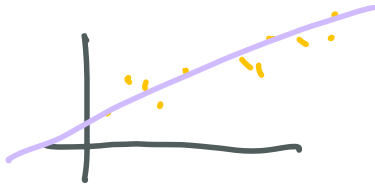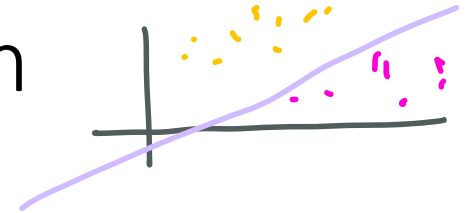# Neural Networks

# Outline

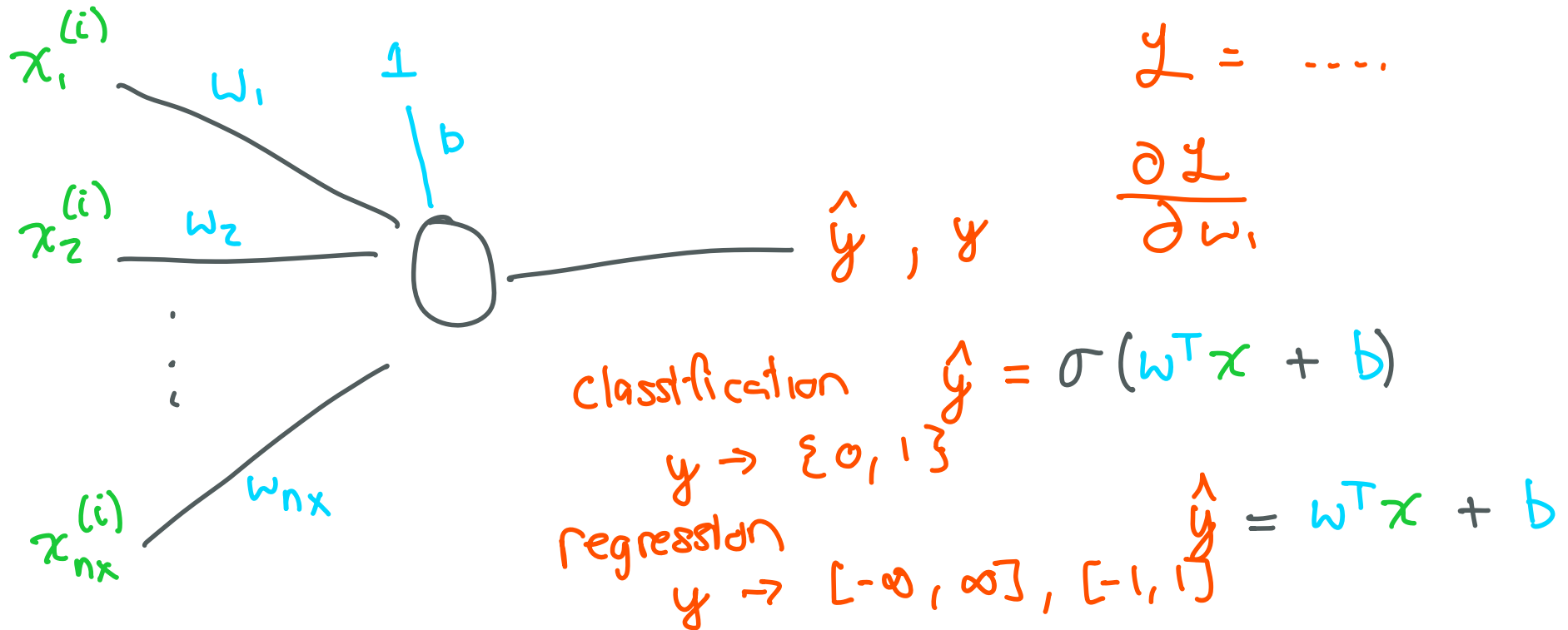- Questions about projects?
- Recap of a singular neuron model
- Notation and terminology
- Compute graphs
- Optimization
- Backpropagation

