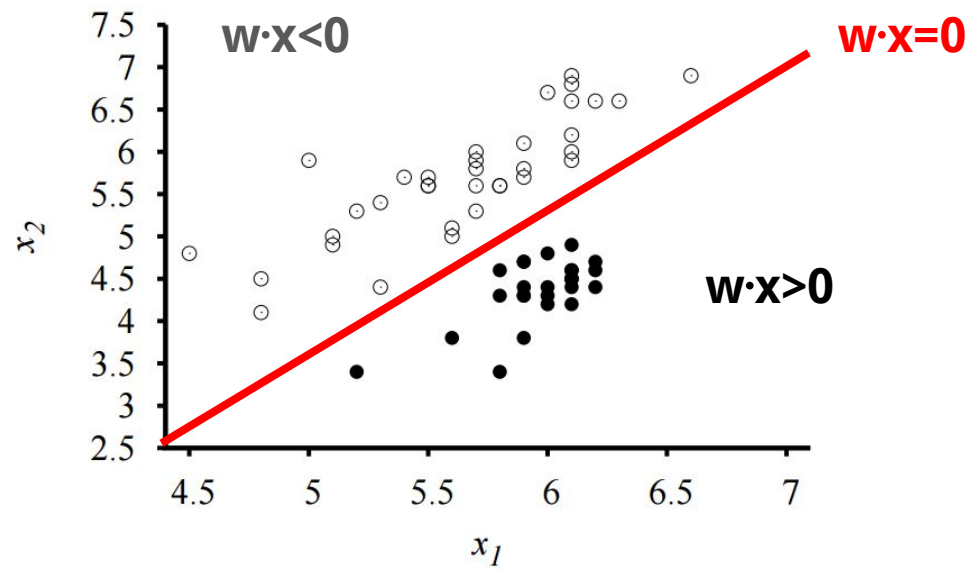


# 4700: Neural Networks (continued)

[Some slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

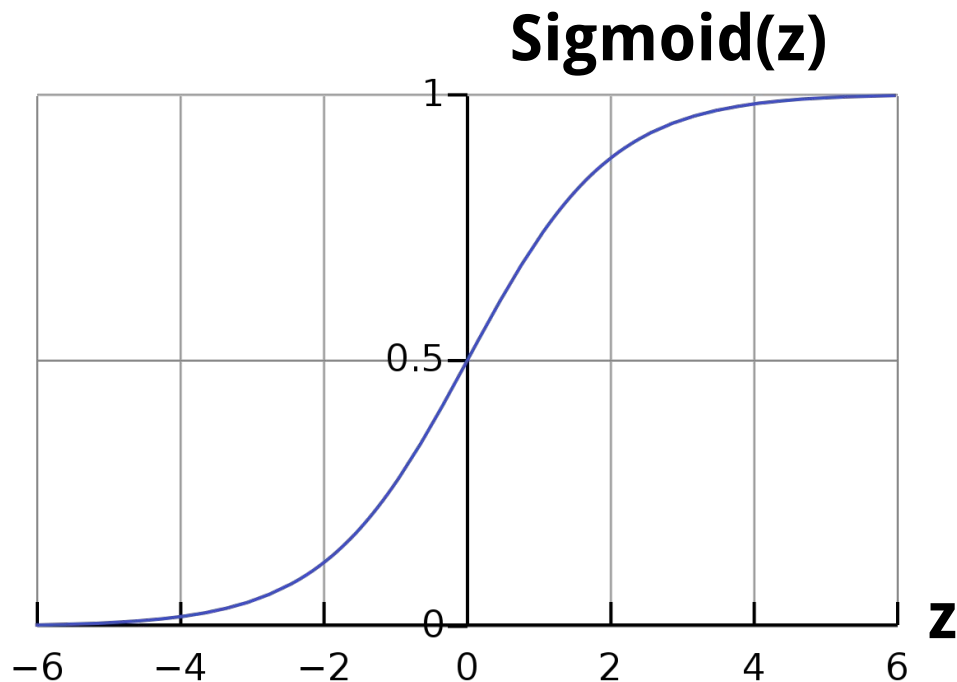


**Sigmoid(0)=0.5**

**As  $z \rightarrow +\infty$ , Sigmoid( $z$ )->1**

**As  $z \rightarrow -\infty$ , Sigmoid( $z$ )->0**

$$\text{Sigmoid}(z) = \frac{1}{1 + e^{-z}}$$



# Logistic Regression

$$P(y=+1 \mid x, w) = \text{Sigmoid}(w \cdot x)$$

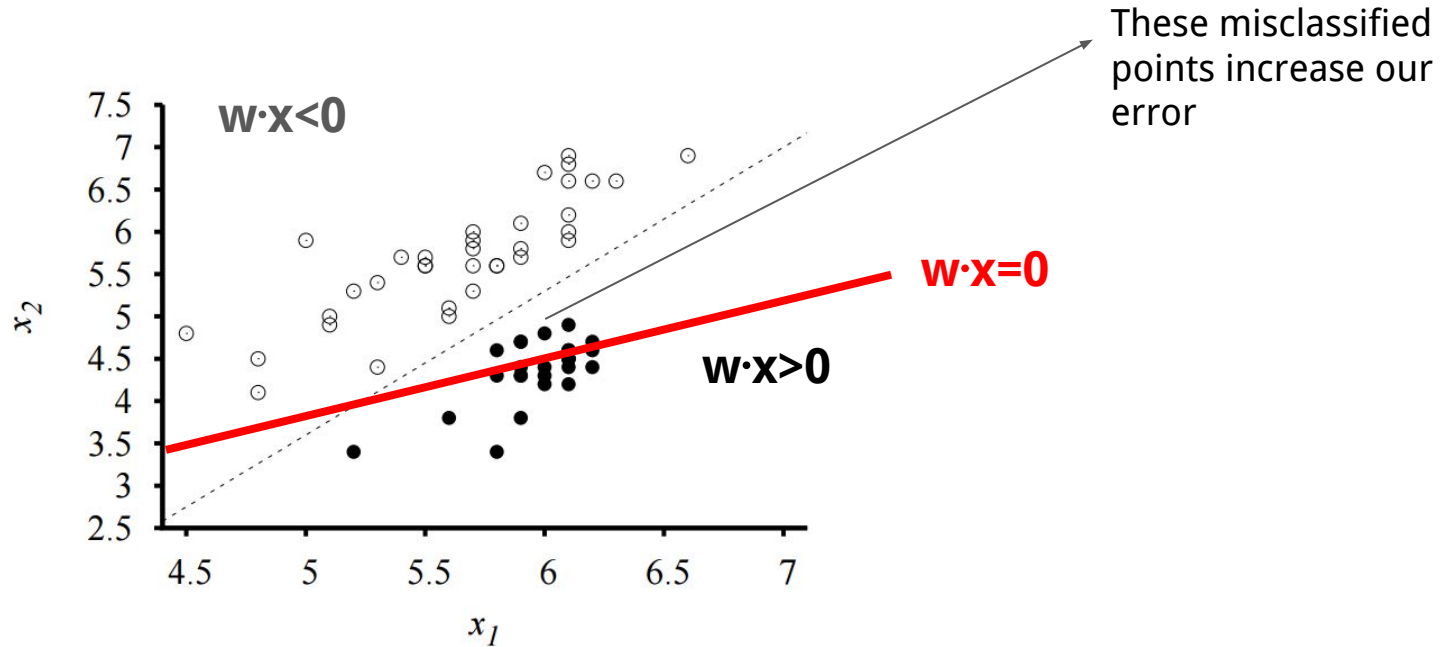
Best  $w$ ?

