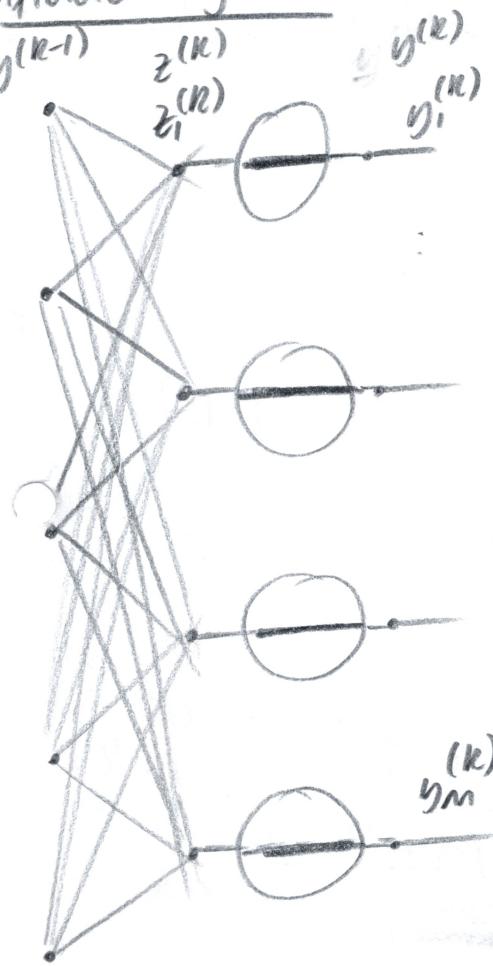


4.5 Backprop calc - review

(*) clarify distinction between scalar and vector act.

Influence diagrams (*)



scalar activation

- each z_i influences
one y_i

- so in diagram

$z_i^{(k)}$ only influences $y_i^{(k)}$

(*) no. outputs in layer $y^{(k)}$

does not have to equal no. inputs in layer $z^{(k)}$

scalar activation

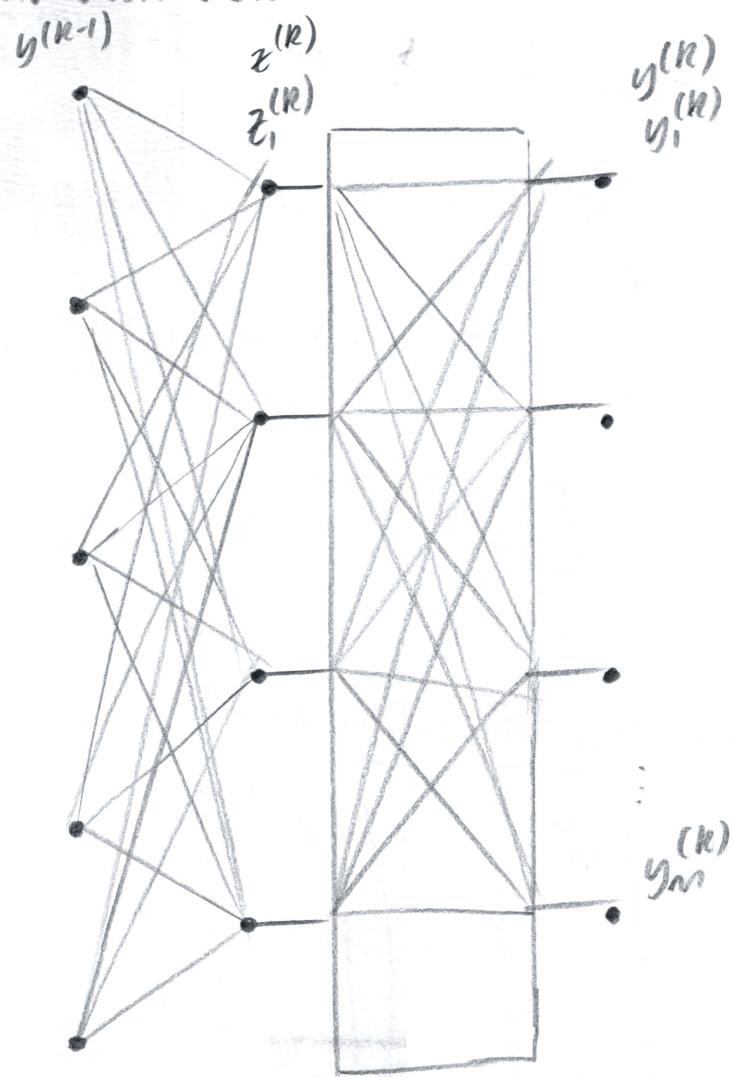
$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \frac{\partial \text{Div}}{\partial y_i^{(k)}} \frac{dy_i^{(k)}}{dz_i^{(k)}}$$