- This will be our most math heavy week
- Next week we'll rely on PyTorch to compute all derivatives
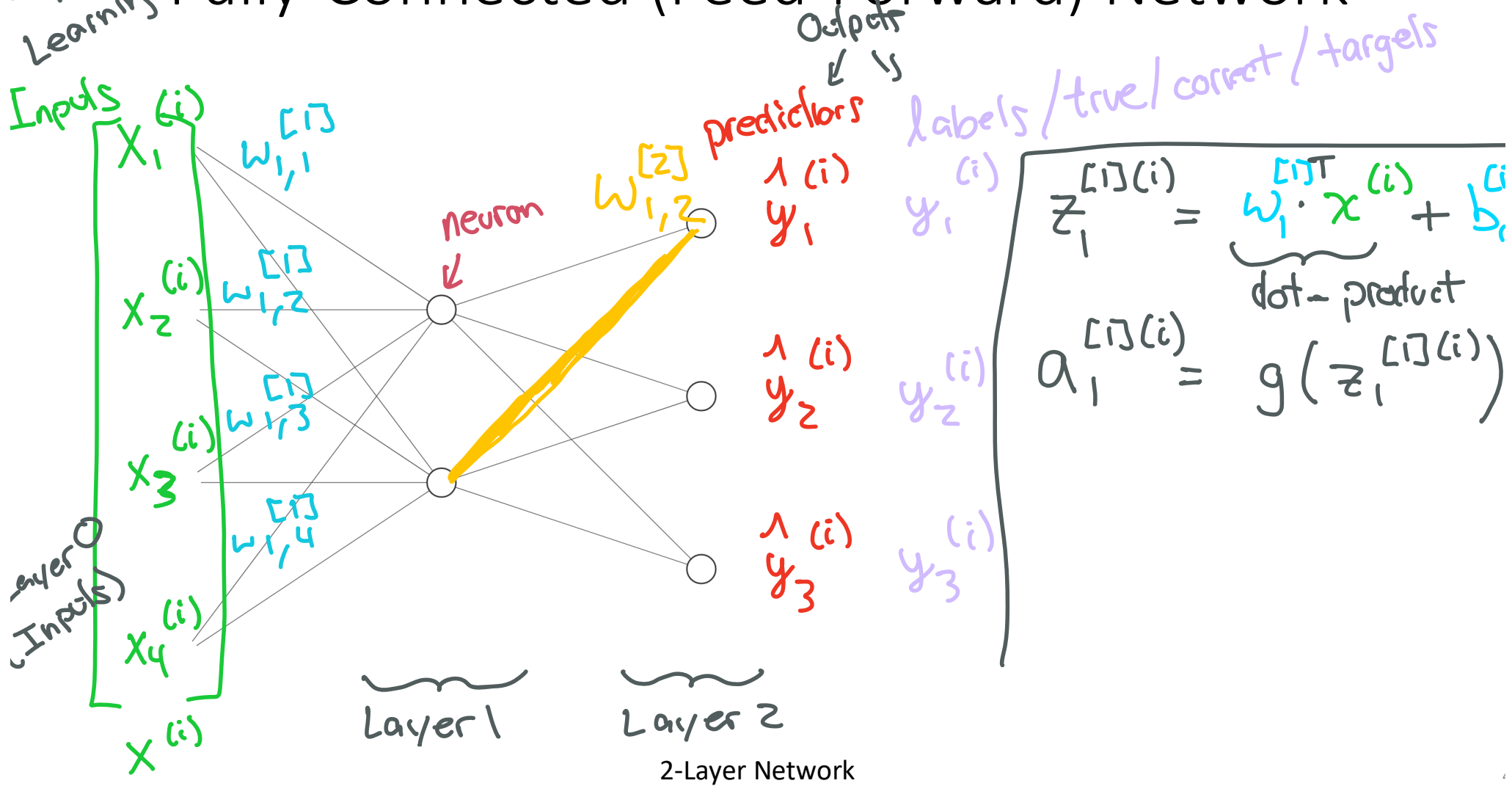
# Recap: A Single Neuron

- Take five minutes to draw
  - Whatever will help you remember (no correct or incorrect drawings)
  - You'll keep a running drawing log the rest of the semester



$x_1^{(i)}$   $w_1$   $1$   $b$

$x_2^{(i)}$   $w_2$

$\vdots$

$x_{n_x}^{(i)}$   $w_{n_x}$

$\hat{y}, y$

$y = \dots$

$\dfrac{\partial L}{\partial w_1}$

classification   $\hat{y} = \sigma(w^T x + b)$

$y \rightarrow \{0, 1\}$

regression

$y \rightarrow [-\infty, \infty], [-1, 1]$   $\hat{y} = w^T x + b$

# Fully-Connected (Feed-Forward) Network

Superv
Learning

Inputs $X_1^{(i)}$

$W_{1,1}^{[1]}$

$x_2^{(i)}$  $W_{1,2}^{[1]}$

$x_3^{(i)}$  $W_{1,3}^{[1]}$

$W_{1,4}^{[1]}$

Layer 0
(Inputs)

$X_4^{(i)}$

$X^{(i)}$

neuron

Output

predictors

$W_{1,2}^{[2]}$

$\hat{y}_1^{(i)}$

$\hat{y}_2^{(i)}$

$\hat{y}_3^{(i)}$

labels / true / correct / targets

$y_1^{(i)}$

$y_2^{(i)}$

$y_3^{(i)}$

$$z_1^{[1](i)} = \underbrace{w_1^{[1]T} \cdot x^{(i)}}_{\text{dot-product}} + b_1^{(i}$$

$$a_1^{[1](i)} = g\left(z_1^{[1](i)}\right)$$

Layer 1     Layer 2

2-Layer Network

# Fully-Connected (Feed-Forward) Network



$$a_1^{[0]} = x_1$$
$$a_2^{[0]} = x_2$$
$$x_3$$
$$x_4$$

$$a_1^{[0]}$$
$$a_2^{[0]}$$
$$a_3^{[0]}$$
$$a_4^{[0]}$$

$$z^{[1]} = w^{[1]T} a^{[0]} + b^{[1]} \quad a^{[1]} = g(z^{[1]})$$
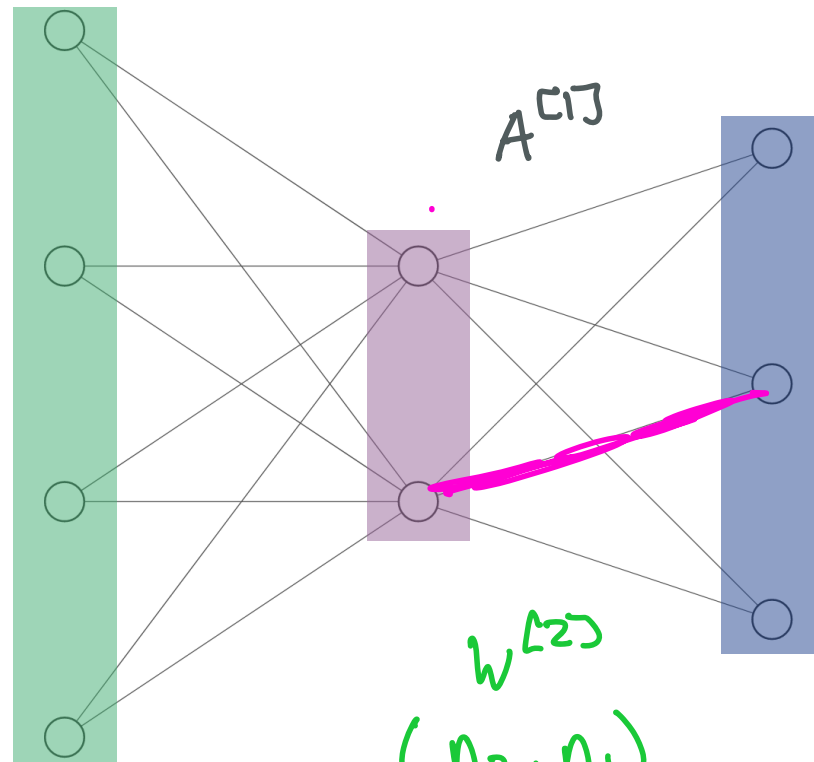
$$a^{[1]} =$$

