4.3.1
layer $\ell$    $1 \le \ell \le L$ :

$$\xrightarrow[\in R^{M_{\ell-1}}]{\frac{x_{\ell-1}}{}} \boxed{\begin{array}{c} \underline{\underline{W}}_\ell \\ \underline{b}_\ell \\ \phi_\ell \end{array}} \xrightarrow{} x_\ell \in R^{M_\ell}$$

$$\underline{\underline{W}}_\ell \in R^{M_\ell \times M_{\ell-1}}$$

$$\underline{b}_\ell \in R^{M_\ell}$$

$$\underline{a}_\ell = \underline{\underline{W}}_\ell \cdot \underline{x}_{\ell-1} + \underline{b}_\ell \in R^{M_\ell}$$

$$\underline{x}_\ell = \phi_\ell(\underline{a}_\ell) \in R^{M_\ell}$$

$M_\ell(M_{\ell-1} + 1)$ unknown parameters to be learnt by training eg by minimising cost fⁿ.

• No shortcuts are possible i.e no direct connection b/w $x_{\ell-1}$ and $x_\ell$

• No ~~feeb~~ feedbacks i.e. $\underline{x}_\ell$ ~~and~~ dont feed data or i/p to
$$\underline{x}_{\ell-1}$$

Network:

$$\underline{x}_L = f(\underline{x}_0 ; \underline{\theta}) : \mathbb{R}^{M_0} \to \mathbb{R}^{M_L}$$

•) parameter vector

$$\underline{\theta} = \begin{bmatrix} \text{vec}(\underline{\underline{W}}_1) \\ \underline{b}_1 \\ \vdots \\ \text{vec}(\underline{\underline{W}}_L) \\ \underline{b}_L \end{bmatrix} \in \mathbb{R}^{N_p} \quad, \text{ learnt from data during}$$

training.

$N_p$ : No of parameters

$$= \sum_{l=1}^{L} M_l \times (M_{l-1} + 1)$$

No of multiplications $N_x = \sum_{i=1}^{L} M_l M_{l-1} \approx N_p \to$ computational complexity

" "  additions  $\approx N_p$

•) No of layers $L$ and no of neurons in each layer

$\{M_1 \dots M_L\}$

•) activation f^n: $\{\phi_1 \dots \phi_L\}$

$\left.\begin{array}{l}\\ \\ \end{array}\right\}$ hyper parameters
– chosen by you
see 6.7

## First neural network in history:

(liner) perceptron by Rosenblatt, 1957

$$y = \underline{w}^T \bullet \underline{x} + b$$

\*) no • non-linear activation function $\phi()$

\*) no hidden layers

$\Rightarrow$ Only for linear tasks

## 4.4 Activation function

Mild requirements on $\phi_L()$:
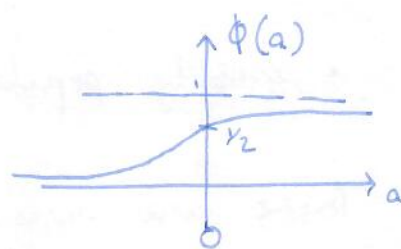
\*) non-linear in general $\rightarrow$ fundamental

\*) smooth, differentiable $\rightarrow$ for training

\*) simple calculation $\rightarrow$ low computational complexity

### Sigmoid function

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



- $0 < \phi(a) < 1 \,\hat{=}\, $ probability
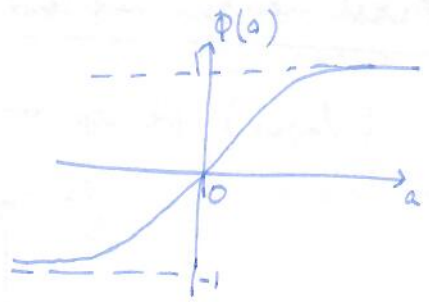
- Symmetrical $\phi(-a) = 1 - \phi(a)$

- derivative

$$\frac{d\phi(a)}{da} = \frac{e^{-a}}{(1 + e^{-a})^2} = \phi(a)\phi(-a) = \phi(a)\left[1 - \phi(a)\right] \in (0, 1)$$

- widely used in conventional NN.

