

What

8. 提前结束

a. 在每次验证集误差有所改善后, 我们保存着模型参数的副本, 当训练集误差经过时, 我们返回这些参数而不是最新的参数。

图 7.1 令 n 为评估间隔的步数

θ 为初始值

$\theta \leftarrow \theta_0$

$i \leftarrow 0$

$j \leftarrow 0$

$u \leftarrow u$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$i^* \leftarrow i$

while $j < p$ do

运行 n 步, 更新 θ .

$i \leftarrow i + n$

$u' \leftarrow \text{ValidErr}(\theta)$

if $u' < u$

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$u \leftarrow u'$

else

$j \leftarrow j + 1$

最佳 θ^* 与 i^*

图 7.2 令 $X^{(train)}$ 和 $y^{(train)}$ 为训练集。

将 $X^{(train)}$ 和 $y^{(train)}$ 分为

$X^{(subtrain)}$, $y^{(subtrain)}$ 与 $X^{(valid)}$, $y^{(valid)}$

从 θ 开始, 训练 7.1 直到 i^* ,

将 θ 再次设为随机值。

在 $X^{(train)}$, $y^{(train)}$ 上训练 j 。

图 7.3

z

运行 7.1 直到 θ

$z \leftarrow J(\theta, X^{(subtrain)}, y^{(subtrain)})$

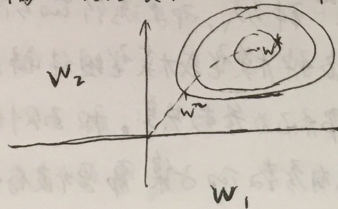
while $J(\theta, X^{(valid)}, y^{(valid)}) > z$ do

在 $X^{(train)}$, $y^{(train)}$ 上训练 n 步。

why 有效

b. 提前终止为何具有正则化效果:

我们可以将 $\frac{1}{\lambda}$ 作为有效容量的度量。假设梯度有界, 限制迭代的次数和训练样本数能够限制从 0 到达的参数空间的大小, 在这个意义上, $\frac{1}{\lambda}$ 的效果就好像是权重衰减系数的倒数。



实线表示负对数似然;

虚线为SGD所经过的轨迹;

提前终止在 w 处停止。

$$\bar{J}(0) = J(w^*) - \frac{1}{2}(w - w^*)^T H (w - w^*)$$

$$w^{(0)} = 0, \quad w^{(T)} = w^{(T-1)} - \frac{1}{\lambda} \nabla_w \bar{J}(w^{(T-1)})$$

$$w^{(T)} - w^* = (I - \frac{1}{\lambda} H) (w^{(T-1)} - w^*)$$

$$w^{(T)} - w^* = (I - \frac{1}{\lambda} Q \Lambda Q^T) (w^{(T-1)} - w^*)$$

假定 $w^{(0)} = 0$, 并且足够小使 $|1 - \frac{1}{\lambda} \lambda_i| < 1$, 经过 T 次后:

$$Q^T w^{(T)} = [I - (I - \frac{1}{\lambda} \Lambda)^T] Q^T w^*$$

$$\text{而 } L^2 \text{ 中有 } Q^T \tilde{w} = [I - (\lambda + \alpha I)^{-1} \alpha] Q^T w^*$$

比较两式有当 $\alpha, \frac{1}{\lambda}, T$ 满足 $(I - \frac{1}{\lambda} \Lambda)^T = (\lambda + \alpha I)^{-1} \alpha$ 时

L^2 正则化与提前终止等价。

$$T \approx \frac{1}{\frac{1}{\lambda} \alpha}, \quad \alpha \approx \frac{1}{T \frac{1}{\lambda}}$$

也就是说, 在这些假设下, 训练次数 T 与 L^2 参数成反比作用。

$T \frac{1}{\lambda}$ 的倒数与权重的衰减系数作用类似。