2-Layer Network

# Fully-Connected (Feed-Forward) Network
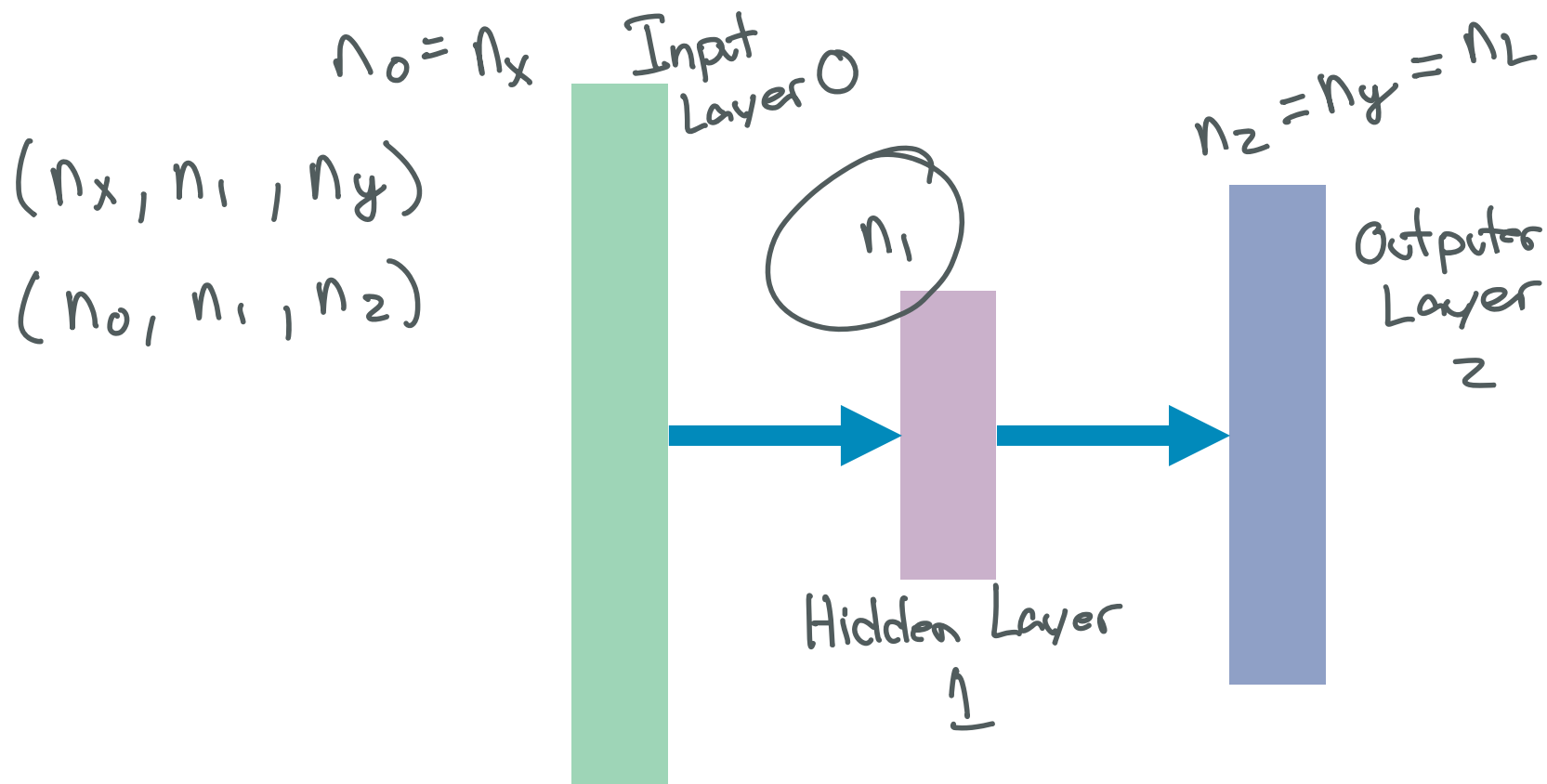


$A^{[0]} = X$

$A^{[0]}$

$\hat{Y} = A^{[2]}$

$A^{[1]}$

$W^{[2]}$

$(n_2, n_1)$

$\dfrac{\partial \mathcal{L}}{\partial W^{[2]}}$

$(n_2, n_1)$

# Fully-Connected (Feed-Forward) Network

$n_0 = n_x$ Input Layer 0

$n_z = n_y = n_2$

$(n_x, n_1, n_y)$

$(n_0, n_1, n_2)$

$n_1$

Output Layer 2

Hidden Layer 1

Simplified (easier to draw) diagram

# Vectorized Equations

16,000 examples

$Z_1^{[1](i)}$

$Z_2^{[1](i)}$

$\vdots$

$Z_{100}^{[1](i)}$

these are scalars

$Z^{[1]}$

matrix

$(100, \cancel{16})$

$n_1$

$(100, 16{,}000)$

$(16{,}000, 100)$

$(N, n_1)$

$Z_1^{[1](i)} = w_1^{[1]T} \cdot x^{(i)} + b_1$

$$Z^{[1]} = \underset{(N, n_x)}{\bigcirc{X}} \; \underset{(n_x, n_1)}{W^{[1]T}} + \underset{(1, n_1)}{b^{[1]T}}$$