- 
- Maximum likelihood estimation:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

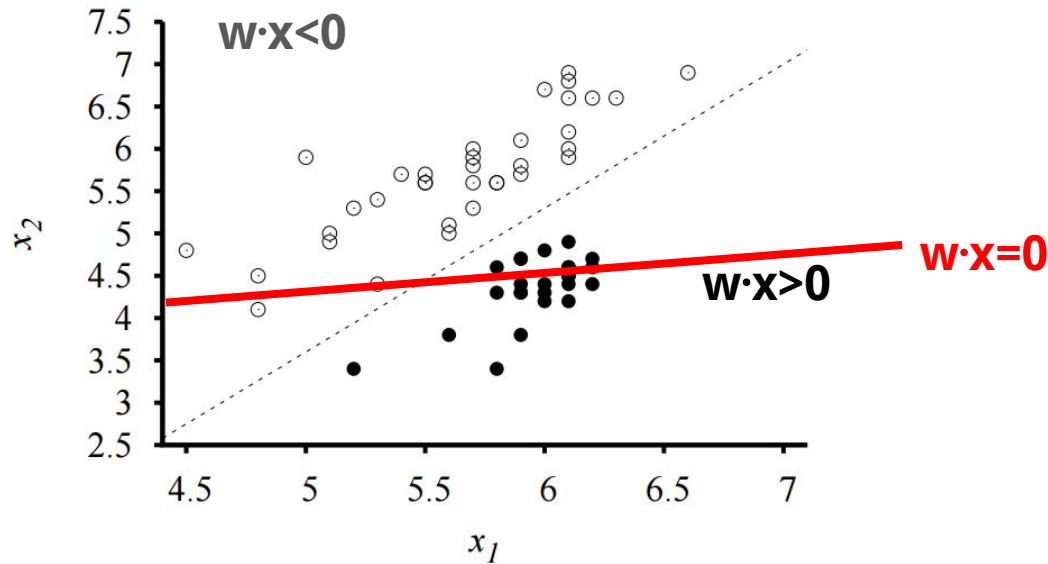


# Classification - What error are we minimizing?



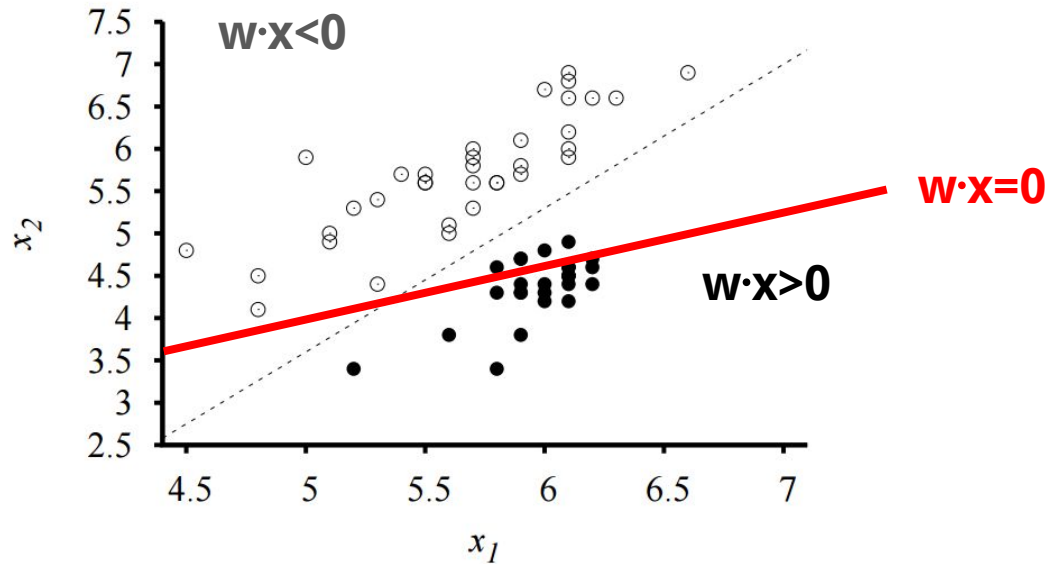
# Classification - What error are we minimizing?

Yikes, even worse!



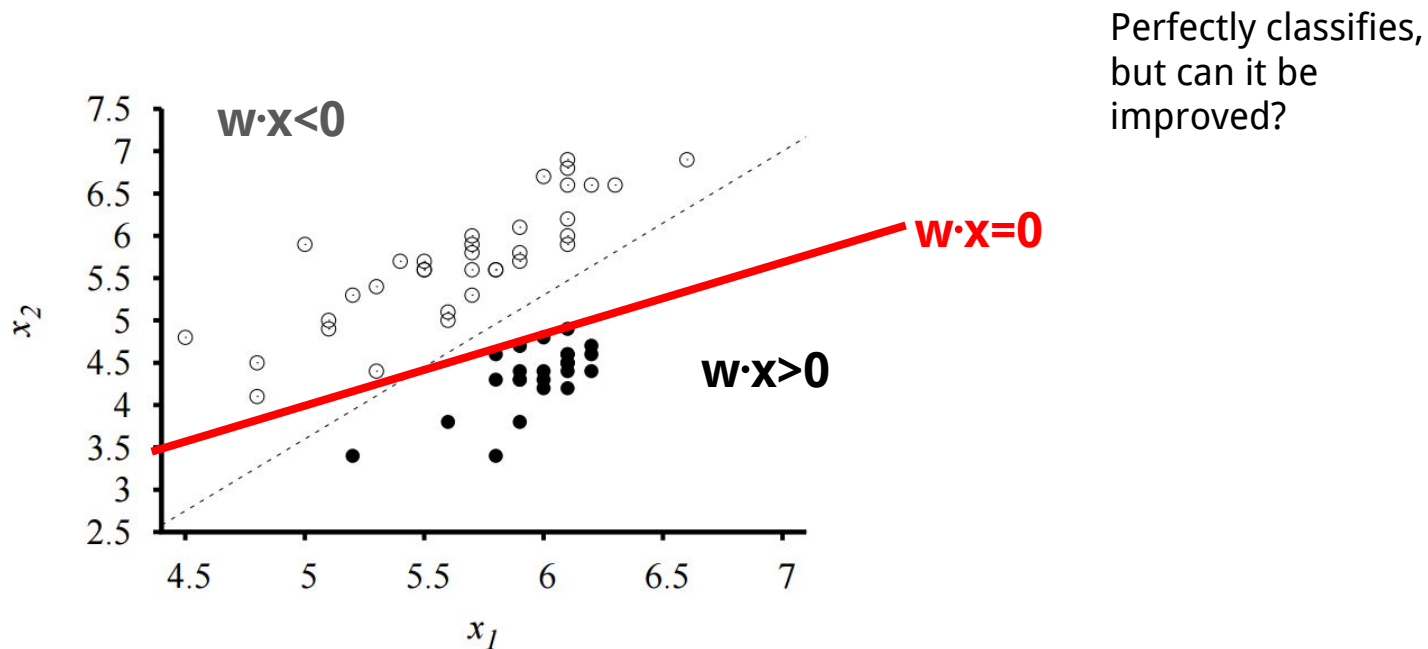
# Classification - What error are we minimizing?

