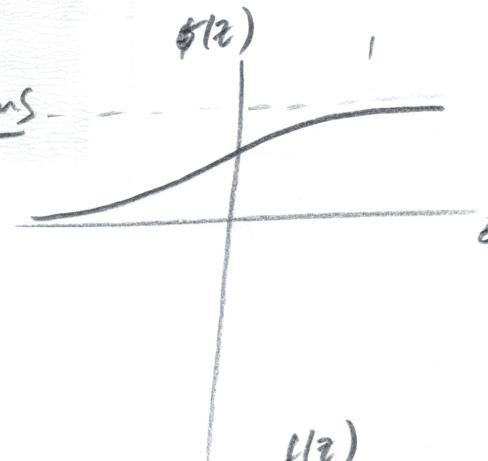


## L4 - Backpropagation - review

(\*) Activation functions (x) will on activation functions

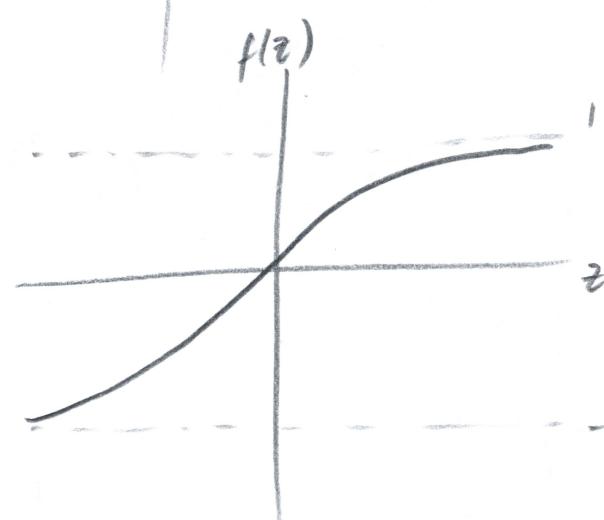
Sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$



Tanh  $f(z) = \tanh(z)$

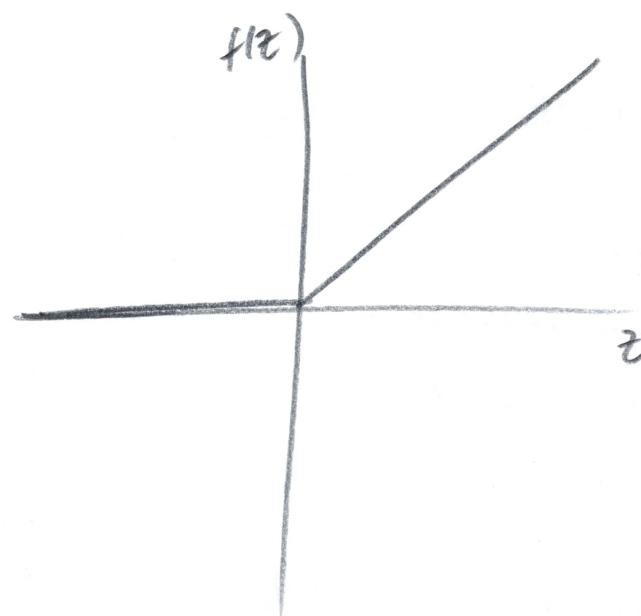
$$f'(z) = 1 - f(z)^2 = 1 - \tanh^2(z)$$



$$= \max(0, z)$$

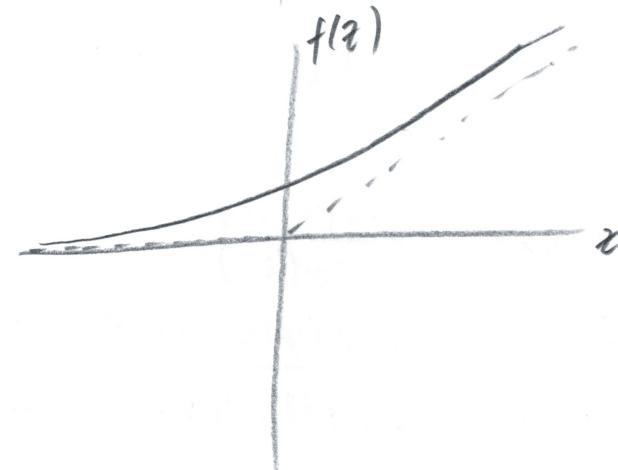
ReLU  $f(z) = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$

