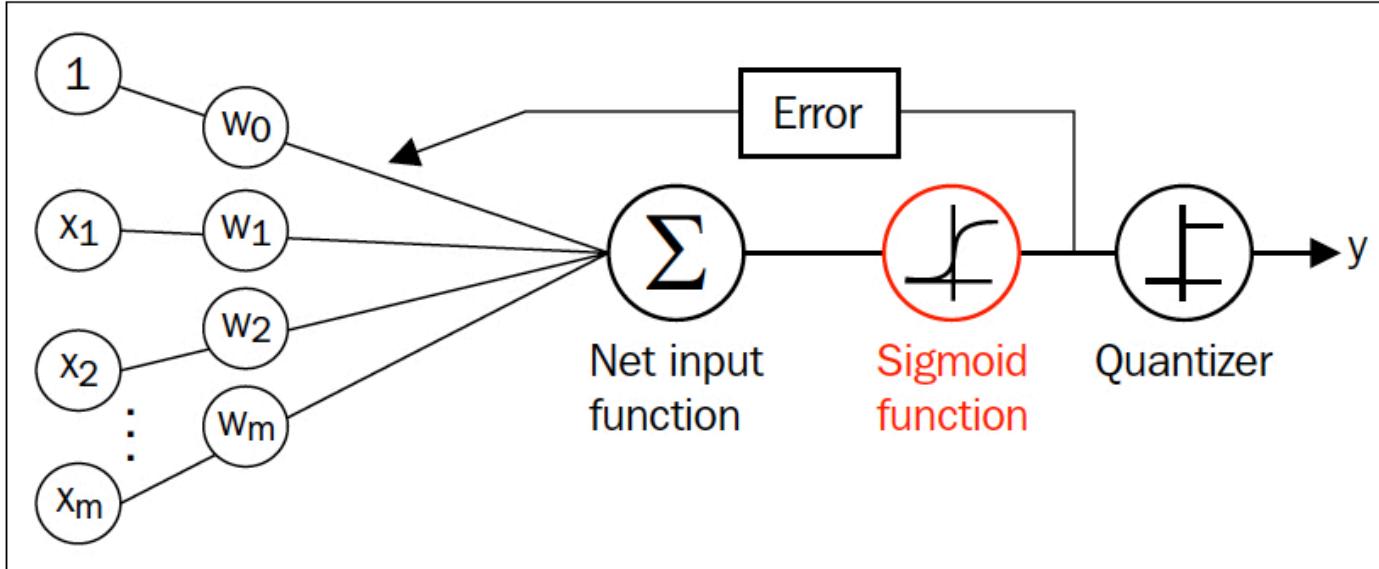


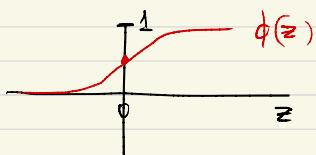
Logistic Regression



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.

condition 19

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

$$p = P(y=1 | x)$$

then

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

defines the activation function

in LR. Since

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

Summary

Define:

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

Odds ratio

p: probability positive outcome
of event

Ex: Event with class label $y = 1$

then LG:

$$O(p(y=1|\mathbf{x})) = \mathbf{w}^T \mathbf{x}$$

Note:

$$z = \ln \frac{x}{1-x} \rightarrow x = \frac{1}{1+e^{-z}}$$

12

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

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

Consider $N=1$ (single sample)¹³,
then

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

Note that the quantizer in LR

is given by

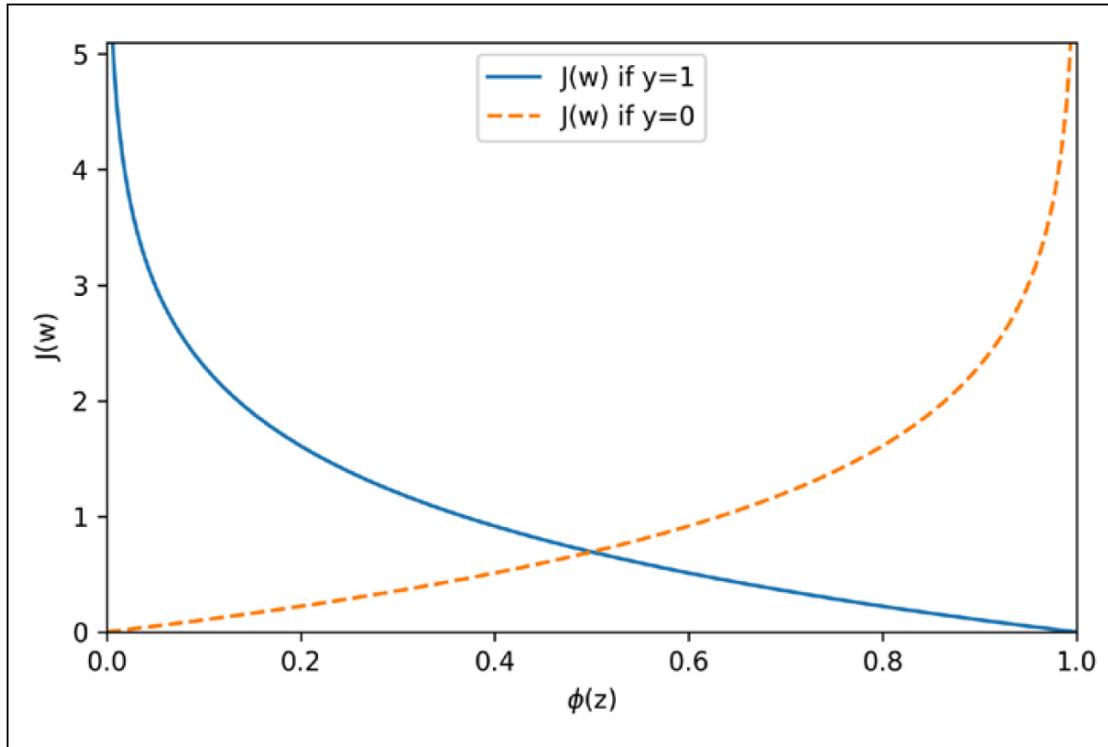
$$q(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

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

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

In summary, we penalize wrong predictions ($\phi(z)=0$ when $y=1$
and $\phi(z)=1$ when $y=0$) by
very large costs.

