

Softmax Regression.

Given training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, where $y^{(i)} \in \{1, 2, \dots, k\}$.

$$h_{\theta}(x^{(i)}) = \begin{bmatrix} p(y^{(i)}=1|x^{(i)}; \theta) \\ p(y^{(i)}=2|x^{(i)}; \theta) \\ \vdots \\ p(y^{(i)}=k|x^{(i)}; \theta) \end{bmatrix}_{k \times 1} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T \cdot x^{(i)}}} \cdot \begin{bmatrix} e^{\theta_1^T \cdot x^{(i)}} \\ e^{\theta_2^T \cdot x^{(i)}} \\ \vdots \\ e^{\theta_k^T \cdot x^{(i)}} \end{bmatrix}_{k \times 1}$$

↓
Normalize the distribution to 1.

$$X = \begin{bmatrix} \dots & x^{(1)T} & \dots \\ \dots & x^{(2)T} & \dots \\ \vdots & \vdots & \vdots \\ \dots & x^{(m)T} & \dots \end{bmatrix}_{m \times (n+1)}, \quad \theta = \begin{bmatrix} \dots & \theta_1^T & \dots \\ \dots & \theta_2^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & \theta_k^T & \dots \end{bmatrix}_{k \times (n+1)}$$

Cost Function:

$$J(\theta) = -\frac{1}{m} \cdot \left[\sum_{i=1}^m \sum_{j=1}^k \underbrace{1\{y^{(i)}=j\}}_{\substack{\downarrow \\ \text{indicator function} \\ 1\{\text{true statement}\}=1 \\ 1\{\text{false statement}\}=0}} \cdot \log \frac{e^{\theta_j^T \cdot x^{(i)}}}{\sum_{k=1}^k e^{\theta_k^T \cdot x^{(i)}}} \right] + \frac{\lambda}{2} \cdot \sum_{i=1}^k \sum_{j=0}^n \theta_{ij}^2$$

$$\begin{aligned} \nabla_{\theta_j} J(\theta) &= -\frac{1}{m} \cdot \sum_{i=1}^m \left[x^{(i)} \cdot (1\{y^{(i)}=j\} - p(y^{(i)}=j|x^{(i)}; \theta)) \right] + \lambda \cdot \theta_j \\ &= \frac{e^{\theta_j^T \cdot x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T \cdot x^{(i)}}} \end{aligned}$$

Softmax Regression vs k Binary classifiers.

If the k classes are mutually exclusive \rightarrow softmax classifier
otherwise, k separate classifiers are preferred.

Sparse function in matlab.

$S = \text{sparse}(i, j, v)$ generates a sparse matrix S from the triplets i, j , and v such that $S(i(k), j(k)) = v(k)$. The $\max(i)$ -by- $\max(j)$ output matrix has space allotted for nonzero elements.

e.g: $y = \begin{bmatrix} 3 \\ 5 \\ 1 \end{bmatrix}_{3 \times 1}$, $\text{sparse}(y, 1:3, 1) \Rightarrow \begin{matrix} (3,1) = 1 \\ (5,2) = 1 \\ (1,3) = 1 \end{matrix}$
↓
set the target value to 1.

full function in matlab: convert sparse matrix to full matrix.

e.g: $\text{full}(\text{sparse}(y, 1:3, 1)) = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

$\text{groundTruth} = \text{full}(\text{sparse}(\text{labels}, 1:\text{numCases}, 1));$

Implementation of Softmax Regression.

1. Cost Function and gradient.

i). Set 1, get the max number in each column.

ii). Set 2, get the max number in each row

$M = \text{bsxfun}(@\text{minus}, O \times \text{data}, \text{max}(O \times \text{data}, [], 1));$

O is $k \times (nH)$, data is $(nH) \times m$, $O \times \text{data}$ is $k \times m$.

In this function, M = each element in each column in $O \times \text{data}$ - max element in each column.

$M = \exp(M);$

→ sum of each column in M .

$P = \text{bsxfun}(@\text{rdivide}, M, \text{sum}(M));$

In this function, P = each element in each column in M divides by each column element in M , respectively, therefore, normalize the probability.

$\text{cost} = \frac{-1}{\text{numCases}} \cdot \text{groundTruth}(:,)^T \cdot \log(P(:,)) + \frac{\text{lambda}}{2} \cdot \text{sum}(\text{theta}(:,)^2);$

$\text{theta grad} = \frac{-1}{\text{numCases}} \cdot (\text{groundTruth} - P) \cdot \text{data}^T + \text{lambda} \cdot \text{theta};$