$(N, n_1)$

$$\begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_{n_x}^{(1)} \\ & & & \\ x_1^{(N)} & x_2^{(N)} & \cdots & x_{n_x}^{(N)} \end{bmatrix}$$

$N$

$\longleftarrow \quad n_x \quad \longrightarrow$

# Vectorized Equations( for any layer)

$W^{[l]} \to (n_l, n_{l-1})$

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

$(N, n_l) \quad \underbrace{(N, n_{l-1})(n_{l-1}, n_l)}_{(N, n_l)} \quad (1, n_l) \leftarrow$ broadcasted to $(N, n_l)$

$+$

$Z^{[1]} = A^{[1-1]} W^{[1]T} + b^{[1]T} \quad A^{[1]} = g(Z^{[1]})$

$Z^{[2]} = A^{[2-1]} W^{[2]T} + b^{[2]T} \quad A^{[2]} = g(Z^{[2]}) = \hat{y}$
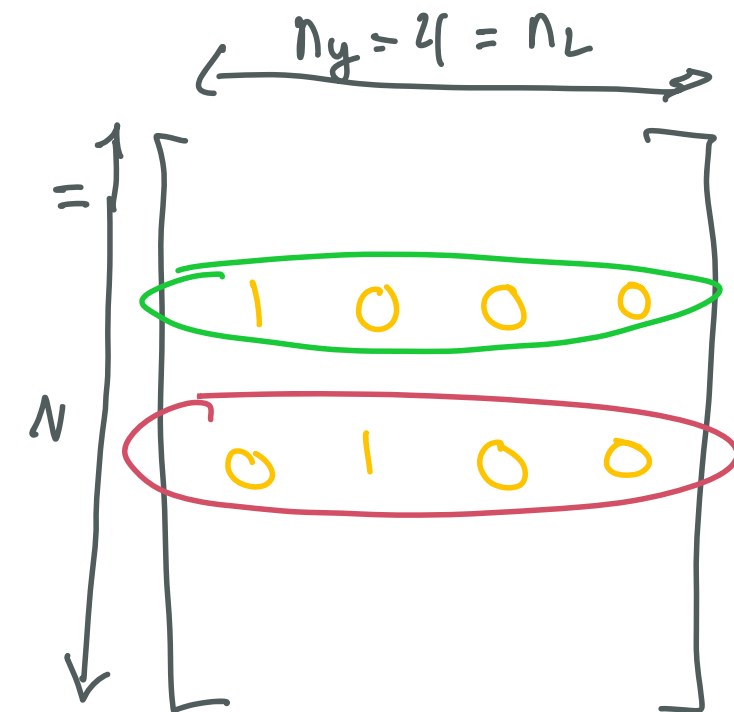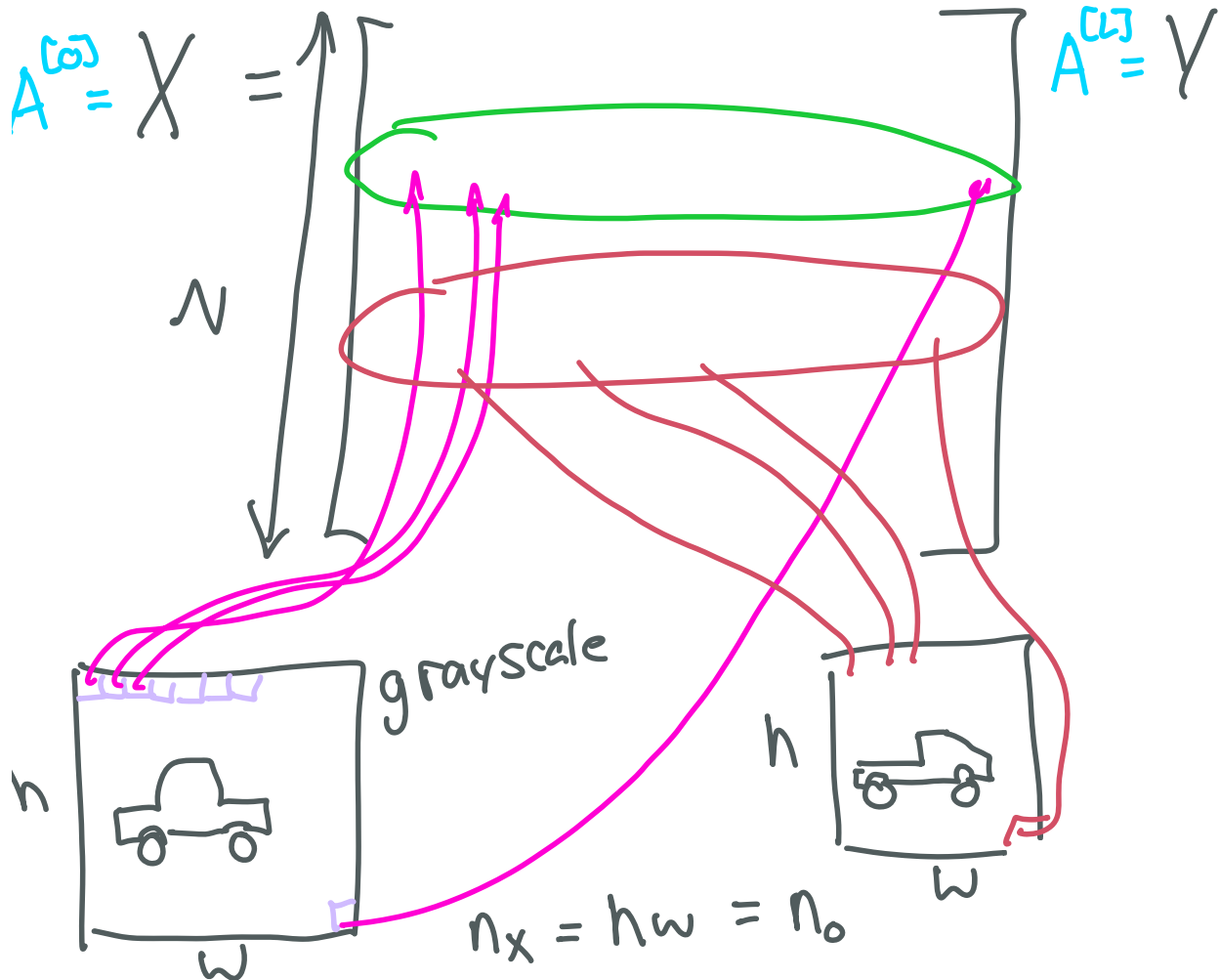
element-wise

