

Multiclass Classification $y = \{0, \dots, K-1\}$

MLE to estimate pmf $p(y|\vec{x})$

$$D = ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n))$$

(\vec{X}_i, Y_i) are i.i.d. with $P_{\vec{X}, Y}, P_{\vec{X}, Y}$ $\vec{\alpha}^*(\vec{x}_i)$

Assume $Y_i | \vec{X}_i = \vec{x}_i \sim \text{Categorical}(\alpha_0^*(\vec{x}_i), \dots, \alpha_{K-1}^*(\vec{x}_i))$

$$p(y=0 | \vec{x}_i) = \alpha_0^*(\vec{x}_i) \quad \sum_{y=0}^{K-1} \alpha_y^*(\vec{x}_i) = 1$$

$$\vdots$$
$$p(y=K-1 | \vec{x}_i) = \alpha_{K-1}^*(\vec{x}_i)$$

$$\sigma(\vec{z}) = (\sigma_0(\vec{z}), \dots, \sigma_{K-1}(\vec{z})) \in [0, 1]^K \quad \text{"softmax"}$$

$$\sigma_y(z_0, \dots, z_{K-1}) = \frac{\exp(z_y)}{\sum_{q=0}^{K-1} \exp(z_q)}$$

$$p(y|\vec{x}) = \sigma_y(\vec{x}^T \vec{w}_0^*, \dots, \vec{x}^T \vec{w}_{K-1}^*)$$

$$\alpha^*(\vec{x}_i) = \sigma(\vec{x}_i^T \vec{w}_0^*, \dots, \vec{x}_i^T \vec{w}_{K-1}^*)$$

$$\vec{w}_{MLE, 0}, \dots, \vec{w}_{MLE, K-1}$$

$$= \arg \min_{\vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1}} - \sum_{i=1}^n \log \left(p(y_i | X_i, \vec{w}_0, \dots, \vec{w}_{K-1}) \right)$$

$$= \arg \min_{\vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1}} - \sum_{i=1}^n \log \left(\sigma_{y_i}(\vec{X}_i^T \vec{w}_0, \dots, \vec{X}_i^T \vec{w}_{K-1}) \right)$$

$$= \arg \min_{\vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1}} - \sum_{i=1}^n \log \left(\frac{\exp(\vec{X}_i^T \vec{w}_{y_i})}{\sum_{q=0}^{K-1} \exp(\vec{X}_i^T \vec{w}_q)} \right) \quad \log\left(\frac{a}{b}\right) = \log(a) - \log(b)$$

$$= \arg \min_{\vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1}} - \sum_{i=1}^n \left[\log(\exp(\vec{X}_i^T \vec{w}_{y_i})) - \log\left(\sum_{q=0}^{K-1} \exp(\vec{X}_i^T \vec{w}_q)\right) \right]$$

$$= \arg \min_{\vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1}} - \sum_{i=1}^n \left[\underbrace{\vec{X}_i^T \vec{w}_{y_i}}_{= g(\vec{w}_0, \dots, \vec{w}_{K-1})} - \underbrace{\log\left(\sum_{q=0}^{K-1} \exp(\vec{X}_i^T \vec{w}_q)\right)}_{= h_i(r_i)} \right] \quad \text{Convex}$$

$$u_{iq} = \vec{X}_i^T \vec{w}_q, \quad v_{iq} = \exp(u_{iq}), \quad r_i = \sum_{q=0}^{K-1} \exp(u_{iq})$$

$$\vec{w}_y = (w_{y_0}, w_{y_1}, \dots, w_{y_d})^T \in \mathbb{R}^{d+1} \quad \text{for all } y \in \mathcal{Y}$$

$$\frac{\partial g}{\partial w_{y_j}}(\vec{w}_0, \dots, \vec{w}_{K-1}) = - \sum_{i=1}^n \left[\frac{\partial u_{iy_i}}{\partial w_{y_j}} - \frac{dh_i}{dr_i} \sum_{q=0}^{K-1} \frac{dv_{iq}}{du_{iq}} \frac{\partial u_{iq}}{\partial w_{y_j}} \right]$$

$$\frac{dh_i}{dr_i} = \frac{1}{r_i}, \quad \frac{dv_{iq}}{du_{iq}} = \exp(u_{iq}), \quad \frac{\partial u_{iq}}{\partial w_{y_j}} = \begin{cases} X_{ij} & \text{if } y=q \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial u_{iy_i}}{\partial w_{y_j}} = \mathbb{I}_{\{y_i\}}(y) X_{ij}$$

$$= \mathbb{I}_{\{y_i\}}(y) X_{ij}$$

Indicator function: $\mathbb{I}_{\mathcal{Y}}$

$$\mathbb{I}_{\mathcal{Y}}(z) = \begin{cases} 1 & \text{if } z \in \mathcal{Y} \\ 0 & \text{otherwise} \end{cases}$$

