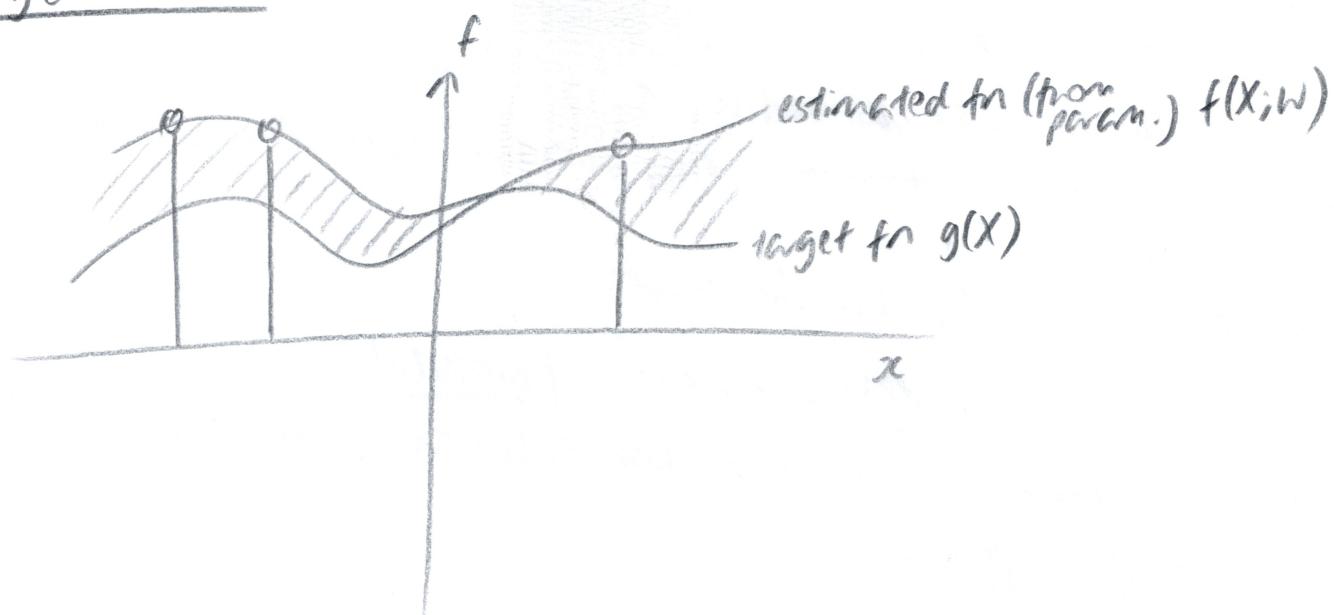


14 - BackpropagationBR: continue training from L3 (will later demarcate)L3 - Training (continued)

- BR:



- BR:- we want to estimate parameters  $w$  to minimise the area between  $f(x; w)$  and  $g(x)$  (or some function of it)
- we require entire function  $g(x)$  to compute the integral
$$\int_X \text{div}(f(x; w), g(x)) P(x) dx$$
  - $g(x)$  unknown, therefore obtain samples  $(x_i, d_i)$   
(i.e. values of  $g(x_i)$ )
  - Rather than minimise area; minimise average of 'errors' as a proxy
  - Quantified by divergence function (~~either~~  $\geq 0$ )
  - Require divergence to be differentiable
  - minimise empirical average error over training data

optimisation

- High school calculus refresher
- distinguish minimum, maximum, turning points
- $f''(x) < 0$  at maxima
- $f''(x) > 0$  at minima

- none of this accounts for concavity/convexity

- writical point  $\rightarrow$  local minima, maxima, inflection points.

-  $f''(x) = 0$  at inflection points  $\rightarrow$  require further investigation

## Multivariate optimisation

- find a location where derivative is 0

- small perturbations will not in any direction will not change the value of function.

- gradient: vector of partial derivatives

- function increases fastest in direction of gradient

- — " decreases —" in opposite direction of gradient

(\*) gradient 2nd property:

- slice the function  $f(x)$  at any height

- edge is a level set / contour

- gradient is orthogonal to the level set (remember Bishop!)

