

Section B: Backprop

Sunday, October 9, 2022 8:44 PM

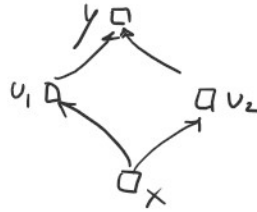
Chain Rule

Def #1: $y = f(u)$
 $u = g(x)$



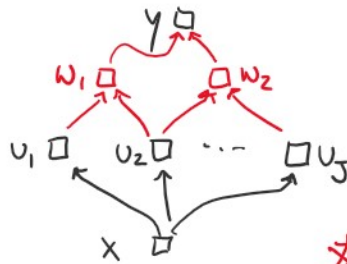
$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

Def #2: $y = f(u_1, u_2)$
 $u_1 = g_1(x)$
 $u_2 = g_2(x)$



$$\frac{dy}{dx} = \frac{dy}{du_1} \frac{du_1}{dx} + \frac{dy}{du_2} \frac{du_2}{dx}$$

Def #3: $y = f(u)$
 $u = g(x)$



$$\frac{dy}{dx} = \sum_{j=1}^J \frac{dy}{du_j} \frac{du_j}{dx}$$

★ holds for any such intermediate vars. w_i

Backprop Ex #1

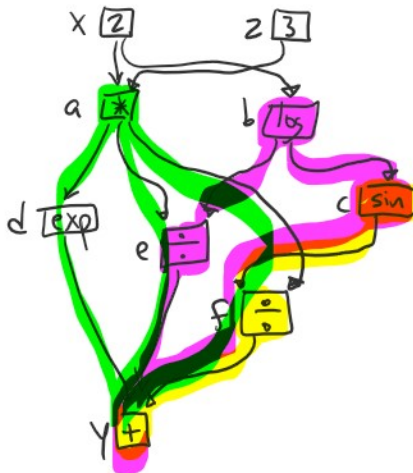
$$y = f(x, z) = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

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

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

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

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

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

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

$$g_a = \frac{dy}{da} = \frac{dy}{dd} \frac{dd}{da} + \frac{dy}{de} \frac{de}{da} + \frac{dy}{df} \frac{df}{da}$$

Updates for Backprop

Updates for Backprop

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

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

Efficient b/c

- reuse of forward comp.
- reuse of backward comp.

$$g_a = \frac{dy}{da} = \frac{dy}{dd} \frac{dd}{da} + \frac{dy}{de} \frac{de}{da} + \frac{dy}{df} \frac{df}{da}$$

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

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

$$g_z = (g_a)(x)$$

Neural Network Training

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

Iterate until convergence:

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

$$g_{\alpha^{(1)}} = \nabla_{\alpha^{(1)}} J^{(i)}(\theta)$$

$$g_{\alpha^{(2)}} = \nabla_{\alpha^{(2)}} J^{(i)}(\theta)$$

$$g_{\beta} = \nabla_{\beta} J^{(i)}(\theta)$$

Aside:

$$\nabla_{\vec{a}, \vec{b}} J(\vec{a}, \vec{b}) = \begin{bmatrix} \partial J(\vec{a}, \vec{b}) / \partial a_1 \\ \vdots \\ \partial J(\vec{a}, \vec{b}) / \partial a_k \end{bmatrix}$$

$$J^{(i)}(\theta) = \ell(h_{\theta}(\vec{x}^{(i)}), y^{(i)})$$

- Update parameters

$$\alpha^{(1)} \leftarrow \alpha^{(1)} - \delta g_{\alpha^{(1)}}$$

$$\alpha^{(2)} \leftarrow \alpha^{(2)} - \delta g_{\alpha^{(2)}}$$

$$\beta \leftarrow \beta - \delta g_{\beta}$$

left out intercept terms

$\vec{\beta}$ $\vec{z}^{(2)}$ $\vec{z}^{(1)}$

Backprop Ex#2 for NN

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

Loss fn. $J = \ell(y, y^*) = -(y^* \log(\hat{y}) + (1-y^*) \log(1-\hat{y}))$

③ Training ex. (\vec{x}, y^*)

Forward Comp.

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

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

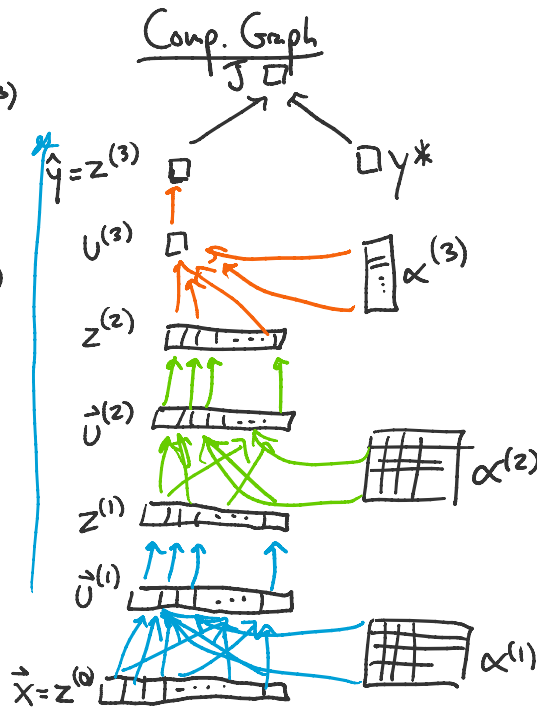
for $i=1, 2, 3$:

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

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

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

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



Backward Comp.

$g_y = [1]$

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

for $i=3, 2, 1$:

$g_{\vec{u}^{(i)}} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$

$g_{z^{(i-1)}} = \dots$

$g_{\alpha^{(i)}} = \dots$

$g_x = g_{z^{(0)}}$

HWS

Vector Chain Rule

$$\underbrace{\frac{dy}{d\vec{x}}}_{P \times 1} = \underbrace{\left(\underbrace{\left(\frac{dy}{d\vec{u}} \right)^T}_{N \times 1} \underbrace{\left(\frac{d\vec{u}}{d\vec{x}} \right)^T}_{P \times N} \right)^T}_{1 \times P}}_{P \times 1}$$

$$= \frac{d\vec{u}}{d\vec{x}} \frac{dy}{d\vec{u}}$$