

Generalized softmax regression: MLE



$$P(\mathbf{y} \mid \mathbf{x}; \theta) = \frac{e^{f(\mathbf{x}; \theta)}}{\sum_k e^{f^{(k)}(\mathbf{x}; \theta)}}$$

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim P_D} [\log P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \sum_{i=1}^K y_i \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \langle \mathbf{y}, \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right] \rangle$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \sum_{i=1}^K \delta_{y_i=1} \cdot \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i=1}^K \left[\delta_{y_i=1} \log \sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}] - \delta_{y_i=1} f(\mathbf{x}; \boldsymbol{\theta})^{(i)} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \log \sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}] - \sum_{i=1}^K \delta_{y_i=1} f(\mathbf{x}; \boldsymbol{\theta})^{(i)}$$

$$\nabla_{f^{(k)}} RHS = \frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(k)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} - \delta_{y_k=1}$$

$$= \operatorname{softmax}(f(\mathbf{x}; \boldsymbol{\theta})^{(k)}) - \delta_{y_k=1}$$

softmax function:

ML

L

E



Gradients:

The gradients w.r.t. $f(\cdot)$ get propagated to the parameters, θ , via the chain rule. This gradient is valid for any distribution in the exponential family, not just the softmax.



Generalized softmax regression: MLE

Softmax function: $P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}) = \frac{e^{f(\mathbf{x}; \boldsymbol{\theta})}}{\sum_k e^{f^{(k)}(\mathbf{x}; \boldsymbol{\theta})}}$

The gradients w.r.t. $f(\cdot)$ get propagated to the parameters, $\boldsymbol{\theta}$, via the chain rule. This gradient is valid for any distribution in the exponential family, not just the softmax.

$$\text{MLE: } \hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim P_D} [\log P(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \sum_{i=1}^K y_i \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \langle \mathbf{y}, \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right] \rangle$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} - \sum_{i=1}^K \delta_{y_i=1} \cdot \log \left[\frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(i)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i=1}^K \left[\delta_{y_i=1} \log \sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}] - \delta_{y_i=1} f(\mathbf{x}; \boldsymbol{\theta})^{(i)} \right]$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \log \sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}] - \sum_{i=1}^K \delta_{y_i=1} f(\mathbf{x}; \boldsymbol{\theta})^{(i)}$$

Gradients: $\nabla_{f^{(k)}} \text{RHS} = \frac{\exp[f(\mathbf{x}; \boldsymbol{\theta})^{(k)}]}{\sum_j \exp[f(\mathbf{x}; \boldsymbol{\theta})^{(j)}]} - \delta_{y_k=1}$

$$= \text{softmax}(f(\mathbf{x}; \boldsymbol{\theta})^{(k)}) - \delta_{y_k=1}$$

Gradient Descent for softmax regression

mapping: $f: \mathbf{x} \rightarrow \mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$

where $\mathbf{x} \in \mathbb{R}^N$

$\mathbf{y}, \mathbf{b}, f(\cdot) \in \mathbb{R}^K$

$\mathbf{x} \in \mathbb{R}^N$

$\mathbf{W} \in \mathbb{R}^{K \times N}$

gradients: $\nabla_{W_{k,i}} \mathbf{L}_{CE} = x_i \left(\frac{\exp[\mathbf{w}_k \mathbf{x} + b_k]}{\sum_j \exp[\mathbf{w}_j \mathbf{x} + b_j]} - \delta_{1,y_k} \right)$

gradient update rule:

