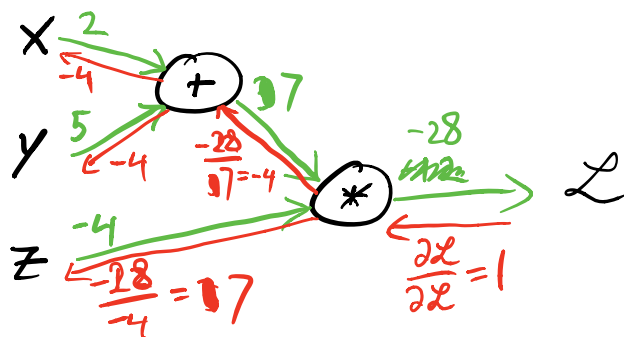


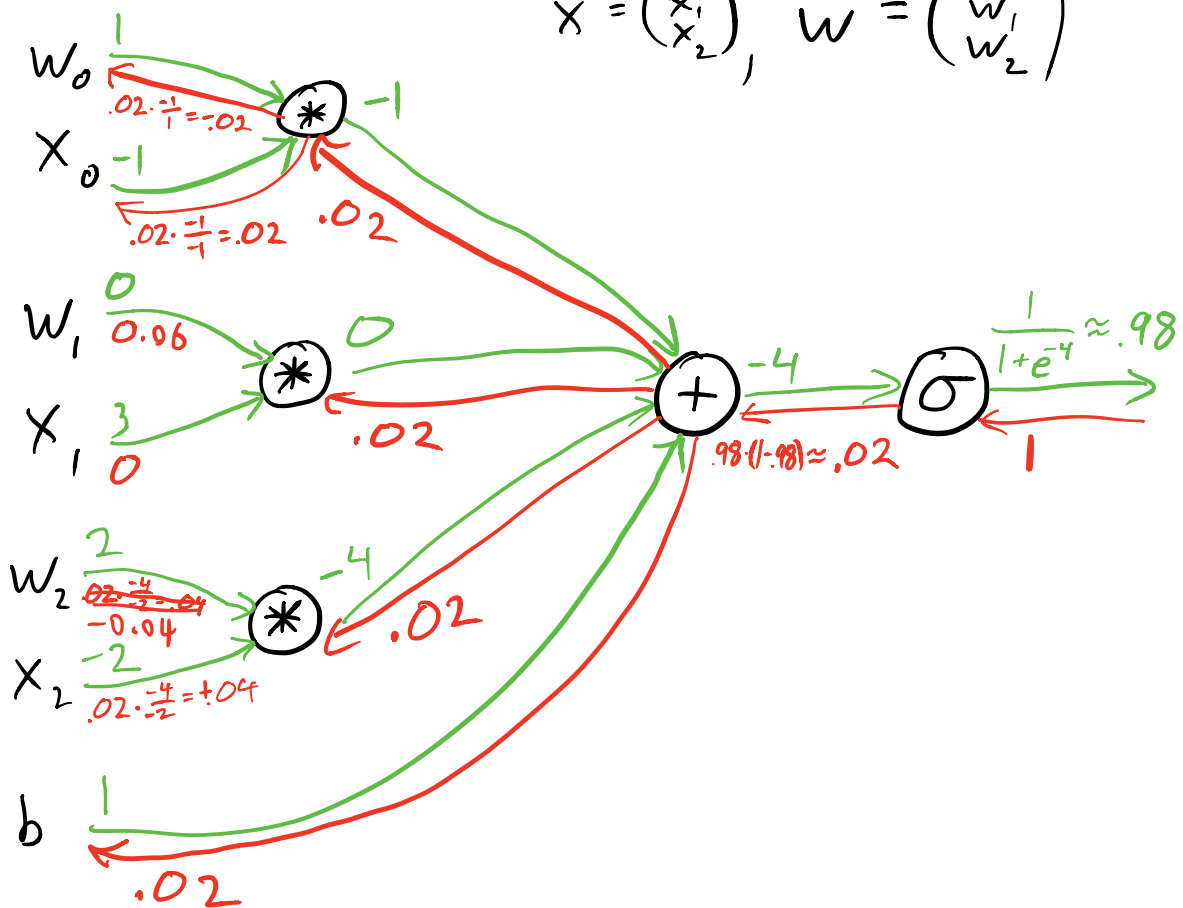
Scalar Operations

+	$z = \sum_i x_i \Rightarrow \frac{\partial z}{\partial x_i} = 1$ $\mathcal{L} = f(z)$ $\frac{\partial \mathcal{L}}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial z}$	<p>"distribute flow unchanged"</p>
*	$z = \prod_i x_i \Rightarrow \frac{\partial z}{\partial x_i} = \frac{z}{x_i}$ $\mathcal{L} = f(z)$ $\frac{\partial \mathcal{L}}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{z}{x_i}$	<p>"split flow & scale by out/in"</p>
max & min	$z = \max(x_1, x_2, \dots, x_n) \Rightarrow \frac{\partial z}{\partial x_i} = \mathbb{1}[z == x_i]$ $\mathcal{L} = f(z)$ $\frac{\partial \mathcal{L}}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial z} \cdot \mathbb{1}[x_i == z]$	<p>"direct flow to max/min input"</p>