$$= - \sum_{i=1}^n \left[\mathbb{I}_{\{y_i\}}(y) X_{ij} - \frac{1}{\sum_{q=0}^{K-1} \exp(\vec{X}_i^T \vec{w}_q)} \sum_{q=0}^{K-1} \left(\exp(\vec{X}_i^T \vec{w}_q) \mathbb{I}_{\{y_i\}}(y) X_{ij} \right) \right]$$

$$= - \sum_{i=1}^n \left[\mathbb{I}_{\{y_i\}}(y) X_{ij} - X_{ij} \frac{\exp(\vec{X}_i^T \vec{w}_y)}{\sum_{q=0}^{K-1} \exp(\vec{X}_i^T \vec{w}_q)} \right]$$

$$= - \sum_{i=1}^n \left(\mathbb{I}_{\{y_i\}}(y) - \sigma_y(\vec{X}_i^T \vec{w}_0, \dots, \vec{X}_i^T \vec{w}_{K-1}) \right) X_{ij}$$

$$\frac{\partial g}{\partial w_{y_j}}(\vec{w}_0, \dots, \vec{w}_{K-1}) = \sum_{i=1}^n \left(\sigma_y(\vec{X}_i^T \vec{w}_0, \dots, \vec{X}_i^T \vec{w}_{K-1}) - \mathbb{I}_{\{y_i\}}(y) \right) X_{ij}$$

for all $y \in \mathcal{Y}$, $j \in \{0, \dots, d\}$

No closed form solution

Use gradient descent

$$\nabla_{\vec{w}_y} g(\vec{w}_0, \dots, \vec{w}_{K-1}) = \left(\frac{\partial g}{\partial w_{y_0}}, \dots, \frac{\partial g}{\partial w_{y_d}} \right)^T \in \mathbb{R}^{d+1} \quad \text{for } y \in \mathcal{Y}$$

$$\vec{w}_y^{(t+1)} = \vec{w}_y^{(t)} - \eta^{(t)} \nabla_{\vec{w}_y} g(\vec{w}_0^{(t)}, \dots, \vec{w}_{K-1}^{(t)})$$

Multiclass Classification Learner:

$$A(D) = \hat{f}_{\text{Mul}}$$

$$f_{\text{Bayes}}(\vec{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y | \vec{x})$$

$$\approx \operatorname{argmax}_{y \in \mathcal{Y}} p(y | \vec{x}, \vec{w}_{\text{MLE},0}, \dots, \vec{w}_{\text{MLE},K-1})$$

$$= \operatorname{argmax}_{y \in \mathcal{Y}} \sigma_y(\vec{x}; \vec{w}_{\text{MLE},0}, \dots, \vec{x}; \vec{w}_{\text{MLE},K-1})$$

$$= \hat{f}_{\text{Mul}}$$

Softmax Regression

$\mathbf{y}_{\text{soft}} = [0, 1]^K$ representing values of
 $(p(y=0|\vec{x}), \dots, p(y=K-1|\vec{x}))$

Labels: $K=3$, $\vec{y}_0 = (1, 0, 0) \Leftrightarrow$ class 0 from before
 $= (0, 1, 0) \Leftrightarrow$ class 1
 $= (0, 0, 1) \Leftrightarrow$ class 2

$$\mathcal{L}(\hat{\vec{y}}, \vec{y}) = - \sum_{q=0}^{K-1} \left[y_q \log(\hat{y}_q) \right] \text{ "multiclass cross-entropy loss"}$$

$$\mathcal{F} = \left\{ f \mid f: \mathbb{R}^{d+1} \rightarrow [0, 1]^K \text{ and } f(\vec{x}) = \sigma(\vec{x}^T \vec{w}_0, \dots, \vec{x}^T \vec{w}_{K-1}) \right. \\ \left. \text{where } \vec{w}_0, \dots, \vec{w}_{K-1} \in \mathbb{R}^{d+1} \right\}$$

Learner: $\mathcal{A}(D) = \hat{f}_{\text{ERM}}$

$$\hat{f}_{\text{ERM}} = \arg \min_{f \in \mathcal{F}} \hat{L}(f)$$

$$= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n - \sum_{q=0}^{K-1} \left[y_{iq} \log(f_q(\vec{x}_i)) \right]$$

$$= \arg \min_{f \in \mathcal{F}} - \sum_{i=1}^n \sum_{q=0}^{K-1} \left[y_{iq} \log(f_q(\vec{x}_i)) \right]$$

$$= f_{\text{MLE}}$$

$$f_{MLE}(\vec{x}) = \sigma(\vec{x}^T \vec{w}_{MLE,0}, \dots, \vec{x}^T \vec{w}_{MLE,k-1})$$

