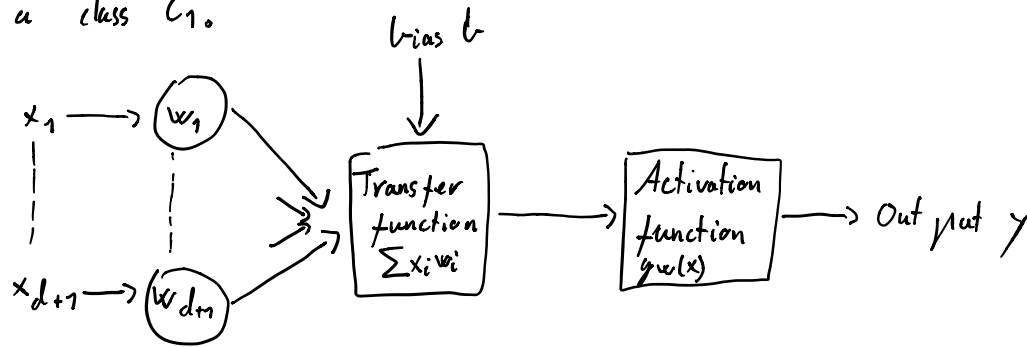


# 1 Logistic Regression

Logistic regression is a tool for binary classification that "uses" one neuron to determine the probability that input  $x$  belongs in a class  $C_1$ .



$$w_i \text{ have : } P(x \in C_1 | x) = g_w(x) = \frac{1}{1 + e^{-w^T x}}$$

$$P(x \in C_2 | x) = 1 - P(x \in C_1 | x) = 1 - g_w(x)$$

$$g_w(x) \in [0, 1]$$

$$\text{Loss function } E(w) = -\frac{1}{N} \sum_{n=1}^N t^n \ln(y^n) + (1 - t^n) \ln(1 - t^n)$$

measures how well our hypothesis function  $g_w(x) = y^n$  does on  $N$  data points.  $t^n$  = target value for example  $n$

Goal:  $t_n = y^n$  for all  $n$ .

\* Softmax regression is a generalization of logistic regression to a multi-class classification.

1.1

Show that  $-\frac{\partial E^n(w)}{\partial w_j} = (t^n - y^n) x_j^n$

We have  $E^n(w) = -t^n \ln(g_w^n) + (1-t^n) \ln(1-g_w^n)$

$$\Rightarrow \frac{\partial E^n(w)}{\partial w_j} = -t^n \frac{1}{g_w^n} \frac{\partial g_w^n}{\partial w_j} \quad \text{Applies chain rule.}$$

Using the hint:  $\frac{\partial E^n(w)}{\partial w_j} = -t^n \frac{1}{\cancel{g_w^n}} x_j^n \cancel{g_w^n} (1-g_w^n)$

$$-\frac{\partial E^n(w)}{\partial w_j} = t^n x_j^n (1-g_w^n) = (t^n - y^n) x_j^n$$

1.2

We have for the Softmax regression cost function

$$E = - \sum_{n=1}^N \sum_{k \in 1}^C t_{k'}^n \ln(y_{k'}^n) + t_k^n \ln(y_k^n), \quad k' \neq k$$

where:  $y_k^n = \frac{e^{a_k^n}}{\sum_{k'} e^{a_{k'}^n}}$  is the Softmax function with

$a_k^n = w_k^T x^n$  called net input to  $y_k$  and

$y_k^n$  is the probability that  $x$  is a member of class  $k$ . Note  $\sum_k y_k^n = 1$

We have two cases:  $k \neq k'$  and  $k = k'$ ,

hence we get:

Forsak 3:

$$E^n(w) = - \overbrace{\sum_{K'=1}^C t_{K'}^n \ln(y_{K'}^n)}^{\text{Case } K' \neq K} + \overbrace{t_K^n \ln(y_K^n)}^{\text{Case } K'=K}$$

$$= - \sum_{K'=1}^C t_{K'}^n \left( \ln(e^{w_{K'}^T x^n}) - \ln\left(\sum_{K'} e^{w_{K'}^T x^n}\right) \right) + t_K^n \ln\left(\frac{e^{w_K^T x^n}}{\sum_{K'} e^{w_{K'}^T x^n}}\right)$$

$$\frac{\partial E^n(w)}{\partial w_{kj}} = - \sum_{K'=1}^C t_{K'}^n \left( \left( \frac{1}{e^{w_{K'}^T x^n}} \cdot e^{w_{K'}^T x^n} \cdot x_j^n \right) - \left( \frac{1}{\sum e^{w_{K'}^T x^n}} e^{w_K^T x^n} x_j^n \right) \right) +$$

$$t_K^n \left( \left( \frac{1}{e^{w_K^T x^n}} \cdot e^{w_K^T x^n} \cdot x_j^n \right) - \left( \frac{1}{\sum e^{w_{K'}^T x^n}} e^{w_K^T x^n} x_j^n \right) \right)$$

$$\frac{\partial E^n(w)}{\partial w_{kj}} = - \sum_{K'=1}^C t_{K'}^n \left( x_j^n y_K^n \right) + t_K^n x_j^n (1 - y_K^n)$$

$$- \frac{\partial E^n(w)}{\partial w_{kj}} = x_j^n y_K^n \sum_{K'=1}^C t_{K'}^n + t_K^n x_j^n - x_j^n y_K^n$$

$$= x_j^n y_K^n \left( \sum_{K'=1}^C t_{K'}^n - 1 \right) + t_K^n x_j^n$$

$$= \underline{\underline{x_j^n (t_K^n - y_K^n)}}$$

### Task 2.2

$$L(w) = \|w_{ij}\|^2 = \sum_{ij} w_{ij}^2$$

$$\frac{\partial L(w)}{\partial w_{ij}} = 2w_{ij}$$

Gradient descent with regularization then becomes:

$$dw = \frac{dE(w)}{dw} + \lambda \frac{dL(w)}{dw}$$