Example: $N = 1$



Summary

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

$$L(\boldsymbol{w}) = P(\mathbf{y} | \mathbf{x}; \boldsymbol{w}) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \boldsymbol{w}) = \prod_{i=1}^n \left(\phi(z^{(i)}) \right)^{y^{(i)}} \left(1 - \phi(z^{(i)}) \right)^{1-y^{(i)}}$$

success failure

$$l(\boldsymbol{w}) = \log L(\boldsymbol{w}) = \sum_{i=1}^n \left[y^{(i)} \log \left(\phi(z^{(i)}) \right) + (1 - y^{(i)}) \log \left(1 - \phi(z^{(i)}) \right) \right]$$

$$J(\boldsymbol{w}) = \sum_{i=1}^n \left[-y^{(i)} \log \left(\phi(z^{(i)}) \right) - (1 - y^{(i)}) \log \left(1 - \phi(z^{(i)}) \right) \right]$$

Mathematical problem

Input

$$\mathbf{x} = (x_0, \dots, x_m)^T$$

Weights

$$\mathbf{w} = (w_0, \dots, w_m)^T$$

data samples: $i = 1, \dots, N$

$$\min_{\mathbf{w}} J(\mathbf{w})$$

$$J(\mathbf{w}) = \sum_{i=1}^n \left[-y^{(i)} \log(\phi(z^{(i)})) - (1 - y^{(i)}) \log(1 - \phi(z^{(i)})) \right]$$

Using gradient descent, we
again update

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

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

$$\begin{aligned} -(\nabla E)_j &= \frac{\partial}{\partial w_j} \left(y \ln \phi(z) + (1-y) \ln (1-\phi(z)) \right) \\ &\quad \text{(assume for simplicity that } N=1 \text{)} \\ &= \left(\frac{y}{\phi(z)} - \frac{(1-y)}{1-\phi(z)} \right) \frac{\partial}{\partial w_j} \phi(z) \end{aligned}$$

$$\begin{aligned} \text{use : } \phi'(z) &= \frac{d\phi}{dz} = \left(\frac{1}{1+e^{-z}} \right)' = \\ &= \frac{+e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \\ &= \phi(z)(1-\phi(z)) \end{aligned}$$

hence $\frac{\partial}{\partial w_j} \phi(z) = \phi'(z) \frac{\partial}{\partial w_j} z$ [15]

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

So the update of the weights is

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

similar to the gradient descent rule
of the Adaline

Summary: gradient descend

Use: $\phi'(z) = \phi(z)(1 - \phi(z))$

$$(\Delta \mathbf{w}_{k+1})_j = \eta \sum_i (y^i - \phi(z^i)) x_j^i$$

1b

D. Overfitting

Overfitting is the problem that a model (Perception, 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

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

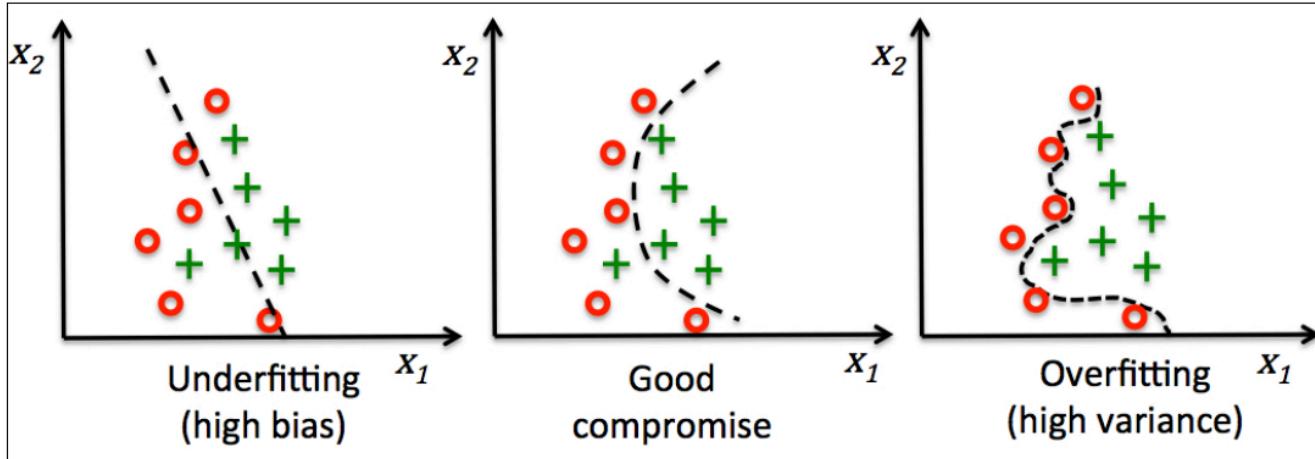
17

which adds a term

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

to the gradient of the cost function.

Overfitting & Regularization



$$J(\boldsymbol{w}) = C \left[\sum_{i=1}^n \left(-y^{(i)} \log(\phi(z^{(i)})) - (1-y^{(i)}) \log(1-\phi(z^{(i)})) \right) \right] + \frac{1}{2} \|\boldsymbol{w}\|^2$$

Exercise B3