$$A^{[l]} = g(\overset{\downarrow}{Z}^{[l]}) = \begin{bmatrix} g(z_{1,1}^{[l]}) & g(z_{1,2}^{[l]}) \cdots \\ \vdots & \ddots \\ g(z_{n_l,1}^{[l]}) & \cdots \end{bmatrix}$$

$(N, n_l) \qquad (N, n_l)$

# MNIST Dataset Example

- MNIST includes 60,000 training images
- Each image is grayscale and 28x28 pixels in size
- Each output is a one-hot encoding of the digits 0 through 9

$28 \cdot 28 = 784$

- What is the shape of $X$?

$$(N, n_x) \Rightarrow (60,000, 784)$$

- What is the shape of $Y$?

$$(N, n_y) \Rightarrow (60,000, 10)$$

1-byte

$(60,000, 28, 28)$

How much memory?

# MNIST Neural Network

- What is the shape of $Z^{[1]}$?

$$(N, n_x) \Rightarrow (N, n_1) \Rightarrow (60{,}000, 17)$$
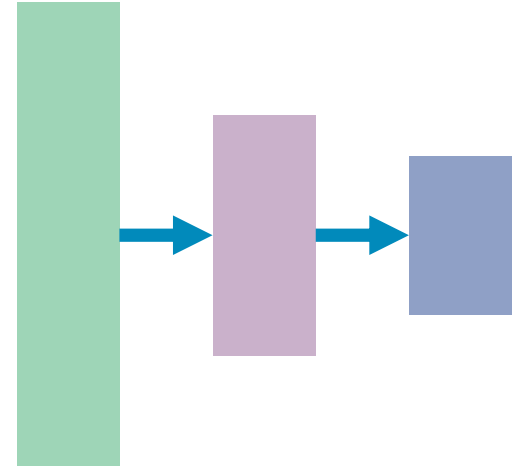
- What is the shape of $A^{[1]}$?

$$(60{,}000, 17)$$

- What is the shape of $Z^{[2]}$?

$$(60{,}000, 10)$$

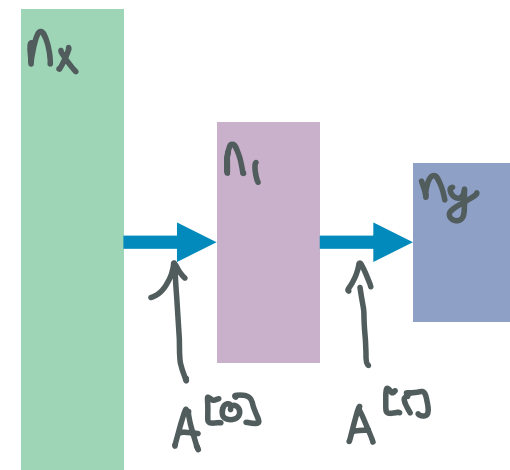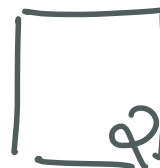- What is the shape of $A^{[2]}$?

$$(60{,}000, 10)$$

# MNIST Neural Network

$n_x$

$n_1$

$n_y$

- Imagine we have a two-layer network
- The hidden layer has 17 neurons

$A^{[0]}$   $A^{[1]}$

- What is the shape of $W^{[1]}$?

$$(n_l, n_{l-1}) \rightarrow (n_1, n_0) \rightarrow (17, 784)$$

- What is the shape of $b^{[1]}$?