σ $\frac{1}{1+e^{-x}}$	$z = \sigma(x) \Rightarrow \frac{\partial z}{\partial x} = z(1-z)$ $\mathcal{L} = f(z)$ $\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial z} \cdot z \cdot (1-z)$	<p>"attenuate flow increasingly as output saturates"</p>

1. $\mathcal{L} = (x+y)z$



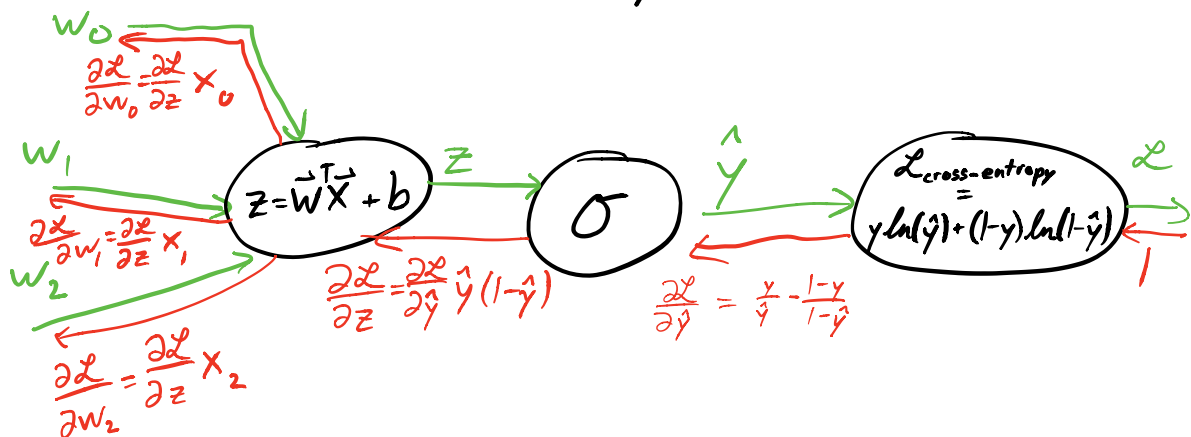
2. $\mathcal{L} = \sigma(\vec{w}^T \vec{x} + b)$ ($\vec{x}, \vec{w} \in \mathbb{R}^{3 \times 1}$, $b \in \mathbb{R}$)
 $\vec{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$, $\vec{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$



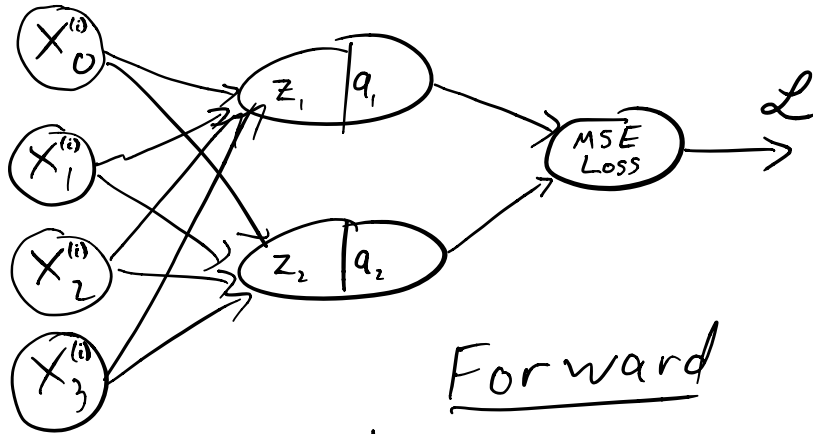
Logistic Regression (Batch Size = 1)

given: x_i inputs, y output
 $\vec{x} \in \mathbb{R}^{n \times 1}$, $\vec{w} \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$

