# Generalized softmax regression

- Generic softmax regression adopts the underlying model mapping y <- x in logistic regression, but extends it to >2 dimensions

- The softmax function is more broadly useful for other (not necessarily linear) mapping functions, too! It comes in very handy for any function that's at least once differentiable w.r.t. parameters that are being learned. Here we'll just say there is some function           that's differentiable with respect to some parameters,    . Note, below represents a vector of probabilities, one for each class.

$$P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{e^{f(\mathbf{x}; \boldsymbol{\theta})}}{\sum_i e^{f^{(i)}(\mathbf{x}; \boldsymbol{\theta})}}$$

# Softmax function:

$$f : \mathbf{x} \to \mathbf{y}$$

*θ*

# Generalized softmax regression

- Generic softmax regression adopts the underlying model mapping y <- x in logistic regression, but extends it to >2 dimensions

- The softmax function is more broadly useful for other (not necessarily linear) mapping functions, too! It comes in very handy for any function that's at least once differentiable w.r.t. parameters that are being learned. Here we'll just say there is some function $f : \mathbf{x} \to \mathbf{y}$ that's differentiable with respect to some parameters, $\boldsymbol{\theta}$. Note, below represents a vector of probabilities, one for each class.

$$\text{Softmax function:} \quad P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{e^{f(\mathbf{x}; \boldsymbol{\theta})}}{\sum_i e^{f^{(i)}(\mathbf{x}; \boldsymbol{\theta})}}$$

# Generalized softmax regression: MLE

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

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

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

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

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

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

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

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

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