$$(n_l, 1) \rightarrow (n_1, 1) \rightarrow (17, 1)$$
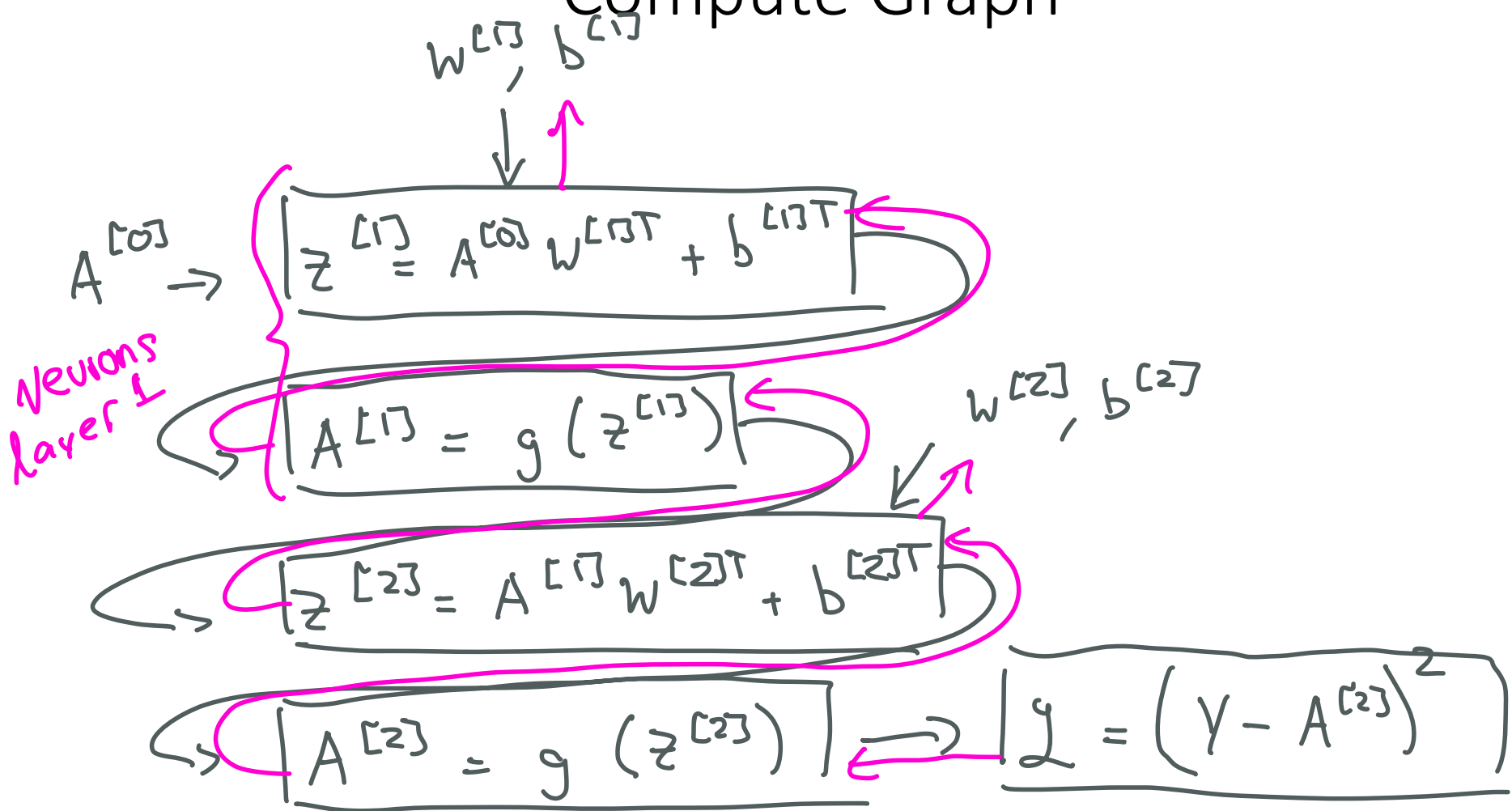
- What is the shape of $W^{[2]}$?

$$(n_l, n_{l-1}) \rightarrow (n_2, n_1) \rightarrow (10, 17)$$

- What is the shape of $b^{[2]}$?

$$(n_l, n_{l-1}) \rightarrow (n_2, 1) \rightarrow (10, 1)$$

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

# Compute Graph

$$W^{[1]}, b^{[1]}$$

$$A^{[0]} \rightarrow \quad Z^{[1]} = A^{[0]} W^{[1]T} + b^{[1]T}$$

Neurons
layer 1

$$A^{[1]} = g\left(Z^{[1]}\right)$$

$$W^{[2]}, b^{[2]}$$

$$Z^{[2]} = A^{[1]} W^{[2]T} + b^{[2]T}$$

$$A^{[2]} = g\left(Z^{[2]}\right)$$

$$\mathcal{L} = \left(Y - A^{[2]}\right)^2$$

# Optimization with Binary Cross Entropy Loss

$$\mathcal{L}(\hat{Y}, Y) = -\| Y \log \hat{Y} + (1-Y) \log (1- \hat{Y})\|$$

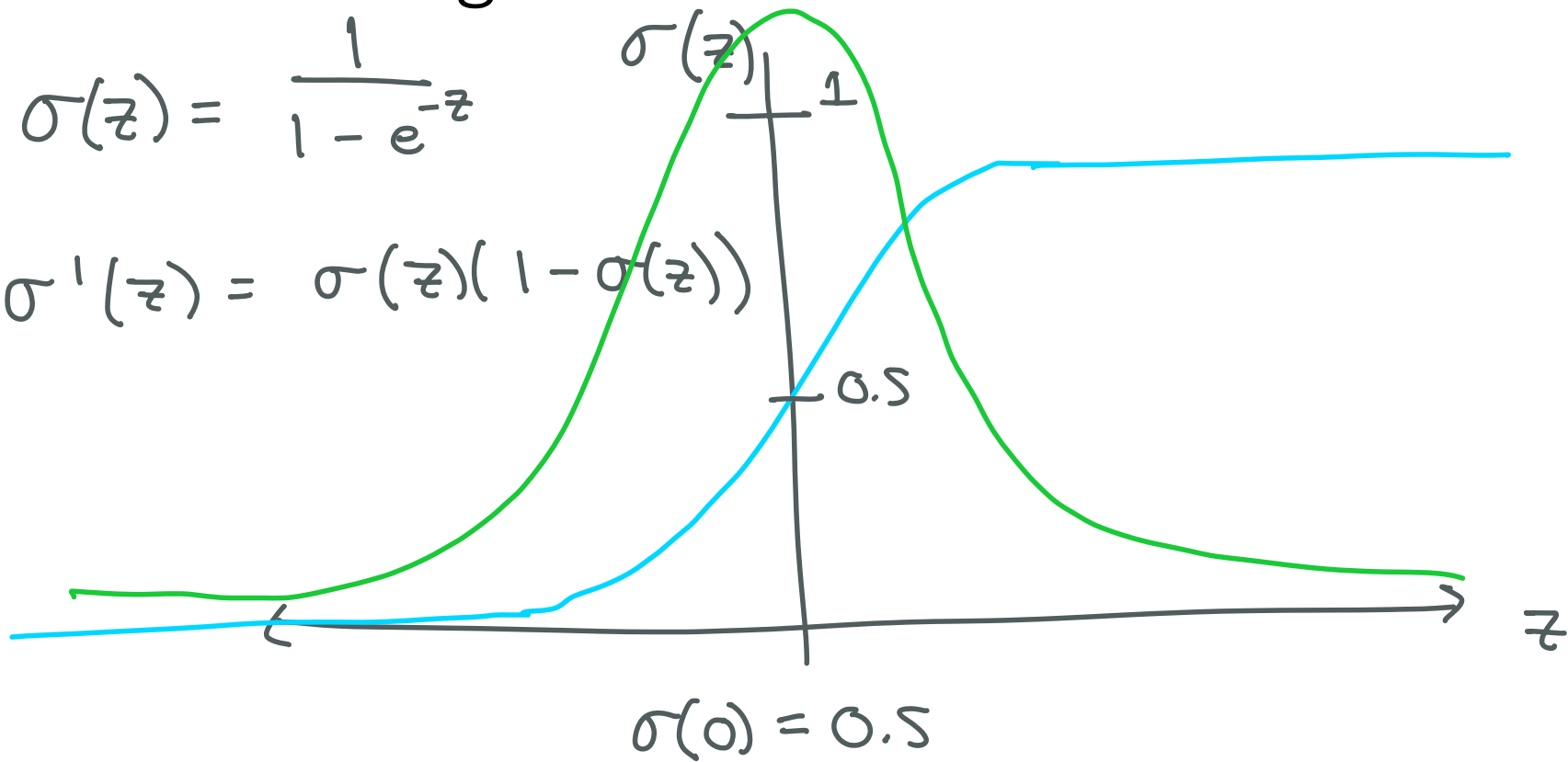| $\hat{y}$ | $y$ | $\log \hat{y}$ | $\log(1-\hat{y})$ | $\mathcal{L}$ |
|---|---|---|---|---|
| 0.1 ☺ 0 | | -1 | -0.046 | -0.046 |
| 0.1 ☹ 1 | | -1 | -0.046 | -1 |
| 0.9 ☹ 0 | | -0.046 | -1 | -1 |
| 0.9 ☺ 1 | | -0.046 | -1 | -0.046 |