A bit better



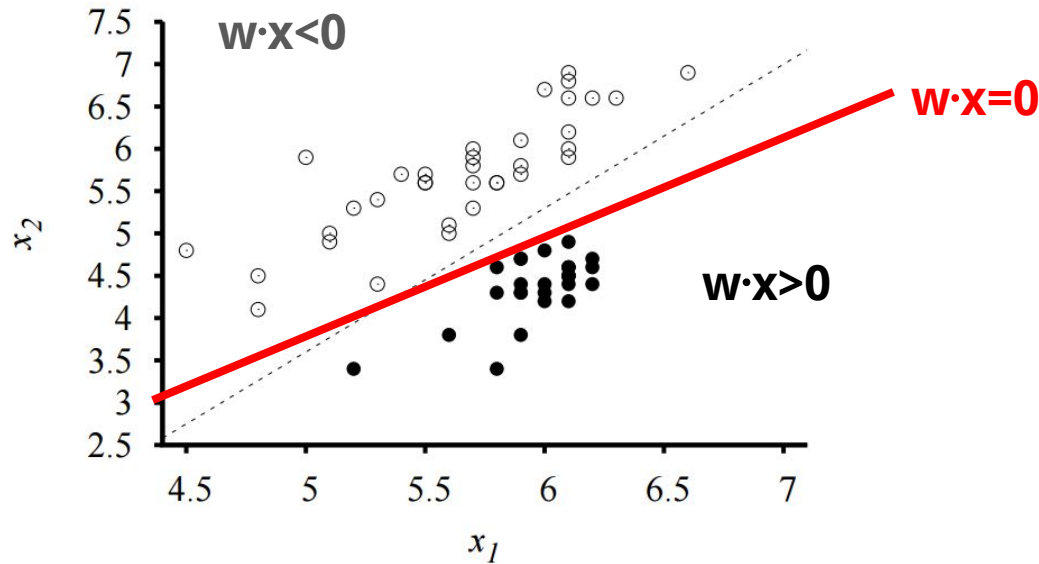


# Classification - What error are we minimizing?



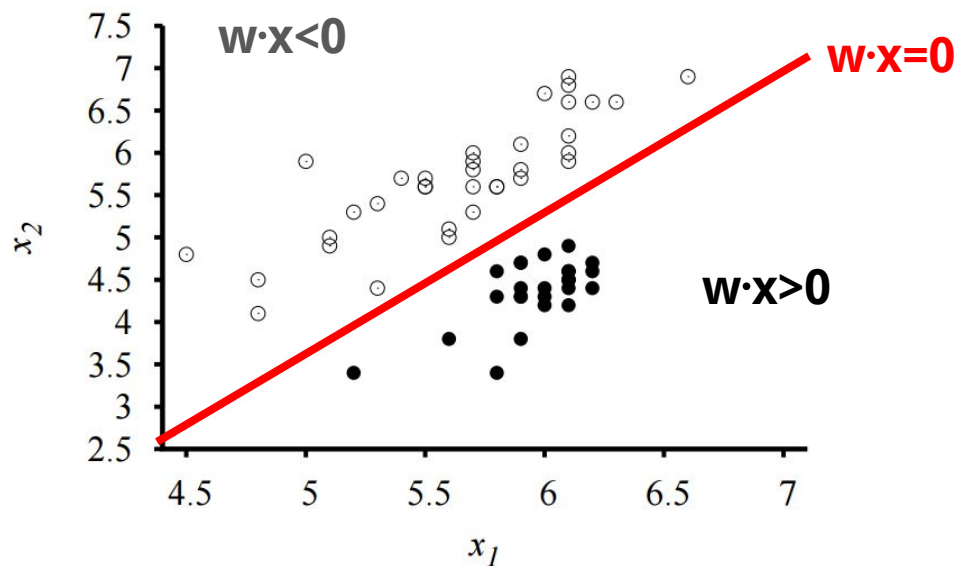
# Classification - What error are we minimizing?

Getting better...



## Classification - What error are we minimizing?

Perfect!



# Local search for weight vector

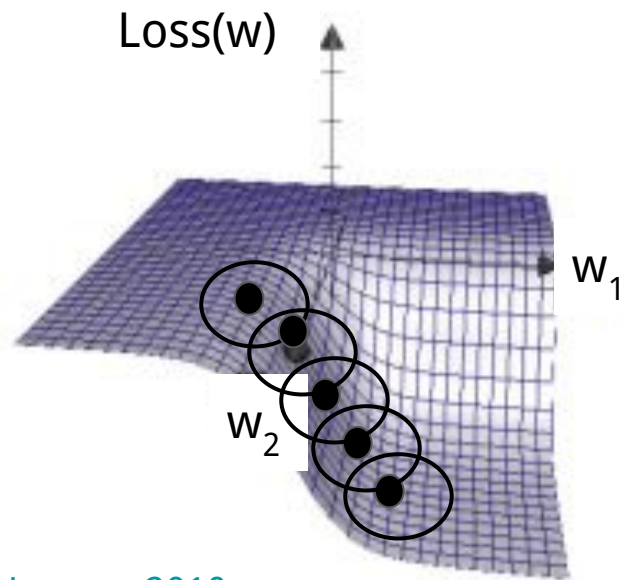


Figure from [Chaudhuri & Solar-Lezama 2010](#)

# Gradient Ascent

- Perform update in uphill direction for each coordinate
- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate
- E.g., consider:  $g(w_1, w_2)$

▪ Updates:

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

▪ Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$

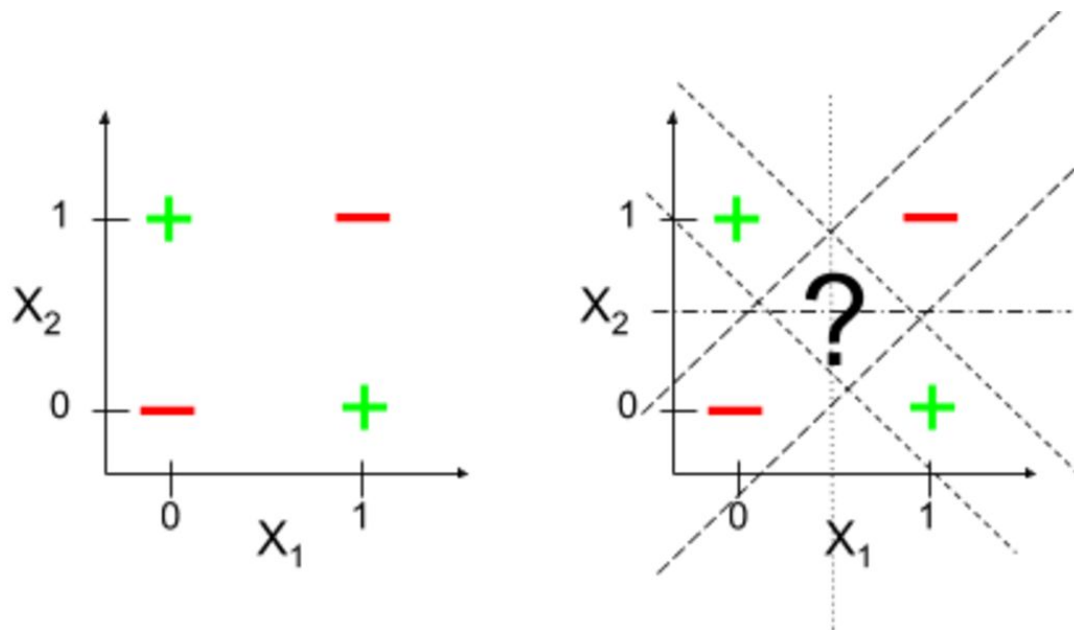
with:  $\nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix}$  = gradient