- derivative of error wrt
mpnt is product

vector activation

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} \frac{dy_j^{(k)}}{\partial z_i^{(k)}}$$

- derivative of error wrt
input - sum of partial deriv
regardless of no. outputs $y_j^{(k)}$



vector activation

- each z_i influences
(multiple) all y_1, y_2, \dots, y_m

- in diagram, each $z_i^{(k)}$ influences
(multiple) all $y_1^{(k)}, y_2^{(k)}, \dots, y_m^{(k)}$

(*) softmax \rightarrow forward HWI pt 1 very confusing (this will provide clarity)

$$y_i^{(k)} = \frac{\exp(z_i^{(k)})}{\sum_j \exp(z_j^{(k)})} \quad \frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} \frac{\partial y_j^{(k)}}{\partial z_i^{(k)}}$$

\rightarrow see 24.5 step notes
for derivatn.

$$\frac{\partial y_j^{(k)}}{\partial z_i^{(k)}} = \begin{cases} y_i^{(k)}(1-y_i^{(k)}) & \text{if } i=j \text{ (non-diag influence path)} \\ -y_i^{(k)}y_j^{(k)} & \text{if } i \neq j \text{ (diagonal infl. path)} \end{cases}$$

$$\frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} y_i^{(k)} (\delta_{ij} - y_j^{(k)}) \quad \begin{array}{l} \text{(using compression of not} \\ \text{using Kronecker deltas)} \end{array} \quad (I)$$

$\delta_{ij} = 1 \text{ if } i=j; \delta_{ij} = 0 \text{ if } i \neq j$

(*) vector activations (vector form)

$$y^{(k)} = \begin{bmatrix} y_1^{(k)} \\ \vdots \\ y_M^{(k)} \end{bmatrix} = f \left(\begin{bmatrix} z_1^{(k)} \\ \vdots \\ z_D^{(k)} \end{bmatrix} \right)$$

- M - no. of output units in layer k (width of layer k)
- D - no. of pre-activation affine combos in layer k

(*) examine how using vector activation for softmax differs from stand backprop scalar form pseudocode
(only change in output act. heading)
i.e. use (I) in output layer part of pseudocode

(*) subgradients

- Multipledef.
 - If $f: \mathcal{U} \rightarrow \mathbb{R}$ is a real-valued convex function defined on a convex open set in the Euclidean space \mathbb{R}^n ; a vector v is called a subgradient at a point x_0 in \mathcal{U} if for any x in \mathcal{U} one has
- $$f(x) - f(x_0) \geq v^T(x - x_0)$$

Multivariable exp. (in 2D space) $I \subset \mathbb{R}$

- At points; function not $f: I \rightarrow \mathbb{R}$
differentiable

- For any x_0 in the domain off; one can draw a line which goes through point $(x_0, f(x_0))$

- This point β is greater either touching or below the graph of f .

