



# 模式识别与机器学习

---

## 4. 16、 Softmax判据的学习



# 待学习的参数

## 学什么

$$\text{Softmax: } \max_i \frac{\exp(\mathbf{w}_i^T \mathbf{x} + w_{0i})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x} + w_{0j})}$$

- 给定训练样本，学习 $K$ 组参数 $\{\mathbf{w}_i, w_{0i}\}_{i=1,2,\dots,K}$ 。

# 训练样本

## 训练样本

- 给定 $K$ 个类别，共 $N$ 个标定过的训练样本：

$$\mathcal{X} = \{(\mathbf{x}_1, \mathbf{t}_1), (\mathbf{x}_2, \mathbf{t}_2) \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$$

- 输出真值采用one-hot方式：

- ✓ 假设 $\mathbf{x}_n \in C_i$ ，则

- ✓ 为了每一项概率分布表述方便， $\mathbf{t}_n$ 也表达为集合形式：

$$\mathbf{t}_n : \{t_n^i\}_{i=1,2,\dots,K}$$

$$\mathbf{t}_n = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \longleftarrow i\text{-th}$$



# 目标函数

## Softmax模型输出的概率分布

- 给定单个样本 $\mathbf{x}$ , one-hot形式的输出标签 $\{l_i\}_{i=1,\dots,K}$ 符合多项分布 (试验次数 $N=1$ )。
- ✓ 分布参数 $p_i$ : 模型输出属于 $C_i$ 类的后验概率 $z_i$ , 其为待学习参数。

$$p(\{l_i\}|\mathbf{x}) = \prod_{i=1}^K z_i^{l_i}, \text{ where } z_i = p(C_i|\mathbf{x}) = \frac{\exp(\mathbf{w}_i^T \mathbf{x} + w_{0i})}{\sum_{j=1}^K \exp(\mathbf{w}_j^T \mathbf{x} + w_{0j})}$$

### 多项分布 (Multinomial)

- 随机变量 $x$ 有 $K$ 个互斥的取值, 每个取值对应的概率为 $p_i$ ,  $\sum_{i=1}^K p_i = 1$ 。
- 该随机变量 $x$ 试验 $N$ 次, 第 $i$ 个值的次数为 $N_i$  ( $\sum_{i=1}^K N_i = N$ ), 则试验次数的概率联合分布为多项分布:

$$\text{Multi}(N_1, N_2, \dots, N_K; \{p_i\}) = N! \prod_{i=1}^K \frac{p_i^{N_i}}{N_i!}$$

- 若 $N=1$ , 则多项分布表达为:

$$\prod_{i=1}^K p_i^{N_i}$$



# 目标函数：最大似然估计

## 似然函数

- 如果参数是最优的，意味着对大部分样本 $(\mathbf{x}_n, \mathbf{t}_n)$ 而言，输出概率 $p(\{\mathbf{t}_n^i\}|\mathbf{x}_n)$ 应该是较大的（无论 $\mathbf{t}_n^i$ 取值是1还是0）。
- 因此，使用**最大似然估计**：针对所有训练样本 $\mathcal{X}$ ，最大化输出标签分布的似然函数，以此求得参数 $\{\mathbf{w}_i, w_{0i}\}_{i=1,2,\dots,K}$ 的最优值。
- 似然函数为所有训练样本输出概率的乘积。
- 所以，目标函数表达为：

$$\max L(\{\mathbf{w}_i, w_{0i}\}|\mathcal{X}) = \max \prod_{n=1}^N p(\{\mathbf{t}_n^i\}|\mathbf{x}_n)$$



# 目标函数：最大似然估计

## 目标函数

$$\max L(\{\mathbf{w}_i, w_{0i}\}|\mathcal{X}) = \max \prod_{n=1}^N p(\{t_n^i\}|\mathbf{x}_n)$$

$$p(\{t_i\}|\mathbf{x}) = \prod_{i=1}^K z_i^{t_i}$$

- 对似然函数求取log再取反，得到目标函数：

$$\begin{aligned} \min J &= -\log L \\ &= -\log \prod_{n=1}^N \prod_{i=1}^K (z_n^i)^{t_n^i} = -\sum_{n=1}^N \sum_{i=1}^K t_n^i \log z_n^i \end{aligned}$$