Forward	Backward
$Z = \vec{w}^T \vec{x} + b$	$\frac{\partial \mathcal{L}}{\partial \vec{w}} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial \vec{w}}$ $\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b}$ $= \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) \hat{y}(1-\hat{y}) \vec{x}$ $= \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) \hat{y}(1-\hat{y})$
$\hat{y} = \sigma(z)$	$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z}$ $= \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) (\hat{y}(1-\hat{y}))$
$\mathcal{L} = y \ln(\hat{y}) + (1-y) \ln(1-\hat{y})$	$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}$



2-Layer MLP (Batch Size=3)



$$1. \quad Z = W_1 X + b_1 u$$

$$X = \begin{bmatrix} | & | & | \\ \text{example} & \text{example} & \text{example} \\ | & | & | \end{bmatrix} = \begin{bmatrix} \text{--- feature 1 ---} \\ \text{--- feature 2 ---} \\ \text{--- feature 3 ---} \\ \text{--- feature 4 ---} \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & X_1^{(2)} & X_1^{(3)} \\ X_2^{(1)} & X_2^{(2)} & X_2^{(3)} \\ X_3^{(1)} & X_3^{(2)} & X_3^{(3)} \\ X_4^{(1)} & X_4^{(2)} & X_4^{(3)} \end{bmatrix}$$

$$W_1 = \begin{bmatrix} | & | & | & | \\ \text{weights} & \text{weights} & \text{weights} & \text{weights} \\ \text{of} & \text{of} & \text{of} & \text{of} \\ \text{feature} & \text{feature} & \text{feature} & \text{feature} \\ | & | & | & | \end{bmatrix} = \begin{bmatrix} \text{--- weights of upper neuron ---} \\ \text{--- weights of lower neuron ---} \end{bmatrix}$$

$$b_1 = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \quad u = [1 \ 1 \ 1]$$

$[2 \times 1]$ 1×3

$$Z = \begin{bmatrix} \text{--- preactivation in upper neuron ---} \\ \text{--- preactivation in lower neuron ---} \end{bmatrix}$$

$[2 \times 3]$

$$= \begin{bmatrix} | & | & | \\ \text{example} & \text{example} & \text{example} \\ \text{preactivations} & \text{preactivations} & \text{preactivations} \\ | & | & | \end{bmatrix}$$

$$2. \quad A = \sigma(Z)$$

$$A = \begin{bmatrix} \text{activation of upper neuron} \\ \text{activation of lower neuron} \end{bmatrix}$$

$$[2 \times 3] = \begin{bmatrix} \text{example 1} & \text{example 2} & \text{example 3} \\ \text{activations} & \text{activations} & \text{activations} \end{bmatrix}$$

$$3. \quad \hat{y} = w_2 A + b_2 u$$

$$w_2 = \begin{bmatrix} \text{coefficient of upper activation} & \text{coefficient of lower activation} \end{bmatrix}$$

$$[1 \times 2]$$

$$b_2 \in \mathbb{R}, \quad u = [1 \quad 1 \quad 1]$$

$$1 \times 3$$

$$\hat{y} = \begin{bmatrix} \text{example 1 prediction} & \text{example 2 prediction} & \text{example 3 prediction} \end{bmatrix}$$

$$[1 \times 3]$$

$$4. \quad \mathcal{L} = \frac{1}{3} \|\hat{y} - y\|^2$$

$$\mathcal{L} \in \mathbb{R}$$

$$y = \begin{bmatrix} \text{ex. 1 ground truth} & \text{ex. 2 ground truth} & \text{ex. 3 ground truth} \end{bmatrix}$$

$$[1 \times 3]$$

Backward

$$1. \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{2}{3} (\hat{y} - y)$$

$[1 \times 3]$

$$2. \frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2} = \frac{2}{3} (\hat{y} - y) A^T$$

$[1 \times 2]$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b_2} = \frac{2}{3} (\hat{y} - y) u^T$$

$[1 \times 1]$

$$\frac{\partial \mathcal{L}}{\partial A} = \frac{2}{3} w_2^T (\hat{y} - y)$$

$[2 \times 3]$

$$3. \frac{\partial \mathcal{L}}{\partial Z} = \frac{\partial \mathcal{L}}{\partial A} \odot \sigma'(Z) = \left(\frac{2}{3} w_2^T (\hat{y} - y) \right) \odot A \odot (1 - A)$$

$[2 \times 3]$

$$4. \frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial Z} \frac{\partial Z}{\partial w_1} = \frac{2}{3} \left((w_2^T (\hat{y} - y)) \odot A \odot (I - A) \right) X^T$$

[2x4]

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial Z} \frac{\partial Z}{\partial b_1} = \frac{2}{3} \left((w_2^T (\hat{y} - y)) \odot A \odot (I - A) \right) u^T$$

[2x1]

In general ...