① ②  $\rightarrow$  clarify on this  $\Rightarrow$  especially when using  
Lagrange multipliers

- recall orthogonal - perpendicular

(\*) Hessian-matrix of 2nd derivatives of Multivariate  $\rightarrow$  opt

- generalisations of 2nd deriv for scalar fn  $\rightarrow$  greatest increase/decrease of

(\*) Hessian eigenvectors give directions of function

eigenvalues give direction of curvature

① ②  $\rightarrow$  clarify

- eigenvalue - 0 - inflection  
+ve - minimum } analog to 2nd deriv  
-ve - maximum

(\*) unconstrained max of Multivariate scalar fn

1. solve  $\nabla_X f(X) = 0$

2. compute  $\nabla_X^2 f(X)$  at candidate soln.; evaluate

If Hessian  $\nabla^2$  is  $S_n^{++}$  (positive def), eigenvalues are all tve  
⇒ local minima.

(W) (B) → clarify

If Hessian  $\nabla^2$  is  $S_n^{-}$  (negative def), eigenvalues are negative and

⇒ local maxima

(\*) Multiple dimensions → maximum in one-direction  $\Rightarrow$  derivative minimum in another  $\Rightarrow$  some  $\nabla^2$  eigenvalues tve, some -ve  
horseshoe

- often non-convexity forbids such solutions analytically

- use iterative solutions

(\*) or is it that local minima are not necessarily global minima (W) (A)

- iterative solutions:- (optimisation of  $f(\underline{x})$ )

i) initial guess  $\underline{x}_0$

• direction to step in

ii) update guess

• how big steps must be

iii) stop when  $f(\underline{x})$  no longer decreases  $\rightarrow \underline{x}^*$

- gradient descent

(W) (A) - clarify

logic; check  
you completely  
understand

- for maximum:-  $\underline{x}^{k+1} = \underline{x}^k + \eta^k \nabla f(\underline{x}^k)^T$

- for minimum:-  $\underline{x}^{k+1} = \underline{x}^k - \eta^k \nabla f(\underline{x}^k)^T$

(\*) You must make sure you understand the logic of gradient descent

impeachably (W) (A) (P) (B) - (D)

(\*) Relation between step size, derivative (magnitude)

convergence of gradient descent:-

INCOMPLETE  
IF YOU DO  
NOT  
REFRESH

- results for convex/non-convex functions

→ 14-Backprop

- back to DL context: given  $(\underline{x}_1, d_1), \dots, (\underline{x}_T, d_T)$

$$\underset{w}{\operatorname{argmin}} \text{loss}(w) = \frac{1}{T} \sum_i \text{div}(f(\underline{x}_i; w), d_i)$$

BR kept this deliberately vague intentionally:-  
allows explicit consideration of following issues

Q1) What are input-output pairs  $(x_i, d_i)$

Q2) What is  $f(\cdot)$  and parameters  $W$

Q3) What is divergence  $\text{div}(\cdot)$

•  $f(\cdot)$  - MLP directed / NN

- directed  $\rightarrow$  no feedback loops (1-way info flow)

(subseq.)

- layered  $\rightarrow$  NN

- each layer only gets input from previous outputs to later layers

- some disagreement about terminology (remember Bishop)

- generally; an individual neuron:-

- A diff. activation applied to an affine input combo

$$y = f\left(\sum w_i x_i + b\right)$$

(\*) Assume this otherwise specified  
in course i.e. always affine

- Any differentiable fn  $f(x_1, \dots, x_N; W)$

(for gradient  
descent)

- Examine effects of parameter perturbation on output

- G.D. - If I modify param; how much does div/error  
change (itself fn of output)

(\*) However, can use different activations

WAF: Examine / familiarise with activation functions and derivatives (\*)

• New peculiarity - not differentiable at 0

BR: vector activations

(\*) specific activations

WAF: clarity on what it means here

- softmax activation (as vector activation)

