

Backpropagation

(Chain Rule on Graphs)

Chain Rule:

$$h(x) = f(g(x))$$

$$\frac{dh}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$= f'(g(x))g'(x)$$

Ex: $h(x) = e^{x^2}$

$$f(u) = e^u \quad f'(u) = e^u$$

$$g(x) = x^2 \quad g'(x) = 2x$$

$$h'(x) = f'(g(x))g'(x)$$

$$= e^{g(x)} \cdot g'(x) = e^{x^2} \cdot 2x$$

$$h(x) = f(g(t(x)))$$

$$\frac{dh}{dx} = h'(x) = f'(g(t(x))) \cdot g'(t(x)) \cdot t'(x)$$

$$h(x) = e^{(2x+3)^2}$$

$$f(x) = e^x \quad f'(x) = e^x$$

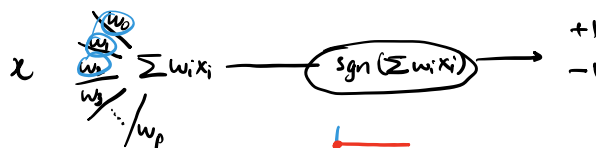
$$g(x) = x^2 \quad g'(x) = 2x$$

$$t(x) = 2x+3 \quad t'(x) = 2$$

$$h'(x) = e^{(2x+3)^2} \cdot 2(2x+3) \cdot 2$$

$$= 6(2x+3)e^{(2x+3)^2}$$

Perceptrons:



Gradient Descent:

$$\min_x f(x)$$

$$x^{(t+1)} = x^{(t)} - \eta f'(x^{(t)})$$

(logistic)

Sigmoid unit/perceptron

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

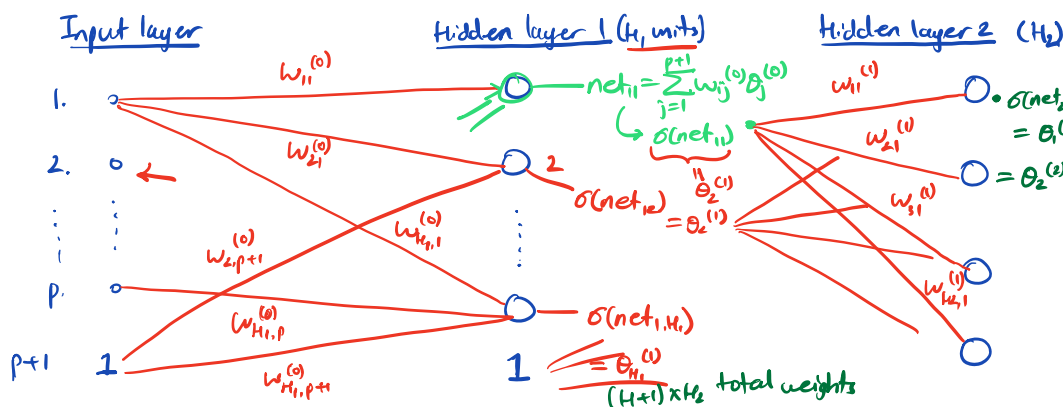
$$\sigma'(x) = \sigma(x)(1-\sigma(x))$$



→ ReLU/tanh

MLP / Neural Network:

$$D = \{(x, y) : x_i \in \mathbb{R}^p, y_i \in \mathbb{R}^k\}$$



$$x_i = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$$

$w_{ji}^{(k)}$
= connecting
ith input to
jth node
in kth layer

$$x_i = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$$

$(p+1) \times H_1$ total weights

Hidden layer 2

0

0

⋮

0

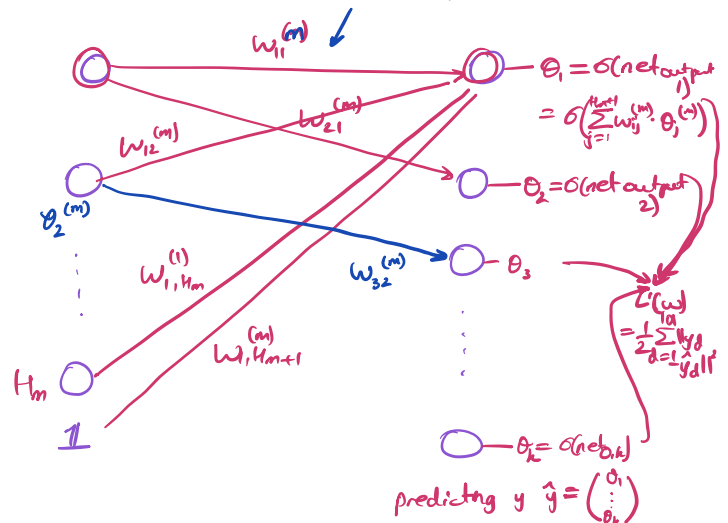
1

$i=2$

$w_{32}^{(2)}$

Hidden layer M (H_M units)