- The slope of such a line (there are/may be many) is the subderivative

- A subgradient is a generalisation to functions of multiple variables

(*) claims I can't fully understand:-

i) subgradient is direction in which function is guaranteed to increase

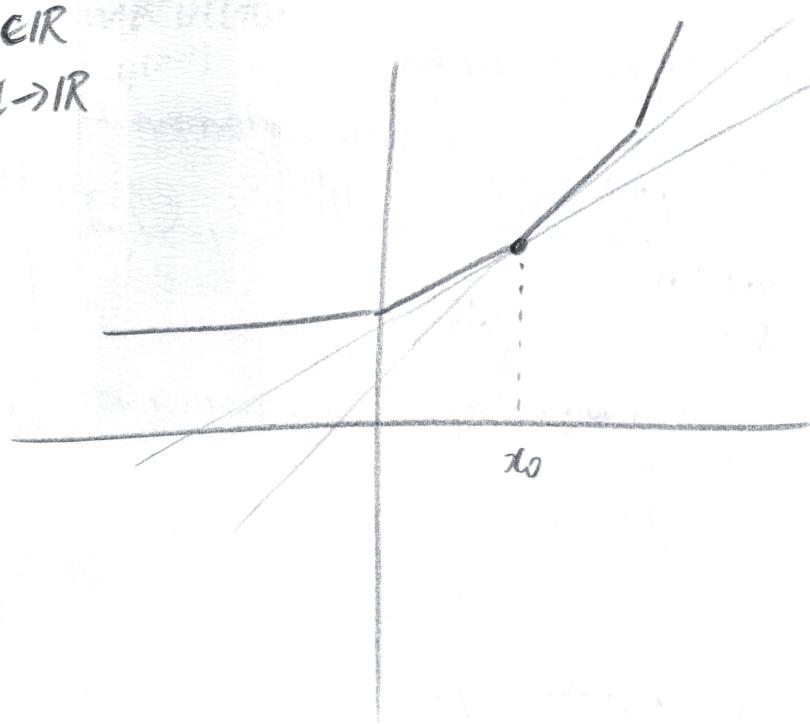
ii) if function differentiable at x_0 ; subgradient is gradient

(*) subgradients used for ReLU prime

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

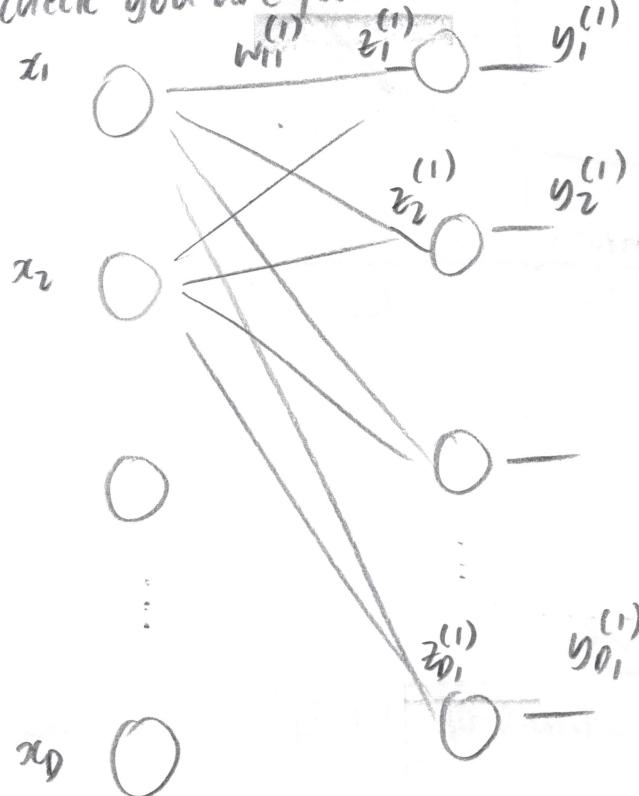
(*) subgradients + max \rightarrow revisit with CNNs

• Note (*) Note where the use of vector activations for non output layers (or output layers) affects the code for backprop



