

Backprop Ex#1

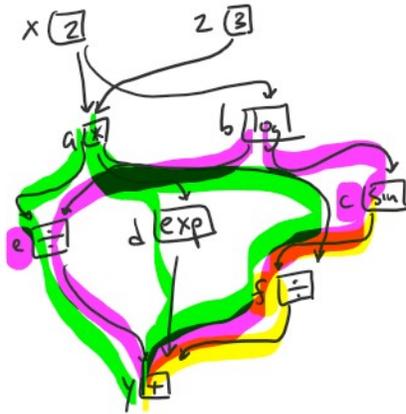
$$y = f(x, z) = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Forward Computation

Given $x=2, z=3$

- $a = xz$
- $b = \log(x)$
- $c = \sin(b)$
- $d = \exp(a)$
- $e = a/b$
- $f = c/a$
- $y = d + e + f$

Computation Graph



Backward Computation

$$g_y = \frac{dy}{dy} = 1$$

$$g_f = \frac{dy}{df} = 1, g_e = \frac{dy}{de} = 1, g_d = \frac{dy}{dd} = 1$$

$$g_c = \frac{dy}{dc} = \frac{dy}{df} \frac{df}{dc} = (g_f) \left(\frac{1}{a} \right)$$

$$g_b = \frac{dy}{db} = \frac{dy}{de} \frac{de}{db} + \frac{dy}{dc} \frac{dc}{db}$$

$$= (g_e) \left(-\frac{a}{b^2} \right) + (g_c) (\cos(b))$$

$$g_a = \frac{dy}{da} = \frac{dy}{de} \frac{de}{da} + \frac{dy}{db} \frac{db}{da} + \frac{dy}{dc} \frac{dc}{da}$$

$$= (g_e) \left(\frac{1}{b} \right) + (g_d) (\exp(a)) + (g_f) \left(\frac{c}{a^2} \right)$$

$$g_x = \frac{dy}{dx} = (g_a)(z) + (g_b) \left(\frac{1}{x} \right)$$

$$g_z = \frac{dy}{dz} = (g_a)(x)$$

Updates for Backprop

$$g_x = \frac{dy}{dx} = \sum_{k=1}^K \frac{dy}{du_k} \frac{du_k}{dx}$$

$$= \sum_{k=1}^K (g_{u_k}) \left(\frac{du_k}{dx} \right)$$

Efficient b/c of...

- █ reuse in forward comp.
- █ reuse in backward comp.

Neural Network Training

- Consider a 2-hidden layer NN
- params we $\vec{\Theta} = [\alpha^{(1)}, \alpha^{(2)}, \beta]$
- SGD Training:

Iterate until convergence:

- ① Sample $i \in \{1, \dots, N\}$:
- ② Compute gradient by backprop:

Background:

$$\nabla_{\vec{a}, \vec{b}} J(\vec{a}, \vec{b}) = \nabla_{\vec{a}, \vec{b}} J(\vec{a}, \vec{b})$$

$$\nabla_{\vec{a}} J(\vec{a}, \vec{b}) = \begin{bmatrix} dJ/da_1 \\ dJ/da_2 \\ \vdots \\ dJ/da_k \end{bmatrix} \quad |a| = k$$

- Sample $i \in \{1, \dots, N\}$
- Compute gradient by backprop:

$$\left. \begin{aligned} g_{\alpha^{(1)}} &= \nabla_{\alpha^{(1)}} J^{(i)}(\vec{\theta}) \\ g_{\alpha^{(2)}} &= \nabla_{\alpha^{(2)}} J^{(i)}(\vec{\theta}) \\ g_{\beta} &= \nabla_{\beta} J^{(i)}(\vec{\theta}) \end{aligned} \right\} J^{(i)}(\vec{\theta}) = \ell(h_{\theta}(\vec{x}^{(i)}), y^{(i)})$$

- Update parameters:

$$\begin{aligned} \alpha^{(1)} &\leftarrow \alpha^{(1)} - \gamma g_{\alpha^{(1)}} \\ \alpha^{(2)} &\leftarrow \alpha^{(2)} - \gamma g_{\alpha^{(2)}} \\ \beta &\leftarrow \beta - \gamma g_{\beta} \end{aligned}$$