- 该目标函数其实就是交叉熵的表达式。



# 目标函数：交叉熵解释

## Softmax模型输出的概率分布

- 针对单个样本  $\mathbf{x}_n$  , 可以得到 softmax 预测输出的概率分布  $p(\{l_n^i\}_{i=1,\dots,K}|\mathbf{x}_n)$ 、输出真值的概率分布  $q(\{l_n^i\}_{i=1,\dots,K}|\mathbf{x}_n)$  。
- $p$  和  $q$  都是多项分布 (试验次数  $N=1$ ) , 分布情况如下表:

给定样本 $\mathbf{x}_n$	模型预测输出的概率分布 $p(\{l_n^i\}_{i=1,\dots,K} \mathbf{x}_n)$	输出真值的概率分布 $q(\{l_n^i\}_{i=1,\dots,K} \mathbf{x}_n)$
属于每个类 $C_i$ 的概率	$z_n^i$	$t_n^i$



# 目标函数：交叉熵解释

## 交叉熵

- 给定单个样本 $x_n$ ，希望softmax模型预测输出的概率分布 $p(l_n|x_n)$ 符合输出真值的概率分布 $q(l_n|x_n)$ 。
  - ✓ 使用交叉熵来估计这两种分布的差异程度：

$$H(p, q) = - \sum_{i=1}^K t_n^i \log z_n^i$$

- 给定 $N$ 个训练样本，把每个样本的交叉熵求和，得到目标函数：

$$\min_{w, w_0} - \sum_{n=1}^N \sum_{i=1}^K t_n^i \log z_n^i$$





# 目标函数优化

## 对参数 $\mathbf{w}_k$ 求偏导

$$J(\{\mathbf{w}_i, w_{0i}\}) = - \sum_{n=1}^N \sum_{i=1}^K t_n^i \log z_n^i$$

$$\text{where } z_n^i = \frac{\exp(y_n^i)}{\sum_{j=1}^K \exp(y_n^j)}, \quad y_n^i = \mathbf{w}_i^T \mathbf{x}_n + w_{0i}$$

- 对于任意参数（为区分求和索引 $i$ ，记做 $\mathbf{w}_k$ ），其只与对应的 $y_n^k$ 有关，但所有的 $\{y_n^k\}_{k=1,\dots,K}$ 决定了一个 $z_n^i$ ，所有的 $\{z_n^i\}_{i=1,\dots,K,n=1,\dots,N}$ 决定了 $J$ 。
- 因此，基于链式法则，目标函数对任意参数 $\mathbf{w}_k$ 的偏导表达为：

$$\frac{\partial J}{\partial \mathbf{w}_k} = - \sum_{n=1}^N \sum_{i=1}^K \frac{\partial J}{\partial z_n^i} \frac{\partial z_n^i}{\partial y_n^k} \frac{\partial y_n^k}{\partial \mathbf{w}_k}$$



# 目标函数优化

## 对参数 $w_k$ 求偏导

$$J(\{\mathbf{w}_i, w_{0i}\}) = - \sum_{n=1}^N \sum_{i=1}^K t_n^i \log z_n^i$$

where  $z_n^i = \frac{\exp(y_n^i)}{\sum_{j=1}^K \exp(y_n^j)}$ ,  $y_n^i = \mathbf{w}_i^T \mathbf{x}_n + w_{0i}$

$$\frac{\partial J}{\partial \mathbf{w}_k} = - \sum_{n=1}^N \sum_{i=1}^K \frac{\partial J}{\partial z_n^i} \frac{\partial z_n^i}{\partial y_n^k} \frac{\partial y_n^k}{\partial \mathbf{w}_k}$$

$$= - \sum_{n=1}^N \sum_{i=1}^K t_n^i (\delta_{ik} - z_n^k) \mathbf{x}_n$$

$$= \sum_{n=1}^N (z_n^k - t_n^k) \mathbf{x}_n$$

$$\begin{aligned} \sum_{i=1}^K t_n^i \delta_{ik} &= t_n^k \\ \sum_{i=1}^K t_n^i z_n^k &= z_n^k \end{aligned}$$