(*) NN vector formulation

- check you are familiar with notation:-



$$W_R = \begin{bmatrix} (k) & (k) & & (k) \\ w_{11} & w_{21} & \dots & w_{D_{k-1}1} \\ (k) & (k) & \dots & w_{D_k-12} \\ w_{12} & w_{22} & \dots & \\ \vdots & & & \\ (k) & (k) & & (k) \\ w_{1D_k} & w_{2D_k} & \dots & w_{D_{k-1}D_k} \end{bmatrix}$$

(*) Each column vector in W_R contains the weights from all units in $(k-1)^{th}$ layer to 1st unit of k^{th} layer

- W_R - weights matrix from $(k-1)$ to k^{th} layer

(*) Note f is a vector function of a vector input
(*) Vector notation for relevant transformations in k^{th} layer of a NN.

(*) forward pass pseudocode (vector form)

• set $y_0 = \underline{x}$

• $\underline{y} = y_N$

• for layer $k=1, \dots, N$

$$\cdot \underline{z}_k = W_k y_{k-1} + b_k$$

$$\cdot y_k = f_k(\underline{z}_k)$$

- note output is a composition of affine, activations layer by layer.
- see 2 representations in ut.5 notes

(*) Jacobians matrices of vector-valued functions are matrices containing all first-order partial derivatives ⊗

- vector-valued fn: $f: \mathbb{R}^D \rightarrow \mathbb{R}^M$ (in context of ut.5 notes not.)

$$\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_D \end{pmatrix} \quad \underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_M \end{pmatrix} = f\left(\begin{pmatrix} z_1 \\ \vdots \\ z_D \end{pmatrix}\right)$$

- BR defines counterintuitively (see ut.5 notes you made)

$J_{\underline{y}}(\underline{z})$ - Jacobian of \underline{z} wrt \underline{x} (!) (does he change this later?)

(*) Note Jacobians can describe

- i) scalar activations (diagonal Jacobian)
- ii) vector activations (fully populated Jacobian)

(*) vector calc. notation using gradients, Jacobians

(*) Affine fn: $\underline{z} = \underline{W}\underline{y} + \underline{b}$

$$J_{\underline{z}}(\underline{y}) = \underline{W} \quad \checkmark$$

(*) chain rule our vector functions of vector inputs $y = f(g(\underline{z})) \quad \underline{z} = g(\underline{x}) \quad y = f(\underline{z})$
 $J_{\underline{y}}(\underline{x}) = J_{\underline{y}}(\underline{z}) J_{\underline{z}}(\underline{x}) \quad \checkmark$

(*) chain rule combining Jacobians and gradients $D = f(g(\underline{x})) \quad \underline{z} = g(\underline{x}) \quad D = f(\underline{z})$
 $\nabla_{\underline{x}} D = \nabla_{\underline{z}} D J_{\underline{z}}(\underline{x}) \quad \checkmark$

- scalar funcs of vector inputs

(*) scalar
fns of affine

$$D = f(Wy + b) \quad \underline{z} = Wy + b \quad D = f(\underline{z})$$

fns:-

$$\nabla_y D = \nabla_{\underline{z}}(D) J_{\underline{z}}(y) = \nabla_{\underline{z}} D W$$

$$(*) \nabla_b D = \nabla_{\underline{z}} D, \quad \nabla_w D = y \nabla_{\underline{z}}(D)$$

(*) gradient descent algorithm pseudocode (vector form)

- Initialise weights W_1, W_2, \dots, W_R
biases

• Do :

- for every layer R , compute:-

$$\cdot \nabla_{W_R} Err = \frac{1}{T} \sum_t \nabla_{W_R} \text{Div}(y_t, d_t)$$

$$\cdot W_R = W_R - \eta \nabla_{W_R} Err^T$$

(compute gradient of error
as average of gradients
or mini divergence
over all instances.)

- update param.

