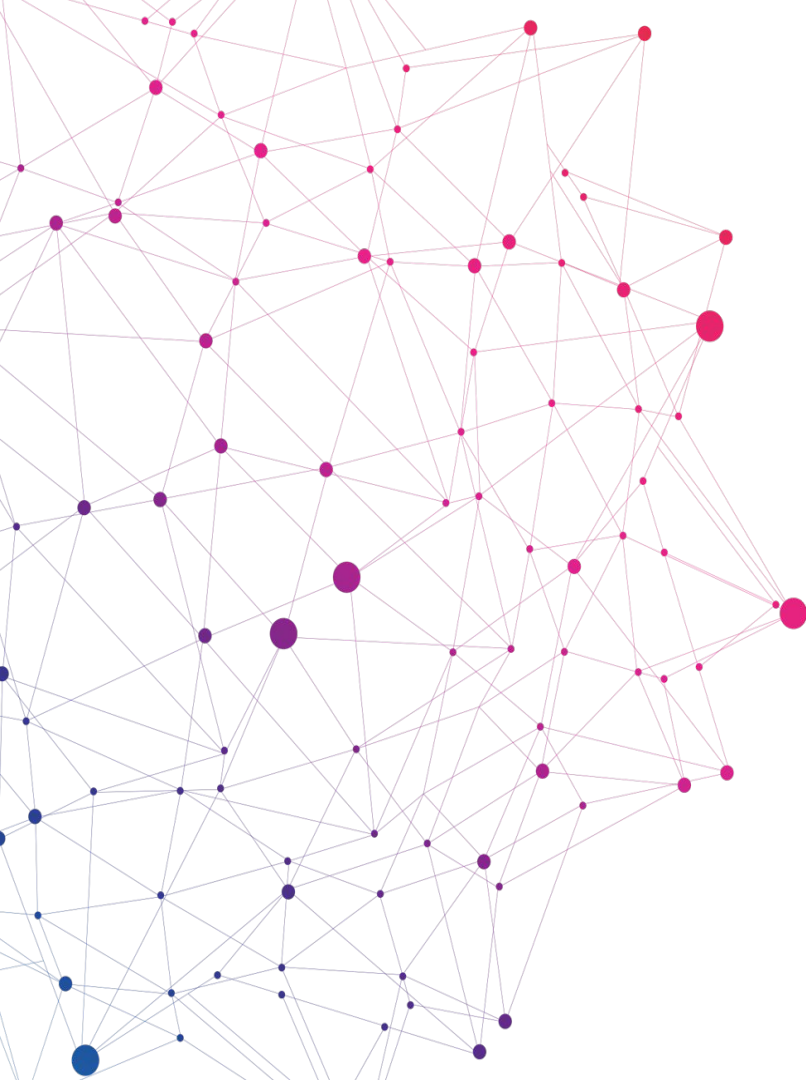


# UVA CS 6316: Machine Learning

## Lecture 15: Neural Network / Deep Learning Basics



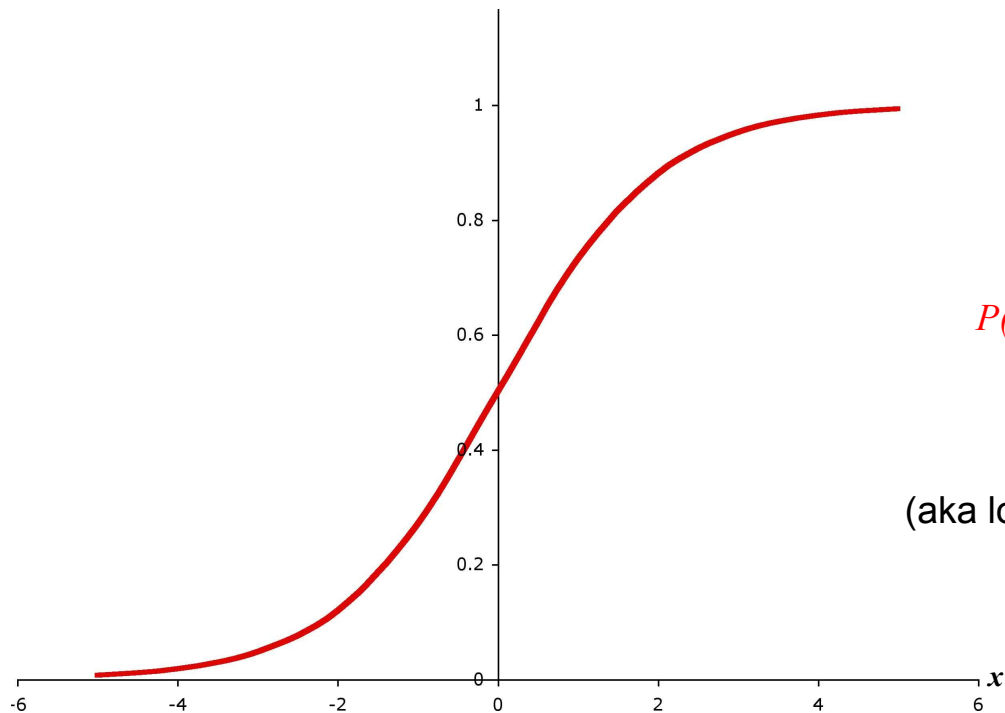
# Intro to Neural Networks and Deep Learning

Jack Lanchantin, Dr. Yanjun Qi



# Neurons

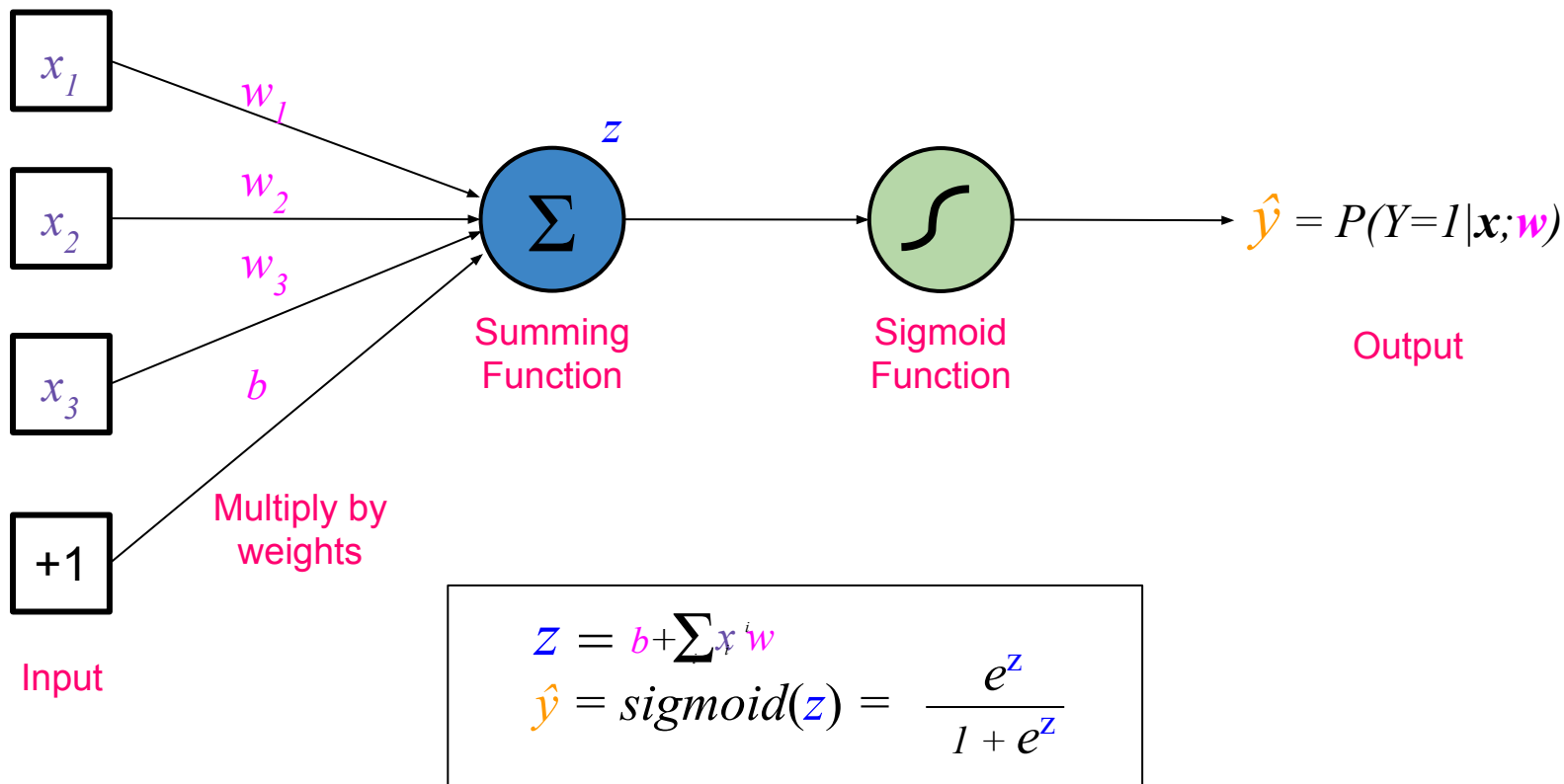
# Logistic Regression



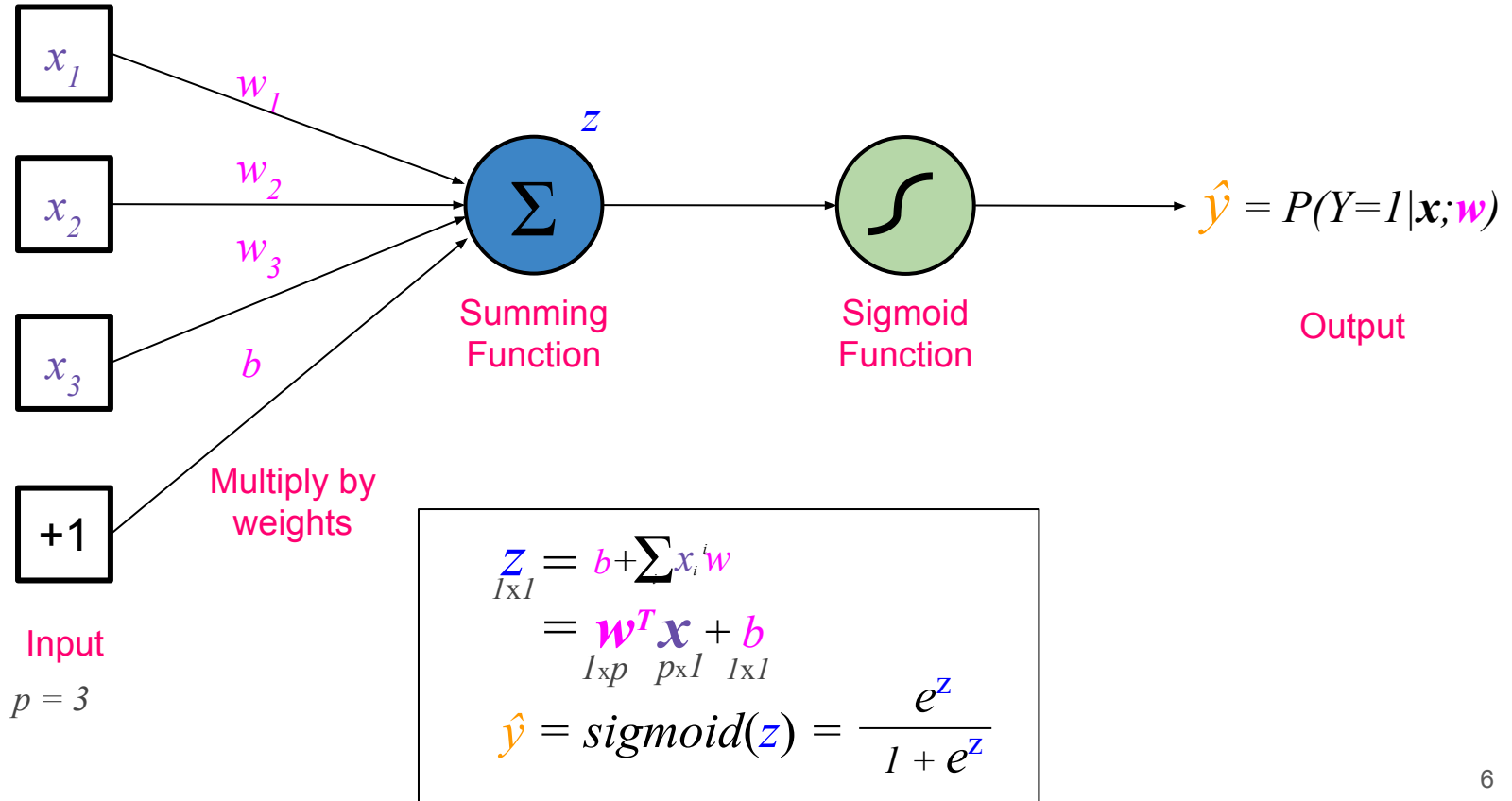
$$P(Y=1|\mathbf{x}) = \frac{e^{w\mathbf{x}+b}}{1 + e^{w\mathbf{x}+b}}$$

**Sigmoid Function**  
(aka logistic, logit, "S", soft-step)

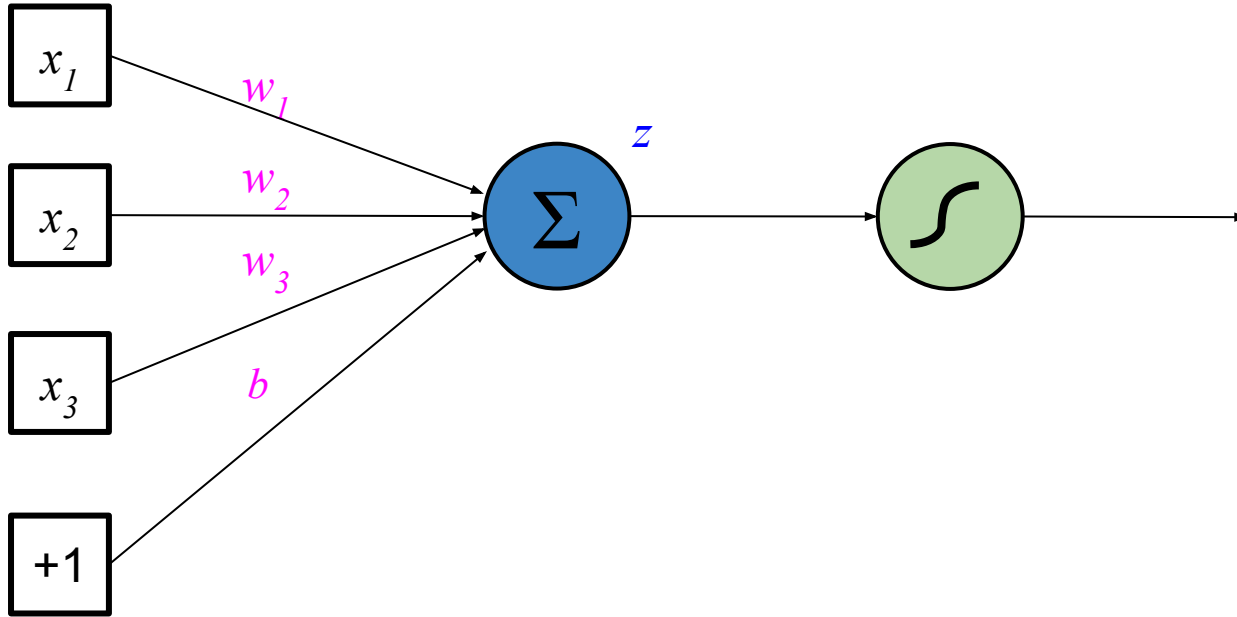
# Logistic Regression



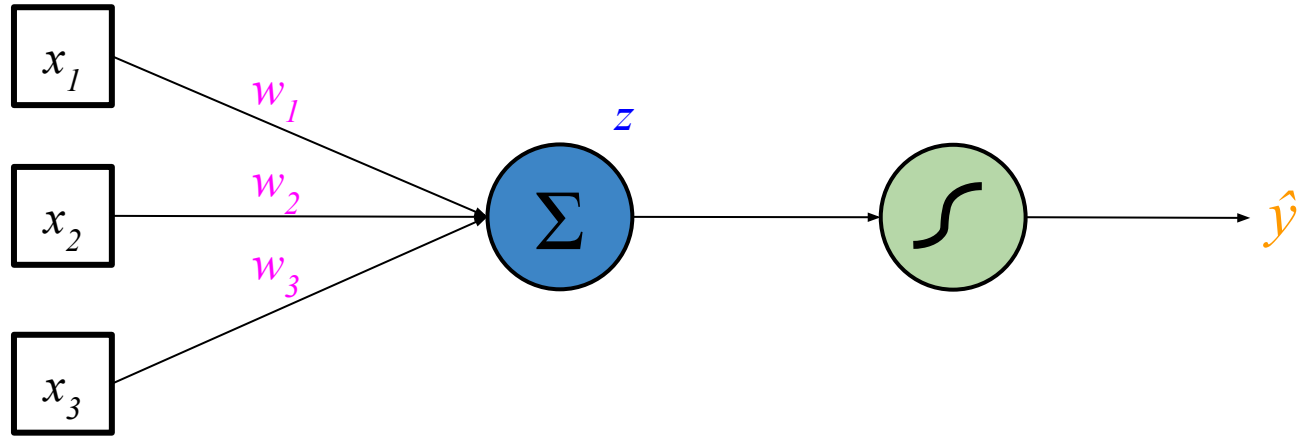
# Logistic Regression



# “Neuron” or “Perceptron”



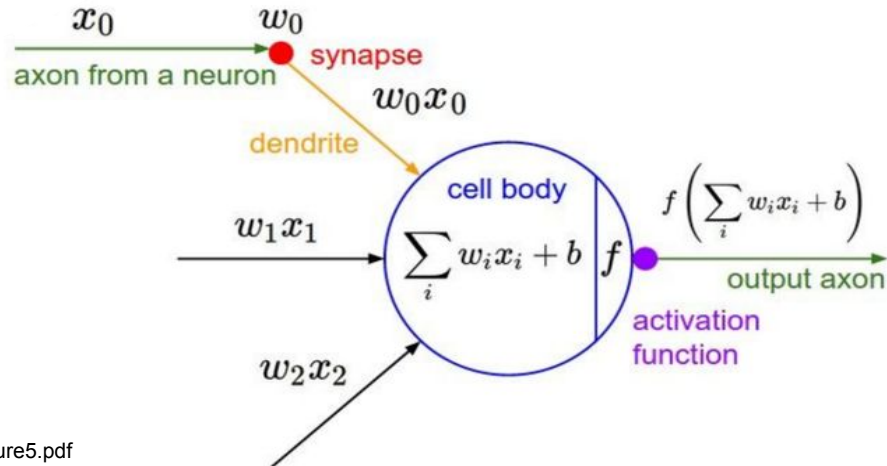
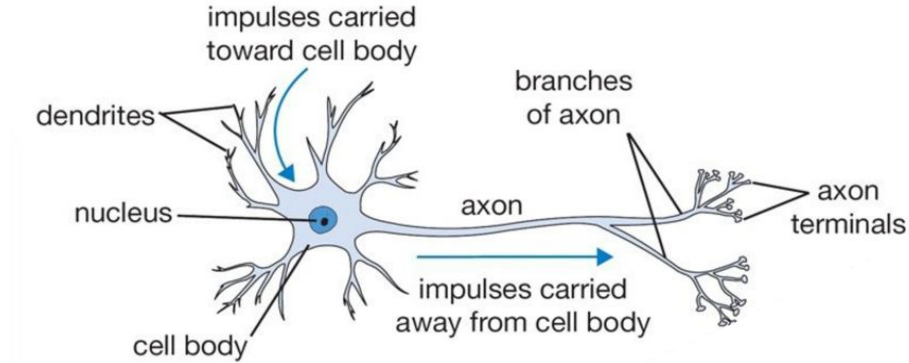
# “Neuron” or “Perceptron”