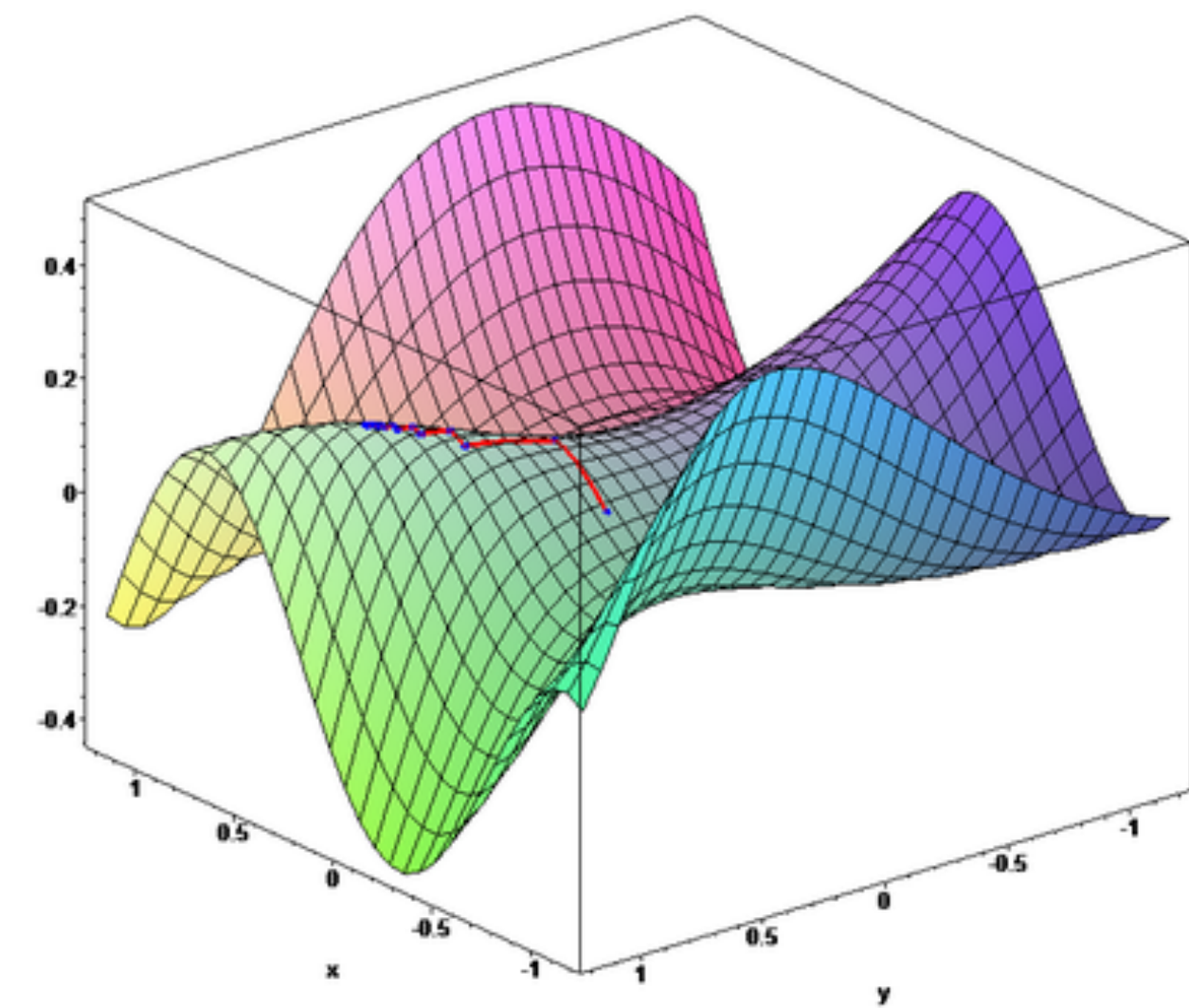
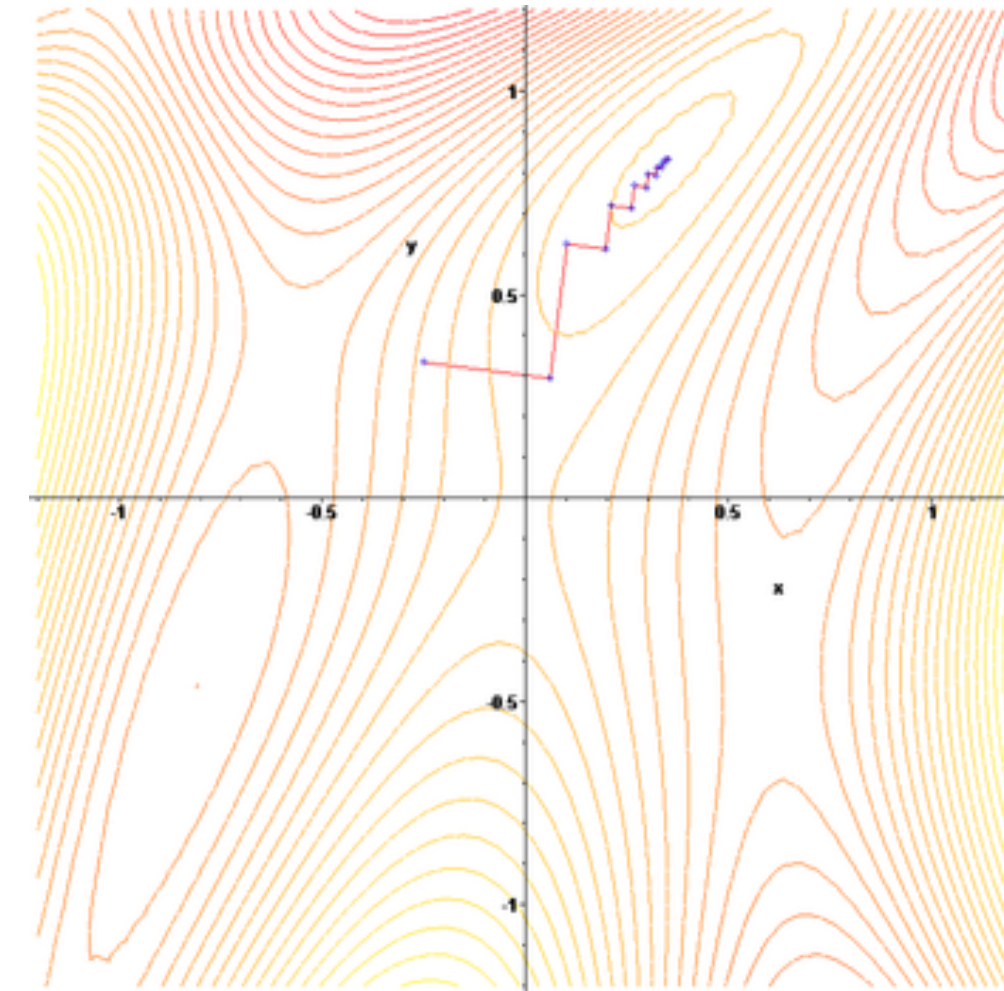
while not converged do : {

$\forall k, j$ do : {

$W_{k,j} := W_{k,j} - \eta \nabla_{W_{k,j}} L_{CE}$

 }

}



- Source: Wikipedia