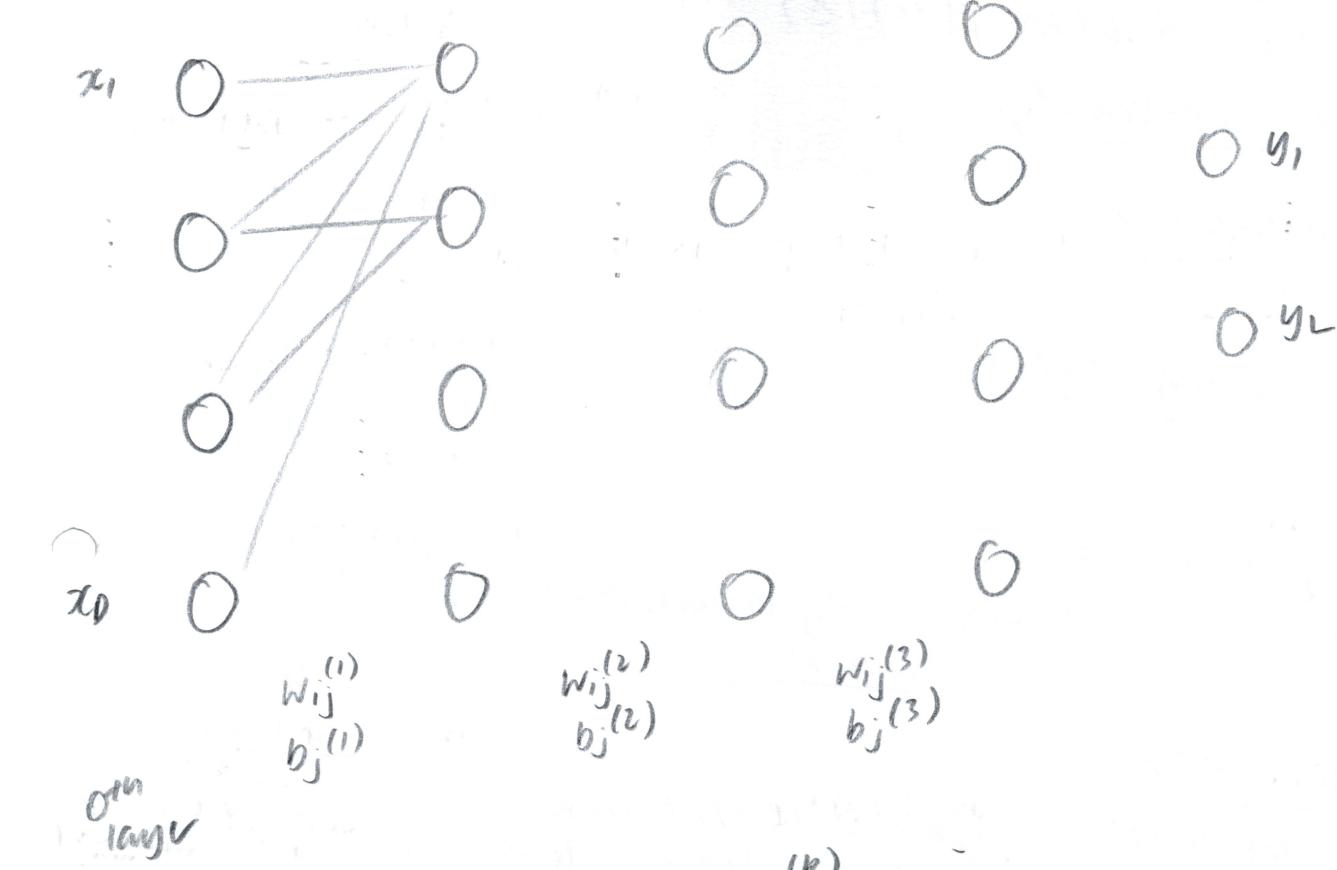
(\*) Nature of activ - change one param  $\Rightarrow$   
all large outputs will change

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

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

vector activation is generalisation of scalar activations  
(layer) (individual neurons)

• notation:  $\oplus$   $\otimes$   $(*)$



- output of  $i^{\text{th}}$  perceptron of  $k^{\text{th}}$  layer -  $y_i^{(k)}$