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



Softplus  $f(z) = \ln(1+e^z)$

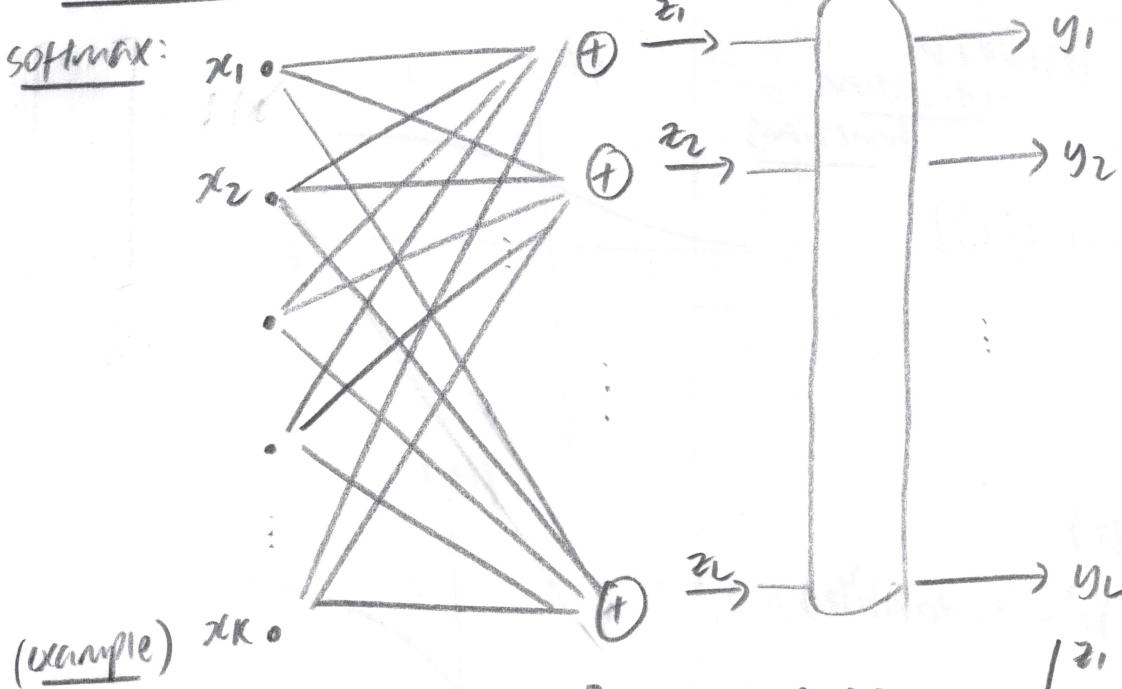
$$f'(z) = \frac{1}{1+e^{-z}}$$



(smooth  
approximation  
to ReLU)

## (\*) vector activations

- softmax:



(example)  $x_{K,0}$

$$z_i = \sum_j w_{ji} x_j + b_i$$

$$y_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

These can  
be stacked

$$\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_L \end{pmatrix}$$

$$y = \begin{pmatrix} \frac{\exp(z_1)}{\sum_j \exp(z_j)} \\ \vdots \\ \frac{\exp(z_L)}{\sum_j \exp(z_j)} \end{pmatrix}$$

(\*) Natural affinity of output unit and individual error/loss/divergence

- Recall you can choose whether to endow with p.o.b. mstep./est. principle
- And that you can instead choose to frame in terms of minimisation
- Prob. interpretation:-
- conditional distn:-  $p(y|z; \theta) \rightarrow ML$

(\*) BR doesn't use consistent, nor memorable, nor standard notation

↳ look elsewhere if you want formalism

(\*) Essentially, BR's 'error' refers to the average training error

(empirical risk minimisation)

(\*) Also; this is a sample average of error/divergence/loss over each training example; that is there is additive structure

⑦

for the relationship between error minimisation, functional form, probabilistic interpretation (ML estimation) and info-theoretic relations; see Bishop for former; and Goodfellow for latter.

### (\*) Real-valued output (vectors)

- Notationally, corresponds to Bishop (2006) - S.5.2
- There is a summation over training instances; and over each component of the target vector in this form:-  
(multiple target variables)
- BR: " $\text{div}(\underline{y}, \underline{d}) = \frac{1}{2} \|\underline{y} - \underline{d}\|_2^2 = \frac{1}{2} \sum_i (y_i - d_i)^2$ "

$\underline{y}$  - vector valued  
 $\underline{d}$  - vector

- over each component  
of target vector

AND "Err(w) =  $\frac{1}{T} \sum_{i=1}^T \text{div}(f(x_i; w), d_i)$ "