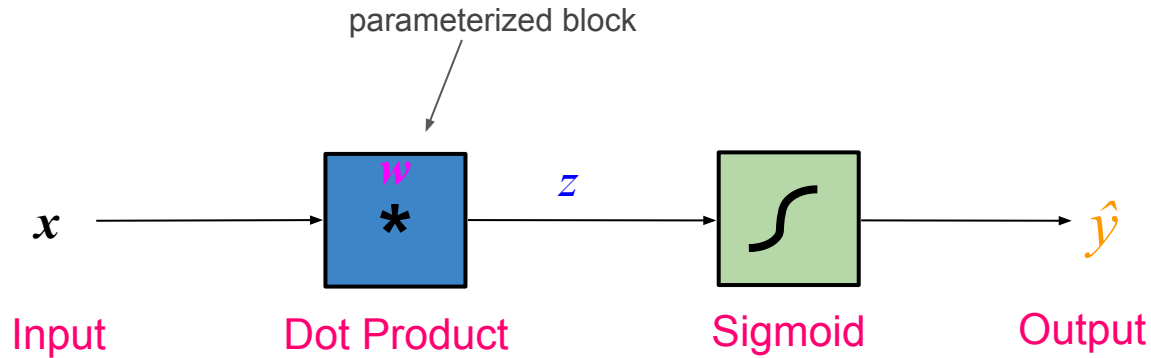
From here on, we leave out bias for simplicity



# Neurons



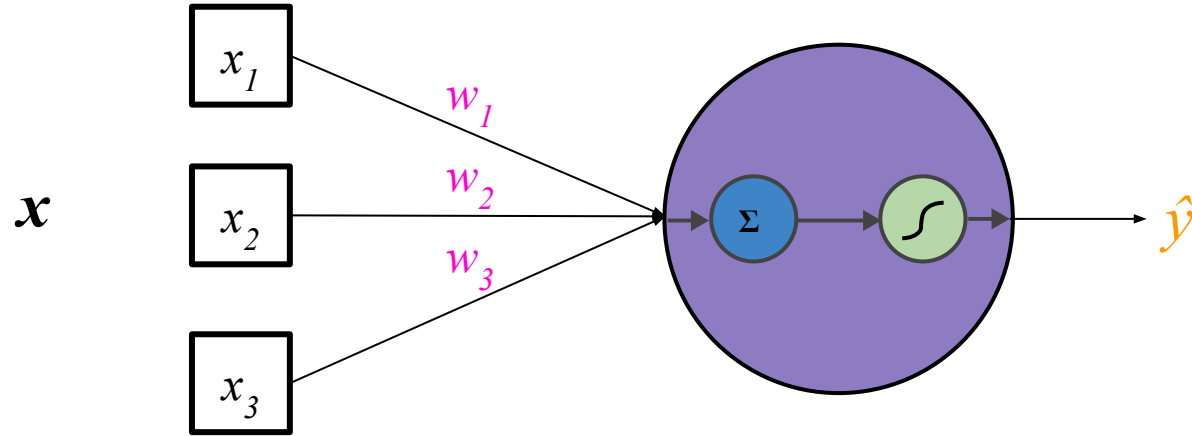
# “Block View” of a Neuron



$$\begin{matrix} z & = & w^T x \\ 1 \times 1 & & 1 \times p \quad p \times 1 \end{matrix}$$

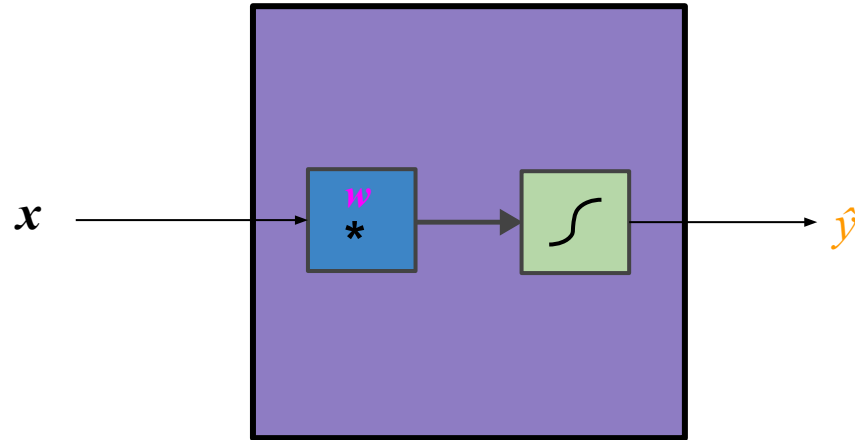
$$\hat{y} = \text{sigmoid}(z) = \frac{e^z}{1 + e^z}$$

# Neuron Representation



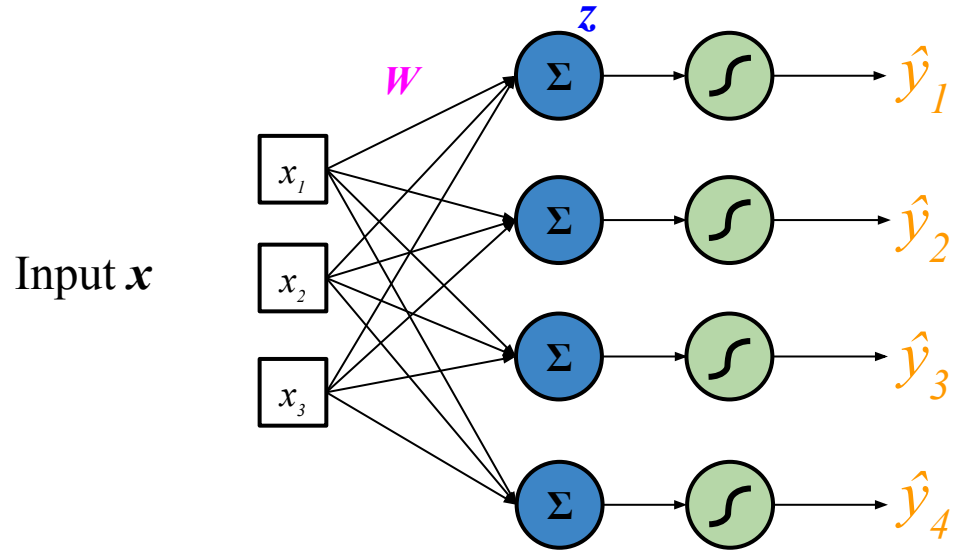
The linear transformation and nonlinearity together is typically considered a single neuron

# Neuron Representation

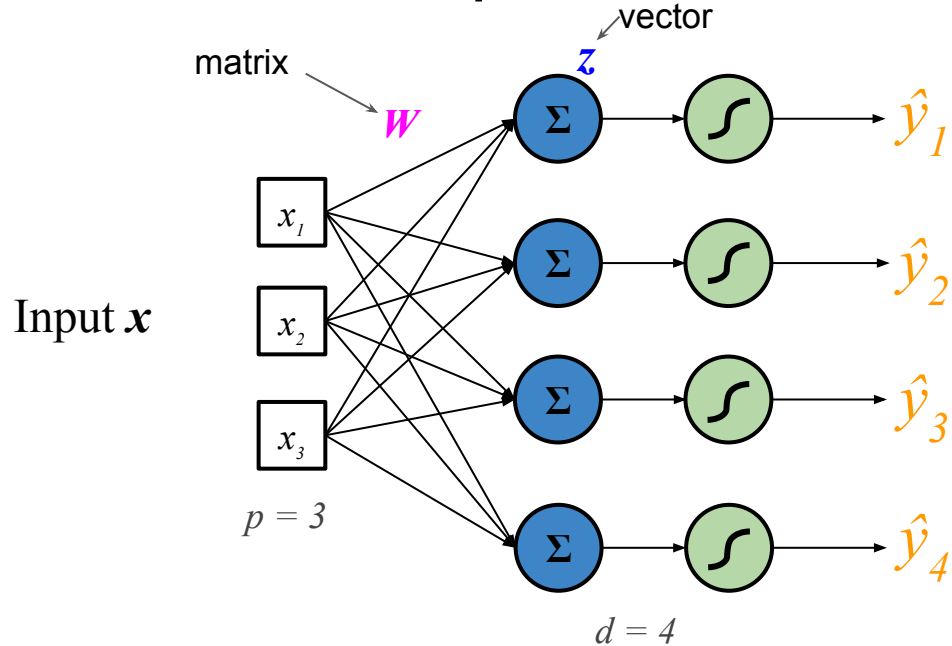


# Neural Networks

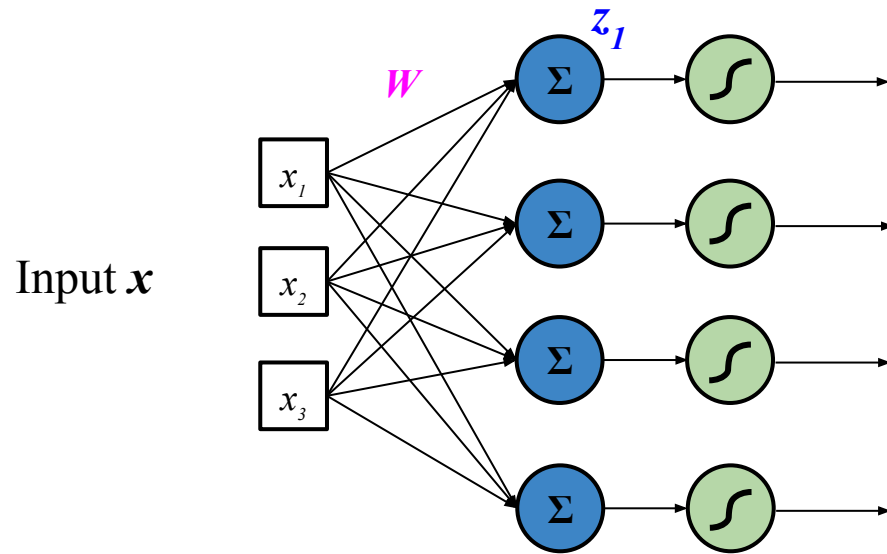
# Multiple Neurons



# Multiple Neurons



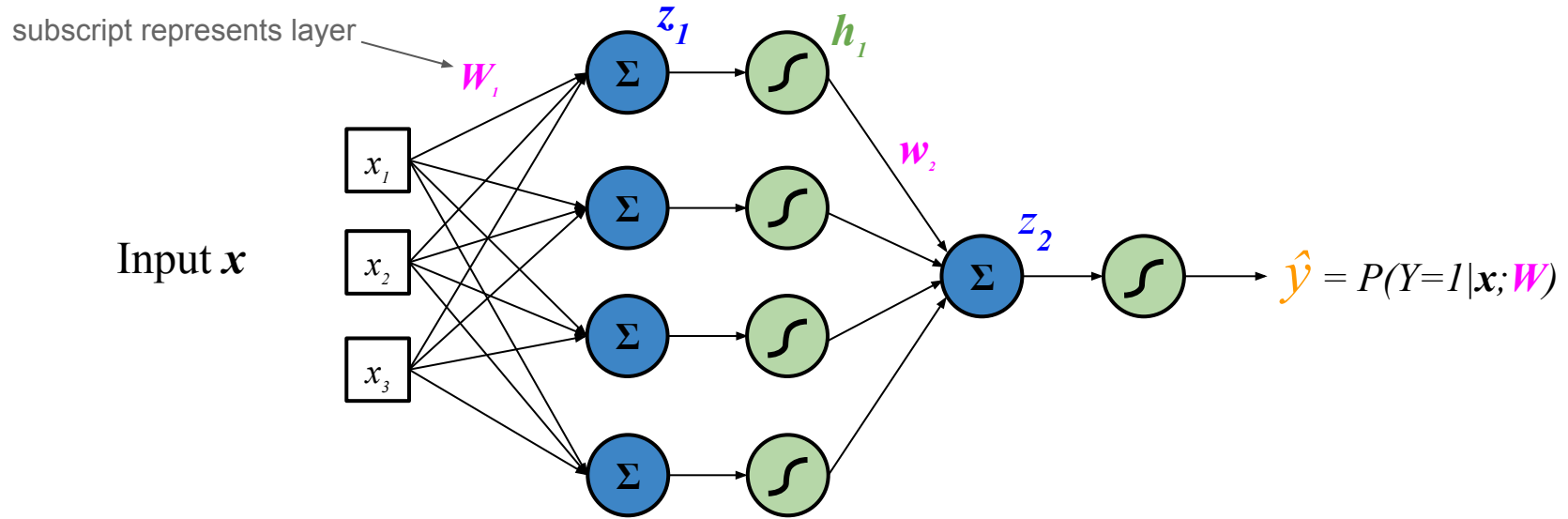
$$\begin{aligned} \underset{d \times 1}{\mathbf{z}} &= \underset{d \times p}{\mathbf{W}^T} \underset{p \times 1}{\mathbf{x}} \\ \underset{d \times 1}{\hat{\mathbf{y}}} &= \underset{d \times 1}{\text{sigmoid}(\mathbf{z})} = \frac{e^{\mathbf{z}}}{1 + e^{\mathbf{z}}} \end{aligned}$$



