# Soft Max Regression w/ NN

**✳** Output layer's activation function
will be the soft-max function

$$\sigma(z_i) = \frac{e^{z_i}}{\sum_{k=1}^{n} e^{z_k}} \qquad \text{where } n = n\_outputs$$
$$= n\_categories$$

From regression,

$$p(y = i \mid x; \theta) = \frac{e^{\theta_i^T x}}{\sum_{j=1}^{n} e^{\theta_j^T x}} \qquad \text{where } \theta_k \sim \text{regression parameters}$$

It is common to take $\theta_n = 0$  $\left( \exists \ (n-1) \text{ inpendent probabilities} \right)$

Binary data: $p(y = i \mid x; \theta) = \dfrac{e^{\theta_i^T x}}{1 + e^{\theta_i^T x}} = \dfrac{1}{1 + e^{-\theta_i^T x}}$

$\sim$ recovers logistic regression

# Maximum likelihood

Take $p(y=i \mid x; \theta) = \dfrac{e^{z_i^L}}{\sum_{k=1}^{\hat{n}} e^{z_i^L}}$

where $\theta \sim$ weights & biases. To determine

$\theta$, use maximum (log) likelihood. In other words,

take the cost to negative log likelihood.

First, whats the likelihood we got this data

$D = \left\{ (x^{(i)}, y^{(i)}) \mid i = 1, \dots, N \right\}$ where $y^{(i)}$'s $\sim$ labels

$$P_\theta(D) = \prod_{i=1}^{N} P(y = y^{(i)} \mid x^{(i)}; \theta)$$

$$P(y \mid x; \theta) = \prod_{i=1}^{\hat{n}} P(y = i \mid x; \theta)^{1\{y=i\}}$$

$$= \prod_{i=1}^{\hat{n}} \left[ \frac{e^{z_i^L}}{\sum e^{z_i^L}} \right]^{1\{y=i\}}$$

$$L_\theta(D) = \sum_{i=1}^{N} \sum_{j=1}^{\hat{n}} \log\left[ \left( \frac{e^{z_j^L}}{\sum e^{z_j^L}} \right)^{1\{y^{(i)}=j\}} \right]$$

# Calculating $\partial_{w_{rs}^{\ell}} L_W(D)$

Write $a_i^{\ell} = s(z_i^{\ell})$ & $z_i^L = \sum_k^{\ell-1} a_k^{\ell-1} w_{ki}^{\ell} + b_i^L$

$$\implies L_\theta(D) = \sum_{i=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} \log s(z_j^L)$$

$$\partial_{w_{rs}^{\ell}} L_W(D) = \sum_{i=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} \frac{1}{s(z_j^L)} \frac{\partial s(z_j^L)}{\partial z_k^L} \frac{\partial z_k^L}{\partial w_{rs}^L}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} \frac{1}{s(z_j^L)} \left[ \frac{e^{z_j} \delta_{jk}}{\sum_0 e^{z_0^L}} - \frac{e^{z_j} e^{z_k}}{(\sum_0 e^{z_0^L})^2} \right] \frac{d}{dw_{rs}^L} \sum_P^{\ell-1} a_p^{\ell-1} w_{pk}^L$$

$$= \sum_{r=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} \frac{s(z_j)(\delta_{jk} - s(z_k))}{s(z_j)} \cdot \sum_P a_p^{\ell-1} \delta_{rp} \delta_{sk}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} (\delta_{jk} - s(z_k)) \cdot a_r^{\ell-1} \delta_{sk}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{n} 1\{y^{(i)}=j\} (\delta_{js} - s(z_s)) a_r^{\ell-1} \qquad \text{say } y^{(i)} = 2$$
$$z1 = 1 !$$

$$= \sum_{i=1}^{N} \left( 1\{y^{(i)} = s\} - s(z_s) \sum_{j=1}^{n} 1\{y^{(i)}=j\} \right) a_r^{\ell-1}$$

$$= \sum_{i=1}^{N} a_r^{\ell-1} \left( 1\{y^{(i)}=s\} - s(z_s) \right)$$

$$\boxed{\frac{\partial L_w(D)}{\partial w_{rs}^L} = \sum_{i=1}^{N} \left( 1\{y^{(i)} = s\} - s(z_s) \right) a_r^{L-1}}$$

Let's take $y = 1, 2, \ldots, K$ for $K$ classes

& let $y \longrightarrow y_i$ where $y_i^c = \delta_{ic}$ for class $c$.

Instead of focusing on $s(z_s)$ producing $1\{y^{(i)} = s\}$,

will use vectors: $\vec{y}^{(i)} - s(\vec{z})$

$$\frac{\partial L_w(D)}{\partial w^L} = \sum_{i=1}^{N} (\vec{a}^L)^T \left( \vec{y}^{(i)} - s(\vec{z}^L | x^{(i)}) \right)$$

Note: $\delta^L = \vec{y} - s(\vec{z}^L)$ !