# What can logistic regression learn?

Lines\*

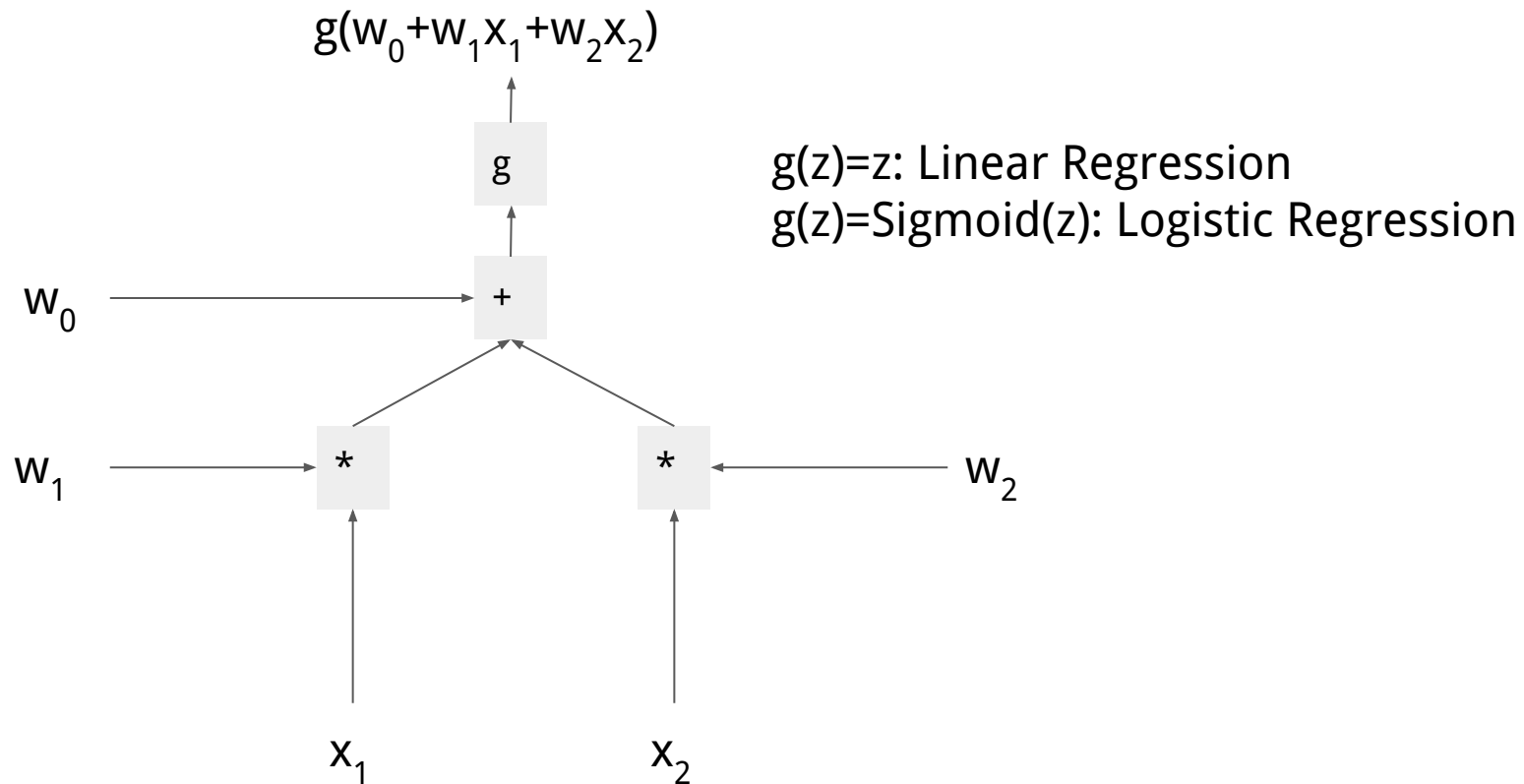
\*in high dimensional spaces

## Richer functions – Can a linear classifier fit this data?



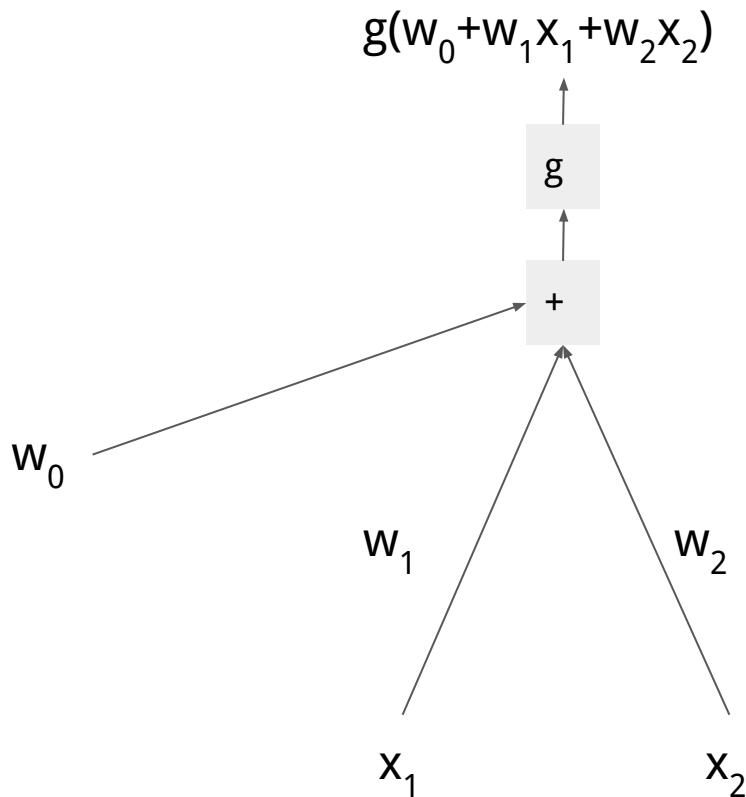
(<https://medium.com/@claude.coulombe/the-revenge-of-perceptron-learning-xor-with-tensorflow-eb52cbdf6c60>)

# A different view of linear models





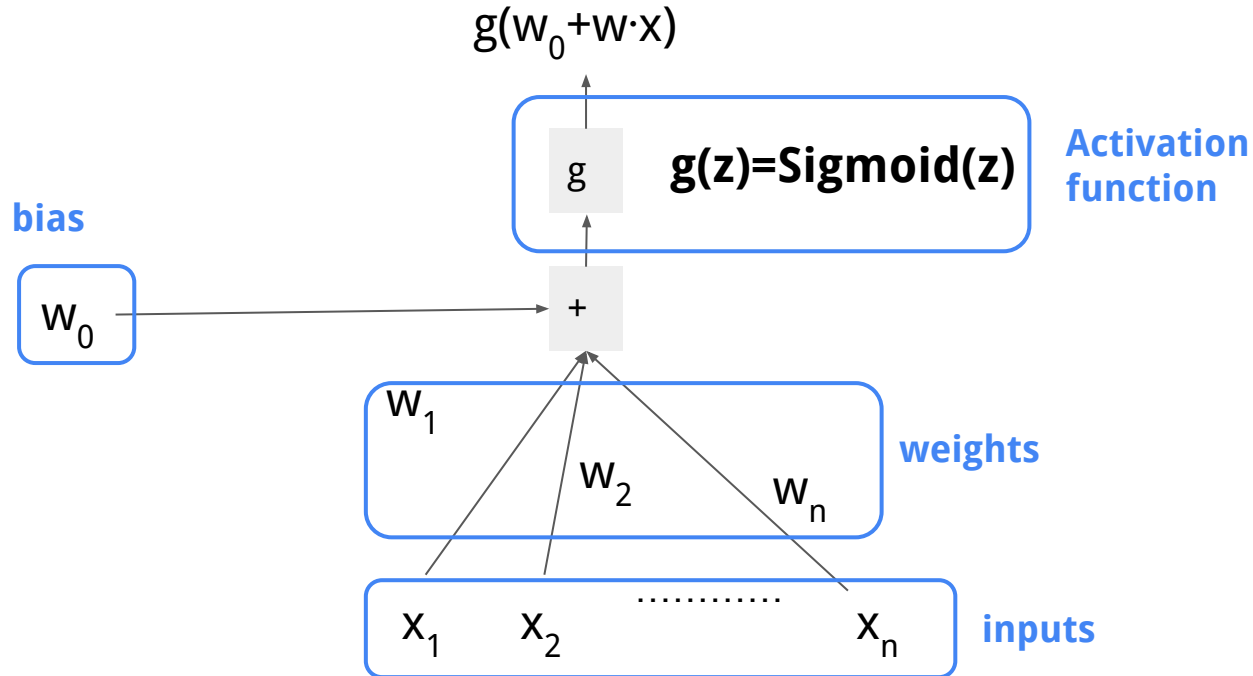
# A different view of linear models



$g(z)=z$ : Linear Regression

$g(z)=\text{Sigmoid}(z)$ : Logistic Regression

# "Neuron"



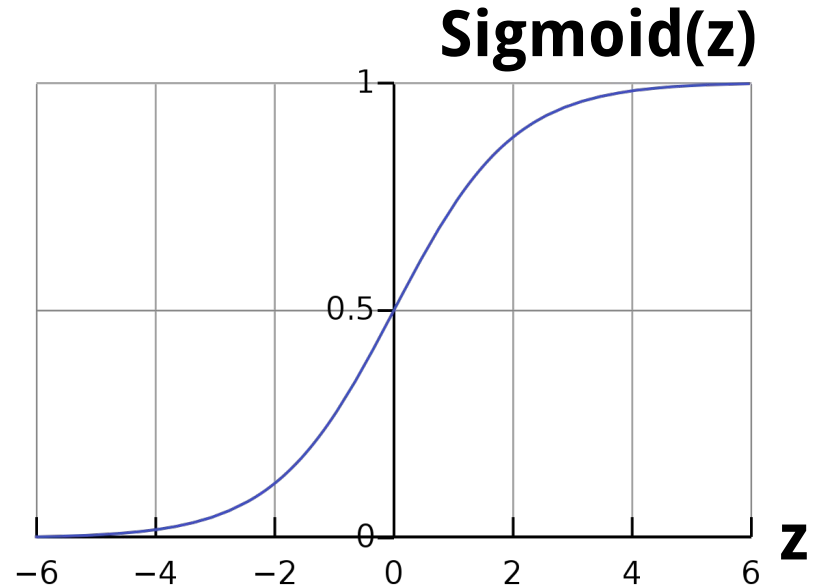
# Nonlinear Logistic Activation +Linear Combination