$\hat{y} = [0, 1]$

$y \in \{0, 1\}$

# Sigmoid Activation Functions



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

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

$\sigma(z)$

1

0.5

z

$\sigma(0) = 0.5$

# Backpropagation $W^{[2]}$

$$\frac{\partial L}{\partial W^{[2]}} = \frac{\partial}{\partial W^{[2]}} - \| Y \log \hat{y} + (1-Y) \log (1-\hat{y}) \|$$

$$= \overset{①}{\frac{\partial L}{\partial \hat{y}}} \cdot \overset{②}{\frac{\partial \hat{y}}{\partial z^{[2]}}} \cdot \overset{③}{\frac{\partial z^{[2]}}{\partial W^{[2]}}} = -\left(\frac{Y}{\hat{y}} - \frac{(1-Y)}{1-\hat{y}}\right) \cdot \hat{y}(1-\hat{y}) A$$

$$= (\hat{y} - Y) A^{[1]}$$

① $\quad \dfrac{\partial L}{\partial \hat{y}} = -\left(\dfrac{Y}{\hat{y}} - \dfrac{(1-Y)}{(1-\hat{y})}\right)$

② $\quad \dfrac{\partial \hat{y}}{\partial z^{[2]}} = \sigma(z^{[2]})(1-\sigma(z^{[2]})) = \hat{y}(1-\hat{y})$

③ $\quad \dfrac{\partial z^{[2]}}{\partial W^{[2]}} = \dfrac{\partial}{\partial W^{[2]}} A^{[1]} W^{[2]T} + b^{[2]T} = A^{[1]}$

# Backpropagation $b^{[2]}$

$$\frac{\partial \mathcal{L}}{\partial b^{[2]}} = \underbrace{\frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[2]}}}_{\text{Already Computed}} \cdot \frac{\partial z^{[2]}}{\partial b^{[2]}} \, \text{③} = (\hat{y} - y)$$

③ $$\frac{\partial z^{[2]}}{\partial b^{[2]}} = \frac{\partial}{\partial b^{[2]}} A^{[1]} W^{[2]T} + b^{[2]T} = \underline{1}$$

# Backpropagation $W^{[1]}$

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \underbrace{\textcircled{1} \; \frac{\partial \mathcal{L}}{\partial \hat{Y}} \cdot \textcircled{2} \; \frac{\partial \hat{Y}}{\partial Z^{[2]}}}_{\substack{\text{Already} \\ \text{Computed} \\ \text{(Back propagate)}}} \cdot \underbrace{\textcircled{3} \; \frac{\partial Z^{[2]}}{\partial A^{[1]}} \cdot \textcircled{4} \; \frac{\partial A^{[1]}}{\partial Z^{[1]}} \cdot \textcircled{5} \; \frac{\partial Z^{[1]}}{\partial W^{[1]}}}$$

## Update for $W^{[1]}$

$$\frac{d}{dx} \log_{10} x = \frac{1}{x} \cdot \frac{1}{\log 10}$$

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial}{\partial W^{[1]}} - \| \ Y \log \hat{Y} + (1-Y) \log(1-\hat{Y}) \|$$

$$= -Y \frac{\partial}{\partial W^{[1]}} \log \hat{Y} - (1-Y) \frac{\partial}{\partial W^{[1]}} \log(1-\hat{Y}) \qquad \textcolor{magenta}{\text{distribet } f \ \log}$$

$$= -Y \frac{1}{\hat{Y}} \frac{\partial}{\partial W^{[1]}} \hat{Y} - (1-Y) \frac{1}{1-\hat{Y}} \frac{\partial}{\partial W^{[1]}} (1 - \hat{Y}) \qquad \textcolor{magenta}{\text{chain rule}}$$

$$= -\frac{Y}{\hat{Y}} \frac{\partial}{\partial W^{[1]}} \sigma(z^{[2]}) + \frac{1-Y}{1-\hat{Y}} \frac{\partial}{\partial W^{[1]}} \sigma(z^{[1]}) \qquad \textcolor{magenta}{\text{regroup}}$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \frac{\partial}{\partial W^{[1]}} \sigma(z^{[2]}) \qquad \textcolor{magenta}{\text{Apply sigmoid derivative}}$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \sigma(z^{[2]}) \left( 1 - \sigma(z^{[2]}) \right) \frac{\partial}{\partial W^{[1]}} z^{[2]} \qquad \textcolor{magenta}{\text{substitute}}$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \hat{Y} (1-\hat{Y}) \frac{\partial}{\partial W^{[1]}} \left( A^{[1]} W^{[2]T} + \cancel{b^{[2]T}} \right)$$