## Hyperbolic tangent  tanh ()

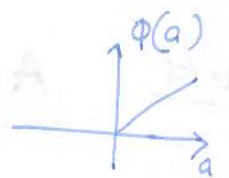$$\phi(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$$= 2\sigma(2a) - 1$$

- like sigmoid, another output range

## Rectifier Linear Unit  (ReLU)

$$\phi(a) = ReLU(a) = \max(a, 0) = \begin{cases} a & a \geqslant 0 \\ 0 & a < 0 \end{cases}$$

- $\hat{=}$ diode ⫠

- Simple calculation

- $\frac{d\phi}{da} = \begin{cases} 1 & a > 0 \\ 0 & a \leqslant 0 \end{cases} = u(a)$ with $u(0) = 0$

- mostly popular in Deep Learning

There are various variants of ReLU like leaky ReLU and Parametric leaky ReLU.

## Softmax

$$\phi(\underline{a}): \underline{a} = [a_i] \in \mathbb{R}^c \longrightarrow \mathbb{R}^c$$

$$\phi(\underline{a}) = \text{softmax}(\underline{a}) = \begin{bmatrix} \phi_1(\underline{a}) \\ \vdots \\ \phi_c(\underline{a}) \end{bmatrix} \text{ with}$$

$$\phi_i(\underline{a}) = \frac{e^{a_i}}{\sum\limits_{j=1}^{c} e^{a_j}} \in (0,1) \quad ; \quad \sum\limits_{i=1}^{c} \phi_i(\underline{a}) = 1$$

- $a_i$ large $\longrightarrow \phi_i(\underline{a})$ close to 1

  $a_i$ " $\longrightarrow \phi_i(\underline{a})$ close to 0

- a non-linear normalization of $\underline{a}$

- used in the o/p layer for classification tasks.

slide 4-8

Special case: $C = 2$

$$\phi_1(\underline{a}) = \frac{e^{a_1}}{e^{a_1} + e^{a_2}} = \frac{1}{1 + e^{-(a_1 - a_2)}} = \sigma(a_1 - a_2)$$

$$\phi_2(\underline{a}) = \frac{e^{a_2}}{e^{a_1} + e^{a_2}} = \frac{1}{1 + e^{-(a_2 - a_1)}} = 1 - \phi_1(\underline{a}) = \sigma(a_2 - a_1)$$

i.e.



one o/p of sigmoid is sufficient for binary classification

Derivative of softmax:

$$\frac{\partial \phi_i(\underline{a})}{\partial a_j} = \ldots = \begin{cases} \phi_i(\underline{a})(1 - \phi_i(\underline{a})) & i = j \\ \\ -\phi_i(\underline{a})\phi_j(\underline{a}) & i \neq j \end{cases}$$
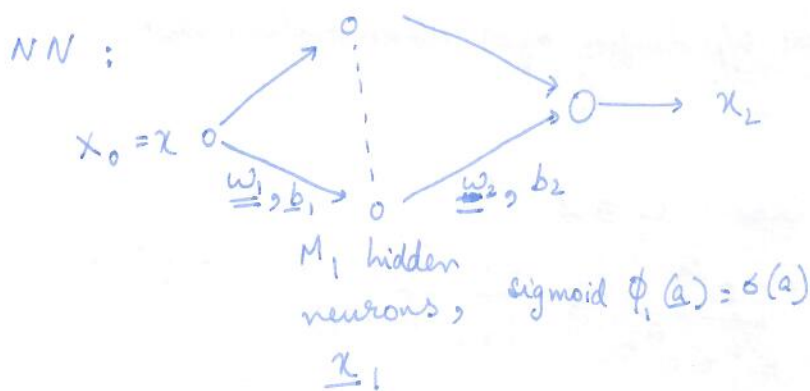
## 4.5  Universal approximation theorem

why can neural network be used every where?
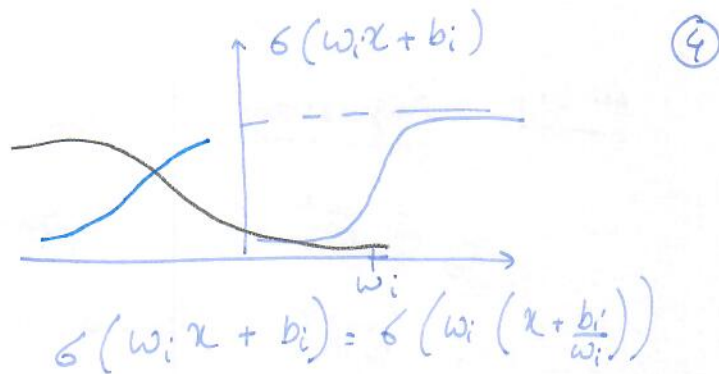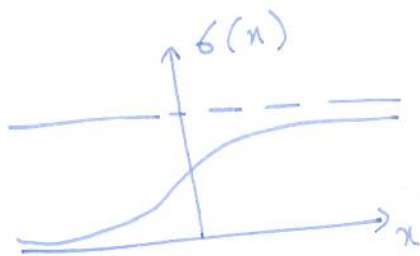
### E4.3  Regression with one hidden layer

