

4.5 - Backpropagation calculation

- continue S.2019 slides 14 from ②
- focus on understanding
- forward pass: present input; compute affine combos z_j and activations
- backward pass: back through net, use forward pass computations to compute derivatives.

Special cases

- assumed some things

- 1. 'Independent' computation
- 2. Affine combination of inputs
- 3. diff. activations

(N): Think in terms
of topological influence
as routes through
a topology

special case

- 1. vector activation vs (scalar activation)

(*) diagram

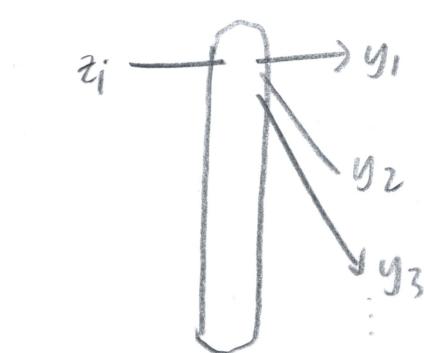
- renders activation multivariate vector function
- Has consequences for 'influence'

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

- Softmax activation is an example of a vector activation

softmax activ.: (as special case vector act.)

$$y_i = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (\text{softmax})$$



$$\frac{dy_i}{dz_i} = \frac{d}{dz_i} e^{z_i} \cdot \frac{1}{\sum_j e^{z_j}} + e^{z_i} \frac{d}{dz_i} \left(\frac{1}{\sum_j e^{z_j}} \right)$$

$$= b \frac{e^{z_i}}{\sum_j e^{z_j}} + \frac{e^{z_i}(-e^{z_i})}{(\sum_j e^{z_j})^2} \rightarrow y_i - y_i^2$$

$$= y_i - \frac{e^{z_i}}{\sum_j e^{z_j}} \cdot \frac{e^{z_i}}{\sum_j e^{z_j}}$$

$$\frac{dy_i}{dz_i} = y_i - y_i^2 \quad (\text{non diagonal path})$$

$$\frac{dy_i}{dz_k} = e^{z_i} \frac{d}{dz_k} \frac{1}{\sum_j e^{z_j}} = \frac{e^{z_i} (-e^{z_k})}{(\sum_j e^{z_j})^2} = -y_i y_k \quad (\text{diagonal path})$$

- can compute derivative of each of y_i s to any z_j
- softmax output activation derivative varies depending on whether diagonal or non-diagonal path.
- use Kronecker delta to compress notation

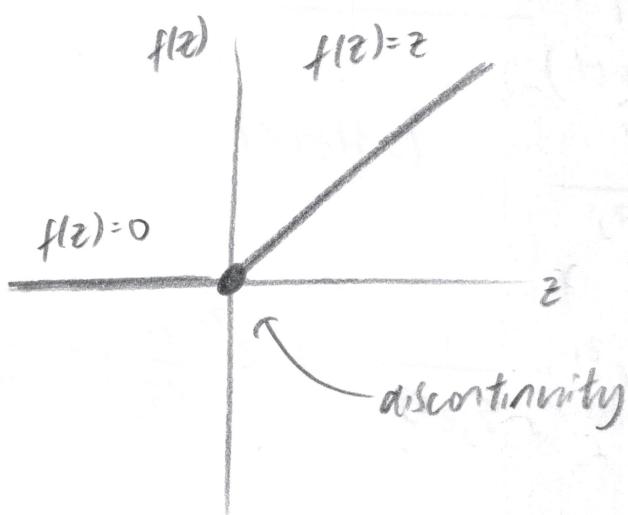
$$\Delta \frac{\partial \text{Div}}{\partial z_i^{(k)}} = \sum_j \frac{\partial \text{Div}}{\partial y_j^{(k)}} y_j^{(k)} (\delta_{ij} - y_j^{(k)})$$

②③④ - multiplicative networks - review

special case 3 - differentiable activ.