$z$ =weighted input to neuron

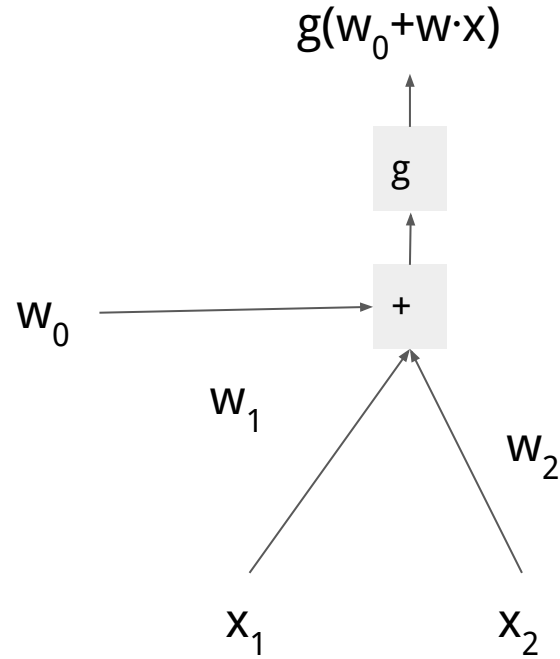
$$z = w_0 + w \cdot x$$

Bias: threshold  
for neuron  
"firing"

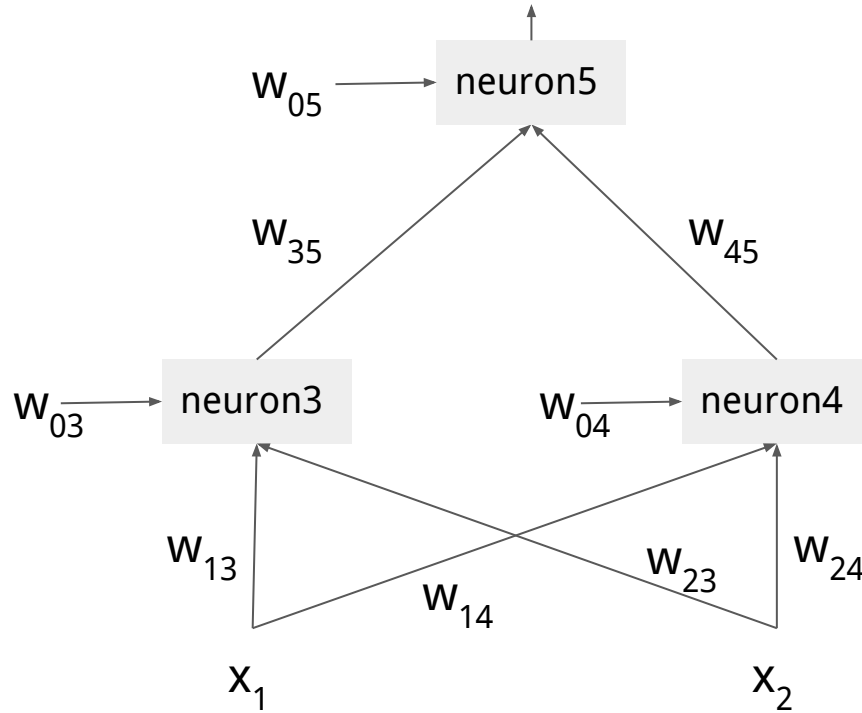
Weights:  
positive/negative?  
excitatory/inhibitory?



# More than one neuron

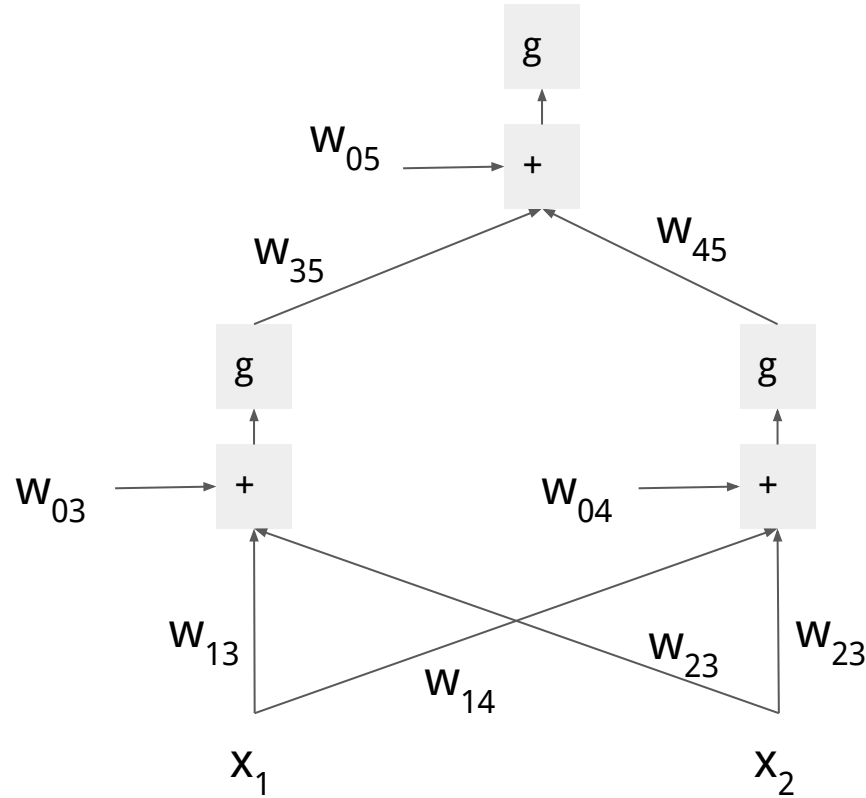


# More than one neuron



Bias for neuron  $i$ :  $w_{0i}$   
Weight from  $i \rightarrow j$ :  $w_{ij}$

# More than one neuron



Bias for neuron  $i$ :  $w_{0i}$   
Weight from  $i \rightarrow j$ :  $w_{ij}$

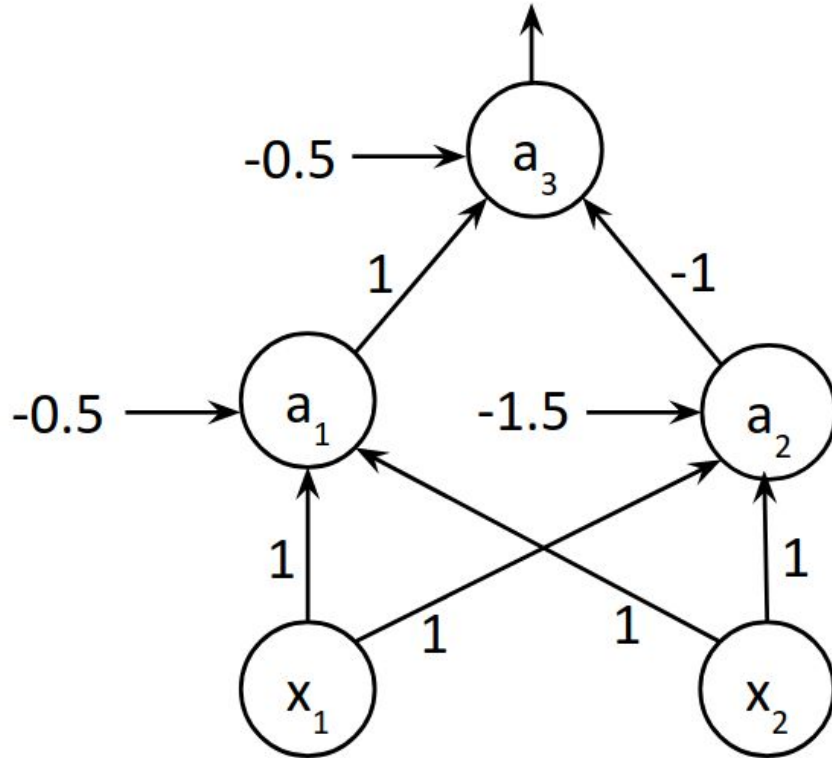
Does it solve xor?