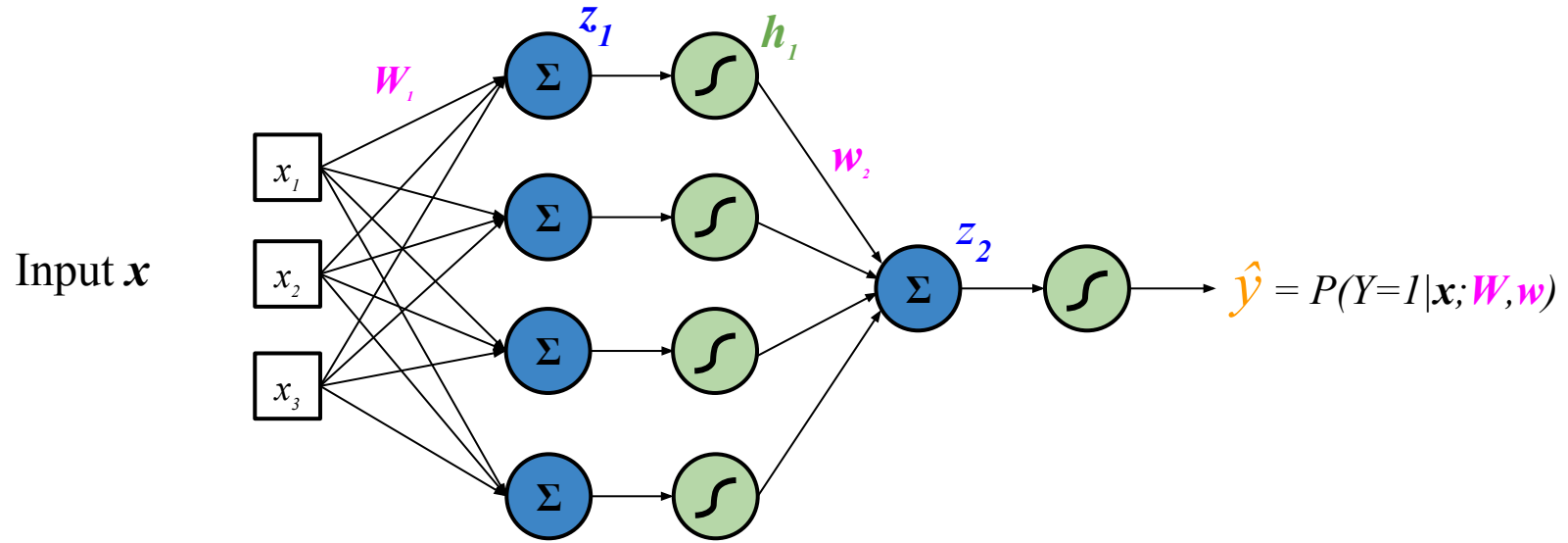
→  $\hat{y} = P(Y=I|\mathbf{x}; W)$



# Neural Network (1 Hidden Layer)

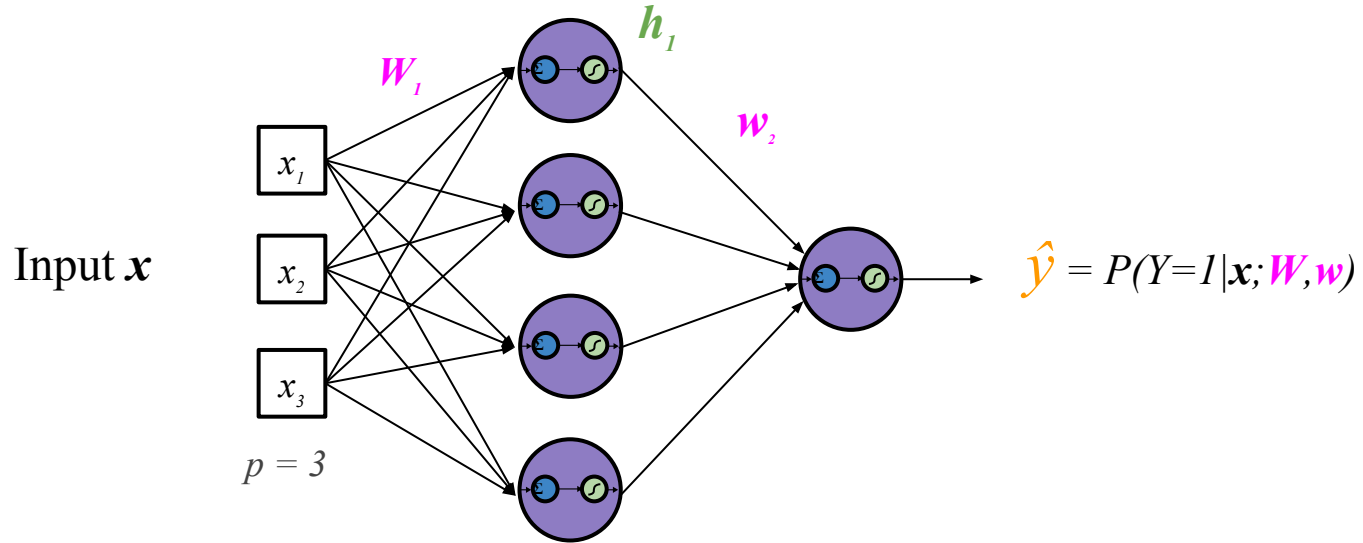


# Neural Network (1 Hidden Layer)



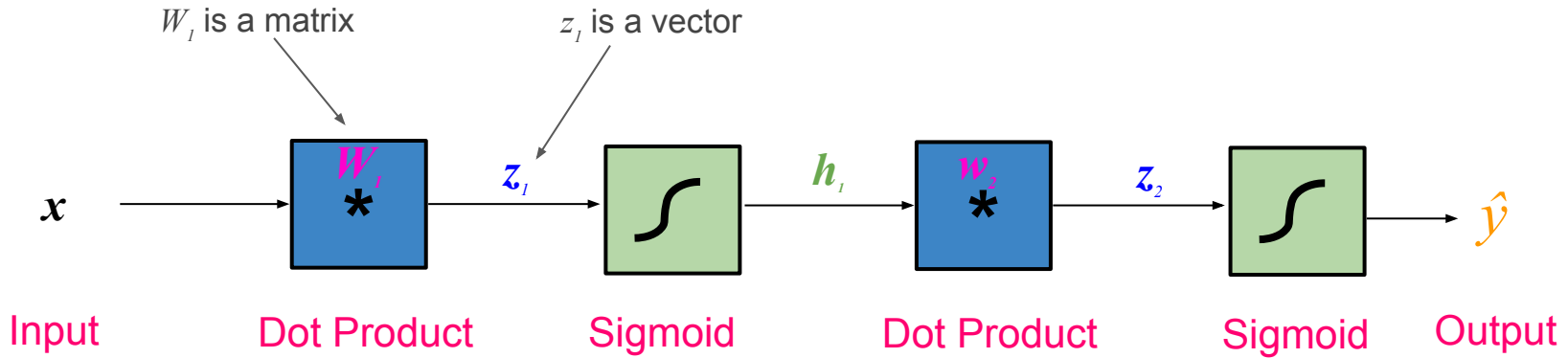
$$\begin{aligned} \underset{dx1}{z_1} &= \underset{d \times p}{W_1^T} \underset{px1}{\mathbf{x}} \\ \underset{dx1}{h_1} &= \underset{dx1}{\text{sigmoid}}(\underset{dx1}{z_1}) \\ \underset{l \times 1}{z_2} &= \underset{l \times p}{w_2^T} \underset{px1}{h_1} \\ \underset{l \times 1}{\hat{y}} &= \underset{l \times 1}{\text{sigmoid}}(\underset{l \times 1}{z_2}) \end{aligned}$$

# Neural Network (1 Hidden Layer)



$$\begin{aligned} \mathbf{z}_1 &= \mathbf{W}_1^T \mathbf{x} \\ \mathbf{z}_1 &_{dx1} \quad \mathbf{W}_1^T \quad \mathbf{x} \\ &_{d \times p} \quad \mathbf{W}_1^T \quad p \times 1 \\ \mathbf{h}_1 &= \text{sigmoid}(\mathbf{z}_1) \\ \mathbf{h}_1 &_{dx1} \quad \text{sigmoid} \quad \mathbf{z}_1 \\ &_{d \times 1} \quad \text{sigmoid} \quad d \times 1 \\ \mathbf{z}_2 &= \mathbf{w}_2^T \mathbf{h}_1 \\ \mathbf{z}_2 &_{l \times 1} \quad \mathbf{w}_2^T \quad \mathbf{h}_1 \\ &_{l \times p_2} \quad \mathbf{w}_2^T \quad p \times 1 \\ \hat{y} &= \text{sigmoid}(\mathbf{z}_2) \\ \hat{y} &_{l \times 1} \quad \text{sigmoid} \quad \mathbf{z}_2 \\ &_{l \times 1} \quad \text{sigmoid} \quad l \times 1 \end{aligned}$$

# Neural Network (1 Hidden Layer)

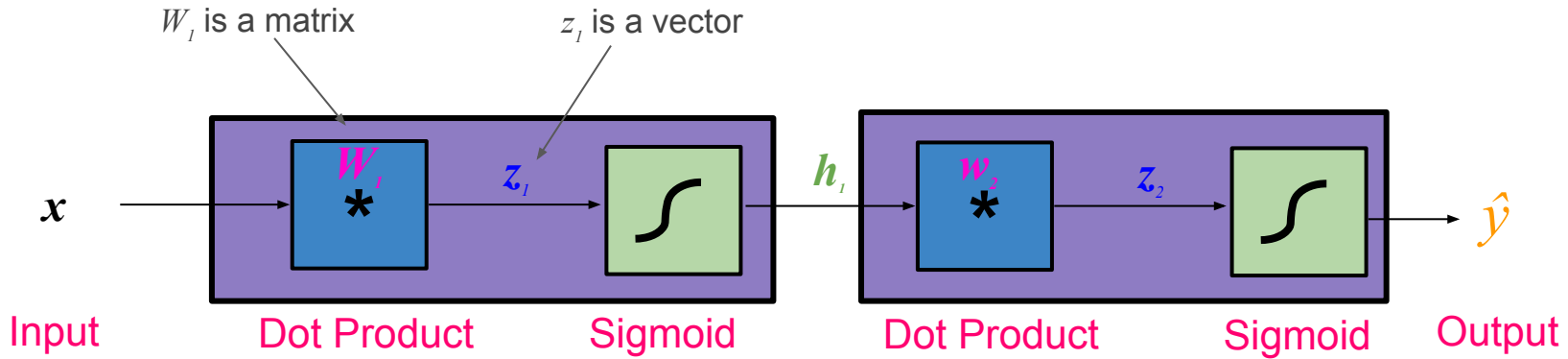


$$\begin{aligned} \mathbf{z}_1 &= \mathbf{W}_1^T \mathbf{x} \\ \mathbf{h}_1 &= \text{sigmoid}(\mathbf{z}_1) \\ \mathbf{z}_2 &= \mathbf{w}_2^T \mathbf{h}_1 \\ \hat{\mathbf{y}} &= \text{sigmoid}(\mathbf{z}_2) \end{aligned}$$

Dimensional annotations for the equations above:

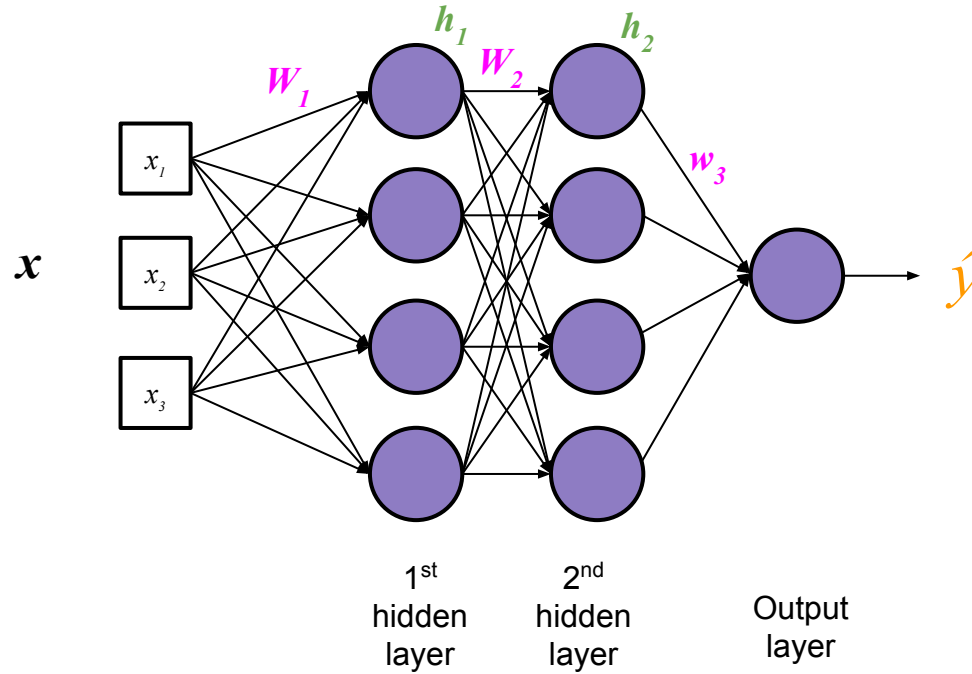
- $\mathbf{z}_1$ :  $dx1$
- $\mathbf{W}_1$ :  $dxp$  (rows),  $px1$  (columns)
- $\mathbf{x}$ :  $px1$
- $\mathbf{h}_1$ :  $dx1$
- $\mathbf{z}_2$ :  $lx1$
- $\mathbf{w}_2$ :  $lpx2$  (rows),  $px1$  (columns)
- $\mathbf{h}_1$ :  $px1$
- $\hat{\mathbf{y}}$ :  $lx1$
- $\mathbf{z}_2$ :  $lx1$

# Neural Network (1 Hidden Layer)

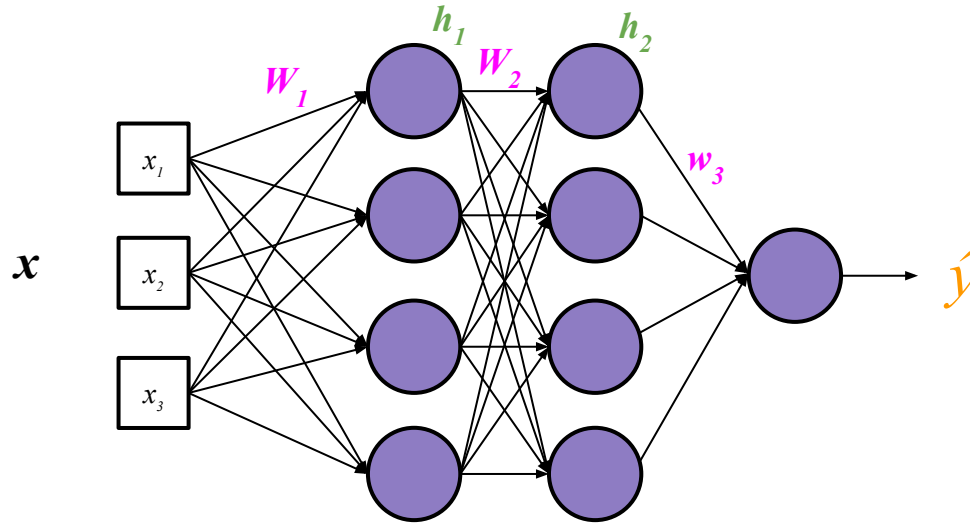


$$\begin{aligned} z_1 &= W_1^T x \\ &\begin{matrix} dx1 \\ dxp & px1 \end{matrix} \\ h_1 &= \text{sigmoid}(z_1) \\ &\begin{matrix} dx1 \\ dx1 \end{matrix} \\ z_2 &= w_2^T h_1 \\ &\begin{matrix} lx1 \\ lxp_2 & px1 \end{matrix} \\ \hat{y} &= \text{sigmoid}(z_2) \\ &\begin{matrix} lx1 \\ lx1 \end{matrix} \end{aligned}$$

# Neural Network (2 Hidden Layers)



# Neural Network (2 Hidden Layers)

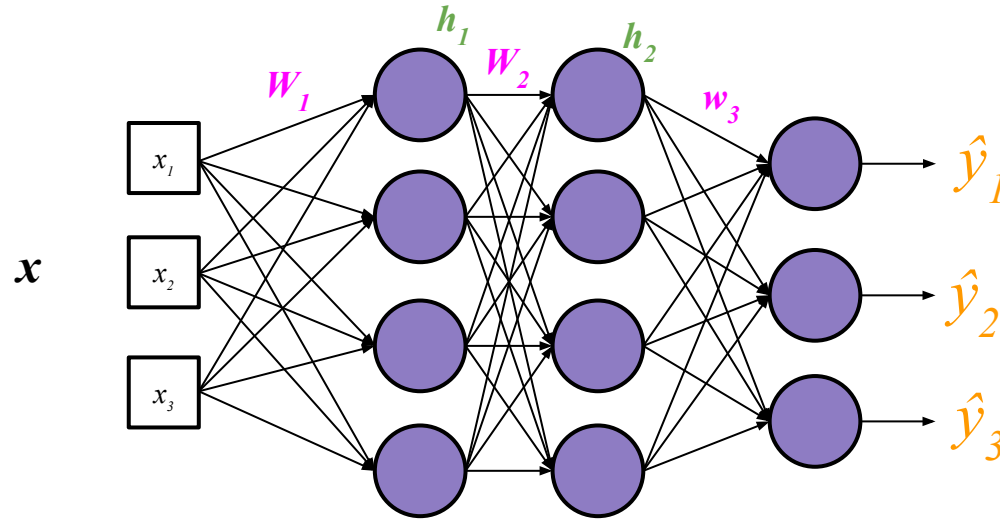


hidden layer 1 output  $\rightarrow$

hidden layer 2 output  $\rightarrow$

$$\begin{aligned} z_1 &= W_1^T x \\ h_1 &= \text{sigmoid}(z_1) \\ z_2 &= W_2^T h_1 \\ h_2 &= \text{sigmoid}(z_2) \\ z_3 &= w_3^T h_2 \\ \hat{y} &= \text{sigmoid}(z_3) \end{aligned}$$

# Multi-Class Neural Network (2 Hidden Layers)



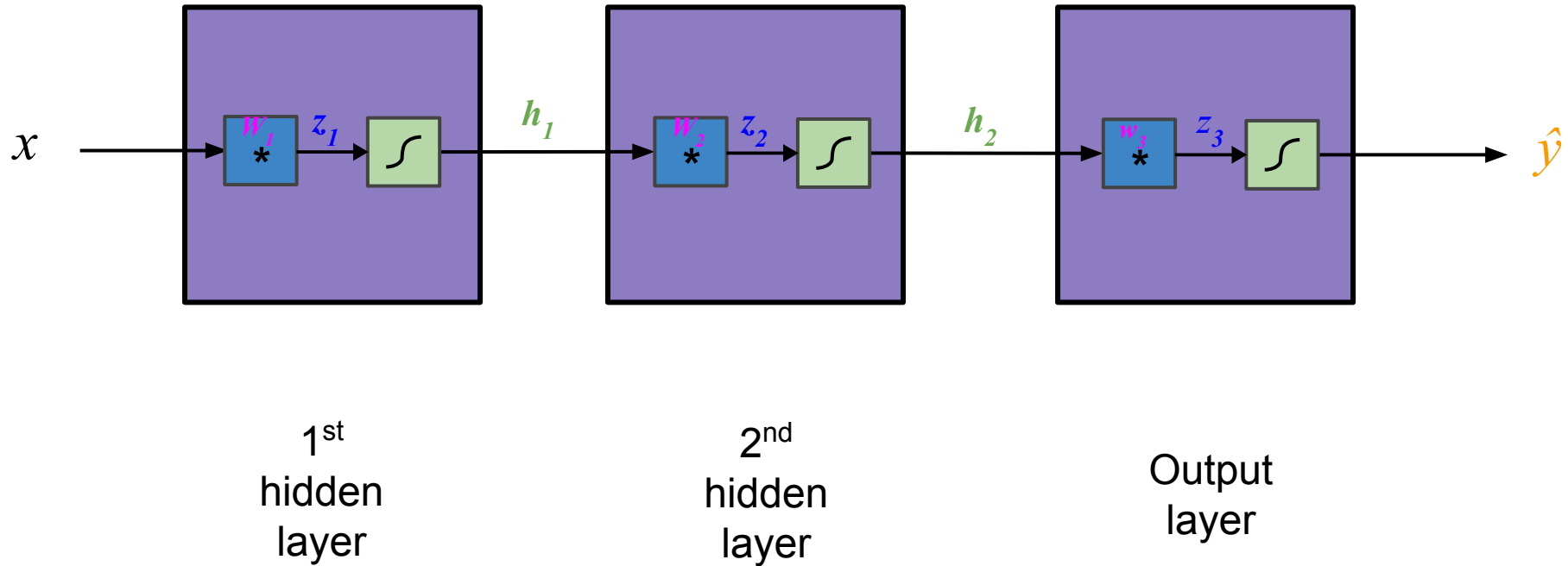
hidden layer 1 output  $\rightarrow$

hidden layer 2 output  $\rightarrow$

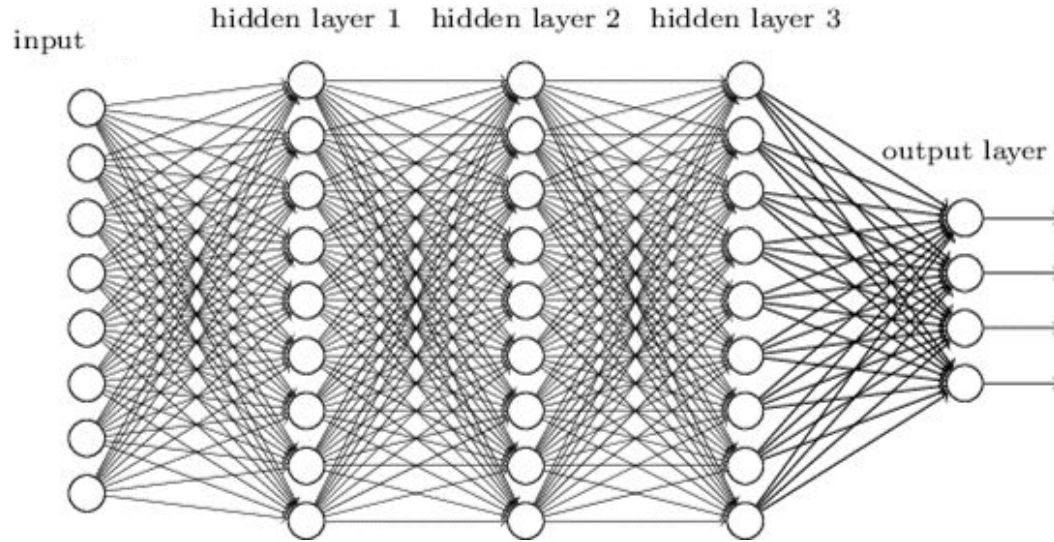
$$\begin{aligned} z_1 &= W_1^T \mathbf{x} \\ h_1 &= \text{sigmoid}(z_1) \\ z_2 &= W_2^T h_1 \\ h_2 &= \text{sigmoid}(z_2) \\ z_3 &= w_3^T h_2 \\ \hat{\mathbf{y}} &= \text{sigmoid}(z_3) \end{aligned}$$