$$= \left( \hat{Y}(1-Y) - Y(1-\hat{Y}) \right) W^{[2]} \frac{\partial}{\partial W^{[1]}} A^{[1]}$$

$$= \left( \hat{Y} - \cancel{\hat{Y}Y} - Y + \cancel{\hat{Y}Y} \right) W^{[2]} \frac{\partial}{\partial W^{[1]}} \sigma(z^{[1]})$$

$$= \left( \hat{Y} - Y \right) W^{[2]} \sigma(z^{[1]})(1 - \sigma(z^{[1]})) \frac{\partial}{\partial W^{[1]}} z^{[1]}$$

$$= \left( \hat{Y} - Y \right) W^{[2]} A^{[1]}(1 - A^{[1]}) \frac{\partial}{\partial W^{[1]}} \left( A^{[0]} W^{[1]T} + \cancel{b^{[1]T}} \right)$$

$$= \left( \hat{Y} - Y \right) W^{[2]} A^{[1]}(1 - A^{[1]}) A^{[0]}$$

$$\frac{\partial \mathcal{L}}{\partial w^{[1]}} = \frac{\partial}{\partial w^{[1]}} - \| Y \log \hat{Y} + (1-Y) \log (1-\hat{Y}) \|$$

$$= -\left( \frac{\partial}{\partial w^{[1]}} (Y \log \hat{Y}) + \frac{\partial}{\partial w^{[1]}} \left( (1-Y) \log(1-\hat{Y}) \right) \right)$$

$$= -\left( Y \frac{1}{\hat{Y}} \frac{\partial}{\partial w^{[1]}} \hat{Y} + (1-Y) \frac{1}{1-\hat{Y}} \frac{\partial}{\partial w^{[1]}} (1 - \hat{Y}) \right)$$

$$= -\left( \frac{Y}{\hat{Y}} - \frac{1-Y}{1-\hat{Y}} \right) \frac{\partial}{\partial w^{[1]}} \hat{Y}$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \frac{\partial}{\partial w^{[1]}} \sigma(z^{[2]})$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \sigma(z^{[2]})(1 - \sigma(z^{[2]})) \frac{\partial}{\partial w^{[1]}} z^{[2]}$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \hat{Y}(1 - \hat{Y}) \frac{\partial}{\partial w^{[1]}} \left( A^{[1]} w^{[2]T} + b^{[2]T} \right)$$

$$= \left( \frac{1-Y}{1-\hat{Y}} - \frac{Y}{\hat{Y}} \right) \hat{Y}(1 - \hat{Y}) w^{[2]} \frac{\partial}{\partial w^{[1]}} A^{[1]}$$

$$= \left( \hat{Y}(1-Y) - Y(1-\hat{Y}) \right) w^{[2]} \frac{\partial}{\partial w^{[1]}} \sigma(z^{[1]})$$

$$= \left( \hat{Y} - Y\hat{Y} - Y + Y\hat{Y} \right) w^{[2]} \sigma(z^{[1]})(1 - \sigma(z^{[1]})) \frac{\partial}{\partial w^{[1]}} z^{[1]}$$

$$= (\hat{Y} - Y) w^{[2]} A^{[1]}(1 - A^{[1]}) \frac{\partial}{\partial w^{[1]}} \left( X w^{[1]T} + b^{[1]T} \right)$$

$$\boxed{= (\hat{Y} - Y) w^{[2]} A^{[1]}(1 - A^{[1]}) X}$$

For e in range (num_epochs)

$$w^{[1]} = w^{[1]} - \alpha \cdot$$

$$X \to \boxed{|A|} \to |A| \to \uparrow$$

$$w^{[1]}, b^{[1]} \qquad w^{[2]}, b^{[2]}$$

$$\frac{d}{dx} \log_{10} X = \frac{1}{X} \frac{1}{\log(10)}$$

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

$$\frac{\partial \mathcal{L}}{\partial W^{[1]}} = \frac{\partial}{\partial W^{[1]}} -\| Y \log \hat{Y} + (1-Y) \log(1-\hat{Y}) \|$$

$$= \frac{\partial}{\partial W^{[1]}} (Y \log \hat{Y}) + \frac{\partial}{\partial W^{[1]}} ((1-Y) \log(1-\hat{Y}))$$

$$Y \log(\hat{Y})$$

$$Y \log(\sigma(z^{[2]}))$$

$$Y \log(\sigma(A^{[1]} W^{[2]T} + b^{[2]T}))$$

$$Y \log(\sigma(\sigma(z^{[1]}) W^{[2]T} + b^{[2]T}))$$

$$\frac{\partial}{\partial W^{[1]}} Y \log(\sigma(\sigma(A^{[0]} W^{[1]T} + b^{[1]T}) W^{[2]T} + b^{[2]T}))$$

FF - Neural Network

$\hat{Y}$

# Backpropagation $b^{[1]}$

Follow the previous slides

# Parameter Updates

$$\frac{\partial \mathcal{L}}{\partial w^{[2]}}$$

For Loop

$$W^{[2]} = w^{[2]} - \alpha \overbrace{(\hat{y} - y) A^{[1]}}$$

$$b^{[2]} = b^{[2]} - \cdots$$

$$W^{[1]} = W^{[1]} - \cdots$$

$$b^{[1]} = b^{[1]} - \cdots$$

$$\alpha = 0.1 \quad 0.01$$

Overfitting

Dropout
Regularization
Batch normalization
RMS Prop

validation

training

y

epochs