- input to network

$$y_i^{(0)} = x_i$$

$$y_i^{(n)} = y_i$$

- output of network

$\otimes$  weight of connection bet  $i^{\text{th}}$  unit of  $(k-1)^{\text{th}}$  layer and  $j^{\text{th}}$  unit of  $k^{\text{th}}$  layer

$w_{ij}^{(k)}$

- bias to  $j^{\text{th}}$  unit of  $k^{\text{th}}$  layer  $b_j^{(k)}$

- input-output pairs  $(x_i, d_i)$

- for  $i^{\text{th}}$  training pair:

$$\underline{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix} \quad \underline{d}_i = \begin{bmatrix} d_{i1} \\ \vdots \\ d_{iL} \end{bmatrix} \quad \underline{y}_i = \begin{bmatrix} y_{i1} \\ \vdots \\ y_{iL} \end{bmatrix}$$

$x_n$  -  $n^{\text{th}}$  training input vector

$d_n$  -  $n^{\text{th}}$  target output vector

$y_n$  -  $n^{\text{th}}$  NN output

- input layer - hidden layer - output layer
- standard ways of representing inputs

### Binary classification

- sigmoid output activation  $p(Y=1|X)$
- i.e. posterior class probability

### Multiclass classification

e.g. 5 classes

- use one hot encoding (vector)

### Multiclass network

- target output  $d$  - one hot
- NN output  $y$  - one hot  $\rightarrow$  ideally binary, but in practice probability

### softmax vector activ.

$$z_i = \sum_j w_{ji}^{(n)} y_j^{(n-1)}$$

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

(\*)  $\rightarrow$  not vectors but can be stacked (see diagram)

- exponentiation, then normalisation
- like probability vector

### Problem statement

#### Binary classification

- well balanced training set

#### Multiclass classification

- one hot representation (10d vector)

- given  $(x_1, d_1), \dots, (x_n, d_n)$ ; estimate  $\underline{w}, \underline{b}$  (weights, biases)

Divergence function  $\text{div}(\cdot)$  - R valued  $\vec{d}$

- for loss to be differentiable, divergence must be diff.  
wrt  $\vec{w}$

-  $L_2$  divergence  $\text{div}(\vec{y}, \vec{d}) = \frac{1}{2} \|\vec{y} - \vec{d}\|^2 = \frac{1}{2} \sum (y_i - d_i)^2$   $\nabla_{\vec{y}} \frac{1}{2} \|\vec{y} - \vec{d}\|^2 = \begin{bmatrix} y_1 - d_1 \\ y_2 - d_2 \end{bmatrix}$

Squared Euclidean distance

Divergence function

- binary class

- Based on KL divergence; use Cross Entropy

Q&A - review + fill in graphs

$$\nabla_{\vec{y}} \text{div}(\vec{y}, \vec{d}) = \left\{ \begin{array}{l} \text{if } y_i = 1 \\ \text{if } y_i = 0 \end{array} \right.$$

convergence (comparison of X-entropy,  $L_2$  divergence)

✓ Q: use gaps to  
fill in notes for  
supplement

- check against  
Bishop

- All filled  
in and  
clarified  
in supplement  
notes

Divergence function

Multiclass class

(\*) for N classes

$$d_1, \dots, d_N \rightarrow \text{1 hot}$$

$$y_1, \dots, y_N \rightarrow \text{probability sum to 1}$$

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

$$\nabla_{\vec{y}} \text{div}(\vec{y}, \vec{d}) =$$

- sometimes instead of one-hot:

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ use } \begin{bmatrix} \epsilon \\ \epsilon \\ 1-(k-1)\epsilon \\ \epsilon \\ \epsilon \end{bmatrix}$$

Q&A - label smoothing - supplement

## Training neural nets through gradient descent

Total training error  $E = \frac{1}{T} \sum_{t=1}^T \text{div}(y_t, d_t)$

Send code

init.  $\{w_{ij}^{(k)}\}$

Do:  
for every layer  $k$ , for all  $i, j$  update:-

$$w_{ij}^{(k)} = w_{ij}^{(k)} - \eta \frac{\partial E}{\partial w_{ij}^{(k)}}$$

until  $E$  has converged

$$\text{loss} = \frac{1}{T} \sum_{t=1}^T \text{div}(y_t, d_t)$$

$$\frac{\partial \text{loss}}{\partial w_{ij}^{(k)}} = \frac{1}{T} \sum_{t=1}^T \frac{\partial \text{div}(y_t, d_t)}{\partial w_{ij}^{(k)}}$$

- continue with his presentation; can harmonise rather than recollect.

