

$$\bullet (Xw - y)^T (Xw - y)$$

↓ L_2 正则

$$(Xw - y)^T (Xw - y) + \frac{1}{2} \alpha w^T w$$

↓

$$\begin{cases} \text{无正则} & w = (X^T X)^{-1} X^T y \\ \text{有正则} & w = (X^T X + \alpha I)^{-1} X^T y \end{cases}$$

不同的仅仅是在对角加 α ，这个矩阵的对角项对应于每个输入特征的值，

L_2 正则化能让学习算法感知到具有较高误差的输入 x ，因此与预测目标值十扰误差较小的特征权重将会收缩。