$$\frac{\partial J}{\partial z_n^i} = \frac{t_n^i}{z_n^i} \quad \delta_{ik} = \begin{cases} 1, & \text{if } i = k \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{\partial z_n^i}{\partial y_n^k} = z_n^i (\delta_{ik} - z_n^k)$$

$$\frac{\partial y_n^k}{\partial \mathbf{w}_k} = \mathbf{x}_n$$

# 目标函数优化

## 对参数 $w_{0k}$ 求偏导

$$J(\{\mathbf{w}_i, w_{0i}\}) = - \sum_{n=1}^N \sum_{i=1}^K t_n^i \log z_n^i$$

where  $z_n^i = \frac{\exp(y_n^i)}{\sum_{j=1}^K \exp(y_n^j)}$ ,  $y_n^i = \mathbf{w}_i^T \mathbf{x}_n + w_{0i}$

- 目标函数对任意 $w_{0k}$ 的偏导为：

$$\begin{aligned} \frac{\partial J}{\partial w_{0k}} &= - \sum_{n=1}^N \sum_{i=1}^K \frac{\partial J}{\partial z_n^i} \frac{\partial z_n^i}{\partial y_n^k} \frac{\partial y_n^k}{\partial w_{0k}} \\ &= - \sum_{n=1}^N \sum_{i=1}^K t_n^i (\delta_{ik} - z_n^k) = \sum_{n=1}^N (z_n^k - t_n^k) \end{aligned}$$



# 梯度分析

## 梯度分析

- 从梯度公式可以看到，第 $i$ 个线性方程的参数 $w_i$ 和  $w_{0i}$ 的更新不仅依赖于第 $i$ 类的样本 $x_n \in C_i$ ，而且还依赖于所有剩余类的样本 $x_n \in C_j, i \neq j$ 。
- 可见，判别式学习是依赖所有类的训练样本来学习参数。

$$\frac{\partial J}{\partial w_i} = \sum_{n=1}^N (z_n^i - t_n^i) x_n$$

$$\frac{\partial J}{\partial w_{0i}} = \sum_{n=1}^N (z_n^i - t_n^i)$$

# 梯度对比

## Softmax vs. Logistic

$$\frac{\partial J}{\partial w_i} = \sum_{n=1}^N (z_n^i - t_n^i) \mathbf{x}_n$$

$$\frac{\partial J(\mathbf{w}, w_0)}{\partial \mathbf{w}} = \sum_{n=1}^N (z_n - t_n) \mathbf{x}_n$$

$$\frac{\partial J}{\partial w_{0i}} = \sum_{n=1}^N (z_n^i - t_n^i)$$

$$\frac{\partial J}{\partial w_0} = \sum_{n=1}^N (z_n - t_n)$$

- Softmax: 针对每个输出类别分别计算梯度值, 但每个参数的梯度值与所有类别样本都相关。

# 目标函数优化

## 参数更新

- 采用梯度下降法更新所有 $\{\mathbf{w}_i, w_{0i}\}$ :
  - ✓ 设当前时刻为 $k$ , 下一个时刻为 $k + 1$
  - ✓  $\eta$ 为更新步长。