# “Block View” Of Multi-Layer Neural Network



# “Deep” Neural Networks (i.e. $> 1$ hidden layer)



# Loss Functions

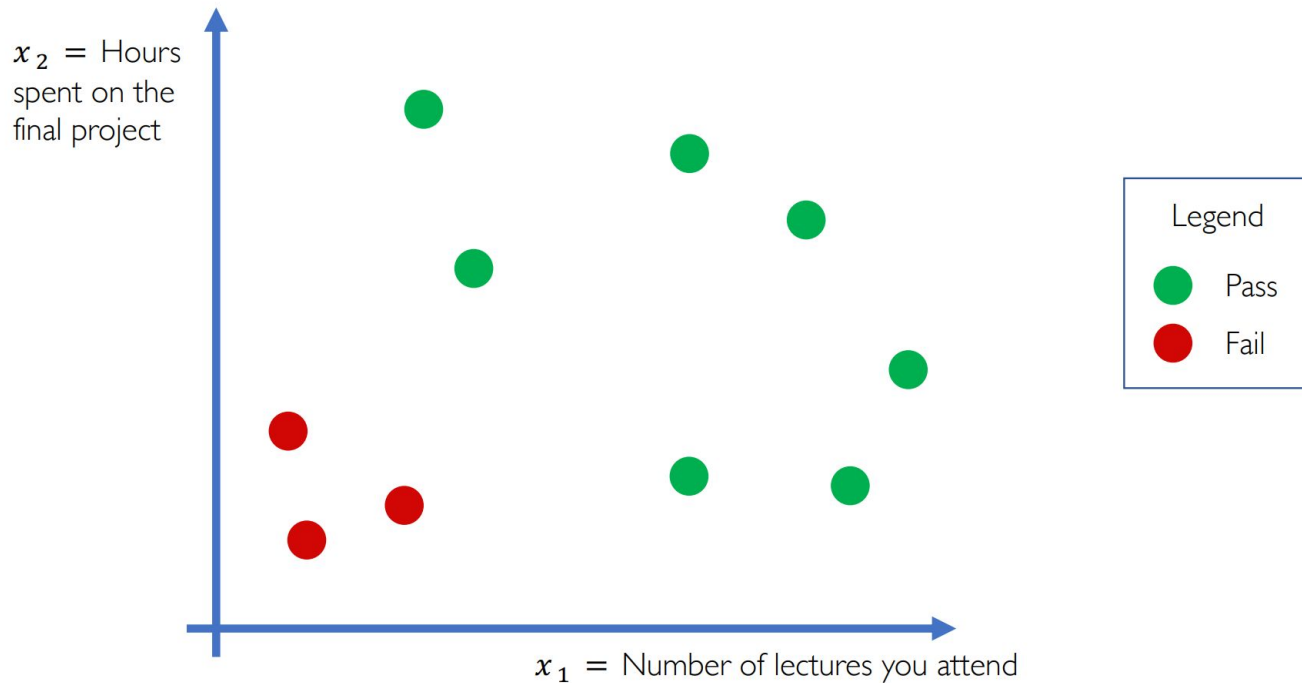
# Toy Example: Will I pass this class?

Let's start with a simple two feature model

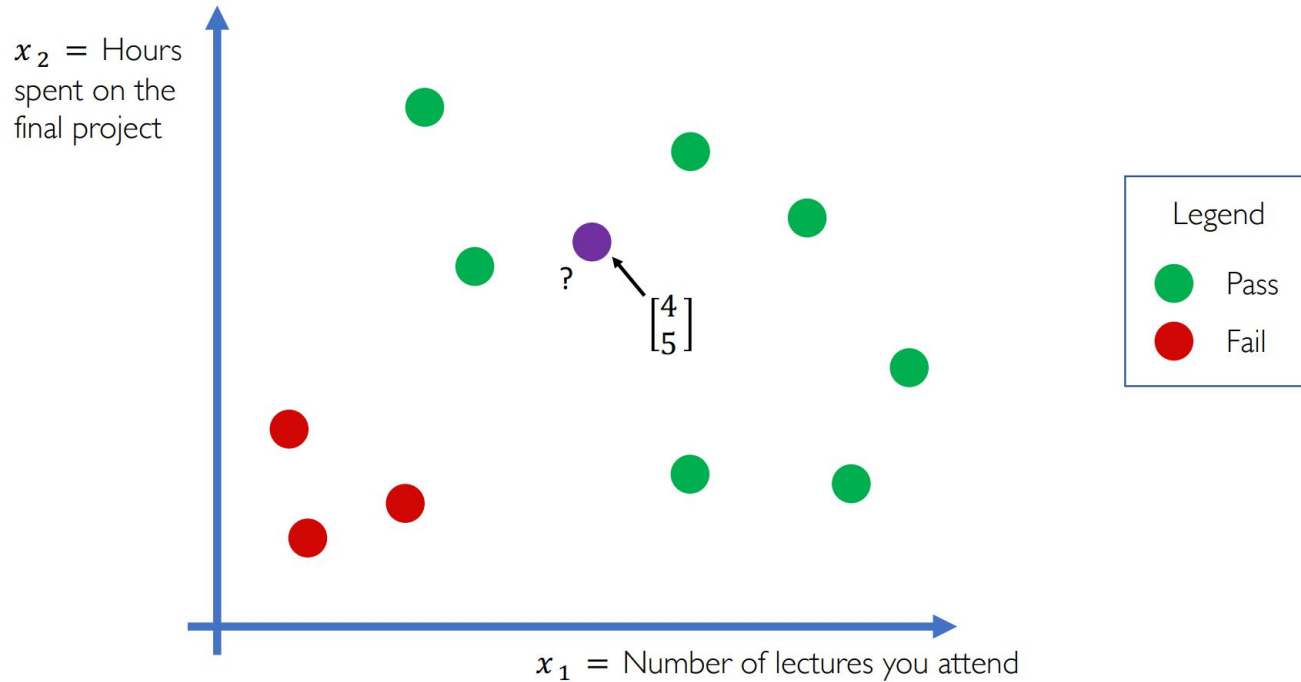
$x_1$  = Number of lectures you attend

$x_2$  = Hours spent on the final project

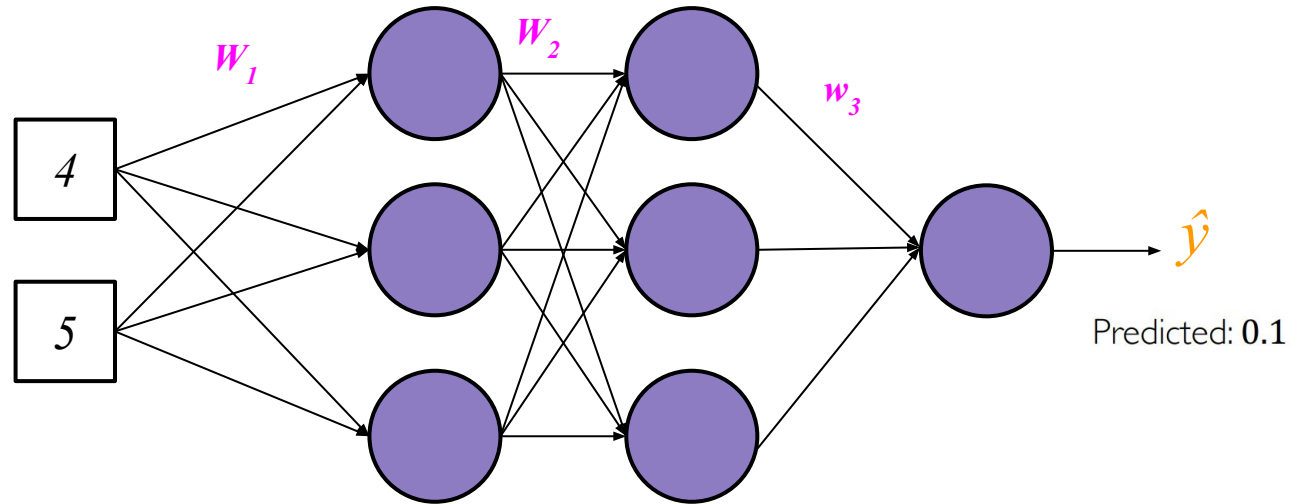
# Toy Example: Will I pass this class?



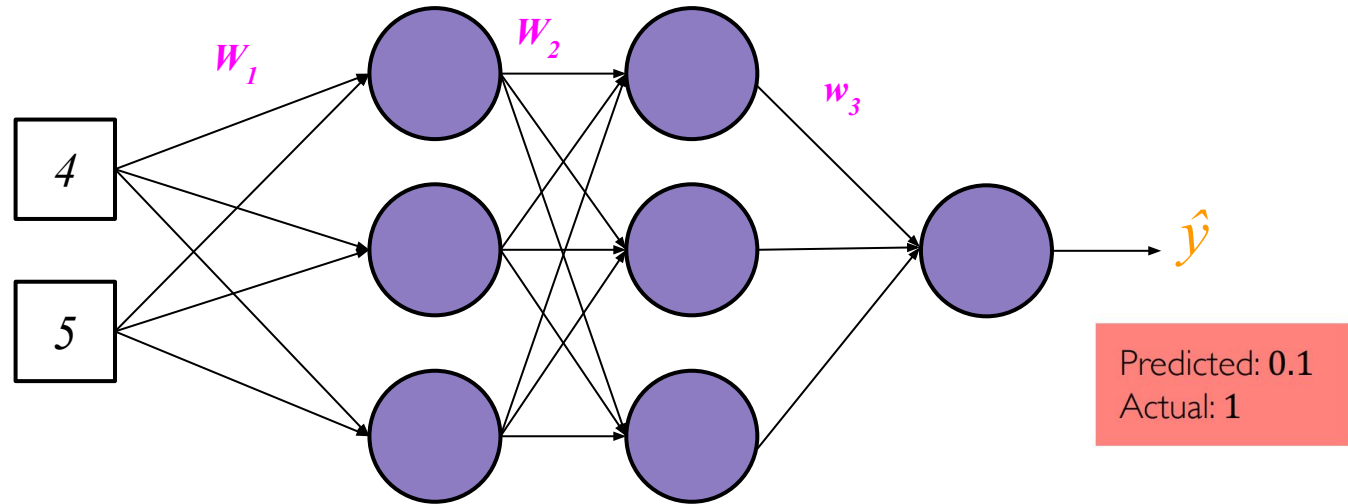
# Toy Example: Will I pass this class?



# Toy Example: Will I pass this class?



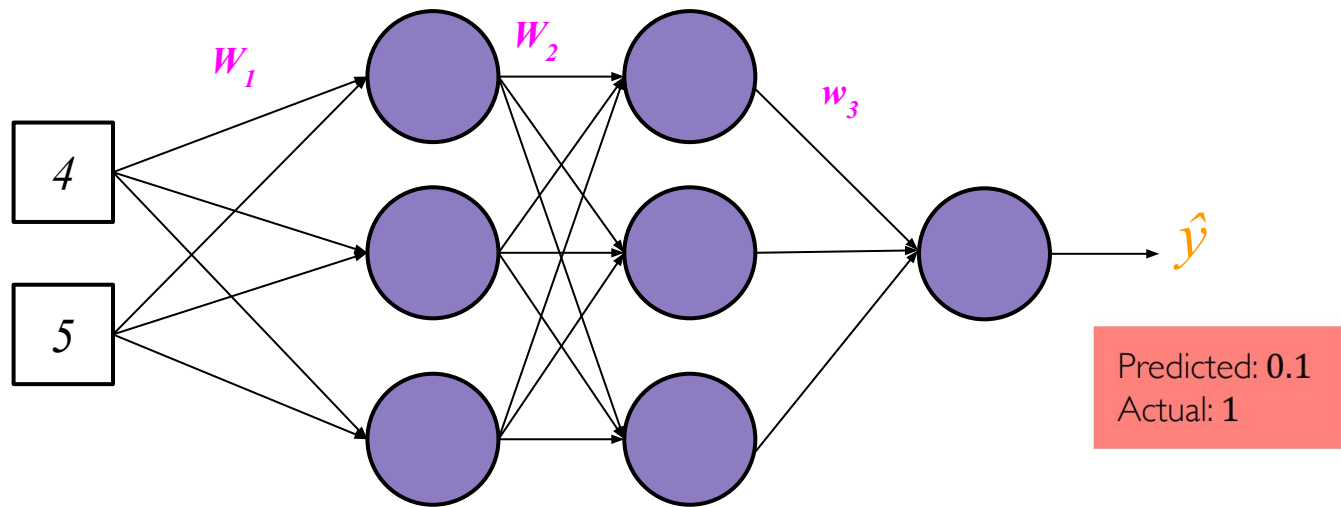
# Toy Example: Will I pass this class?





# Toy Example: Will I pass this class?

The **loss** of our network measures the cost incurred from incorrect predictions



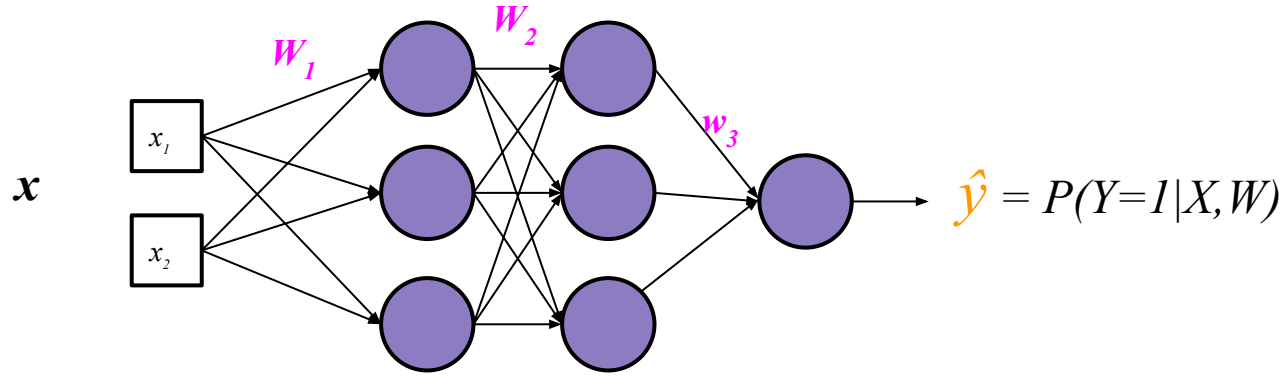
$$\mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

Predicted

Actual

# Binary Classification Loss

Binary Cross Entropy Loss

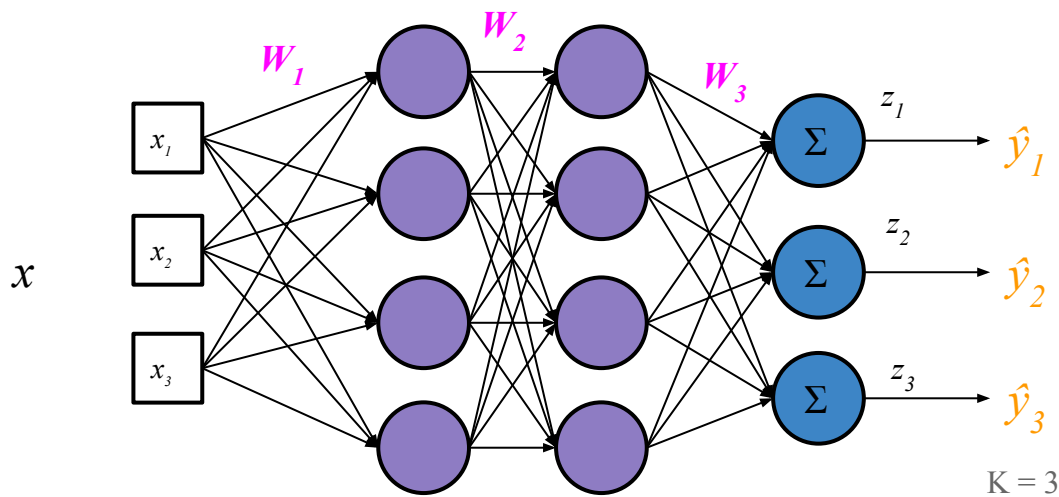


$$E = \text{loss} = -\log P(Y = \hat{y} | \mathbf{X} = \mathbf{x})$$
$$= -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

true output

# Multi-Class Classification Loss

## Cross Entropy Loss



$$\hat{y}_i = \frac{e^{z_i}}{\sum_j e^{z_j}} = P(y_i = I | \mathbf{x})$$

**“Softmax” function.**  
Normalizing function which converts each class output to a probability.

$$\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0.1 \\ 0.7 \\ 0.2 \end{pmatrix}$$