Output layer (k units)



$$L(w) = \frac{1}{2} \sum_{d=1}^{10} \|y_d - \hat{y}_d\|^2$$

$$= \frac{1}{2} \sum_{d=1}^{10} \sum_{k=1}^k (y_{dk} - \hat{\theta}_{dk})^2$$

w = vector of all weights in network:

$$\hat{\theta}^{(1)} = \sigma(w^{(0)} \cdot \theta^{(0)})$$

$$\theta^{(1)} = \begin{pmatrix} \hat{\theta}^{(1)} \\ 1 \end{pmatrix} \in \mathbb{R}^{(H_2+1) \times 1}$$

$$\begin{bmatrix} w_{11}^{(0)} \\ w_{12}^{(0)} \\ \vdots \\ w_{H_1, p+1}^{(0)} \\ w_{11}^{(1)} \\ \vdots \\ w_{11}^{(m)} \end{bmatrix}$$

Aside:

$$W^{(0)} = \begin{bmatrix} w_{11}^{(0)} & w_{12}^{(0)} & \dots & w_{1, p+1}^{(0)} \\ w_{21}^{(0)} & w_{22}^{(0)} & \dots & w_{2, p+1}^{(0)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{H_1, 1}^{(0)} & \dots & \dots & w_{H_1, p+1}^{(0)} \end{bmatrix}$$

$$\theta_1^{(1)} = \sigma\left(\sum_{j=1}^{p+1} w_{1j}^{(0)} \cdot \theta_j^{(0)}\right)$$

$$= \sigma(\text{1st element of } W^{(0)} \cdot \theta^{(0)})$$

$$\theta_2^{(1)} = \sigma\left(\sum_{j=1}^{p+1} w_{2j}^{(0)} \cdot \theta_j^{(0)}\right)$$

$$= \sigma(\text{2nd el. of } W^{(0)} \cdot \theta^{(0)})$$

$$\hat{\theta}^{(1)} = \sigma(W^{(0)} \cdot \theta^{(0)}) = \mathbb{R}^{H_1 \times 1}$$

$\mathbb{R}^{(H_1 \times (p+1))} \cdot \mathbb{R}^{(p+1) \times 1}$

$$\theta^{(1)} = \begin{bmatrix} \hat{\theta}^{(1)} \\ 1 \end{bmatrix} \in \mathbb{R}^{(H_1+1) \times 1}$$

Input

H_1

H_2

H_3

...

$$x = \theta^{(0)}$$

$$\theta^{(1)} = \begin{bmatrix} \sigma(w^{(0)} \theta^{(0)}) \\ 1 \end{bmatrix} \rightarrow \theta^{(2)} = \begin{bmatrix} \sigma(w^{(1)} \cdot \theta^{(1)}) \\ 1 \end{bmatrix} \dots$$

$$\theta = \sigma(w^{(m)} \cdot \theta^{(m)})$$

$$= \sigma(w^{(m)} \cdot \begin{bmatrix} \hat{\theta}^{(m)} \\ 1 \end{bmatrix})$$

$$= \sigma(w^{(m)} \left(\sigma(w^{(m-1)} \left(\sigma(w^{(m-2)} \left(\sigma(w^{(m-3)} \sigma^{(m-2)} \right) \right) \right) \right) \right)$$

$$= \sigma(w^{(n)} \cdot \begin{bmatrix} \sigma(w^{(n-1)} \cdot \theta^{(n-1)}) \\ 1 \end{bmatrix}) \quad \checkmark$$

Train NN:

- GD: $w_{ji}^{(t+1)} = w_{ji}^{(t)} - \eta \left[\frac{\partial L(w)}{\partial w_{ji}^{(t)}} \right]$

- Stochastic GD:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} - \eta \left[\frac{\partial L_d(w)}{\partial w_{ji}^{(t)}} \right]$$