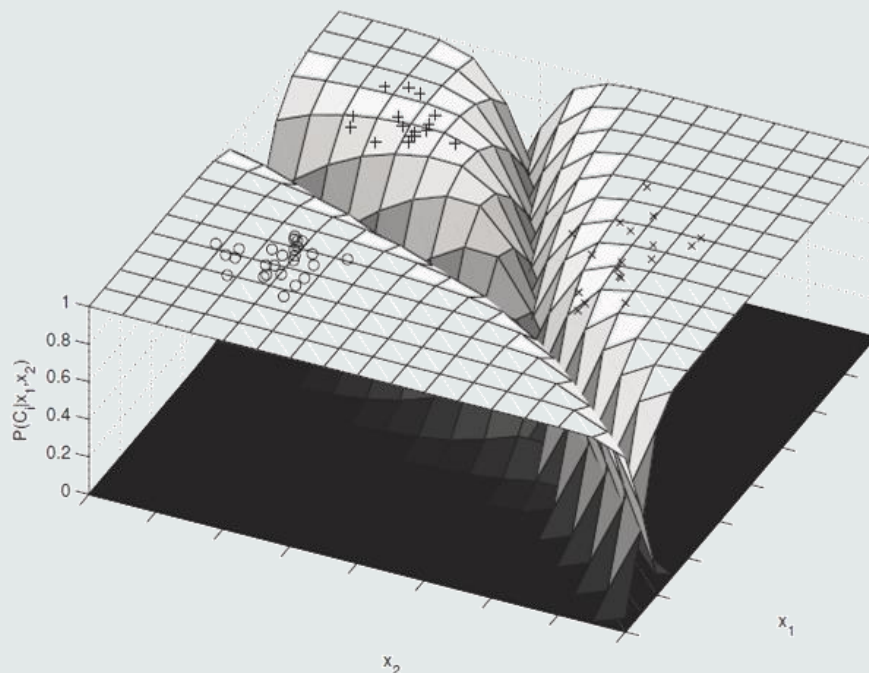
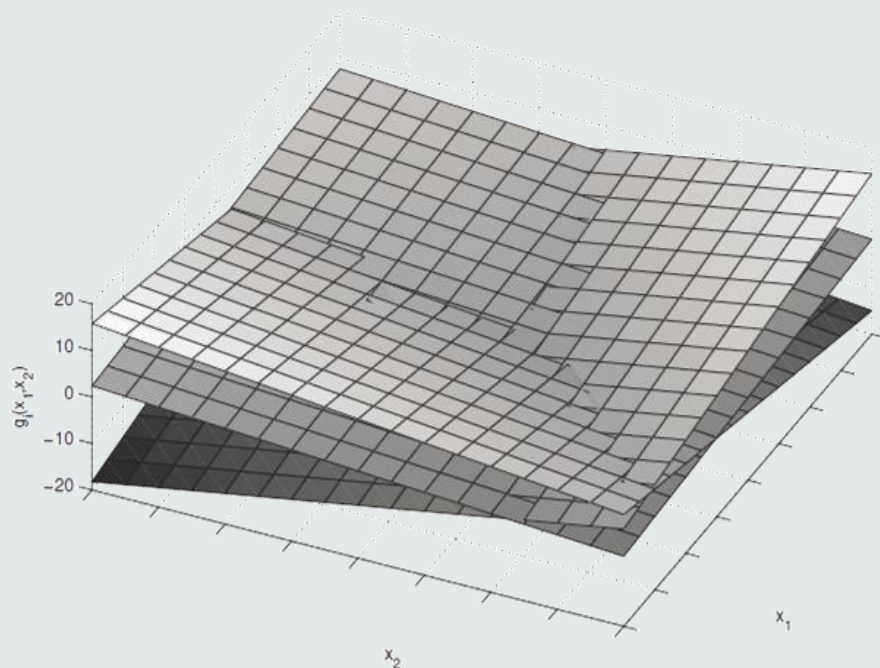
$$\mathbf{w}_i^{(k+1)} = \mathbf{w}_i^{(k)} - \eta^{(k)} \sum_{n=1}^N \left( z_n^i |_{\{\mathbf{w}_i^{(k)}, w_{0i}^{(k)}\}} - t_n^i \right) \mathbf{x}_n$$

$$w_{0i}^{(k+1)} = w_{0i}^{(k)} - \eta^{(k)} \sum_{n=1}^N \left( z_n^i |_{\{\mathbf{w}_i^{(k)}, w_{0i}^{(k)}\}} - t_n^i \right)$$

# Softmax判据学习示例

- $K = 3$
- 左图：3条线性判据  $\mathbf{w}_i^T \mathbf{x} + w_{0i} = 0$
- 右图：3条softmax( $\mathbf{w}_i^T \mathbf{x} + w_{0i}$ )曲线

✓ Softmax输出的非线性形式就是exp函数的非线性形式。





Softmax判据虽然输出非线性，但仍然只能刻画线性分类边界，如何实现非线性分类边界？