(i) Go over Bishop ✓

(x)

- uncomfortable with notation

- see supp. notes

- Adopt

(ii) notation ✓

- calculus refresher

scalar diff fn:-

$$y=f(x) \quad \frac{dy}{dx} \quad \Delta y \approx \frac{dy}{dx} \Delta x \quad \text{small } \Delta x$$

Multivar.

$$y=f(x_1, x_2, \dots, x_M)$$

$$\Delta y \approx \frac{\partial y}{\partial x_1} \Delta x_1 + \frac{\partial y}{\partial x_2} \Delta x_2 + \dots \quad \text{sufficiently small } \Delta x_1, \Delta x_2, \dots, \Delta M$$

$$+ \frac{\partial y}{\partial x_M} \Delta x_M$$

engineering def. of chain rule (use above)

$$y=f(g(x))$$

$$z=g(x) \Rightarrow \Delta z = \frac{dg(x)}{dx} \Delta x$$

$$\frac{dy}{dx} = \frac{\partial f}{\partial g} \frac{dg}{dx}$$

$$y=f(z) \Rightarrow \Delta y = \frac{df(z)}{dz} \Delta z = \frac{df}{dz} \frac{dg(x)}{dx} \Delta x$$

$$y = f(g_1(x), g_2(x), \dots, g_m(x))$$

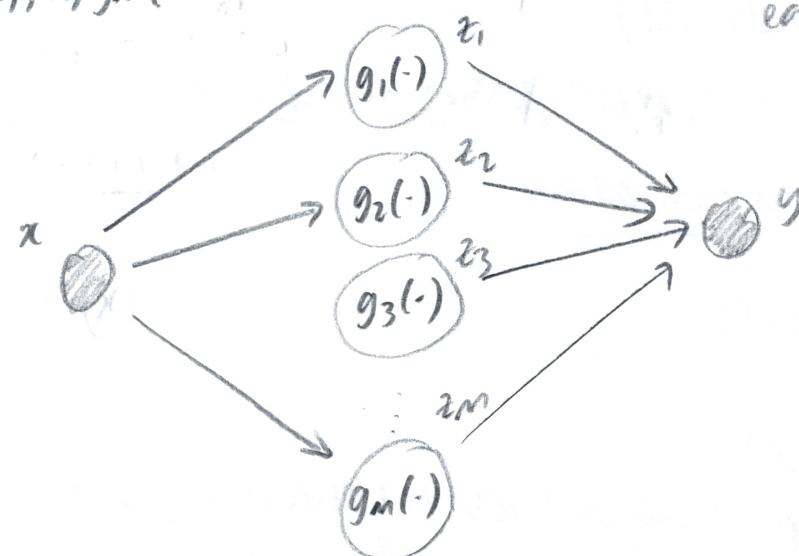
$$\frac{dy}{dx} = \frac{\partial f}{\partial g_1(x)} \frac{dg_1(x)}{dx} + \frac{\partial f}{\partial g_2(x)} \frac{dg_2(x)}{dx} + \dots + \frac{\partial f}{\partial g_m(x)} \frac{dg_m(x)}{dx}$$

- BR: similar illustration

- influence diagram

$$y = f(g_1(x), \dots, g_m(x))$$

(\*)  $x$  affects  $y$  through each of  $g_1, \dots, g_m$



$$(*) \frac{dy}{dx} = \frac{dy}{dz_1} \frac{dz_1}{dx} + \dots + \frac{dy}{dz_m} \frac{dz_m}{dx}$$