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(Z^{[l]}) \quad (\text{where } g \text{ is an elementwise nonlinearity like } \sigma \text{ or ReLU})$$

Let $\delta^{[l]} = \frac{\partial \mathcal{L}}{\partial Z^{[l]}}$, which we can calculate starting at the last layer using $\delta^{[N]} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \odot g'^{[N]}(Z^{[N]})$ and working backwards either by

$$\delta^{[l]} = \left((W^{[l+1]})^T \delta^{[l+1]} \right) \odot g'^{[l]}(Z^{[l]}).$$

Then $\frac{\partial \mathcal{L}}{\partial w^{[l]}} = \delta^{[l]} (A^{[l-1]})^T$

$$\frac{\partial \mathcal{L}}{\partial b^{[l]}} = \delta^{[l]} u^T$$

In our example we had
 $g^{[1]} = \sigma$ and $g^{[2]} = \text{identity function}$

[2]

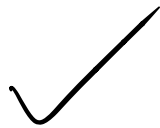
$$\delta^{[2]} = \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{2}{3}(\hat{y} - y), \quad \frac{\partial \mathcal{L}}{\partial w^{[2]}} = \frac{2}{3}(\hat{y} - y)A^T, \quad \frac{\partial \mathcal{L}}{\partial b^{[2]}} = \frac{2}{3}(\hat{y} - y)u^T$$

[1]

$$\delta^{[1]} = ((w^{[2]})^T \cdot \frac{2}{3}(\hat{y} - y)) \odot A \odot (I - A)$$

$$\frac{\partial \mathcal{L}}{\partial w^{[1]}} = \delta^{[1]} (A^{[0]})^T = (((w^{[2]})^T \cdot \frac{2}{3}(\hat{y} - y)) \odot A \odot (I - A)) x^T$$

$$\frac{\partial \mathcal{L}}{\partial b^{[1]}} = \delta^{[1]} (A^{[0]})^T = (((w^{[2]})^T \cdot \frac{2}{3}(\hat{y} - y)) \odot A \odot (I - A)) u^T$$



Optimizers

- Mini-Batch Gradient Descent:

$$W^{[L]} = W^{[L]} - \alpha \frac{\partial \mathcal{L}}{\partial W^{[L]}}$$

"learning rate"

$$b^{[L]} = b^{[L]} - \alpha \frac{\partial \mathcal{L}}{\partial b^{[L]}}$$

(Stochastic Gradient Descent has batch size=1)

- Momentum:

Use exponential weighted average of past gradient updates rather than the raw ones (dampens oscillations)

$$V^{[L]} = \beta V^{[L]} + (1-\beta) \frac{\partial \mathcal{L}}{\partial W^{[L]}}$$
$$W^{[L]} = W^{[L]} - \alpha V^{[L]}$$

- RMSProp:

Divide gradient by exponential weighted average of its magnitude

$$S^{[L]} = \beta S^{[L]} + (1-\beta) \left(\frac{\partial \mathcal{L}}{\partial W^{[L]}} \right)^2 \quad (\text{element-wise square})$$
$$W^{[L]} = W^{[L]} - \alpha \frac{\frac{\partial \mathcal{L}}{\partial W^{[L]}}}{\sqrt{S^{[L]} + \epsilon}} \quad (\epsilon = \text{small non-zero param to avoid dividing by } 0)$$

• Adam:

Combine Momentum & RMS Prop with an additional bias correction

$$V = \frac{\beta_1 V + (1 - \beta_1) \frac{\partial \mathcal{L}}{\partial w}}{(1 - \beta_1)^t}$$

(dividing by $(1 - \beta)^t$
yields less biased
values for small t)

$$S = \frac{\beta_2 S + (1 - \beta_2) \left(\frac{\partial \mathcal{L}}{\partial w} \right)^2}{(1 - \beta_2)^t}$$

$$W = W^- \propto \frac{V}{\sqrt{S} + \epsilon}$$