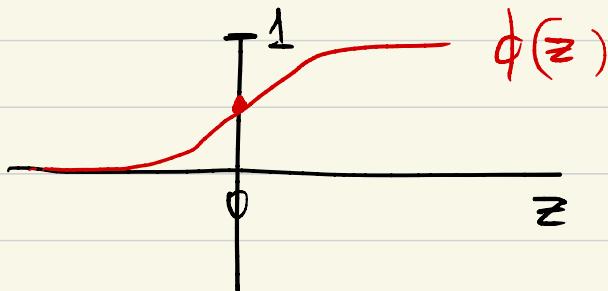


C. Logistic regression (LR)

In LR a sigmoid activation function is used, i.e.

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



To see where this comes from, let p

be the probability of a certain (binary) event (such as a coin),

then define (logit: $[0, 1] \rightarrow \mathbb{R}$)

$$\text{logit}(p) = \ln \frac{p}{1-p}$$

where $p/(1-p)$ is the odds ratio,

Let p now indicate the probability
that a particular sample belongs
to a class (e.g. $y = 1$),
under input data \underline{x} , i.e.

$$P = P(y=1 | \underline{x})$$

then

$$\text{logit}(P) = z = \sum_{j=0}^n w_j x_j$$

defines the activation function

in LR. Since

$$\begin{aligned} z &= \ln \frac{P}{1-P} \rightarrow -z = \ln \left(\frac{1}{P} - 1 \right) \\ \rightarrow e^{-z} &= \frac{1}{P} - 1 \rightarrow P = \frac{1}{1 + e^{-z}}. \end{aligned}$$

Assuming that the samples ($i = 1 \dots, N$) in the training data are independent, then the cost function $L(\underline{w})$ is defined as

$$\begin{aligned} L(\underline{w}) &= P(y | X; \underline{w}) = \\ &= \prod_{i=1}^N P(y^i | X^i, \underline{w}) = \\ &= \prod_{i=1}^N \left(\phi(z^i) \right)^{y^i} \left(1 - \phi(z^i) \right)^{1-y^i} \end{aligned}$$

probability
that feature
gives y^i

not.

In practice, we use $J(\underline{w}) = -\ln L(\underline{w})$

given by

$$J(\underline{w}) = -\sum_{i=1}^N \left\{ y^i \ln \phi(z^i) + (1-y^i) \ln (1-\phi(z^i)) \right\}$$

13

Consider $N=1$ (single sample),

then

$$\mathcal{J}(\underline{w}) = - \left(y \ln \phi(z) + (1-y) \ln(1-\phi(z)) \right)$$

Note that the quantizer in LR

is given by

$$\phi(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

hence $\mathcal{J}(\underline{w}) = -\ln(1-\phi(z)) \quad y=0$

and $\mathcal{J}(\underline{w}) = -\ln \phi(z) \quad y=1$

In summary, we penalize wrong predictions ($\phi(z) = 0$ when $y=1$)

and $\phi(z)=1$ when $y=0$) by very large costs.

Using gradient descent, we again update

$$\Delta w_{k+1} = -\eta \nabla E(w_k)$$

where $(\nabla E)_j$ is now computed as

$$-\nabla E_j = \frac{\partial}{\partial w_j} \left(y \ln \phi(z) + (1-y) \ln (1-\phi(z)) \right)$$

(assume for simplicity
that $N=1$)

$$= \left(\frac{y}{\phi(z)} - \frac{(1-y)}{1-\phi(z)} \right) \frac{\partial}{\partial w_j} \phi(z)$$

use : $\phi'(z) = \frac{d\phi}{dz} = \left(\frac{1}{1+e^{-z}} \right)' =$

$$= \frac{+e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}}$$

$$= \phi(z)(1-\phi(z))$$

15

hence $\frac{\partial}{\partial w_j} \phi(z) = \phi'(z) \frac{\partial}{\partial w_j} z$

$$\begin{aligned}
 (-\nabla \mathcal{E})_j &= \left(\frac{y}{\phi(z)} - \frac{1-y}{1-\phi(z)} \right) \phi(z) (1-\phi(z)) \\
 &\quad * \frac{\partial}{\partial w_j} z \quad , \quad z = \sum_j w_j x_j \\
 &= \left(y (1-\phi(z)) - (1-y) \phi(z) \right) x_j \\
 &= (y - \phi(z)) x_j
 \end{aligned}$$

So the update of the weights is

$$\begin{aligned}
 \Delta w_{k+1} &= -\eta \nabla \mathcal{E}(w) \\
 (\Delta w_{k+1})_j &= \eta \sum_{i=1}^N (y_i - \phi(z_i)) x_j
 \end{aligned}$$

similar to the gradient descend rule
of the Adaline

1b

I. Overfitting

Overfitting is the problem that a model (Perceptron, Adaline or LR) perform well on the training data, but not on the test data. A solution to overfitting is regularization by adding terms to the cost function to avoid large weights. For example, in LR the cost function is written as

$$\begin{aligned} J(\underline{w}) = & -\sum_{i=1}^N \left(y^i \ln \phi(z^i) + (1-y^i) \ln (1-\phi(z^i)) \right) \\ & + \frac{\lambda}{2} \|\underline{w}\|^2 \end{aligned}$$

which adds a term

$$-\nabla \mathcal{E}_j = -\lambda w_j$$

to the gradient of the cost function.