Assumed activations diff. everywhere

ReLU



(*) Subgradients \rightarrow review

① Q: check you understand definition

- ReLU prime \rightarrow use subgradients
to get around discontinuity

- Subgradient and max

- max \rightarrow use subgradients

Q: I am confused by inclusion of averaging factor $\frac{1}{\|x\|_2^2}$ $\frac{1}{T}$ $\frac{1}{N}$ m

actual error

Q: Acknowledges 'average error/loss' rather than actual (?)
→ is this really used?

$$\underline{w} \leftarrow \underline{w} - \eta \nabla_{\underline{w}} \text{Err}$$

Training via backprop algo pseudocode

→ deserves place in supplement notes

move to vector formulation Q: WA3 → review notes

Q: (*) indicated this;
more compact

- simpler conceptually

- fast matrix libraries (e.g. numpy)

(*) $\underline{z}_k = \underline{w}_k \underline{y}_{k-1} + \underline{b}_k$

$$\underline{y}_k = f_k(\underline{z}_k)$$

Q: This is key; everything
before is fine;
but if you want to commit
anything to memory

(*) $\underline{y} = f_N(f_{N-1}(\dots f_2(w_2 f_1(w_1 \underline{x} + \underline{b}) + \underline{b}_2) \dots) + \underline{b}_N)$

init $y_0 = \underline{x}$

$k=1, \dots, N$

$$\underline{z}_k = \underline{w}_k \underline{y}_{k-1} + \underline{b}_k$$

$$y_k = f_k(z_k)$$

output: $\underline{y} = y_N$

(*) $\text{Div}(y, d)$

$$\cdot \text{Div} [f_N(w_N f_{N-1}(\dots f_2(w_2 f_1(w_1 x + b_1) + b_2) \dots) + b_N), d]$$

Jacobians \rightarrow derivative of vector fn. with respect to vector input

$$y = f(z)$$

i.e.

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = f \left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_D \end{bmatrix} \right)$$

perturb-calc rule

$$\Delta y = J_y(z) \Delta z$$

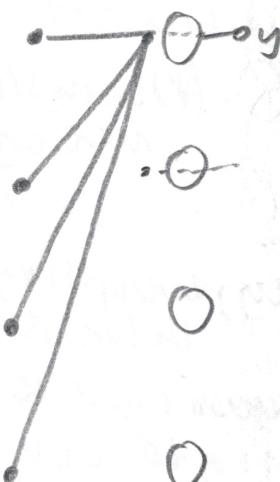
$$J_y(z) = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_D} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & \dots & \frac{\partial y_2}{\partial z_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial z_1} & \frac{\partial y_m}{\partial z_2} & \dots & \frac{\partial y_m}{\partial z_D} \end{bmatrix}$$

- Remember
know

Jacobian for scalar vs for vector activations

(*) notice distinction

scalar activation

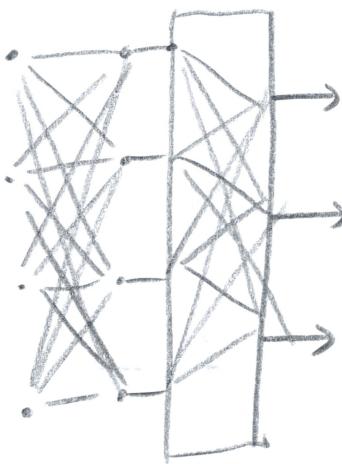


$$J_y(z) = \begin{bmatrix} \frac{dy_1}{dz_1} \\ \vdots \\ \frac{dy_D}{dz_D} \end{bmatrix}$$

$$y_i = f(z_i)$$

$$J_y(z) = \begin{bmatrix} f'(z_1) \\ \vdots \\ f'(z_D) \end{bmatrix}$$

vector activations



$$J_y(z) = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \dots & \frac{\partial y_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial z_1} & \dots & \frac{\partial y_m}{\partial z_n} \end{bmatrix}$$

⑥ Check Jacobian notation (that you understand)
 Remember → highly targeted learning

$J_z(z)$ - Jacobian of z not x (BR notation)

- chain rule applied to Jacobians

- a cosmetic difference

⑦ use Peter/understand BR way of understanding chain rule

w/ \odot - Jacobian / scalar fm of vector m / - familiarity

special cases ⑧

- scalar ⑨ ⑩: dimensionality ??

- dimensionality for higher-order calc:

⑪ requires supplement

(*) Review and ensure you understand
backprop in vector form; including dimensionality considerations

- ①: once you understand
3 (*) calculations in vector-matrix form; repeated a, recursive
application of some formula
-
- backward pass pseudocode
- check
 - most difficult lecture * (BR)
 - Review; review; review
 - BR: vector variant of backprop - you must be fluent.
-
- convergence of SGD }
generalisation
intep of output } next lecture