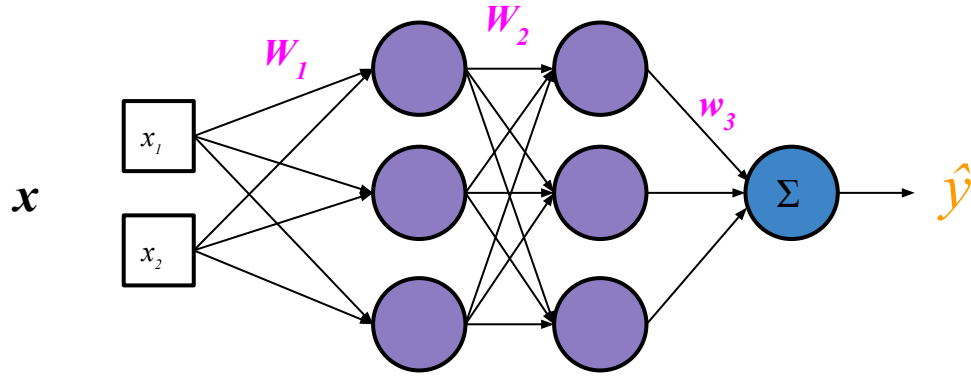
$\mathbf{y}$        $\hat{\mathbf{y}}$

$$E = \text{loss} = - \sum_{j=1 \dots K} y_j \log \hat{y}_j$$

“0” for all except true class

# Regression Loss

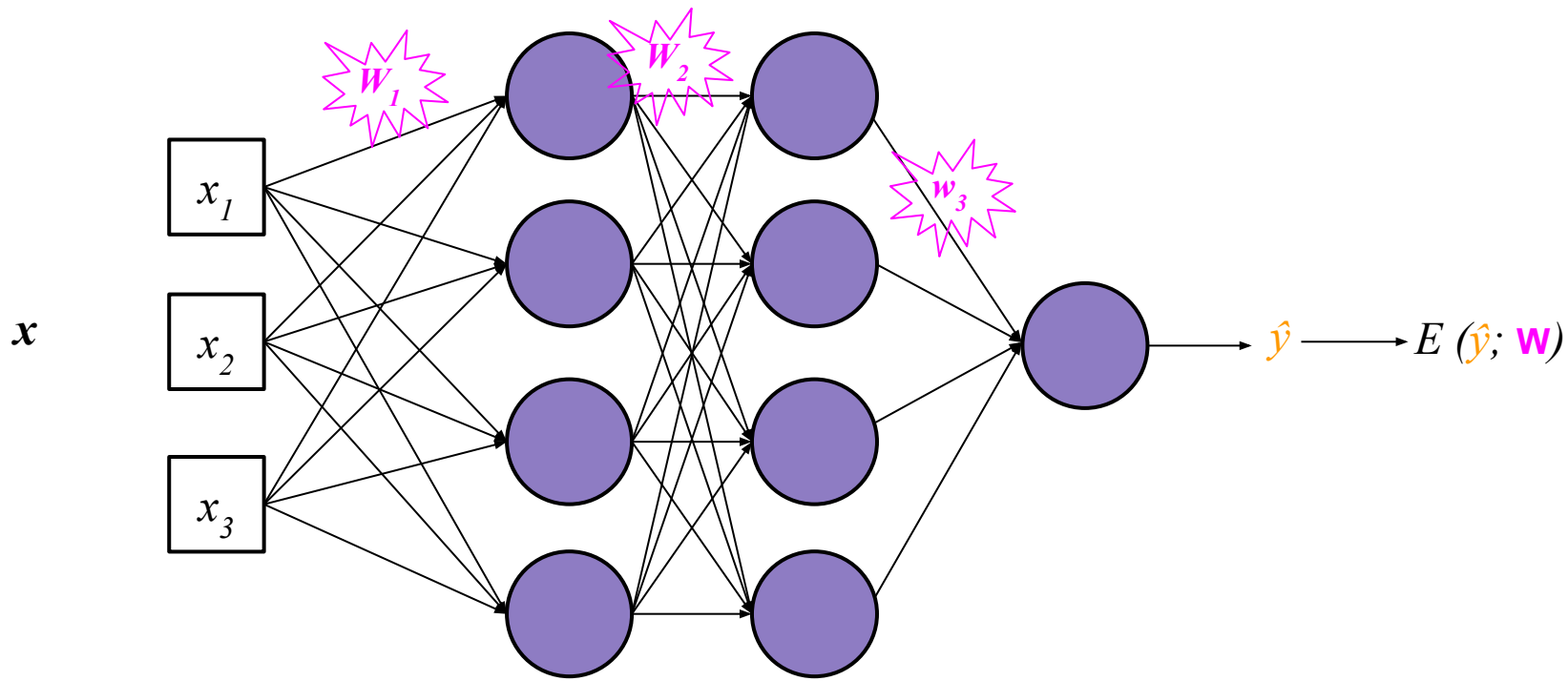
Mean Squared Error Loss



$$E = \text{loss} = \frac{1}{2} (y - \hat{y})^2$$

true output

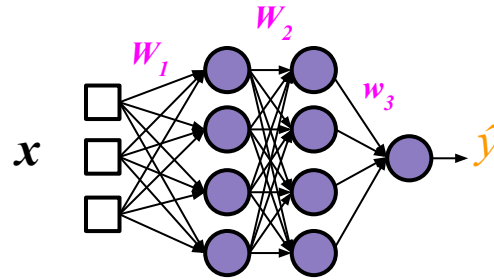
# Backpropagation



We want to find the network weights that **achieve the lowest loss**

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} E(\hat{y}; \mathbf{w})$$

# Training Neural Networks



How do we learn the optimal weights  $W_L$  for our task??

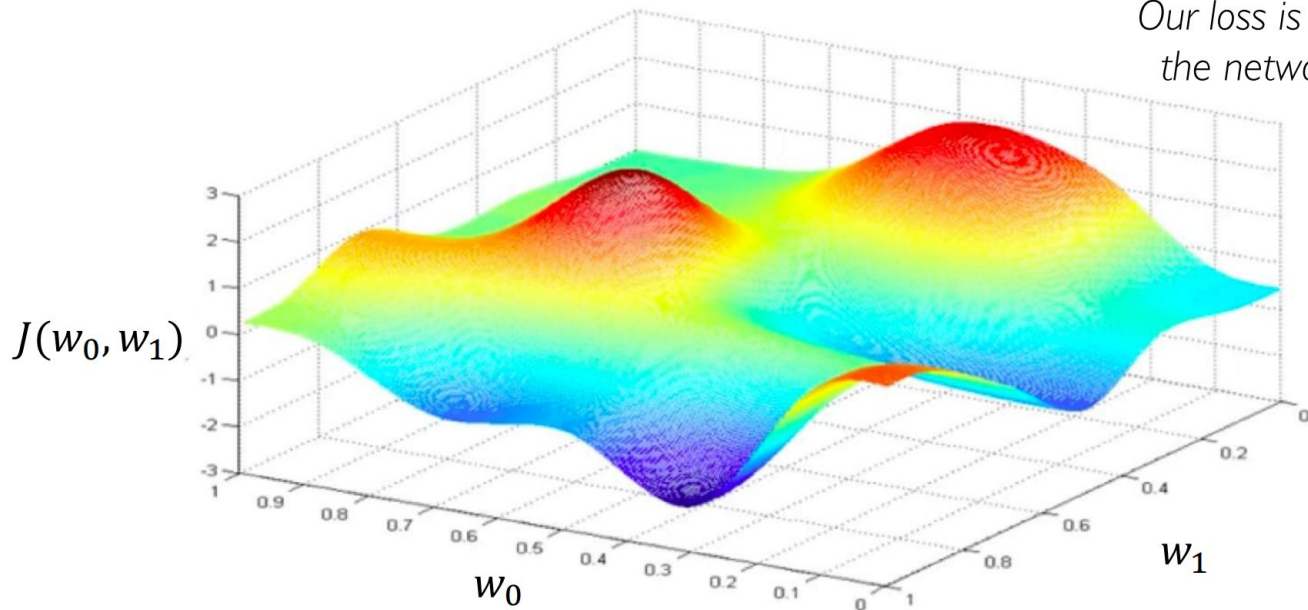
- **Gradient descent:**

$$W_L(t+1) = W_L(t) - \eta \frac{\partial E}{\partial W_L(t)}$$

# Gradient Descent

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} E(\hat{y}; \mathbf{W})$$

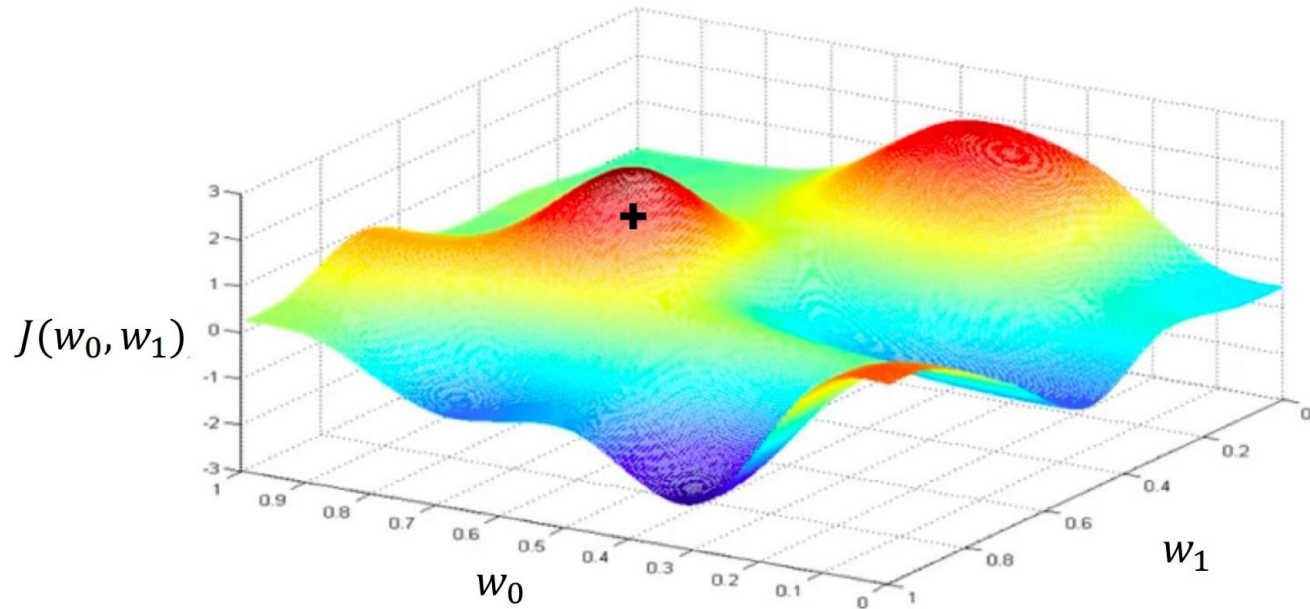
Remember:  
*Our loss is a function of  
the network weights!*





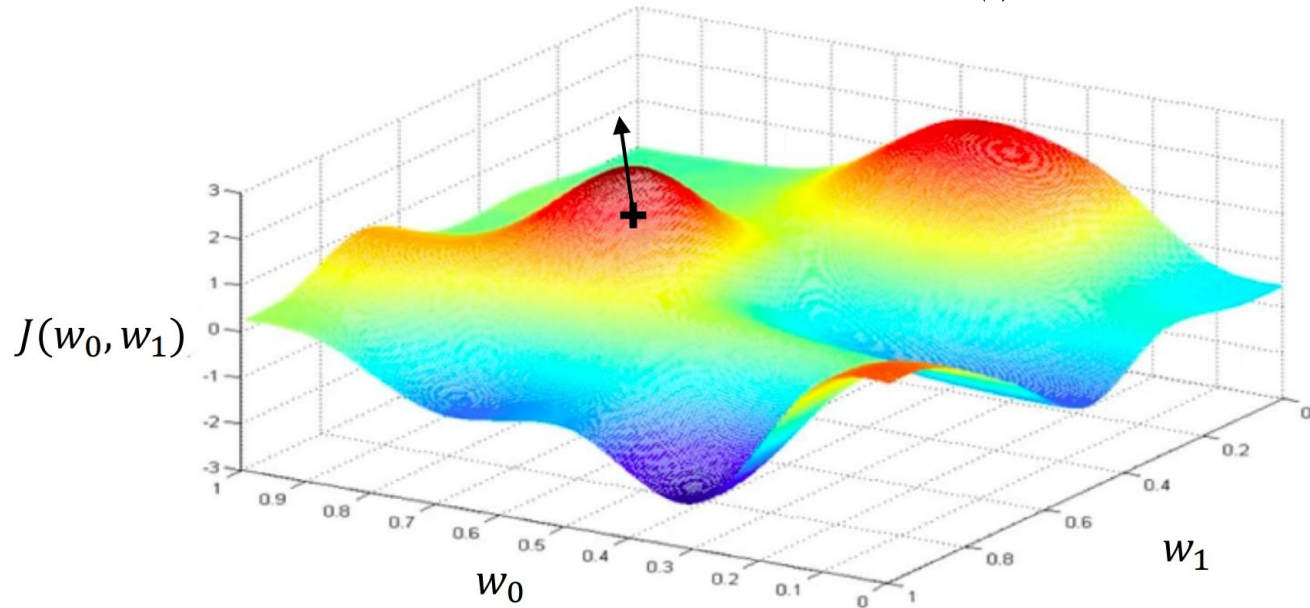
# Gradient Descent

Randomly pick an initial  $(w_0, w_1)$



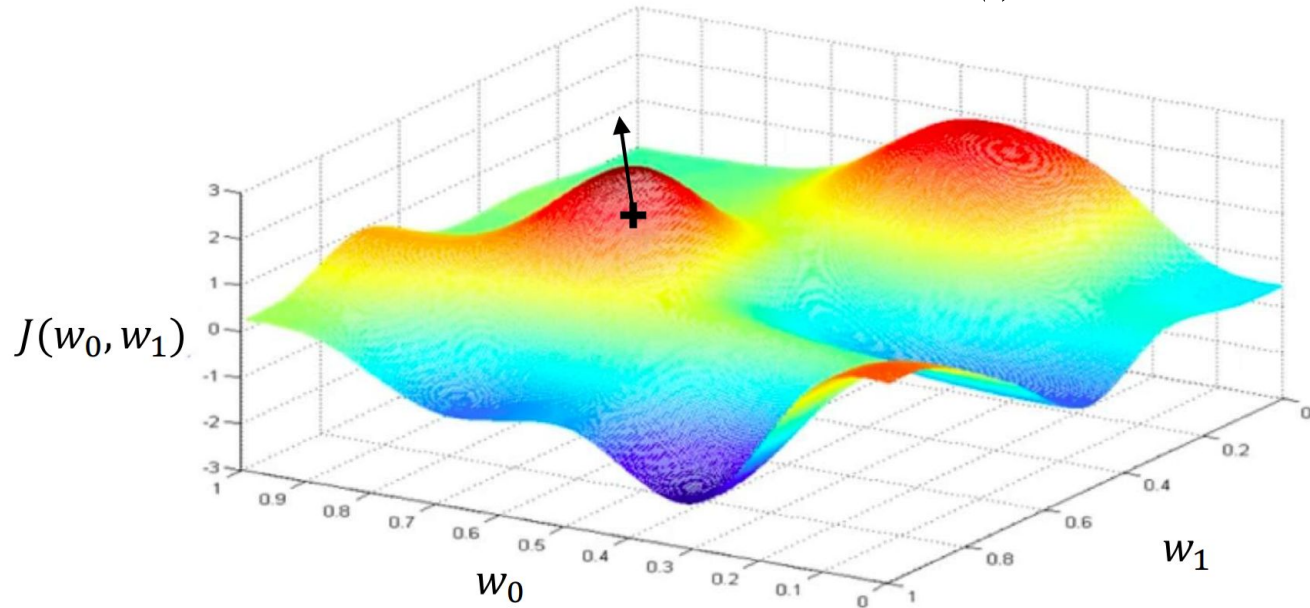
# Gradient Descent

Compute gradient,  $\frac{\partial E}{\partial W(t)}$



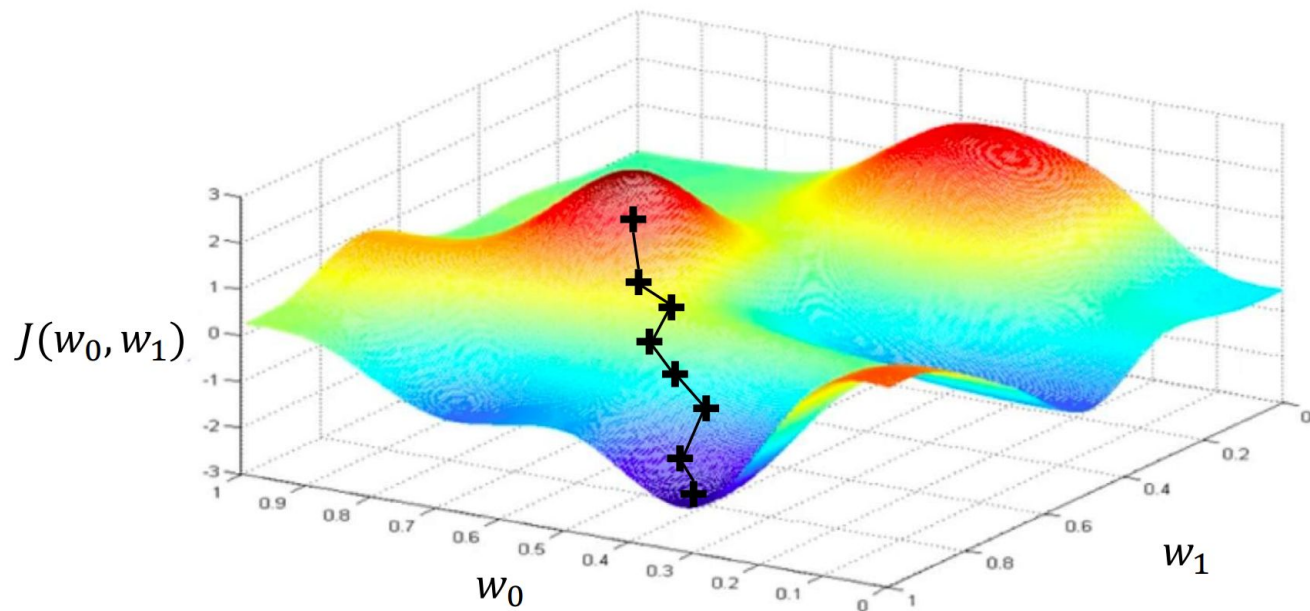
# Gradient Descent

Compute gradient,  $\frac{\partial E}{\partial W(t)}$

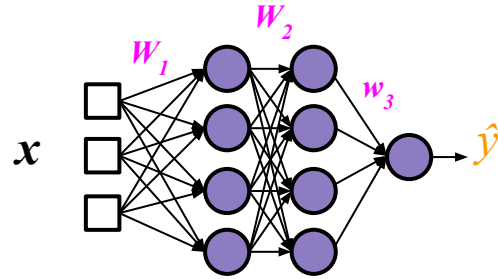


# Gradient Descent

Repeat until convergence



# Training Neural Networks



But how do we get gradients of lower layers (e.g.  $W_1$ ) ?