Recall: $\vec{\theta} \leftarrow \vec{\theta} - \gamma g_{\vec{\theta}}$

Backprop Ex#2: for NN

- Given:
- Dec. fn. $\hat{y} = h_{\theta}(\vec{x}) = \sigma((\alpha^{(3)})^T \sigma((\alpha^{(2)})^T \sigma((\alpha^{(1)})^T \vec{x})))$
 - Loss fn. $J = \ell(\hat{y}, y^*) = -(y^* \log(\hat{y}) + (1-y^*) \log(1-\hat{y}))$
 - Training ex. (\vec{x}, y^*)

Forward Comp.

Given $\vec{x}, y^*, \alpha^{(1)}, \alpha^{(2)}, \alpha^{(3)}$

$$z^{(0)} = \vec{x}$$

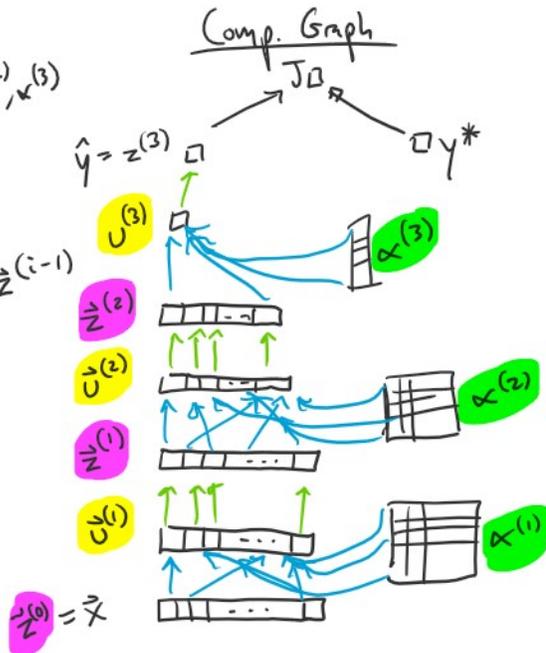
for $i=1, 2, 3$:

$$u^{(i)} = (\alpha^{(i)})^T z^{(i-1)}$$

$$z^{(i)} = \sigma(u^{(i)})$$

$$\hat{y} = z^{(3)}$$

$$J = \ell(\hat{y}, y^*)$$



Backward Comp.

$$g_y = [1]$$

$$g_y = -\left(\frac{y^*}{\hat{y}} - \frac{1-y^*}{1-\hat{y}}\right)$$

for $i=3, 2, 1$:

$$\begin{aligned} g_{z^{(i)}} &= \dots \\ g_{z^{(i-1)}} &= \dots \\ g_{\alpha^{(i)}} &= \dots \end{aligned}$$

$$g_{\vec{x}} = g_{z^{(0)}}$$

HWS

Vector Chain Rule

$$y \in \mathbb{R} \quad \vec{u} \in \mathbb{R}^N \quad \vec{x} \in \mathbb{R}^P$$

$$\frac{\partial y}{\partial \vec{x}} = \left(\left(\frac{\partial y}{\partial \vec{u}} \right)^T \left(\frac{\partial \vec{u}}{\partial \vec{x}} \right)^T \right)^T$$

$$\left[\frac{\partial y}{\partial \vec{x}} \right] = \sum_{n=1}^N \frac{\partial y}{\partial u_n} \frac{\partial u_n}{\partial \vec{x}}$$

$$\begin{aligned} \left[\frac{dy}{dx} \right]_{P \times 1} &= \left(\left[\frac{dy}{d\vec{u}} \right]_{N \times 1}^T \left[\frac{d\vec{u}}{dx} \right]_{P \times N}^T \right)^T \\ &= \frac{d\vec{u}}{dx} \frac{dy}{d\vec{u}} \end{aligned}$$

$$\begin{aligned} \left[\frac{dy}{dx} \right]_p &= \sum_{n=1}^N \frac{dy}{du_n} \frac{du_n}{dx_p} \\ &= \sum_{n=1}^N \left[\frac{dy}{d\vec{u}} \right]_n \left[\frac{d\vec{u}}{dx} \right]_{pn} \end{aligned}$$