- Batch GD

$$L_d(w) = \frac{1}{2} \sum_{d_i \in d} \sum_{k=1}^d (y_{d,i,k} - a_{d,i,k})^2$$

Case 1: $w_{ji} = w_{ji}^{(n)} \rightarrow$ output layer

$$\begin{aligned} \frac{\partial L_d(w)}{\partial w_{ji}^{(n)}} &= \frac{\partial}{\partial w_{ji}^{(n)}} \frac{1}{2} \sum_{k=1}^k (y_k - \underline{a_k})^2 \\ &= \frac{\partial}{\partial w_{ji}^{(n)}} \frac{1}{2} \sum_{k=1}^k (y_k - \sigma(\text{net}_{\text{output},k}))^2 \\ &= \frac{\partial}{\partial w_{ji}^{(n)}} \frac{1}{2} \sum_{k=1}^k (y_k - \sigma(\sum_{\ell=1}^{H_{n+1}} w_{k\ell}^{(n)} \theta_{\ell}^{(n)}))^2 \\ &= \frac{\partial}{\partial w_{ji}^{(n)}} \frac{1}{2} (y_j - \underbrace{\sigma(\sum_{\ell=1}^{H_{n+1}} w_{j\ell}^{(n)} \theta_{\ell}^{(n)})}_{\sigma(\text{net output}, j)})^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial L_d(w)}{\partial w_{ji}^{(n)}} &= \left[\frac{\partial L_d(w)}{\partial \sigma(\text{net}_{o,j})} \right] \left[\frac{\partial \sigma(\text{net}_{o,j})}{\partial \text{net}_{o,j}} \right] \cdot \frac{\partial \text{net}_{o,j}}{\partial w_{ji}^{(n)}} \\ &= -(y_j - \sigma(\text{net}_{o,j})) [\sigma(\text{net}_{o,j}) (1 - \sigma(\text{net}_{o,j}))] [\theta_i^{(n)}] \end{aligned}$$

$$\begin{aligned} \frac{\partial L_d(w)}{\partial w_{ji}^{(n)}} &= -(y_j - \theta_j) \theta_j (1 - \theta_j) \theta_i^{(n)} \\ &= -\delta_{\text{output},j} \theta_i^{(n)} \end{aligned}$$

$$\delta_{\text{output},j} = \frac{-\partial L_d(w)}{\partial \text{net}_{\text{output},j}}$$

SGD:

$$(w_{ji}^{(n)})^{(t+1)} = (w_{ji}^{(n)})^t + \eta \delta_{o,j} \theta_i^{(n)}$$

Linear Reg: $\hat{y} = w^T x$

$$L(w) = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$\frac{\partial L(w)}{\partial w} = 0$$

$$d = \{d_i\} \quad w = (x^T x)^{-1} x^T y$$

$$d = \{d_1, d_2, d_3\}$$

$$\text{Goal: } \frac{\partial L_d(w)}{\partial w_{ji}}$$

$$\frac{\partial L_d(w)}{\partial \sigma(\text{net}_{o,j})} = \frac{\partial}{\partial \sigma(\text{net}_{o,j})} \frac{1}{2} (y_i - \sigma(\text{net}_{o,j}))^2$$

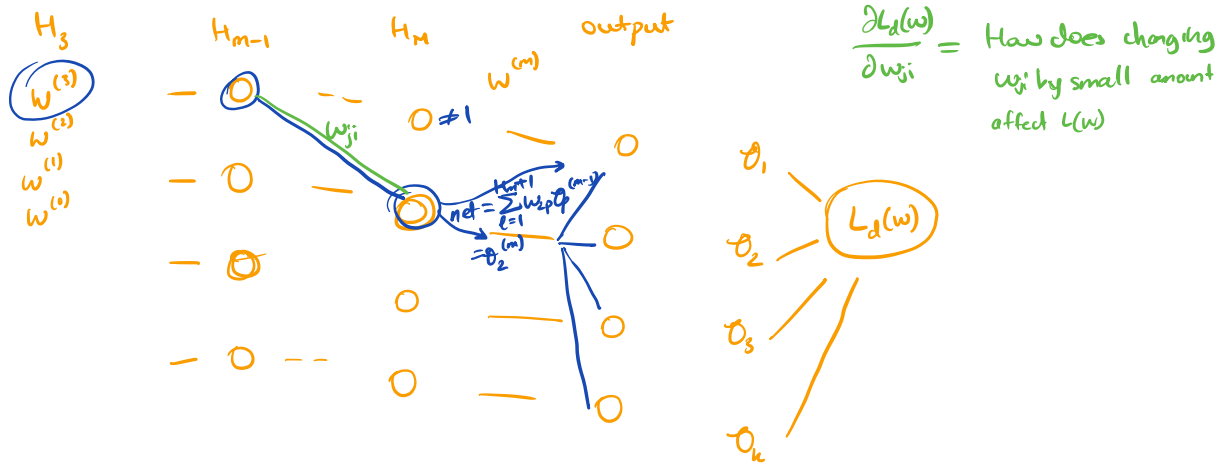
$$\frac{\partial}{\partial x} \frac{1}{2} (y - x)^2 = -(y - x)$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma'(x) = \sigma(x) (1 - \sigma(x))$$