- **Backpropagation!**

- Repeated application of chain rule of calculus
- Locally minimize the objective
- Requires all “blocks” of the network to be differentiable

# Backpropagation Intro

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

# Backpropagation Intro

$$f(x, y, z) = (x + y)z$$

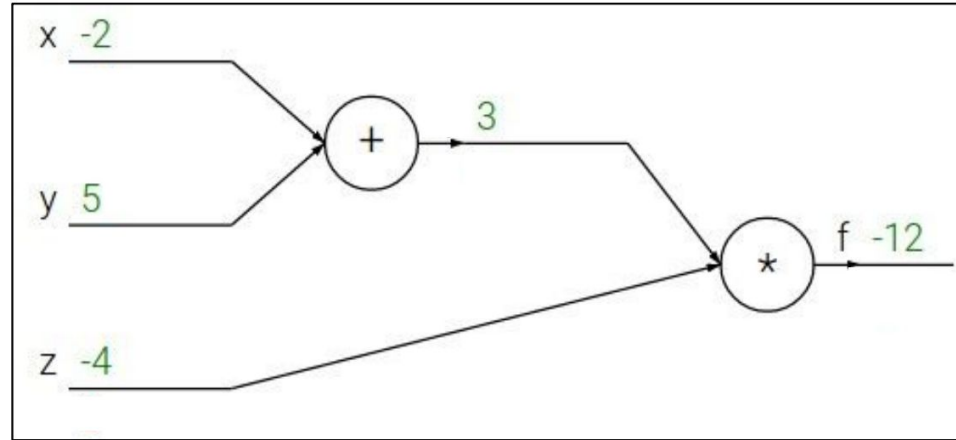
e.g.  $x = -2, y = 5, z = -4$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backpropagation Intro

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$



Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

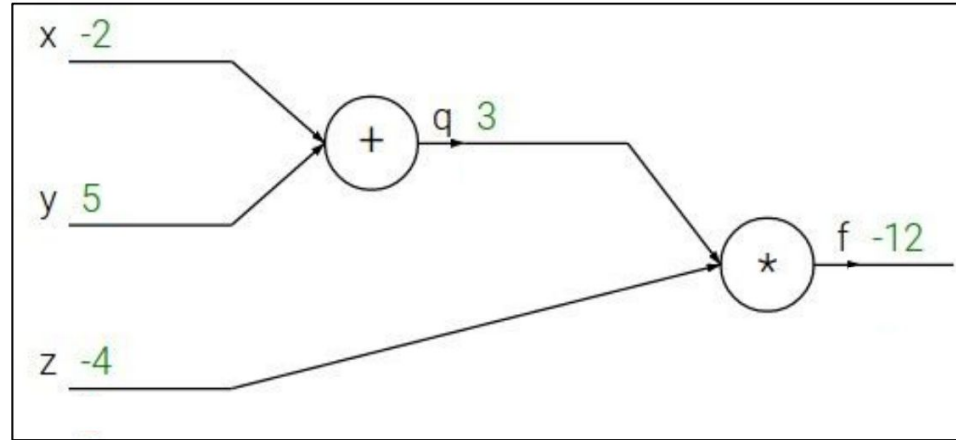


# Backpropagation Intro

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y$$



Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backpropagation Intro

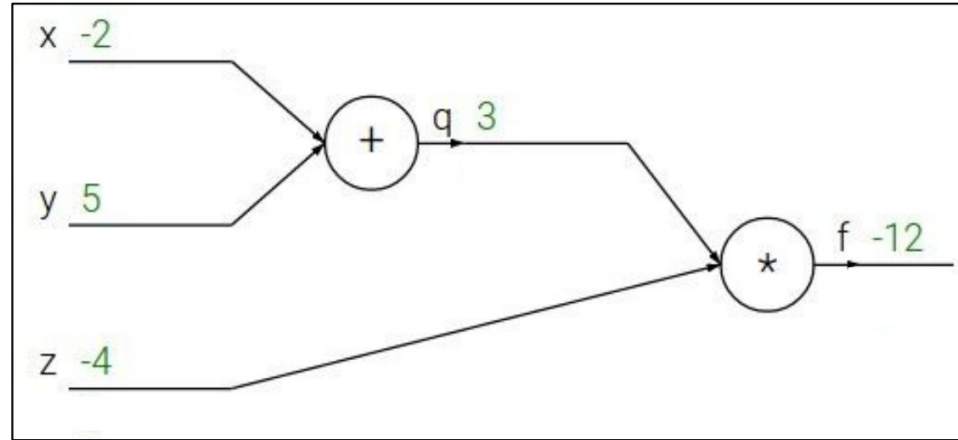
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y$$

$$f = qz$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



# Backpropagation Intro

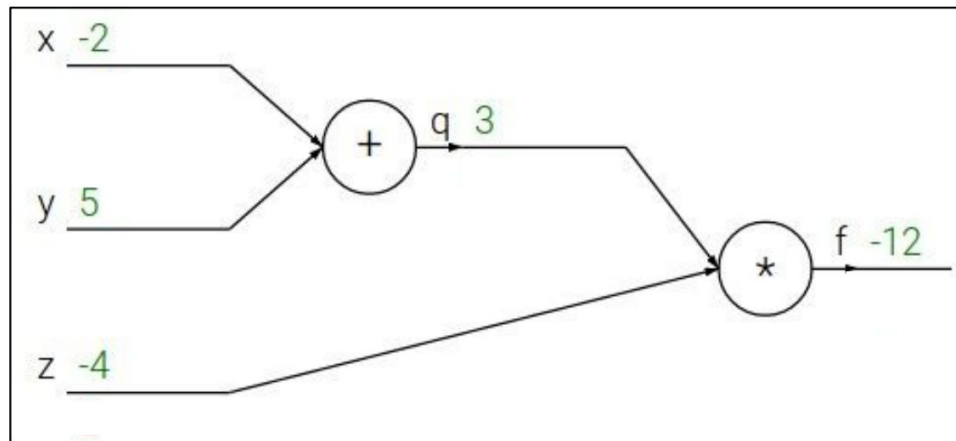
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



# Backpropagation Intro

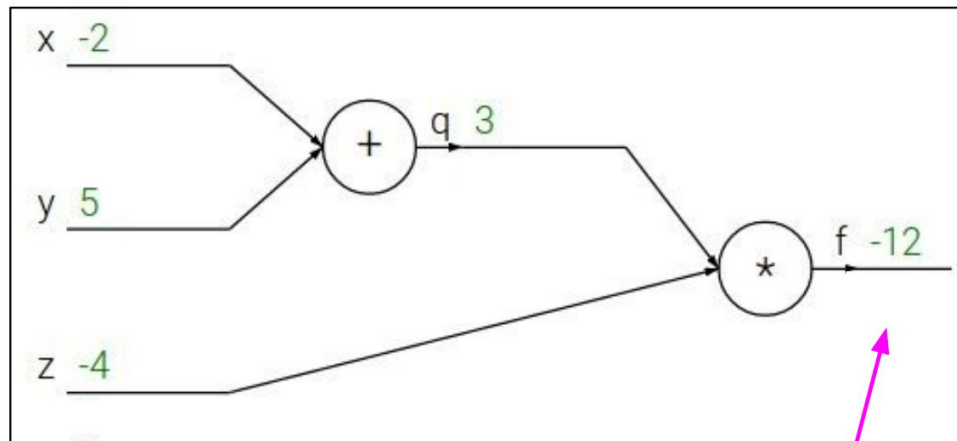
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial f}$$

# Backpropagation Intro

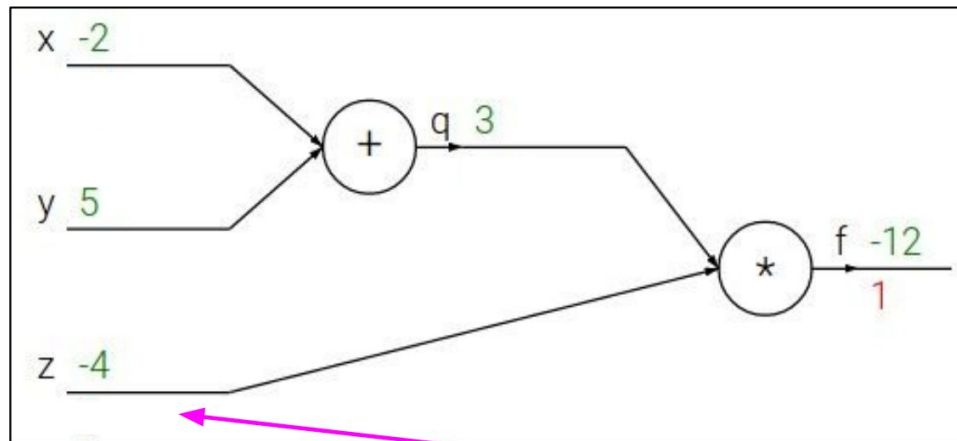
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

# Backpropagation Intro

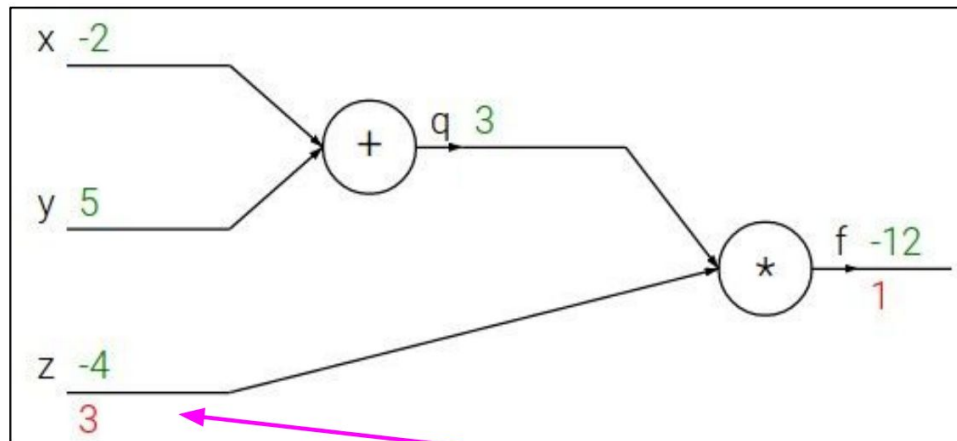
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

# Backpropagation Intro

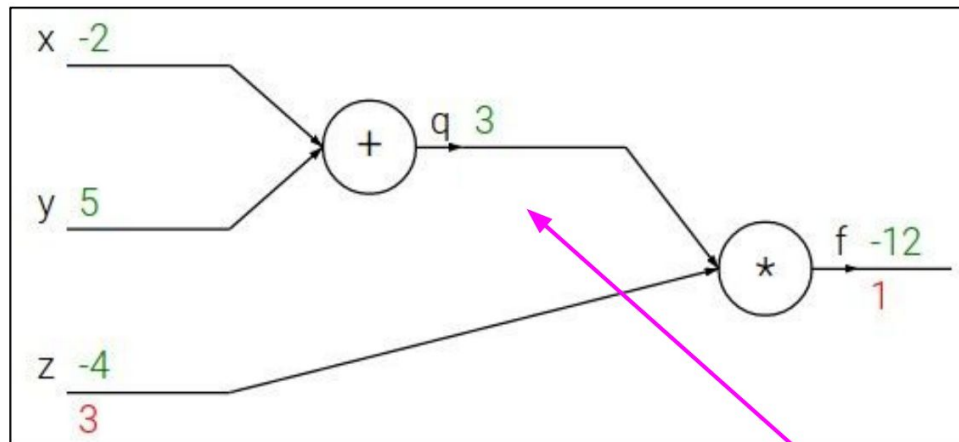
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

# Backpropagation Intro

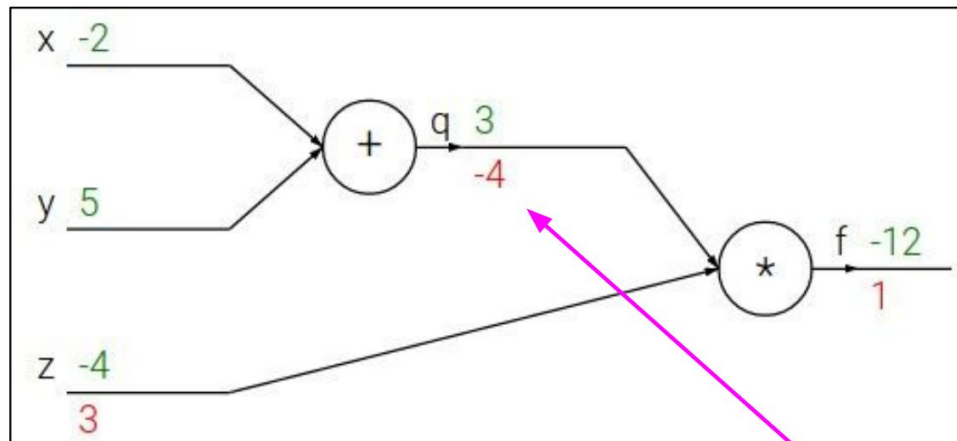
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$



# Backpropagation Intro

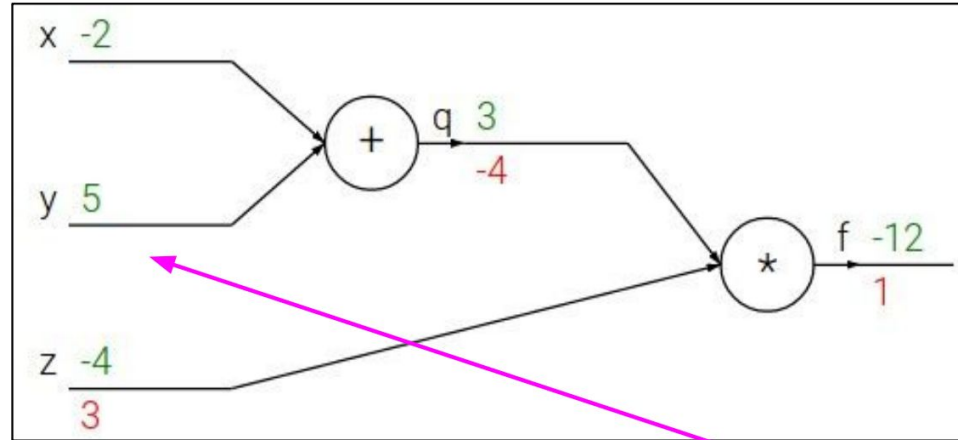
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

# Backpropagation Intro

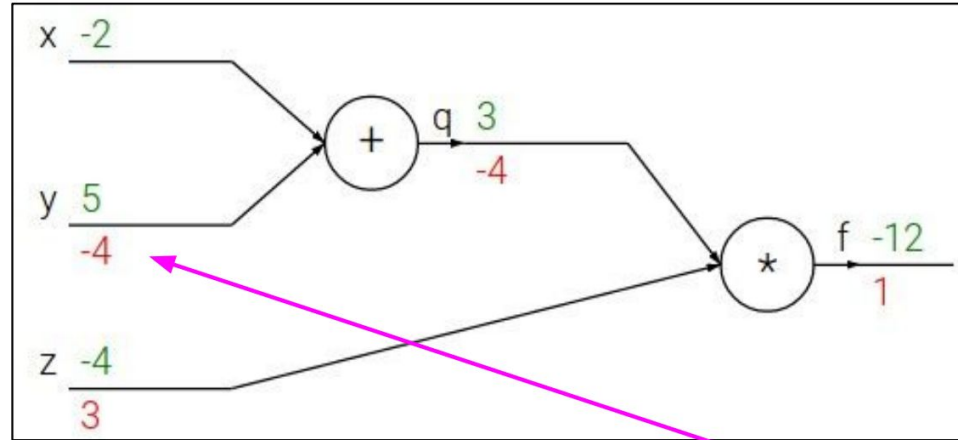
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

# Backpropagation Intro

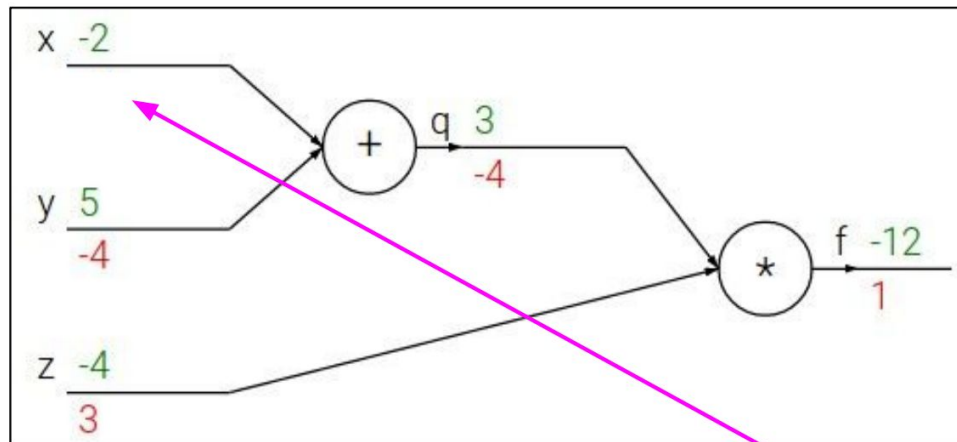
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

# Backpropagation Intro

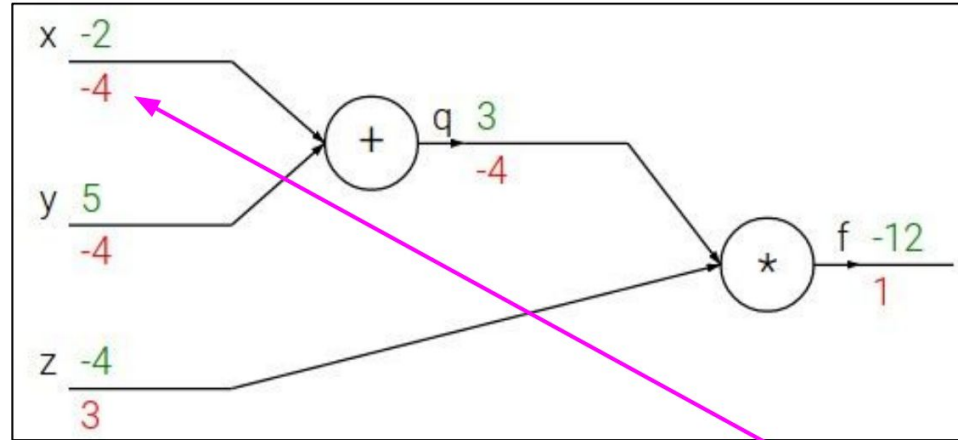
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

# Backpropagation Intro

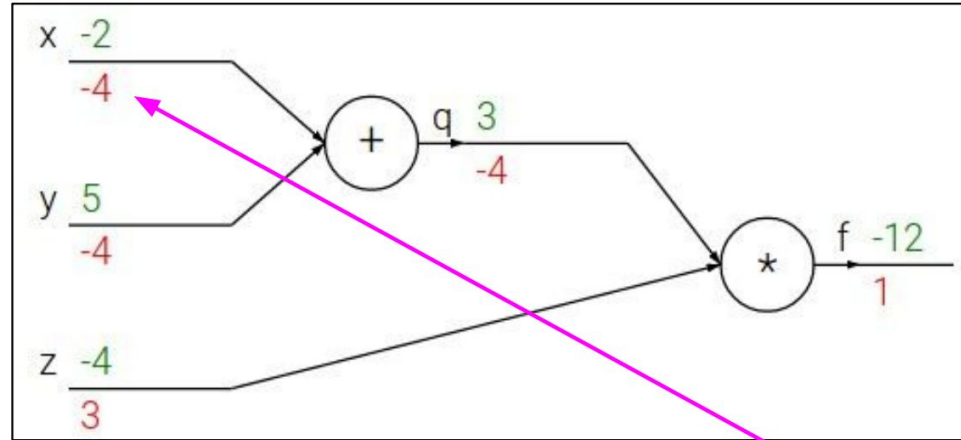
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