- over all training instances

(\*) Also fails to make clear that div and error are also a function of parameters entirely explicitly (or at least does not emphasise)

BS  $\frac{\partial \text{div}(\underline{y}, \underline{d})}{\partial y_i} = (y_i - d_i) \checkmark$

$\nabla_{\underline{y}} \text{div}(\underline{y}, \underline{d}) = \begin{bmatrix} y_1 - d_1 \\ y_2 - d_2 \\ \vdots \end{bmatrix} \checkmark$

Bishop not.:  $E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|_2^2 \quad \frac{\partial E}{\partial w_k} = y_k - t_k \quad \checkmark$

- (\*) Binary class.
- scalar output  $y \in (0,1)$        $d \in (0,1)$
  - $\text{div}(y, d) = -d \log y - (1-d) \log(1-y)$       (cross entropy)
  - scalars
  - $\frac{d \text{div}(y, d)}{dy} = \begin{cases} -\frac{1}{y} & \text{if } d=1 \\ \frac{1}{1-y} & \text{if } d=0 \end{cases}$
  - (see Appendix of HW1)
    - info-theoretic entropies
    - Bernoulli
  - sigmoid activation used in combination with binary class.  
(for outputs) (see Bishop)

- (\*) Multi-class classification
- output:  $\begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix}$       target output:  $\begin{bmatrix} d_1 \\ d_2 \\ \vdots \end{bmatrix}$
  - one-hot encoding
    - with  $c^{\text{th}}$  position  $d_c = 1$  if output is class  $c$ .

$$\text{div}(y, d) = - \sum_i d_i \log y_i = - \log y_c$$

→ a summation over components of vector  $y$  and  $d$

$$\frac{d \text{div}(y, d)}{dy_i} = \begin{cases} -\frac{1}{y_c} & \text{for } c^{\text{th}} \text{ component} \\ 0 & \text{for remaining components} \end{cases}$$

$$\nabla_y \text{div}(y, d) = \begin{bmatrix} 0 \\ \vdots \\ -\frac{1}{y_c} \\ \vdots \\ 0 \end{bmatrix}$$

$$E(w) = - \sum_{n=1}^N \sum_{k=1}^K t_n k \ln y_k(z; w)$$



Bishop note: full summed error over training instances and components of vector outputs

- If considering  $i^{th}$  training instance:

$$E(w) = \sum_{k=1}^N E_k(w)$$

i.e. suppress training instance indexing

$$E_k(w) = -\sum_{n=1}^K t_{kn} \ln y_{kn}$$

### (\*) Gradient descent algorithm (scalar formulation)

• Initialise weights  $\{w_{ij}^{(k)}\}$

• DO:-

•  $\forall i, j, k$  (i.e. inter unit weights in all layers), initialise  $\frac{dE_k}{dw_{ij}^{(k)}} = 0$

•  $\forall t = 1, \dots, T$  (i.e. all training instances)

- For layer  $k$  and  $t$ :

- compute  $\frac{dDiv(Y_t, d_t)}{dw_{ij}^{(k)}}$

-  $\frac{dE_k}{dw_{ij}^{(k)}} += \frac{dDiv(Y_t, d_t)}{dw_{ij}^{(k)}}$

• For every layer  $K$ , and  $t$ :

$$w_{ij}^{(k)} = w_{ij}^{(k)} - \eta \frac{dE_k}{T dw_{ij}^{(k)}}$$

• until  $E_k$  has converged

(\*) Require efficient algorithm for computing  $\frac{dDiv(Y_t, d_t)}{dw_{ij}^{(k)}}$  for all weights

- This is the essence of backpropagation

- That is to find a way of calculating how much individual per training instance loss/error/divergence is responding to changes in your parameters

## (\*) chain rule of calculus + influence diagram

Backprop → forward pass

→ backward pass

- forward pass → fairly clear

- propagate info through network

## Forward pass (pseudocode) (scalar)

- input: D dim vector  $\underline{x} = [x_1, x_j, x_D] \quad j=1, \dots, D$

► Set:-

-  $D_0 = D$  (width of 0<sup>th</sup> layer (input))

-  $y_j^{(0)} = x_j \quad j=1, \dots, D \quad y_0^{(k=1, \dots, N)} = x_0 = 1$  (load 0<sup>th</sup> layer (inputs) and initialise biases)

► for layer  $k=1, \dots, N$

- for  $j=1, \dots, D_k$  ( $D_k$  is size of k<sup>th</sup> layer)

$$\cdot z_j^{(k)} = \sum_{i=0}^{D_{k-1}} w_{i,j}^{(k)} y_i^{(k-1)}$$

$$\cdot y_j^{(k)} = f_k(z_j^{(k)})$$

output:  $\underline{y} = y_j^{(N)} \quad j=1, \dots, D_N$  (final layer outputs)