$$\text{xor}(0,0)=0$$

$$\text{xor}(1,0)=1$$

$$\text{xor}(0,1)=1$$

$$\text{xor}(1,1)=0$$



$$a_1 = \text{Logistic}(-0.5 + 1 \cdot x_1 + 1 \cdot x_2)$$

$$\sim \text{is } x_1 + x_2 - 0.5 > 0$$

$$\sim \text{OR } x_1 \text{ } x_2$$

$$a_2 = \text{Logistic}(-1.5 + 1 \cdot x_1 + 1 \cdot x_2)$$

$$\sim \text{is } x_1 + x_2 - 1.5 > 0$$

$$\sim \text{AND } x_1 \text{ } x_2$$

$$a_3 = \text{Logistic}(-0.5 + 1 \cdot a_1 - 1 \cdot a_2)$$

$$\sim \text{is } a_1 - a_2 - 0.5 > 0$$

$$\sim \text{AND } a_1 \text{ (NOT } a_2)$$



Does it solve the problem of needing feature extractors?

# feature vector recap

Actual thing we want  
to classify



**FEATURE EXTRACTOR**

“feature vector”

$$\mathbf{x}_n \in \mathbb{R}^D$$

$$\mathbf{x}_n = (x_{n1}, x_{n2}, x_{n3}, \dots, x_{nD})$$

*f*

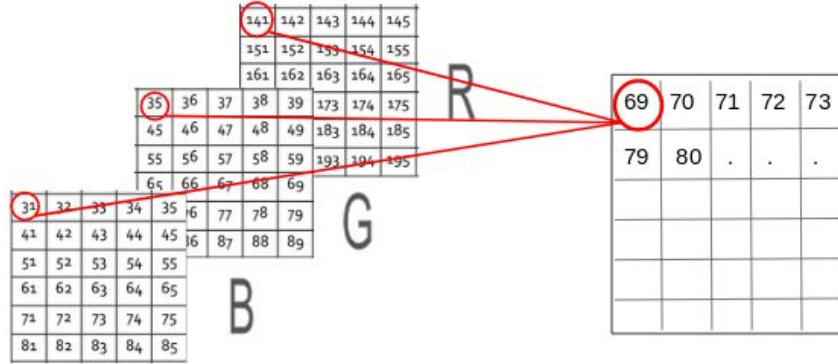
Is it dog or not

*y*

# Example Image Feature Extractor



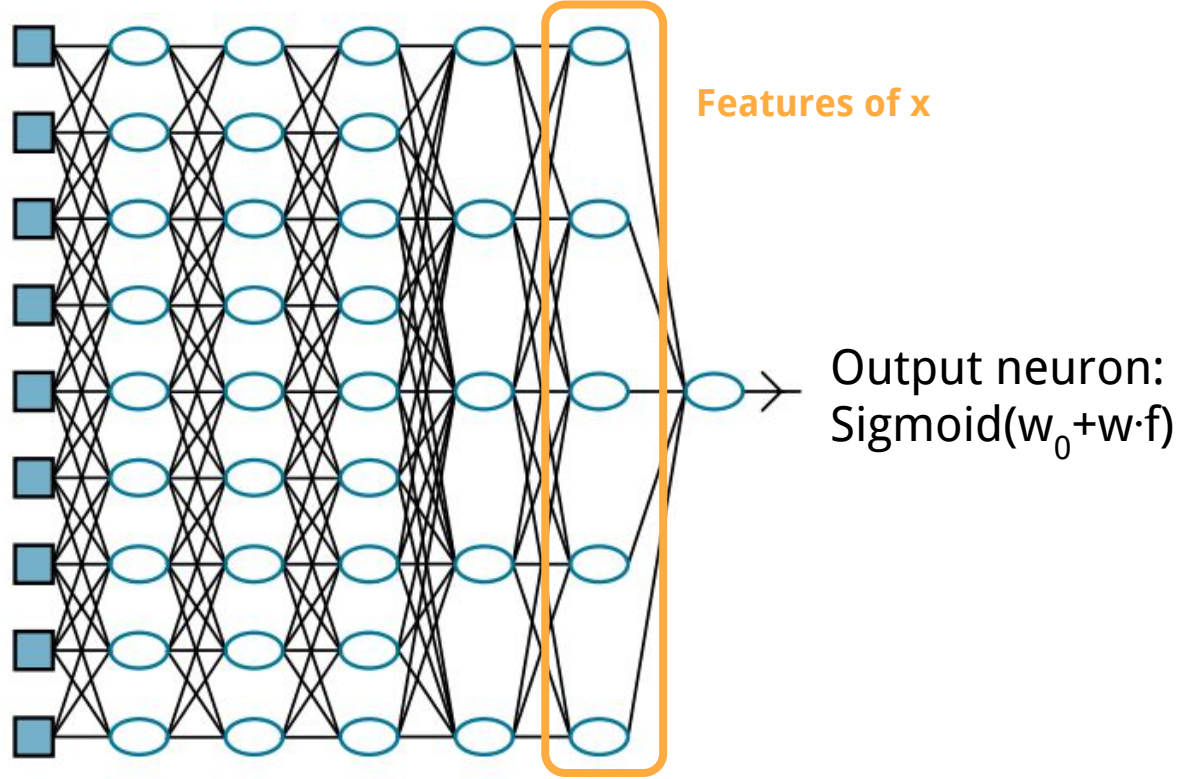
Colour Image



Flatten into  
1d array



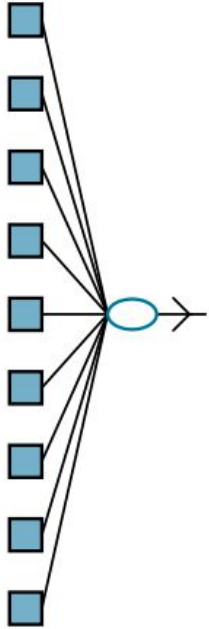
# Neural Net Implicitly Does a Feature Extractor