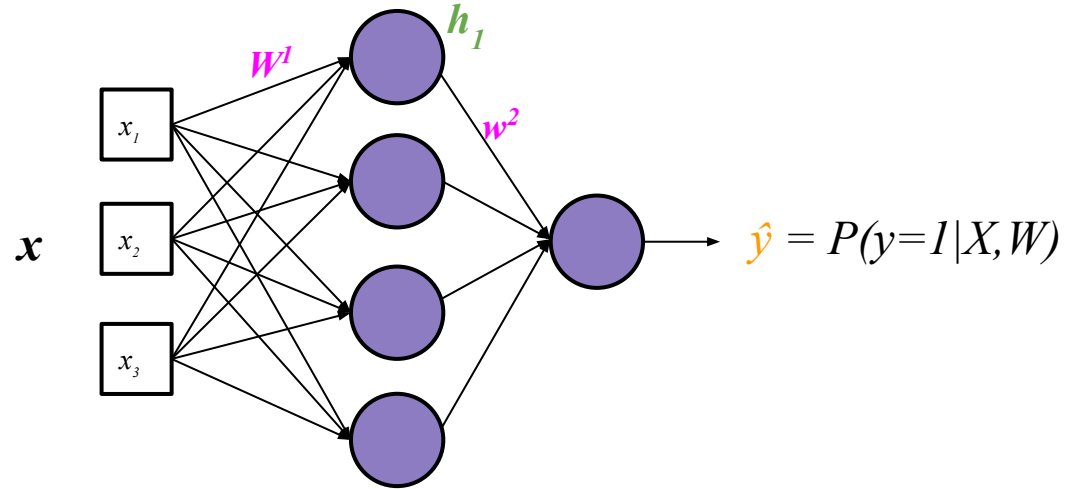


$$\frac{\partial f}{\partial x}$$

**Tells us:** by increasing  $x$  by a scale of 1, we decrease  $f$  by a scale of 4

# Backpropagation

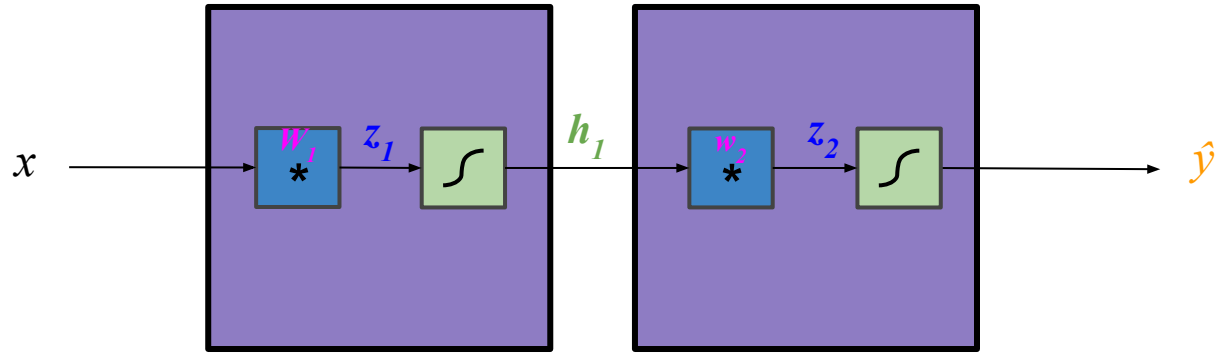
(binary classification example)



Example on 1-hidden layer network for binary classification

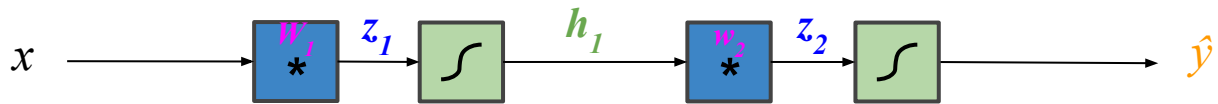
# Backpropagation

(binary classification example)



# Backpropagation

(binary classification example)



$$E = \text{loss} = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

**Gradient Descent  
to Minimize loss:**

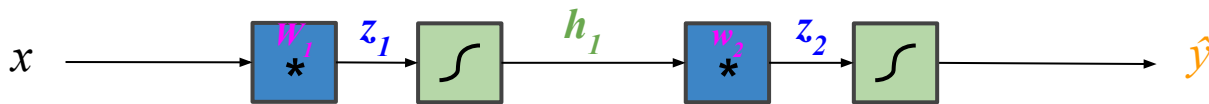
$$\mathbf{w}_2(t + 1) = \mathbf{w}_2(t) - \eta \frac{\partial E}{\partial \mathbf{w}_2(t)}$$
$$W_1(t + 1) = W_1(t) - \eta \frac{\partial E}{\partial W_1(t)}$$

Need to find these!



# Backpropagation

(binary classification example)



$$E = f_4 = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = f_3 = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = f_3 = \mathbf{w}_2^T \mathbf{h}_1$$

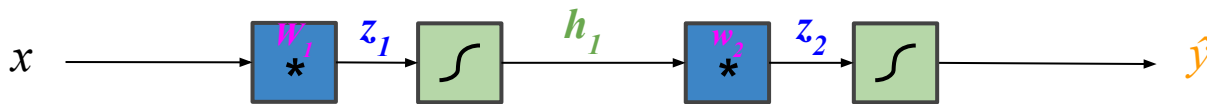
$$\mathbf{h}_1 = f_2 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = f_1 = \mathbf{W}_1^T \mathbf{x}$$

$$E = f_4(f_3(f_2(f_1(x))))$$

# Backpropagation

(binary classification example)



$$E = f_4 = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = f_3 = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = f_3 = \mathbf{w}_2^T \mathbf{h}_1$$

$$\mathbf{h}_1 = f_2 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = f_1 = \mathbf{W}_1^T \mathbf{x}$$

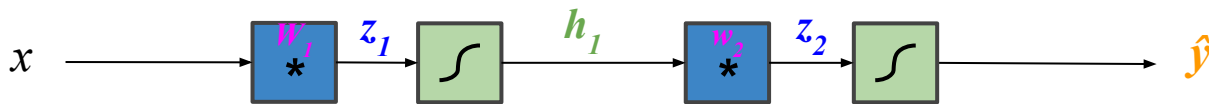
$$\frac{\partial E}{\partial \mathbf{w}_2} = ??$$

$$\frac{\partial E}{\partial \mathbf{W}_1} = ??$$

$$E = f_4(f_3(f_2(f_1(x))))$$

# Backpropagation

(binary classification example)



$$E = f_4 = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = f_3 = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = f_3 = \mathbf{w}_2^T \mathbf{h}_1$$

$$\mathbf{h}_1 = f_2 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = f_1 = \mathbf{W}_1^T \mathbf{x}$$

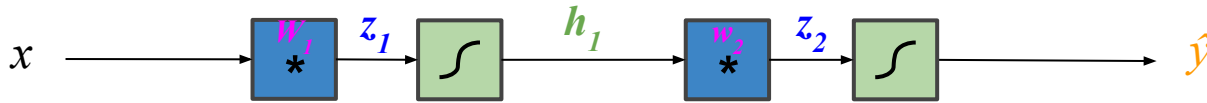
$$\frac{\partial E}{\partial \mathbf{w}_2} = ??$$

$$\frac{\partial E}{\partial \mathbf{W}_1} = ??$$

$E = f_4(f_3(f_2(f_1(x))))$   Exploit the chain rule!

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

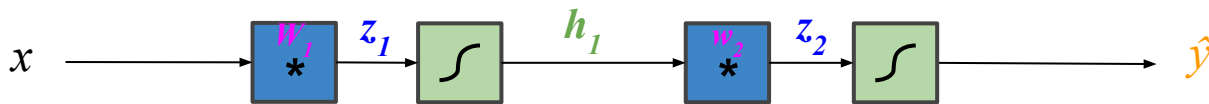
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$\mathbf{z}_1 = W_1^T \mathbf{x}$$

$$\frac{\partial E}{\partial w_2} =$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

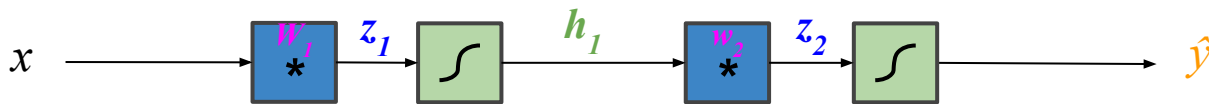
$$\mathbf{z}_1 = \mathbf{W}_1^T \mathbf{x}$$

chain rule

$$\frac{\partial E}{\partial \mathbf{w}_2} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{w}_2}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

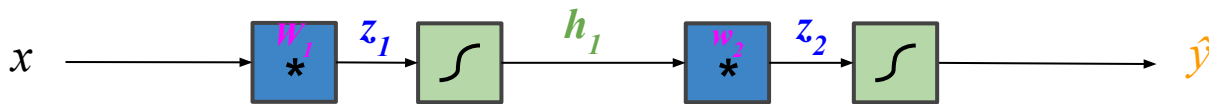
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = \mathbf{W}_1^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{w}_2} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{w}_2} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot \end{aligned}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

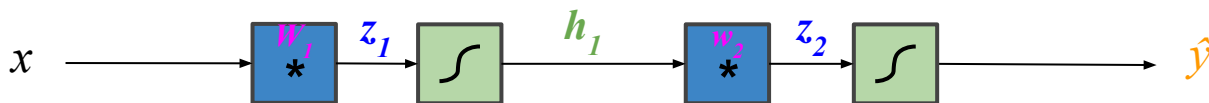
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$\mathbf{z}_1 = W_1^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{w}_2} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{w}_2} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot \left( \frac{e^{z_2}}{1 + e^{z_2}} \left( 1 - \frac{e^{z_2}}{1 + e^{z_2}} \right) \right) \cdot \end{aligned}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

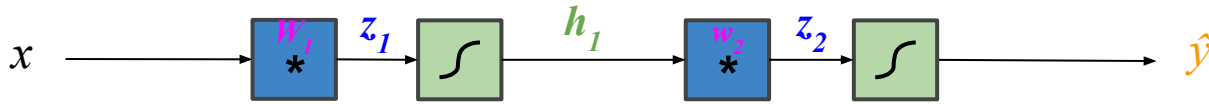
$$\mathbf{z}_1 = W_1^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{w}_2} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{w}_2} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot \left( \frac{e^{z_2}}{1 + e^{z_2}} \left( 1 - \frac{e^{z_2}}{1 + e^{z_2}} \right) \right) \cdot (\mathbf{h}_1) \end{aligned}$$



# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

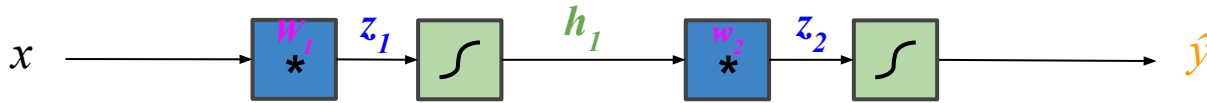
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$\mathbf{z}_1 = W_1^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{w}_2} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{w}_2} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot \left( \frac{e^{z_2}}{1 + e^{z_2}} \left( 1 - \frac{e^{z_2}}{1 + e^{z_2}} \right) \right) \cdot (\mathbf{h}_1) \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot (\hat{y}(1 - \hat{y})) \cdot (\mathbf{h}_1) \end{aligned}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

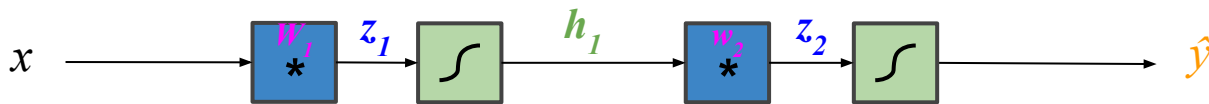
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = \mathbf{W}_1^T \mathbf{x}$$

$$\frac{\partial E}{\partial \mathbf{W}_1} =$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

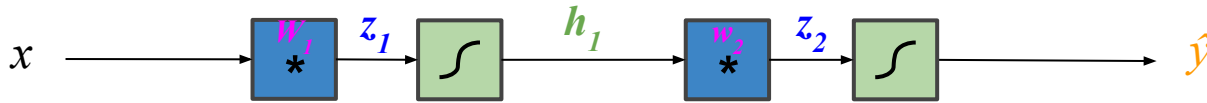
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$\mathbf{z}_1 = \mathbf{W}_1^T \mathbf{x}$$

$$\frac{\partial E}{\partial \mathbf{W}_1} = \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{h}_1} \cdot \frac{\partial \mathbf{h}_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial \mathbf{W}_1}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

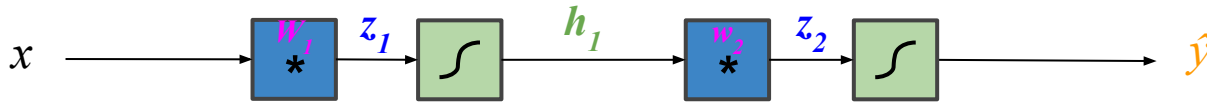
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$z_1 = \mathbf{W}_1^T \mathbf{x}$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{W}_1} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{h}_1} \cdot \frac{\partial \mathbf{h}_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial \mathbf{W}_1} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot (\hat{y}(1 - \hat{y})) \cdot (\mathbf{w}) \cdot (\mathbf{h}_1(1 - \mathbf{h}_1)) \cdot (\mathbf{x}) \end{aligned}$$

# Backpropagation

(binary classification example)



$$E = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\hat{y} = \frac{e^{z_2}}{1 + e^{z_2}}$$

$$z_2 = \mathbf{w}_2^T \mathbf{h}_1$$

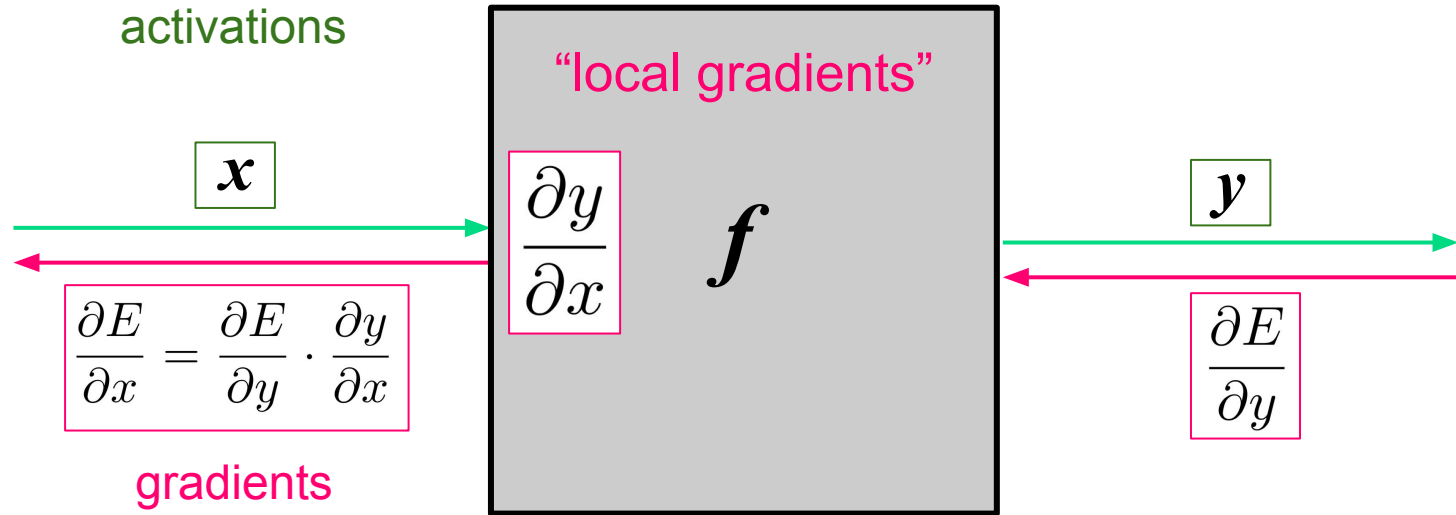
$$\mathbf{h}_1 = \frac{e^{z_1}}{1 + e^{z_1}}$$

$$\mathbf{z}_1 = W_1^T \mathbf{x}$$

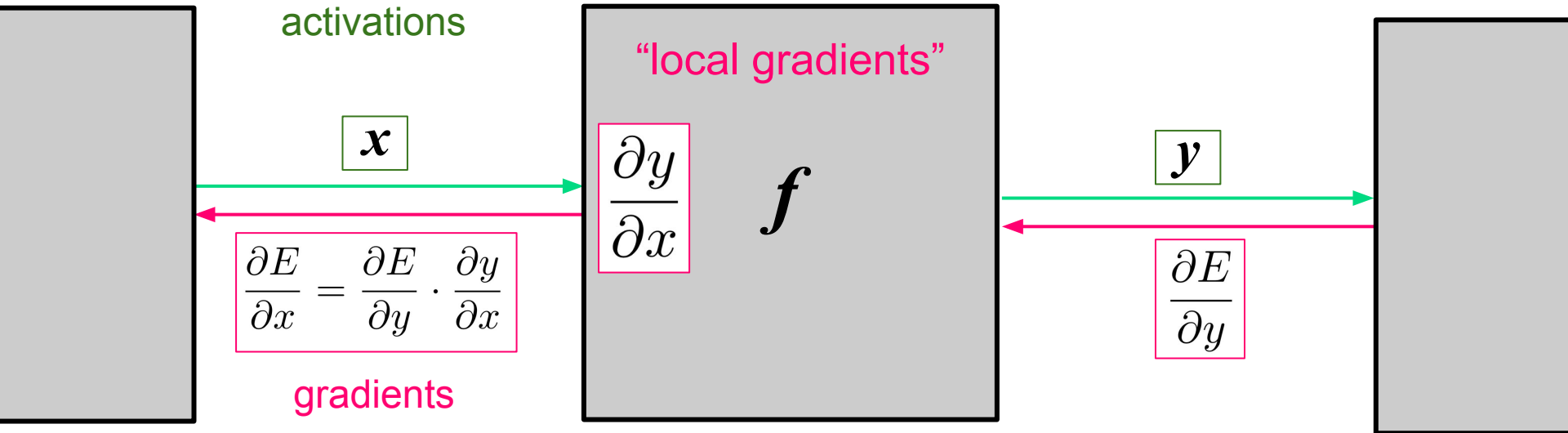
already computed!

$$\begin{aligned} \frac{\partial E}{\partial W_1} &= \frac{\partial E}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_2} \cdot \frac{\partial z_2}{\partial \mathbf{h}_1} \cdot \frac{\partial \mathbf{h}_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial W_1} \\ &= \left( \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \right) \cdot (\hat{y}(1 - \hat{y})) \cdot (\mathbf{w}) \cdot (\mathbf{h}_1(1 - \mathbf{h}_1)) \cdot (\mathbf{x}) \end{aligned}$$

