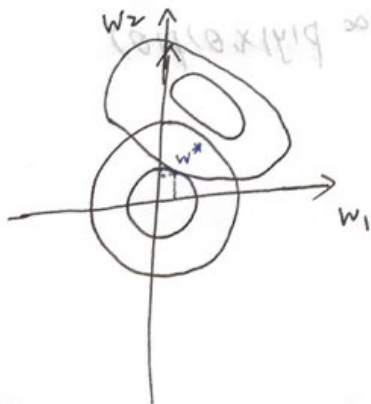


3)  $l_2$  正则化.

同理, 对于同一条等高线来说, 当  $\sum_j w_j^2 \leq C$  与其相切时,

$C$  最小. 且此时  $w^*$  的坐标靠近原点, 不会落在坐标轴上,

所以得到的是  $w^*$  解, 且不为 0.



从贝叶斯(先验/后验)的角度看，和Laplace正则化。

1) Laplace 正则化 ~ Laplace 先验

Laplace 分布:

$$f(x|u, b) = \frac{1}{2b} e^{-\frac{|x-u|}{b}}$$

b 越小, 数据越集中在 u 附近

假设参数  $\theta \sim \text{Laplace}(0, \frac{1}{\lambda})$ , 则:

$$p(\theta_j) = \frac{\lambda}{2} e^{-\lambda |\theta_j|}$$

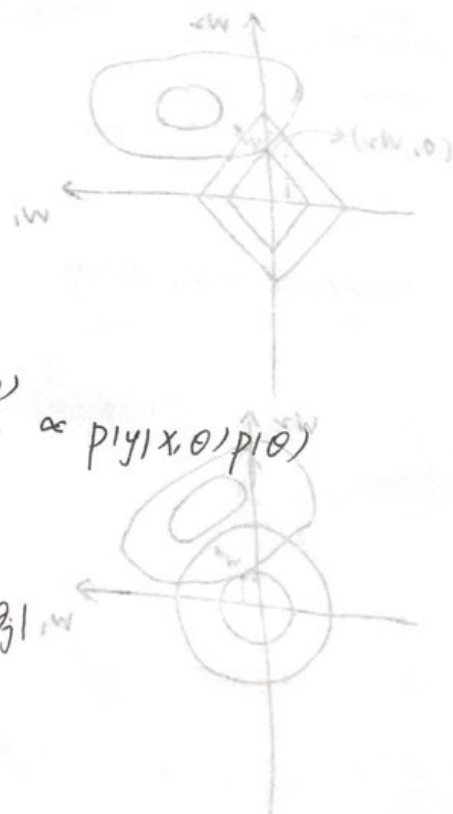
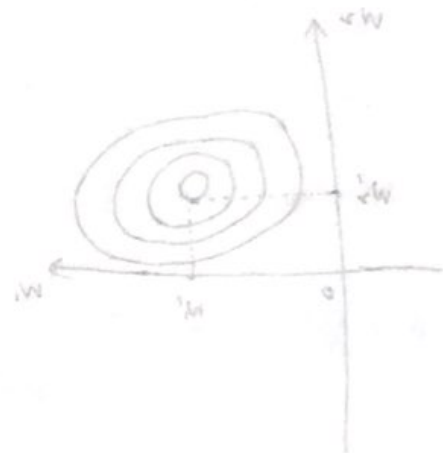
$$\therefore p(\theta) = \prod_{j=1}^M p(\theta_j) = \prod_{j=1}^M \frac{\lambda}{2} e^{-\lambda |\theta_j|}$$

$$\therefore \log p(\theta) = \sum_{j=1}^M \log \frac{\lambda}{2} - \lambda \sum_{j=1}^M |\theta_j|$$

最大后验概率:

$$p(\theta|x, y) = \frac{p(x, y, \theta)}{p(x, y)} = \frac{p(y|x, \theta) p(\theta)}{p(x, y)} \propto p(y|x, \theta) p(\theta)$$

$$\therefore -\log p(\theta|x, y) \approx \sum_{j=1}^N \|y_j - h(x_j)\|^2 + \lambda \sum_{j=1}^M |\theta_j|$$





2)  $l_2$  正则化  $\sim$  高斯先验

假设  $\theta \sim$  高斯分布, 则:

$$p(\theta_j) = \frac{\lambda}{\sqrt{\pi}} e^{-\lambda \theta_j^2}$$

$$\therefore \log p(\theta) = \log \prod_{j=1}^m p(\theta_j) = \sum_{j=1}^m \log \frac{\lambda}{\sqrt{\pi}} - \lambda \sum_{j=1}^m \theta_j^2.$$

$\therefore$  最大后验概率:

$$-\log p(\theta | x, y) \approx \sum_{i=1}^n \|y_i - h_\theta(x_i)\|^2 + \lambda \sum_{j=1}^m \|\theta_j\|^2.$$

$\left\{ \begin{array}{l} \text{Lasso Regression} \Rightarrow l_1 \text{ 正则化 线性回归.} \\ \text{Ridge Regression} \Rightarrow l_2 \text{ 正则化 线性回归.} \end{array} \right.$