(\*) not entirely clear what intermediate forward value comput. are stored - needs clarity

notes + slides cover already

## (\*) Backpropagation

(\*) compute derivatives layer by layer

- (i)  $\frac{\partial \text{Div}(\underline{y}, \underline{d})}{\partial y_i} = \frac{\partial \text{Div}(\underline{y}, \underline{d})}{\partial y_i^{(n)}} \quad \forall i$  (derivative of div not output)

(ii)  $\frac{\partial \text{Div}(\underline{y}, \underline{d})}{\partial z_i^{(N)}}$       (iii)  $\frac{\partial \text{Div}(\underline{y}, \underline{d})}{\partial w_{ij}^{(N)}}$       iv)  $\frac{\partial \text{Div}(\underline{y}, \underline{d})}{\partial y_i^{(N-1)}}$  - bias included in weight

do this for each layer

$N, (N-1), \dots$

(\*)

- (\*) Key is to note  
derivative of each layer is made up of
- i) an element computed in an earlier iteration of algorithm  
(div) (i.e. a later layer)
  - ii) an element computed during the forward pass, which has  
(div) been 'cached'
- 
- OR: Does not make clear  
what values are cached from forward pass
- 
- ② ① - Backpropagation  
(\*) one for one (final layer) - 4 characteristic types (see diagrams  
so you can visualise it)  
(\*) do for  $(N-1)^{th}$  layer to illustrate chain rule summation

$$\frac{\partial \text{div}}{\partial y_i^{(N-1)}} = \sum_j \underbrace{\frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}}}_{= w_{ij}^{(N)}} \cdot \underbrace{\frac{\partial \text{div}}{\partial z_j^{(N)}}}_{\text{already comp.}}$$

as  $z_j^{(N)} = w_{ij}^{(N)} y_i^{(N-1)} + \text{other terms}$  (forward pass?)

$$\Rightarrow \frac{\partial \text{div}}{\partial y_i^{(N-1)}} = \sum_j w_{ij}^{(N)} \frac{\partial \text{div}}{\partial z_j^{(N)}}$$

$$\frac{\partial \text{div}}{\partial z_i^{(N-1)}} = f_{N-1}'(z_i^{(N-1)}) \frac{\partial \text{div}}{\partial y_i^{(N-1)}}$$

w w already comp.  
already comput.  
in forward pass.

$$\frac{\partial \text{div}}{\partial w_{ij}^{(N-1)}} = \frac{\partial z_j^{(N)}}{\partial w_{ij}^{(N-1)}} \cdot \frac{\partial \text{div}}{\partial z_j^{(N)}}$$

$$= y_i^{(N-2)} \frac{\partial \text{div}}{\partial z_j^{(N)}}$$

$$y_0^{(N-2)} \text{ bias} = 1$$

## backward pass pseudocode (scalar form)

- output layer ( $N$ ):  
    - For  $i = 1, \dots, D_N$  (iterate over size of  $N^{\text{th}}$  layer from 1- $D_N$ )
  - $\frac{\partial \text{Div}}{\partial y_i} = \frac{\partial \text{Div}(y, d)}{\partial y_i^{(N)}}$
  - $\frac{\partial \text{Div}}{\partial z_i^{(N)}} = \frac{\partial \text{Div}}{\partial y_i^{(N)}} \frac{\partial y_i^{(N)}}{\partial z_i^{(N)}}$

- for layer  $k = N-1$  to 0: (iterate over layers starting from  $N-1$  in rev.)

- For  $i = 1, \dots, D_k$  (iterate over every unit in layer)

- $\frac{\partial \text{Div}}{\partial y_i^{(k)}} = \sum_j w_{ij}^{(k+1)} \frac{\partial \text{Div}}{\partial z_j^{(k+1)}}$

- $\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \frac{\partial \text{Div}}{\partial y_i^{(k)}} f'_k(z_i^{(k)})$

- $\frac{\partial \text{Div}}{\partial w_{ji}^{(k+1)}} = y_j^{(k)} \frac{\partial \text{Div}}{\partial z_i^{(k+1)}} \quad \text{for } j = 1, \dots, D_{k+1}$

(\*) - still can't entirely see the analogy/similarity between  
backwards and forwards pass

(\*) Backpropagation  $\rightarrow$  derivative of error (difference) is backpropagated  
recursively through network  
"already computed"