True f$^n$ : $f_0(x)$

Given : $x(n)$ and noisy measurements $y(n) = f_0(x(n)) + z(n)$

$\qquad\qquad\qquad\qquad 1 \leq n \leq N$

NN :



$x_0 = x$

$\underline{w}_1, \underline{b}_1 \quad$ $\underline{w}_2, b_2$

$M_1$ hidden neurons, sigmoid $\phi_i(\underline{a}) = \sigma(a)$

$\underline{x}_1$

$$x_2 = f(x; \underline{\theta}) = \underline{w}_2^T \sigma(\underline{w}_1 x + \underline{b}_1) + b_2$$

$$= \sum_{i=1}^{M_1} w_{2,i}\, \underbrace{\sigma(w_{1,i} x + b_{1,i})} + b_2$$

$\qquad\qquad\qquad\qquad M_1$ non-linear basis functions of $x$

$\sigma(x)$

$\sigma(w_i x + b_i)$

$$\sigma(w_i x + b_i) = \sigma\left(w_i\left(x + \frac{b_i}{w_i}\right)\right)$$
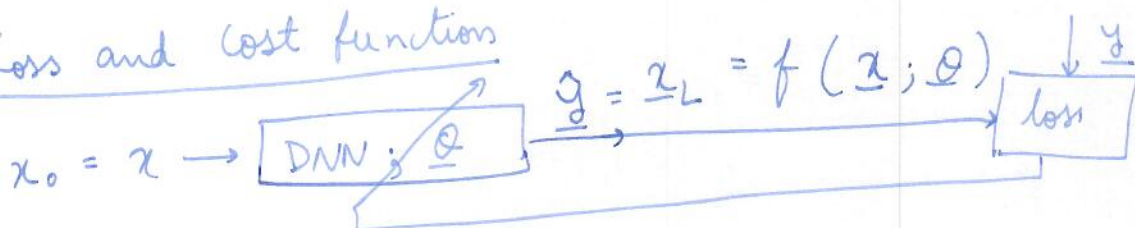
- new centre at $-\frac{b_i}{w_i}$
- new slope $w_i$

Universal approx theorem guarentees that there is a $f^n$ that exists but doesn't tell us more about it. Optimal is found using deep NN.

No theory on which activation $f^n$ is better than the others. It has to be found experimentally.

### 4.6  Loss and cost functions

$$x_0 = x \longrightarrow \boxed{DNN \,;\, \underline{\theta}} \quad \underline{\hat{y}} = \underline{x}_L = f(\underline{x}; \theta) \longrightarrow \boxed{loss} \quad \downarrow \underline{y}$$

ch 3.4
$$\min_{\underline{\theta}} L(\underline{\theta}) = \frac{1}{N} \sum_{n=1}^{N} \ell(\underline{x}(n), \underline{y}(n); \underline{\theta}) \quad \begin{array}{l} \text{overall} \\ \text{cost } f^n \text{ to} \\ \text{be minimized from} \\ D_{train}. \end{array}$$

$$\ell(\underline{x}, \underline{y}; \underline{\theta}) = -\ln g(\underline{y}|\underline{x}; \underline{\theta}) \longleftarrow \text{NLL } \underline{loss} \text{ for one } \underline{x}, \underline{y}$$

$$\underline{g}\,() \xleftarrow{\quad z_0 \quad} DNN$$

## 4.6.1 Regression

$\underline{x} \in \mathbb{R}^d$      random input

$\underline{y} \in \mathbb{R}^c$ : desired "      o/p

Assumption : DNN estimates the mean of $\underline{y}$ i.e.

$$\underline{y} = f(\underline{x}; \underline{\theta}) + \underline{z} \quad ; \quad \underline{z} - \text{noise}$$