- source  $\rightarrow$  destination  
- all paths

- influence diagrams  
give topological way of  
enforcing chain rule

(\*)

$$\frac{d \text{Div}(Y, d)}{d w_{ij}^{(k)}} - \text{compute}$$

(\*) BR: each neuron is 2 computations:- affine  $\rightarrow$  activation ✓  
- computing derivative for a single input      (\*) some rotational  
inconsistencies across slides

- & 2 - affine combination

$y$  - output after activation applied  
to affine combination

(co-influence)

(i): tracing effect of  $\frac{d \text{Div}(Y, d)}{d w_{ij}^{(k)}}$

imagine a route, from  
where you tweak parameter  
to disperse through the topology

(\*) derivative computation requires intermediate and final output values of network in response to input.

scalar formulation  $\rightarrow$  ✓

- forward propagation - very clear

$$y_N \quad y^{(N)} = f_N(z^{(N)}) \quad - \text{vector activation e.g. softmax}$$

(\*) issue with algo on slides (possibly included)

- All  $z_j^{(h)}$  need to be computed before passing through final activation
- As  $y^{(N)} = f_N(z^{(N)})$  is a vector activation

forward pass pseudocode checked.

computing derivatives (backprop)

start from final layer, compute derivatives backward layer by layer

- Backpropagation  $\rightarrow$  see slides

Backpropagation

high level

- illustration for a single training instance (actual = multiple training instances)

$$(*) \text{compute } \frac{\partial \text{Div}(Y, d)}{\partial y_i} = \frac{\partial \text{Div}(Y, d)}{\partial y_i^{(N)}} \quad \forall i$$

$$(*) \text{compute } \frac{\partial \text{Div}}{\partial z_i^{(N)}} = \underbrace{\frac{\partial \text{Div}}{\partial y_i^{(N)}}}_{\text{already computed}} \underbrace{\frac{\partial y_i^{(N)}}{\partial z_i^{(N)}}}_{\rightarrow f'_N(z_i^{(N)})} - \text{derivative of activation function computed in forward pass.}$$

$$\Rightarrow \frac{\partial \text{Div}}{\partial z_i^{(N)}} = f'_N(z_i^{(N)}) \frac{\partial \text{Div}}{\partial y_i^{(N)}} \quad (*)$$

$$\frac{\partial \text{Div}}{\partial w_{ii}^{(N)}} = \underbrace{\frac{\partial z_i^{(N)}}{\partial w_{ii}^{(N)}}}_{\sim} \frac{\partial \text{Div}}{\partial z_i^{(N)}} = y_i^{(N-1)} \frac{\partial \text{Div}}{\partial z_i^{(N)}} \quad (*) \quad (\star) \text{ bias term } y_0^{(N-1)} = 1$$

- already computed in previous backprop step

$$= y_i^{(N-1)} \text{ as } z_i^{(N)} = w_{ii}^{(N)} y_i^{(N-1)} + \text{(other terms)} \quad (\text{computed in forward pass})$$


---


$$\frac{\partial \text{Div}}{\partial y_i^{(N-1)}} = \sum_j \underbrace{\frac{\partial z_j^{(N)}}{\partial y_i^{(N-1)}}}_{\sim} \frac{\partial \text{Div}}{\partial z_j^{(N)}} \quad (*) \quad (\star) \text{ remember topological ordering; how it relates here to the chain rule}$$

- already computed

$\rightarrow = w_{ij}^{(N)}$  because  $z_j^{(N)} = w_{ij}^{(N)} y_i^{(N-1)} + \dots$  other terms not containing  $y_i^{(N-1)}$

$$= \sum_j w_{ij}^{(N)} \frac{\partial \text{Div}}{\partial z_j^{(N)}}$$

### • W A(i): Backpropagation clarity

- the three (\*) types of derivatives summarise the process
- All are instances of the same principle  $\star \star$   
(i.e. chain rule, efficient use of computation)

• Backward pass pseudocode  $\rightarrow$  W A(ii) check you understand

(\*) similarity of forward and backward pass algos.