Comparison to Logistic Regression

If $K=2$ and $\mathcal{Y} = \{0, 1\}$

Assume $Y_i | \vec{X}_i = \vec{x}_i \sim \text{Categorical}(\alpha_0^*(\vec{x}_i), \alpha_1^*(\vec{x}_i))$
 $= \text{Bernoulli}(\alpha_1^*(\vec{x}_i) = \alpha^*(\vec{x}_i))$

$$p(y=0 | \vec{x}_i) = \alpha_0^*(\vec{x}_i), \quad p(y=1 | \vec{x}_i) = \alpha_1^*(\vec{x}_i)$$

$$\alpha_0^*(\vec{x}_i) + \alpha_1^*(\vec{x}_i) = 1 \Rightarrow \alpha_0^*(\vec{x}_i) = 1 - \alpha_1^*(\vec{x}_i)$$

$$\alpha_1^*(\vec{x}) = \sigma_1(\vec{x}^T \vec{w}_0^*, \vec{x}^T \vec{w}_1^*), \quad \alpha_0^*(\vec{x}) = \sigma(\vec{x}^T \vec{w}^*)$$

Goal: Show $\sigma_1(\vec{x}^T \vec{w}_{MLE,0}, \vec{x}^T \vec{w}_{MLE,1}) = \sigma(\vec{x}^T \vec{w}_{MLE})$

$$\sigma_1(z_0, z_1) = \frac{\exp(z_1)}{\exp(z_1) + \exp(z_0)}$$

$$= \frac{\exp(z_1)}{\exp(z_0) + \exp(z_1)} \cdot \frac{\exp(-z_1)}{\exp(-z_1)}$$

$$= \frac{\exp(z_1 - z_0)}{\exp(z_0 - z_0) + \exp(z_1 - z_0)}$$

$$= \frac{1}{\exp(z_0 - z_1) + 1}$$

$$\begin{aligned} \exp(a)\exp(b) \\ = \exp(a+b) \end{aligned}$$

$$\begin{aligned}\sigma_1(\vec{x}^T \vec{w}_0, \vec{x}^T \vec{w}_1) &= \frac{1}{\exp(\vec{x}^T \vec{w}_0 - \vec{x}^T \vec{w}_1) + 1} \\ &= \frac{1}{1 + \exp(-\vec{x}^T (\underbrace{\vec{w}_1 - \vec{w}_0}_{=\vec{w}}))} = \sigma(\vec{x}^T \vec{w})\end{aligned}$$

$$\sigma_0(z_0, z_1) = 1 - \sigma_1(z_0, z_1)$$

Multiclass

$$\vec{w}_{MLE,0}, \vec{w}_{MLE,1}$$

$$= \arg \min_{\vec{w}_0, \vec{w}_1 \in \mathbb{R}^{d_H}} - \sum_{i=1}^n \log(p(y_i | x_i, \vec{w}_0, \vec{w}_1))$$

$$= \arg \min_{\vec{w}_0, \vec{w}_1 \in \mathbb{R}^{d_H}} - \sum_{i=1}^n \log(\sigma_{y_i}(\vec{x}_i^T \vec{w}_0, \vec{x}_i^T \vec{w}_1))$$

$$= \arg \min_{\vec{w}_0, \vec{w}_1 \in \mathbb{R}^{d_H}} - \sum_{i=1}^n \left[y_i \log(\sigma_1(\vec{x}_i^T \vec{w}_0, \vec{x}_i^T \vec{w}_1)) + (1 - y_i) \log(\sigma_0(\vec{x}_i^T \vec{w}_0, \vec{x}_i^T \vec{w}_1)) \right]$$

$$= \arg \min_{\vec{w}_0, \vec{w}_1 \in \mathbb{R}^{d_H}} - \sum_{i=1}^n \left[y_i \log(\sigma(\vec{x}_i^T (\vec{w}_1 - \vec{w}_0))) + (1 - y_i) \log(\sigma_0(\vec{x}_i^T (\vec{w}_1 - \vec{w}_0))) \right]$$

$$= g(\vec{w}_0, \vec{w}_1) = h(\vec{w}) \quad \text{where } \vec{w} = \vec{w}_1 - \vec{w}_0$$

Binary

$$\begin{aligned}\vec{w}_{MLE} &= \arg \min_{\vec{w} \in \mathbb{R}^{d_H}} - \sum_{i=1}^n \left[y_i \log(\sigma(\vec{x}_i^T \vec{w})) + (1-y_i) \log(\sigma_0(\vec{x}_i^T \vec{w})) \right] \\ &= \arg \min_{\vec{w} \in \mathbb{R}^{d_H}} h(\vec{w})\end{aligned}$$

$$\vec{w}_{MLE} = \vec{w}_{MLE,1} - \vec{w}_{MLE,0}$$

$$\Rightarrow \sigma_1(\vec{x}^T \vec{w}_{MLE,0}, \vec{x}^T \vec{w}_{MLE,1}) = \sigma(\vec{x}^T \vec{w}_{MLE})$$