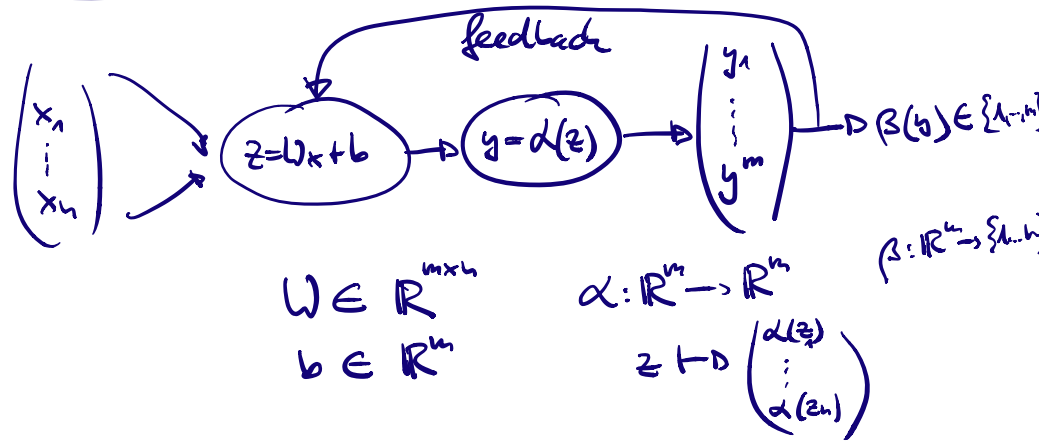


Multiclass Classification  $m \in \mathbb{N}$  classes,  $n \in \mathbb{N}$  inputs



Convention: for all functions  $f: \mathbb{R} \rightarrow \mathbb{R}$   
we def.  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$   
 $x \mapsto \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$

Loss function

$$L(W, b) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \alpha(Wx^{(i)} + b))$$

Update:

$$W \mapsto W^{\text{new}} := W - \eta \frac{\partial L}{\partial W}$$

$$b \mapsto b^{\text{new}} := b - \eta \frac{\partial L}{\partial b}$$

Example: ① Argmax Logistic regression model

Training data  $(x^{(i)}, y^{(i)})_{i=1 \dots n}$   $y^{(i)} \in \{1, \dots, m\}$

$$\alpha(z) = \frac{1}{1 + e^{-z}} \in (0, 1), \quad p(x^{(i)}) = \alpha(W \cdot x^{(i)} + b)$$

$$y = \alpha(Wx + b)$$

$$P(y_j^{(i)} = p_j(x^{(i)})) = p_j(x^{(i)})^{y_j^{(i)}} [1 - p_j(x^{(i)})]^{1 - y_j^{(i)}}$$

$$P(y^{(i)} = p(x^{(i)})) = \prod_{j=1}^m p_j(x^{(i)})^{y_j^{(i)}} [1 - p_j(x^{(i)})]^{1 - y_j^{(i)}}$$

$$P(y^{(i)} = p(x^{(i)})) \text{ for } i=1 \dots n$$

$$= \prod_{i=1}^n \prod_{j=1}^m p_j(x^{(i)})^{y_j^{(i)}} [1 - p_j(x^{(i)})]^{1 - y_j^{(i)}}$$

$$-\log P(y^{(i)} = p(x^{(i)})) \text{ for } i=1 \dots n$$

$$= -\sum_{i=1}^n \sum_{j=1}^m [y_j^{(i)} \log p_j(x^{(i)}) + (1 - y_j^{(i)}) \log (1 - p_j(x^{(i)}))]$$

$$L(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[ y_i^{(j)} \log p(x_i^{(j)}) + (1 - y_i^{(j)}) \log (1 - p(x_i^{(j)})) \right]$$

$\beta(y) :=$  choose element from  $\operatorname{argmax}_{i=1 \dots m} y_i$

$$\operatorname{argmax}_{i=1 \dots m} y_i := \left\{ j=1 \dots m \mid y_j \geq y_i \text{ for all } i=1 \dots m \right\}$$

② Argmax Softmax model

$$\alpha(z) = \left( \frac{e^{z_i}}{\sum_{j=1}^m e^{z_j}} \right)_{i=1 \dots m} \in (0, 1)$$

models the  
argmax by  
putting large  
 $z_i$  on a different  
scale than small  
ones.

everything blue above

$$\begin{aligned} \frac{d\alpha(z)}{dz_i} &= \frac{e^{z_i}}{\sum_{j=1}^m e^{z_j}} - \frac{e^{2z_i}}{\left(\sum_{j=1}^m e^{z_j}\right)^2} \\ &= \alpha(z) (1 - \alpha(z)) \end{aligned}$$

So that Cross-Entropy is a good loss  
function.