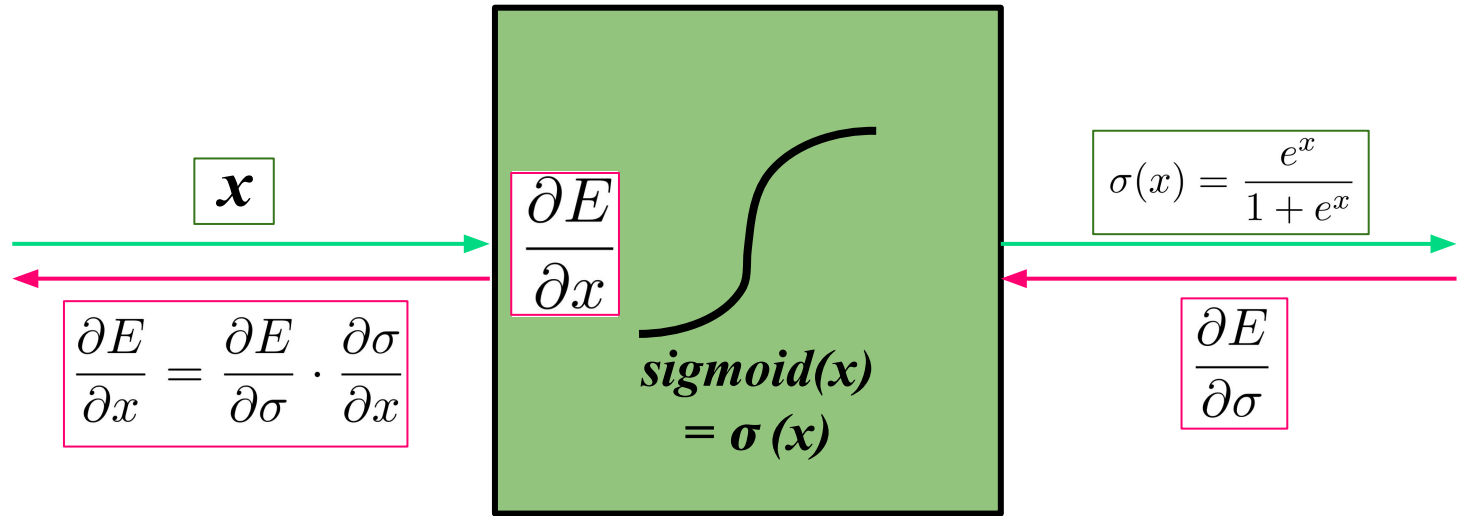
# “Local-ness” of Backpropagation



# “Local-ness” of Backpropagation



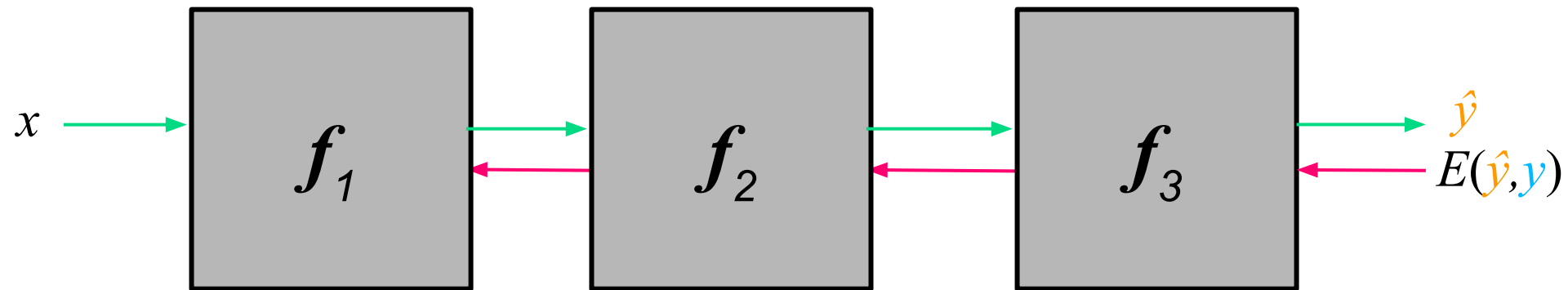
# Example: Sigmoid Block



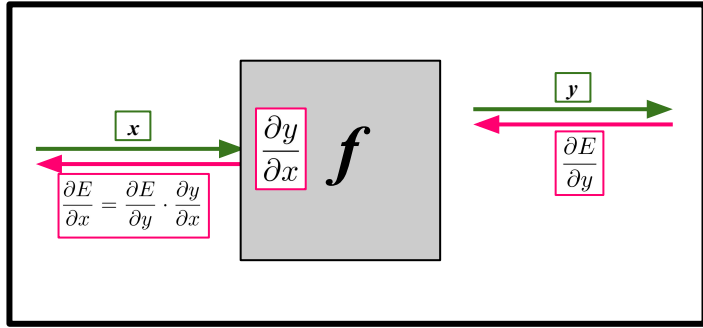


# Deep Learning:

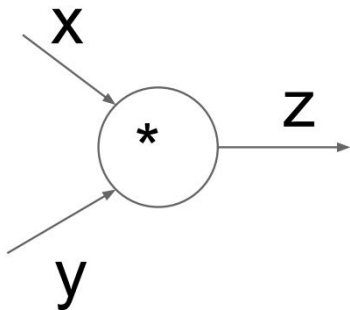
Layers of Differentiable Parameterized Functions (with nonlinearities)



# Building Deep Neural Nets



# Block Example Implementation



(x,y,z are scalars)

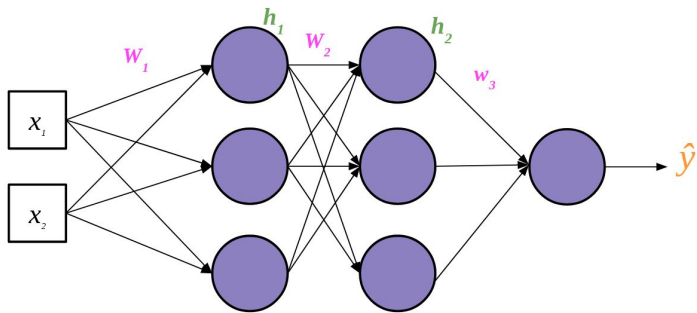
```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        self.x = x # must keep these around!  
        self.y = y  
        return z  
    def backward(dz):  
        dx = self.y * dz # [dz/dx * dL/dz]  
        dy = self.x * dz # [dz/dy * dL/dz]  
        return [dx, dy]
```



# Deep Learning Frameworks



# Pytorch Sample Code



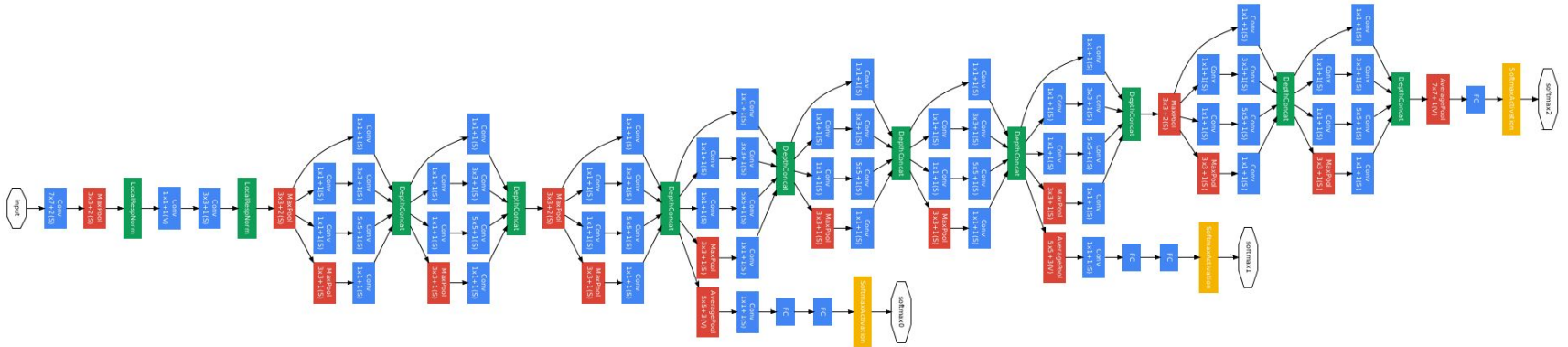
```
import torch.nn as nn
import torch.nn.functional as F

class ThreeLayerNet(torch.nn.Module):
    def __init__(self, d_in, d_hidden, d_out):
        super().__init__()
        self.W1 = nn.Linear(d_in, d_hidden)
        self.W2 = nn.Linear(d_hidden, d_hidden)
        self.w3 = nn.Linear(d_hidden, d_out)
        self.nonlinear = nn.Sigmoid()

    def forward(self, x):
        h1 = self.nonlinear(self.W1(x))
        h2 = self.nonlinear(self.W2(h1))
        y_hat = self.nonlinear(self.w3(h2))
        return y_hat

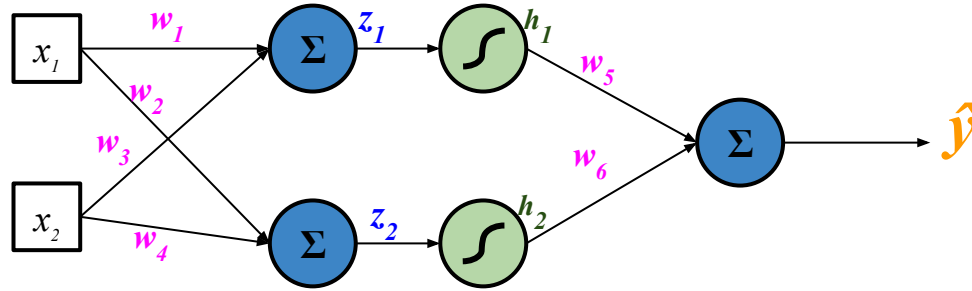
model = ThreeLayerNet(2, 3, 1)
```

# Building Deep Neural Nets



“GoogLeNet” for Object Classification

# Backprop Demo



$f_4$	$E = (y - \hat{y})^2$
$f_3$	$\hat{y} = h_1 w_5 + h_2 w_6$
$f_2$	$h_1 = \frac{\exp(z_1)}{1 + \exp(z_1)}$ $h_2 = \frac{\exp(z_2)}{1 + \exp(z_2)}$
$f_1$	$z_1 = x_1 w_1 + x_2 w_3$ $z_2 = x_1 w_2 + x_2 w_4$

$$w(t+1) = w(t) - \eta \frac{\partial E}{\partial w(t)}$$

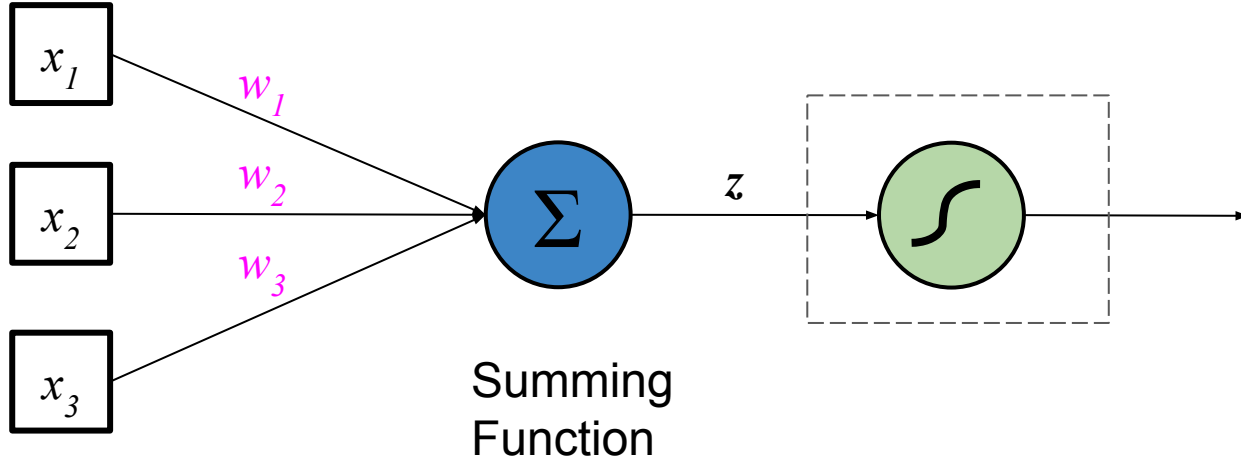
$$\frac{\partial E}{\partial w_i} = ??$$

# Nonlinearity Functions



# Nonlinearity Functions

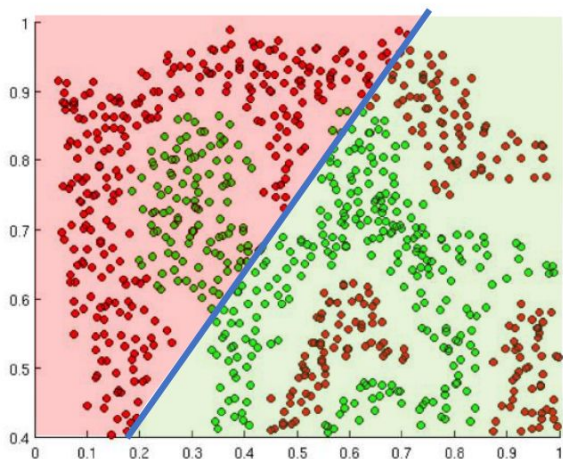
(i.e. transfer or activation functions)



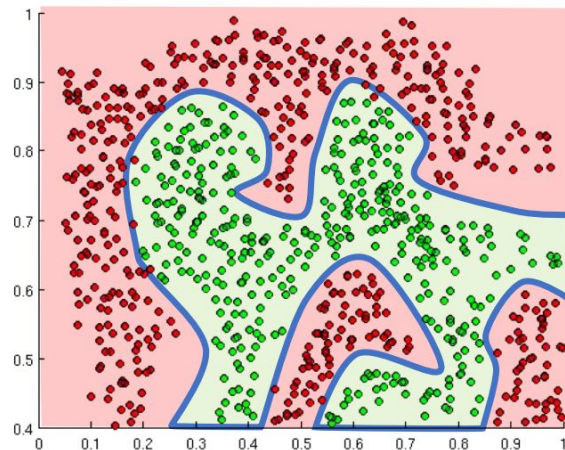
# Nonlinearity Functions

(i.e. transfer or activation functions)

The purpose of activation functions is to **introduce non-linearities** into the network







Linear Activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions





# Nonlinearity Functions

(i.e. transfer or activation functions)

Name	Plot	Equation	Derivative ( w.r.t $x$ )
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectifier (ReLU) <sup>[9]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$





# Nonlinearity Functions

(i.e. transfer or activation functions)

Name	Plot	Equation	Derivative ( w.r.t $x$ )
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectifier (ReLU) <sup>[9]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

# Nonlinearity Functions

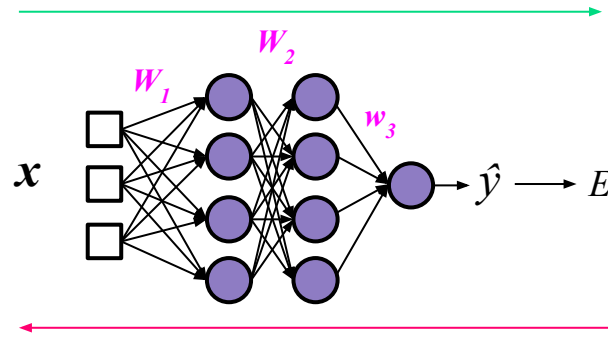
(i.e. transfer or activation functions)

Name	Plot	Equation	Derivative ( w.r.t $x$ )
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectifier (ReLU) <sup>[9]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

usually works best in practice

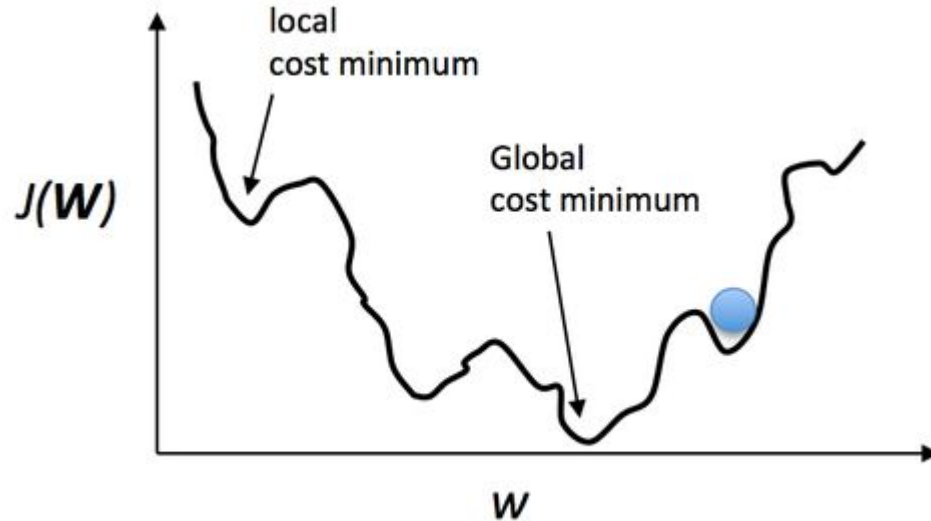
# Neural Nets in Practice

# Neural Net Pipeline



1. Initialize weights
2. For each batch of input  $x$  samples  $S$ :
  - a. Run the network “**Forward**” on  $S$  to compute outputs and loss
  - b. Run the network “**Backward**” using outputs and loss to compute gradients
  - c. Update weights using SGD (or a similar method)
3. Repeat step 2 until loss convergence

# Non-Convexity of Neural Nets

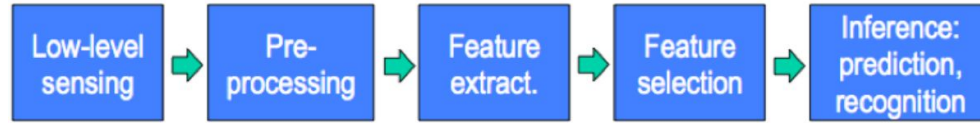


In very high dimensions, there exists many local minimum which are about the same.

Pascanu, et. al. *On the saddle point problem for non-convex optimization* 2014



# Advantage of Neural Nets



## Feature Engineering

- ✓ Most critical for accuracy
- ✓ Account for **most of the computation** for testing
- ✓ Most time-consuming in development cycle
- ✓ Often **hand-craft** and **task dependent** in practice

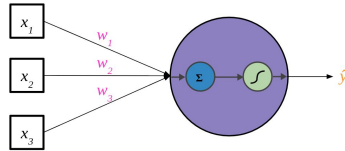
## Feature Learning

- ✓ Easily **adaptable to new** similar tasks
- ✓ Layerwise representation
- ✓ Layer-by-layer unsupervised training
- ✓ Layer-by-layer supervised training

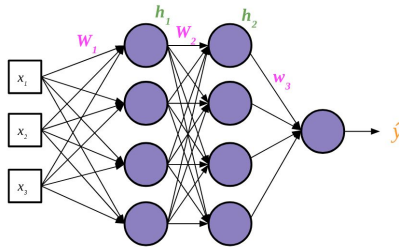
As long as it's fully differentiable, we can train the model to automatically learn features for us.

# Summary

Neurons



Neural Networks



Loss Functions and Backpropagation