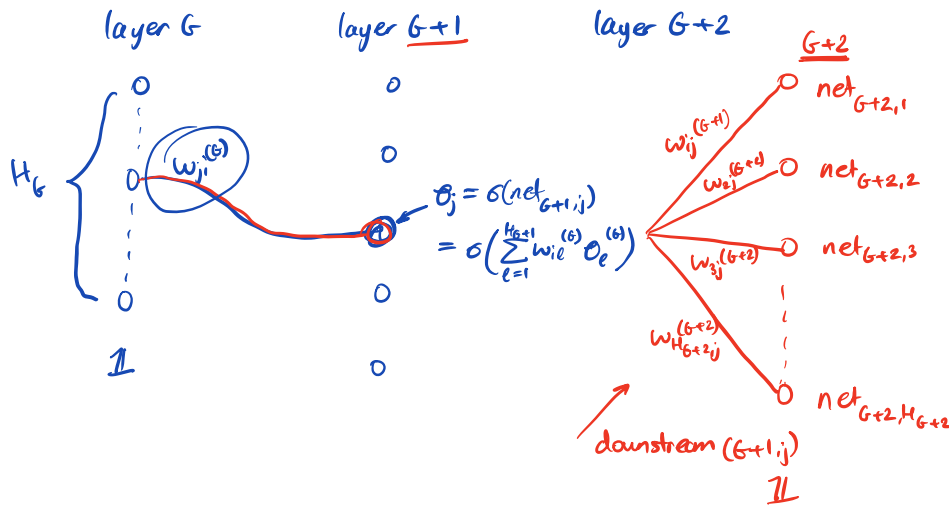
$$\text{net}_{o,j} = \sum_{\ell=1}^{H_{n+1}} w_{j\ell}^{(n)} \theta_{\ell}^{(n)}$$

$$= w_{j1}^{(n)} \theta_1^{(n)} + \underbrace{w_{j2}^{(n)} \theta_2^{(n)}}_{\text{circled}} + \dots + w_{jH_{n+1}}^{(n)} \theta_{H_{n+1}}^{(n)}$$

$$\frac{\partial L_d(w)}{\partial \text{net}_{o,j}} = \frac{\partial L_d}{\partial \sigma(\text{net}_{o,j})} \cdot \frac{\partial \sigma(\text{net}_{o,j})}{\partial \text{net}_{o,j}}$$



Case 2: $w_{ji} = w_{ji}^{(G)}$



$$\begin{aligned} \frac{\partial L_d(w)}{\partial w_{ji}^{(G)}} &= \frac{\partial L_d(w)}{\partial net_{G+1,j}} \times \frac{\partial net_{G+1,j}}{\partial w_{ji}^{(G)}} \\ &= \left[\sum_{z \in \text{downstream}(G+1,j)} \frac{\partial L_d(w)}{\partial net_{G+2,z}} \cdot \frac{\partial net_{G+2,z}}{\partial net_{G+1,j}} \right] \frac{\partial net_{G+1,j}}{\partial w_{ji}^{(G)}} \\ &= \left[\sum_{z \in \text{downstream}(G+1,j)} \delta_{G+2,z} (w_{zj}^{(G+1)} o_j^{(G+1)} (1 - o_j^{(G+1)})) \right] o_j^{(G)} \end{aligned}$$

SGD:

$$\begin{aligned} (w_{ji}^{(G)})^{(t+1)} &= (w_{ji}^{(G)})^{(t)} + \eta o_i^{(G)} \left[\delta_{j,G+1} (1 - o_j^{(G+1)}) \sum_z \delta_{G+2,z} w_{zj}^{(G+1)} \right] \\ &= (w_{ji}^{(G)})^{(t)} + \eta o_i^{(G)} \delta_{j,G+1} \left[\delta_{j,G+1} = \text{effect of } net_{G+1,j} \text{ on loss} \right] \end{aligned}$$