

## 1. $l_1, l_2$ 正则化

优化目标:

在原目标(代价)函数中添加惩罚项,

$$\tilde{J}(w; X, y) = J(w; X, y) + \alpha \Omega(w)$$

### 1. 基于约束条件的最优化

模型复杂度用VC维衡量。VC维与系数 $w$ 个数成线性关系, VC维越大, 模型越复杂。

用范数来限制 $\gamma$ 数。

$$\min_w J(w; X, y)$$

$$\text{s.t. } \|w\|_1 \leq C$$

向量中非0的 $\gamma$ 数

但这个问题 NP 问题, 无法求解

改成要求 $w$ 接近于0, 尽量小, 所以用 $l_1$ 范数近似 $l_0$ 范数。

$$\text{即: } \min J(w; X, y) \quad \text{或} \quad \min J(w; X, y)$$

$$\text{s.t. } \|w\|_1 \leq C \quad \text{s.t. } \|w\|_2 \leq C$$

拉格朗日乘子法:

$$L(w, \alpha) = J(w; X, y) + \alpha(\|w\|_1 - C) \quad \alpha > 0$$

$$L(w, \alpha) = J(w; X, y) + \alpha(\|w\|_2 - C)$$

假设 $\alpha$ 最优解是 $\alpha^*$ , 则最优化等价于

$$\min_w J(w; X, y) + \alpha^* \|w\|_1 \quad \text{或}$$

$$\min_w J(w; X, y) + \alpha^* \|w\|_2$$

这与 $\min \tilde{J}(w; X, y)$ 等价

所以是在原 $J(w; X, y)$ 增加约束条件

$$\begin{cases} \|w\|_1 \leq C \\ \|w\|_2 \leq C \end{cases}$$

## 2. 基于最大后验概率估计

将 $w$ 看作随机变量, 也具有某种分布, 从而

$$p(w|X, y) = \frac{p(y|X, w)p(w)}{p(y)} \propto p(y|X, w)p(w)$$

$$\text{MAP} = \underbrace{\log p(y|X, w)}_{\text{似然函数}} + \underbrace{\log p(w)}_{\substack{\text{增加对} w \text{ 的概率分布} \\ \text{的先验假设}}}$$

收集到 $X, y$ 后, 根据 $w$ 在 $\{X, y\}$ 下的后验概率对 $w$ 进行修正。

① 假设 $w$ 的先验分布是均值为0的高斯分布  
 $w_j \sim N(0, \sigma^2)$

$$\begin{aligned} \log p(w) &= \log \prod_j p(w_j) \\ &= \log \prod_j \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(w_j)^2}{2\sigma^2}} \right] \\ &= -\frac{1}{2\sigma^2} \sum_j w_j^2 + C' \end{aligned}$$

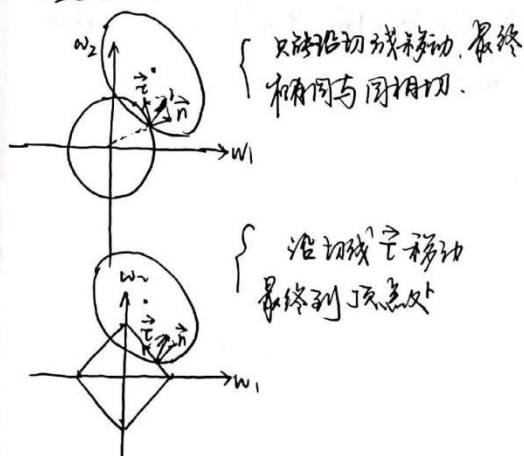
等价于 $l_2$ 正则项

② 假设均值为0, 参数为 $a$ 的拉普拉斯分布  
 $p(w) = \frac{1}{\sqrt{2a}} e^{-\frac{|w|}{a}}$

$$\begin{aligned} \log p(w) &= \log \prod_j \frac{1}{\sqrt{2a}} e^{-\frac{|w_j|}{a}} \\ &= -\frac{1}{a} \sum_j |w_j| + C' \end{aligned}$$

等价于增加 $l_1$ 正则项

3. 直观上



4 理论上

设原目标函数  $J(w)$  最优解为  $w^*$ . 假设  $J(w)$  在  $w^*$  处泰勒展开有

$$\hat{J}(w) = J(w^*) + \frac{1}{2} (w - w^*)^T \underset{\substack{\uparrow \\ \text{Hessian 矩阵}}}{H} (w - w^*)$$

$w^*$  为  $J(w)$  最优解:  $J'(w^*) = 0$

$$\therefore \hat{J}(w) = \frac{1}{2} (w - w^*)^T H (w - w^*)$$

$\therefore \hat{J}(w)$  取得最优值

$$\nabla_w \hat{J}(w) = H (w - w^*) = 0$$

添加 L2 正则项后

$$\begin{aligned} \nabla_w \tilde{J}(w) &= \nabla_w \hat{J}(w) + \nabla_w \Omega(w) \\ &= H (w - w^*) + \alpha w \end{aligned}$$

设最优解为  $\tilde{w}$ , 则有

$$H (\tilde{w} - w^*) + \alpha \tilde{w} = 0$$

$$\therefore \tilde{w} = (H + \alpha I)^{-1} H w^*$$

$H$  是对称阵, 可特征分解  $H = Q \Lambda Q^T$   $Q$  为正交阵

$$\tilde{w} = Q \underbrace{(\Lambda + \alpha I)^{-1} \Lambda}_{\substack{\downarrow \lambda_j \\ \lambda_j + \alpha}} Q^T w^*$$

$\therefore \tilde{w}$  是  $w^*$  在  $H$  的每个特征向量上的分量以  $\frac{\lambda_j}{\lambda_j + \alpha}$  比例被压缩得到.

① 若  $\lambda_j \gg \alpha$ , 则  $w_j^*$  受正则化影响较小

② 若  $\lambda_j \ll \alpha$  则  $w_j^*$  受正则化影响较大.  
收缩到接近于 0 的值. 但  $w_j^* \neq 0 \therefore \tilde{w}_j \neq 0$ . 因而不会稀疏.