• until

err has converged

(*) convergence clarity \rightarrow next supp. notes

(*) NN training algo pseudocode (vector form)

- Initialise all weights, biases $(W_1, b_1, W_2, b_2, \dots, W_R, b_R)$

• Do :-

- $Err = 0$
- $\forall R$ layers, initialise $\nabla_{W_R} Err = 0, \nabla_{b_R} Err = 0$

- $\forall t=1, \dots, T$ instances, comp. (FORWARD)
 - divergence $\text{Div}(y_t, d_t)$

$$\cdot Err += \text{Div}(y_t, d_t)$$

(BACKWARD)

- $\forall R$ layers, comp.

$$(i) - \nabla_{y_R} \text{Div} = \nabla_{z_{R+1}} \text{Div} W_R$$

$$(ii) - \nabla_{z_R} \text{Div} = \nabla_{y_R} \text{Div} J_{y_R}(z_R)$$

$$\text{iii) } \nabla_{W_R} \text{Div}(Y_t, d_t); \nabla_{b_R} \text{Div}(Y_t, d_t)$$

$$\text{iv) } \nabla_{W_R} \text{Err} \triangleq \nabla_{W_R} \text{Div}(Y_t, d_t); \nabla_{b_R} \text{Err} \triangleq \nabla_{b_R} \text{Div}(Y_t, d_t)$$

• VR, update:-

$$w_k = w_k - \frac{\eta}{\gamma} (\nabla_{W_R} \text{Err})^\top; b_k = b_k - \frac{\eta}{\gamma} (\nabla_{b_R} \text{Err})^\top$$

(*) do you understand
comput. flow

• until Err has converged

(*) An individual look at each step of backprop. (view in combo with slides)

- $\nabla_y \text{Div}$ - dimensionality varies on divergence/problem setting?

as a generic way of specifying div, grad, Jacob

- ○ - (i)

$$\nabla_{z_N} \text{Div} = \nabla_y \text{Div} \cdot \nabla_{z_N} Y = \nabla_y \text{Div} J_y(z_N)$$

(ii)

$$\nabla_{y_{N-1}} \text{Div} = \nabla_{z_N} \text{Div} \cdot \nabla_{y_{N-1}} z_N = \nabla_{z_N} \text{Div} W_N$$

intermediate step; store what we're interested in

- - - - (non output layers)

$$\begin{aligned} \nabla_{z_{N-1}} \text{Div} &= \nabla_{y_{N-1}} \text{Div} \cdot \nabla_{z_{N-1}} y_{N-1} \\ &= \nabla_{y_{N-1}} \text{Div} \cdot J_{y_{N-1}}(z_{N-1}) \end{aligned}$$

$$(*) \nabla_{w_N} \text{Div} = y_{N-1} \nabla_{z_N} \text{Div}$$

$$\nabla_{b_N} \text{Div} = \nabla_{z_N} \text{Div}$$

$$\nabla_{y_{N-2}} = \nabla_{z_{N-1}} \text{Div} \cdot \nabla_{y_{N-2}} z_{N-1}$$

(ii)

$$= \nabla_{z_{N-1}} \text{Div} W_{N-1}$$

store

$$\nabla_{w_{N-1}} \text{Div} = y_{N-2} \nabla_{z_{N-1}} \text{Div}$$

- - - - (up to inputs)

$$\nabla_{z_1} \text{Div} = \nabla_{y_1} \text{Div} J_{y_1}(z_1)$$

(*)

$$\nabla_{b_{N-1}} \text{Div} = \nabla_{z_{N-1}} \text{Div}$$

store

$$\nabla_{w_1} \text{Div} = \underline{x} \nabla_{z_1} \text{Div}$$

(iii) note

$$\text{i) note } \nabla_{z_N} Y = J_y(z_N)$$

(gradient of vector fn wrt vector is same as Jacobian)

$$\nabla_{b_1} \text{Div} = \nabla_{z_1} \text{Div}$$

ii) Jacobian diagonal for scalar act.

$$\text{iii) } y_0 = z$$

