

softmax 回归的概率分类 是如何实现 计算损失的?

普通线性回归

函数 $\hat{y} = Xw + b$

损失 L_2 损失 $\frac{1}{2}(\hat{y} - y)^2$ $\frac{1}{2}$ 是为了抵消 2

更新参数 SGD

softmax

↓
用 one-hot 对类别编码.

用交叉熵来衡量损失

↓
衡量两个概率区别

交叉熵

$$H(p, q) = \sum_i -p_i \log(q_i)$$

$$L(y, \hat{y}) = \sum_i -y_i \log(\hat{y}_i)$$

真实值 输出值

由于 one-hot 编码

y_i 只有一个元素为 1

简化

$$L(y, \hat{y}) = -\log \hat{y}_y \quad \text{真实类别的概率}$$

综上

代码中要获取交叉熵结果 反知道 真实类别的概率即可

为什么初始化参数要跟据输入输出来规定矩阵大小?

$$\hat{y} = XW + b$$

↓
权重

X 是一个批量的样本 输入层 $28 \times 28 = 768$ 个特征
输出层 10 个类

batch-size 256

X 大小为 $\frac{256 \times 768}{X}$ 矩阵乘法 $\frac{768 \times 10}{W}$

所以输入输出决定了初始化权重大小。

为什么 batch-size 中的梯度要 loss 求和计算?

1. 因为求梯度消耗资源。

2. loss 求和指求整个 batch-size 的总 loss 对不同参数求偏导。

↓

得到的也是对不同参数的

总 batch-size 梯度

3. 会在 SGD 中更新参数时 除以 batch-size 以使其归一化。