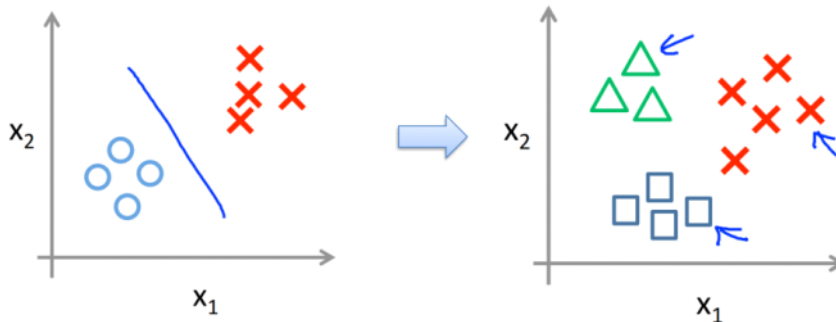


Class 3 Softmax回归

2018年10月17日 17:31

• 简介:

- Softmax回归是一种多分类模型，也叫多分类逻辑回归；
- 在NLP(Natural Language Processing) 中，也被叫做最大熵模型；
- Softmax是一种被广泛使用的分类算法。



• 模型假设:

• 假设

$$p(y = j|x; \theta) = h_j(x) = \frac{e^{\theta_j^T x}}{1 + \sum_{j'=1}^{C-1} e^{\theta_{j'}^T x}}, j = 1, \dots, C-1$$

$$p(y = C|x; \theta) = h_C(x) = \frac{1}{1 + \sum_{j'=1}^{C-1} \exp\{\theta_{j'}^T x\}}$$

• 上面两个式子可合并为简洁形式

$$p(y = j|x; \theta) = h_j(x) = \frac{e^{\theta_j^T x}}{\sum_{j'=1}^C e^{\theta_{j'}^T x}}, j = 1, 2, \dots, C, \text{ where } \theta_C = \vec{0}$$

• 参数

$$\theta_{C \times M}$$

其中, c 为类别数量, M为样本数量

• 最大似然估计:

• 对数似然函数 (多分类问题)

$$\begin{aligned}
 l(\theta) &= \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; \theta) && \text{Softmax Regression} \\
 &= \sum_{i=1}^N \log \prod_{j=1}^C \left(\frac{e^{\theta_j^T x}}{\sum_{j'=1}^C e^{\theta_{j'}^T x}} \right)^{1\{y^{(i)}=j\}} \\
 &= \sum_{i=1}^N \sum_{j=1}^C 1\{y^{(i)} = j\} \log \left(\frac{e^{\theta_j^T x}}{\sum_{j'=1}^C e^{\theta_{j'}^T x}} \right) \\
 &= \sum_{i=1}^N \sum_{j=1}^C 1\{y^{(i)} = j\} \log h_j(x^{(i)})
 \end{aligned}$$

➤ 对于二分类问题，采用逻辑回归

$$l(\theta) = \sum_{i=1}^N y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \quad \text{Logistic Regression}$$

➤ 当Softmax回归模型的类别数为2时，退化为逻辑回归模型。

• 梯度下降优化

• 梯度

$$\frac{\partial l(\theta)}{\partial \theta_k} = \sum_{i=1}^N \boxed{(1\{y^{(i)} = k\} - h_k(x^{(i)})) x^{(i)}} \quad \text{Error} \times \text{Feature}$$

其中：

$$\frac{\partial \log h_j(x)}{\partial \theta_k} = \begin{cases} (1 - h_k(x))x, & j = k \\ -h_k(x)x, & j \neq k \end{cases}$$

$$\begin{aligned}
 \frac{\partial \sum_{j=1}^C 1\{y = j\} \log h_j(x)}{\partial \theta_k} &= \begin{cases} (1 - h_k(x))x, & y = k \\ -h_k(x)x, & y \neq k \end{cases} \\
 &= (1\{y = k\} - h_k(x))x
 \end{aligned}$$

• 梯度下降

$$\begin{aligned}
 \theta_k &:= \theta_k + \alpha \sum_{i=1}^N (1\{y^{(i)} = k\} - h_k(x^{(i)})) x^{(i)} \\
 \text{where } h_k(x) &= \frac{e^{\theta_k^T x}}{\sum_{k'=1}^C e^{\theta_{k'}^T x}}, k = 1, 2, \dots, C
 \end{aligned}$$

• 随机梯度下降

随机选择一个样本，更新参数：

$$\theta_k := \theta_k + \alpha (1\{y = k\} - h_k(x)) x$$