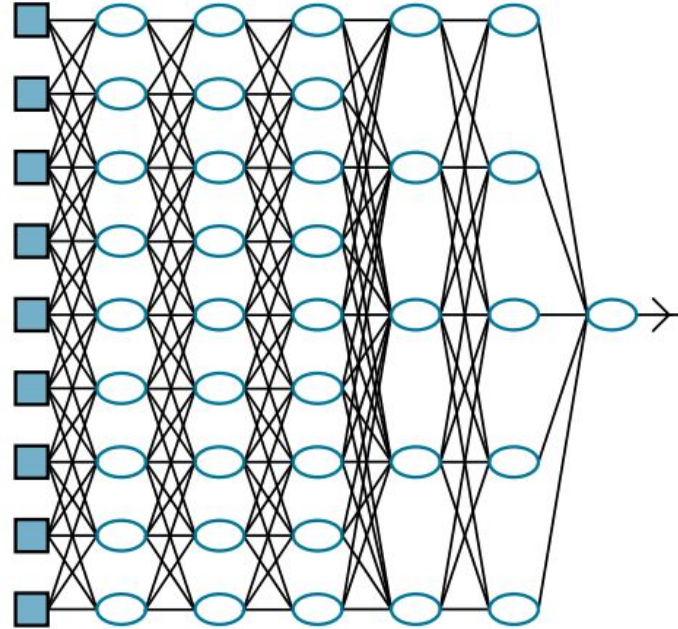
Inputs:  $x$

Second-to-last neurons: features

# Neural Net Structure

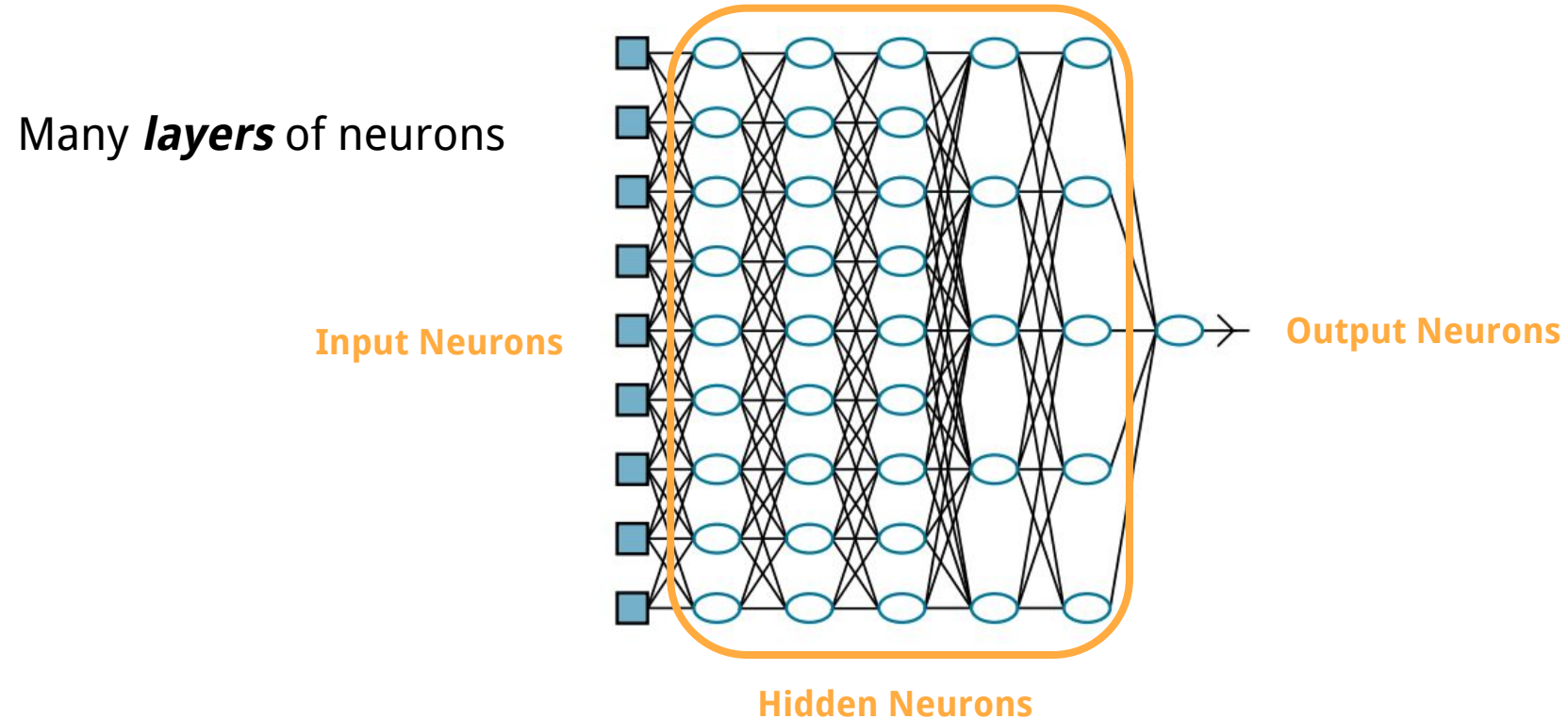


Linear Classifier or Regressor



Neural Network

# Neural Net Structure



# Neural Networks Summary so far

Biologically inspired computation model

1 logistic regression  $\sim$  1 neuron

Wire up many different neurons, arranged in layers

Neural network output depends on:

1. Input, from the training data
2. Weights and biases, which we'll learn

`output = neural_network(input, w)`

Real-valued, can be a single number, a vector of numbers, a 2D array of numbers, ...

Network input

“Parameters” -- a (typically huge) vector of real numbers  
Weights and biases

**Differentiable:**

Can calculate:

$d/dw$  `neural_network(input, w)`



This is a neural network with parameters  $(w_1, w_2)$

$$\begin{aligned} \text{output}_2 &= \text{neural\_network}_2(\text{output}_1, w_2) \\ \text{output}_1 &= \text{neural\_network}_1(\text{input}, w_1) \end{aligned}$$

Real-valued, can be a single number, a vector of numbers, a 2D array of numbers, ...

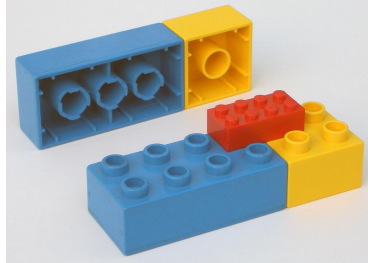
Network input

"Parameters" -- a (typically huge) vector of real numbers  
Weights and biases

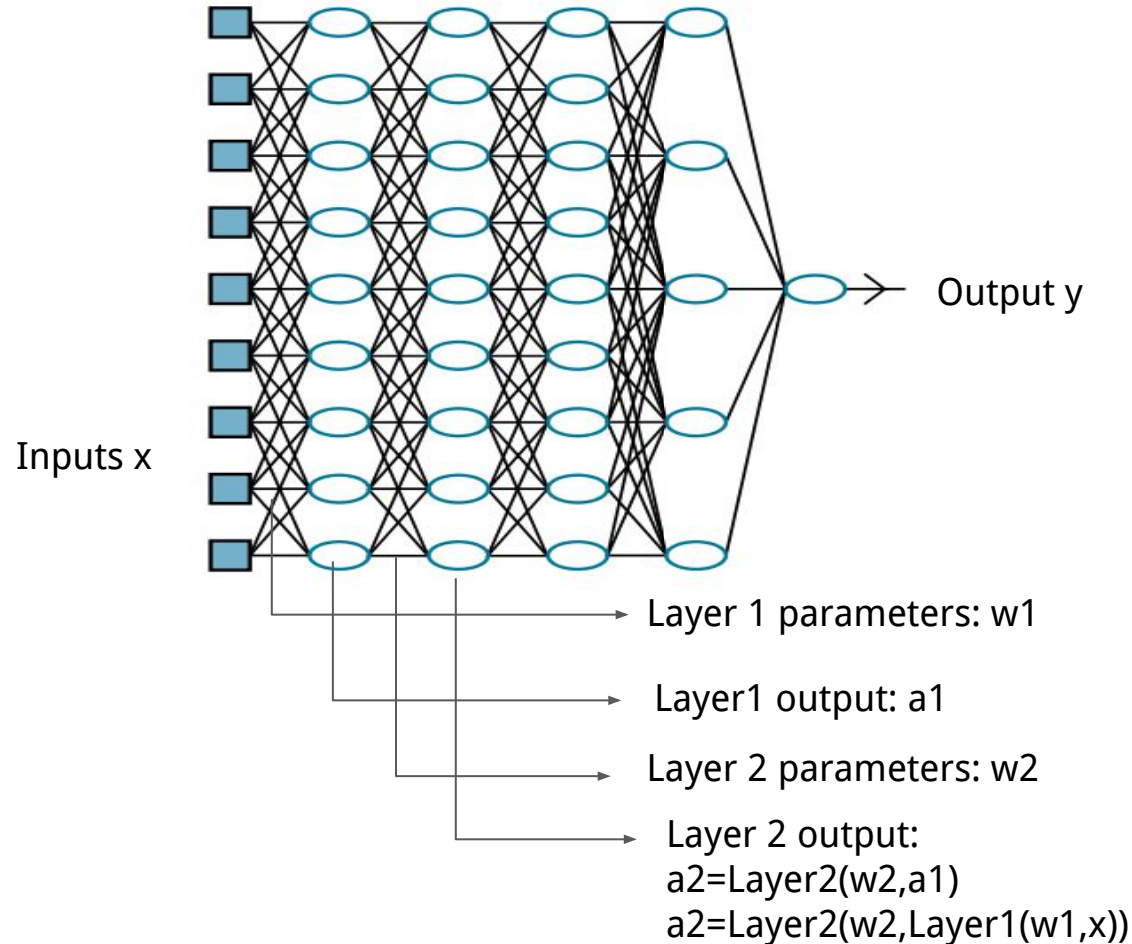
**Differentiable:**

Can calculate:

$d/dw \text{ neural\_network}(\text{input}, w)$



# Layers = Function Composition



# Neural Network Supervised Learning

Training data:  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

Weights (and biases):  $w \in \mathbb{R}^d$

$$w = \underset{w \in \mathbb{R}^D}{\operatorname{argmin}} \underbrace{\sum_n \operatorname{Loss}(\operatorname{neural\_network}(w, x_n), y_n)}$$

$$w = \underset{w \in \mathbb{R}^D}{\operatorname{argmin}} \operatorname{Loss}(w)$$

Loss is a smooth, differentiable function of  $w$   
 $\Rightarrow$  gradient descent

# Gradient Ascent

- consider:  $g(w_1, w_2)$

- Updates:

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

- Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$

$$\text{with: } \nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix}$$

# Gradient Descent Challenges

1. Taking a bunch of derivatives of a complicated function

# Gradient Descent Challenges

1. Taking a bunch of derivatives of a complicated function
2. Gradient descent takes lots of little steps.  
Each step requires looking at the entire training set—which might be huge

# Gradient Descent Challenges

1. Taking a bunch of derivatives of a complicated function  
=> backpropagation, automatic differentiation, pytorch, tensorflow

2. Gradient descent takes lots of little steps.

Each step requires looking at the entire training set—which might be huge

=> stochastic gradient descent, see previous lecture

# Next up

1. Taking a bunch of derivatives of a complicated function

=> backpropagation, automatic differentiation, pytorch, tensorflow

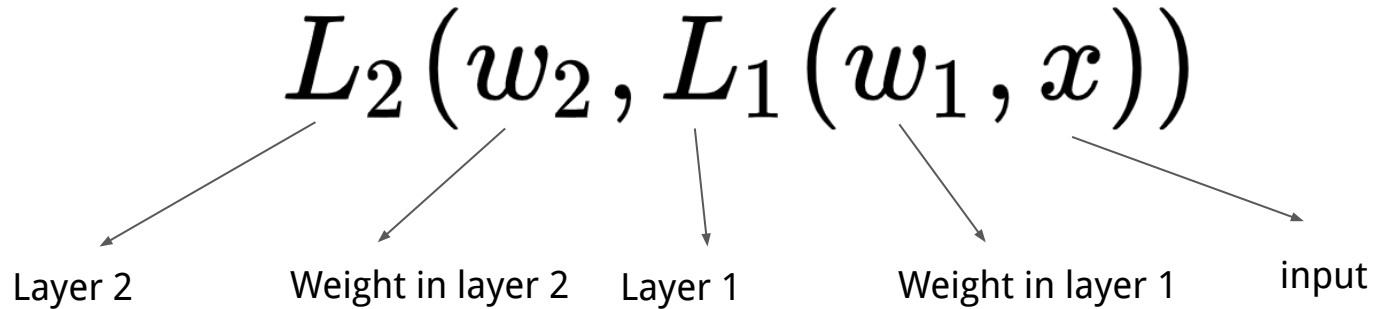
2. Gradient descent takes lots of little steps.

Each step requires looking at the entire training set—which might be huge

=> stochastic gradient descent



# Neural network layers = Function composition



# Neural network layers = Function composition

$$\frac{d}{dw_1} L_2(w_2, L_1(w_1, x))$$

The diagram illustrates the components of the derivative expression  $\frac{d}{dw_1} L_2(w_2, L_1(w_1, x))$ . Arrows point from the following labels to their corresponding parts in the expression:

- Layer 2 points to  $L_2$
- Weight in layer 2 points to  $w_2$
- Layer 1 points to  $L_1$
- Weight in layer 1 points to  $w_1$
- input points to  $x$

# Neural network layers = Function composition

$$\frac{d}{dw_1} L_2(w_2, L_1(w_1, x))$$

Diagram illustrating the components of the function composition:

- Layer 2 (points to  $L_2$ )
- Weight in layer 2 (points to  $w_2$ )
- Layer 1 (points to  $L_1$ )
- Weight in layer 1 (points to  $w_1$ )
- input (points to  $x$ )

$$= \frac{dL_2(w_2, L_1)}{dL_1} \frac{dL_1(w_1, x)}{dw_1}$$

# Modern Neural Network Libraries Do Backpropagation



```
# begin of forward computation
>>> w=torch.rand(6, 6, requires_grad=True)
>>> a3=logistic(w[0,3]+w[1,3]*x1+w[2,3]*x2)
>>> a4=logistic(w[0,4]+w[1,4]*x1+w[2,4]*x2)
>>> a5=logistic(w[0,5]+w[3,5]*a3+w[4,5]*a4)
>>> y=0.7
>>> loss=(a5-y)**2 # end of forward computation
>>> loss.backward() # run back propagation
>>> w.grad[3,5] # dLoss / d w_{35}
tensor(0.0570) # automatic differentiation
```

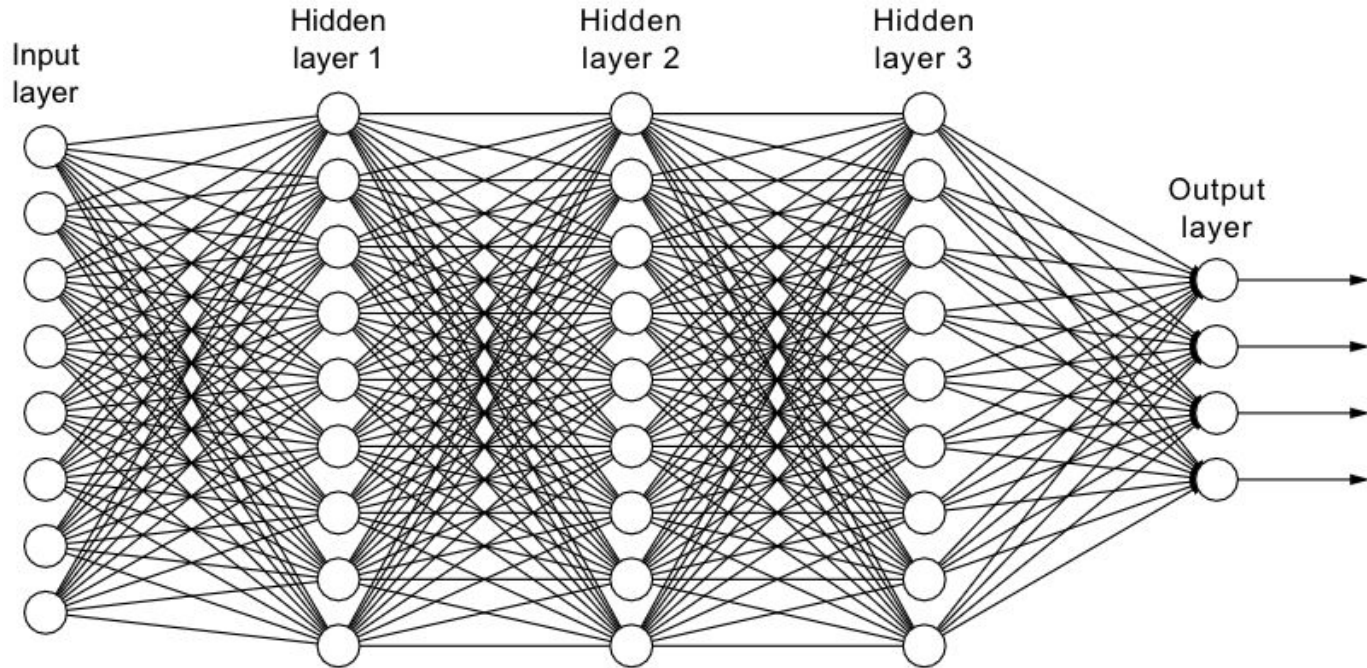
# Learning Autodifferentiation

Implement it! <https://sidsite.com/posts/autodiff/>

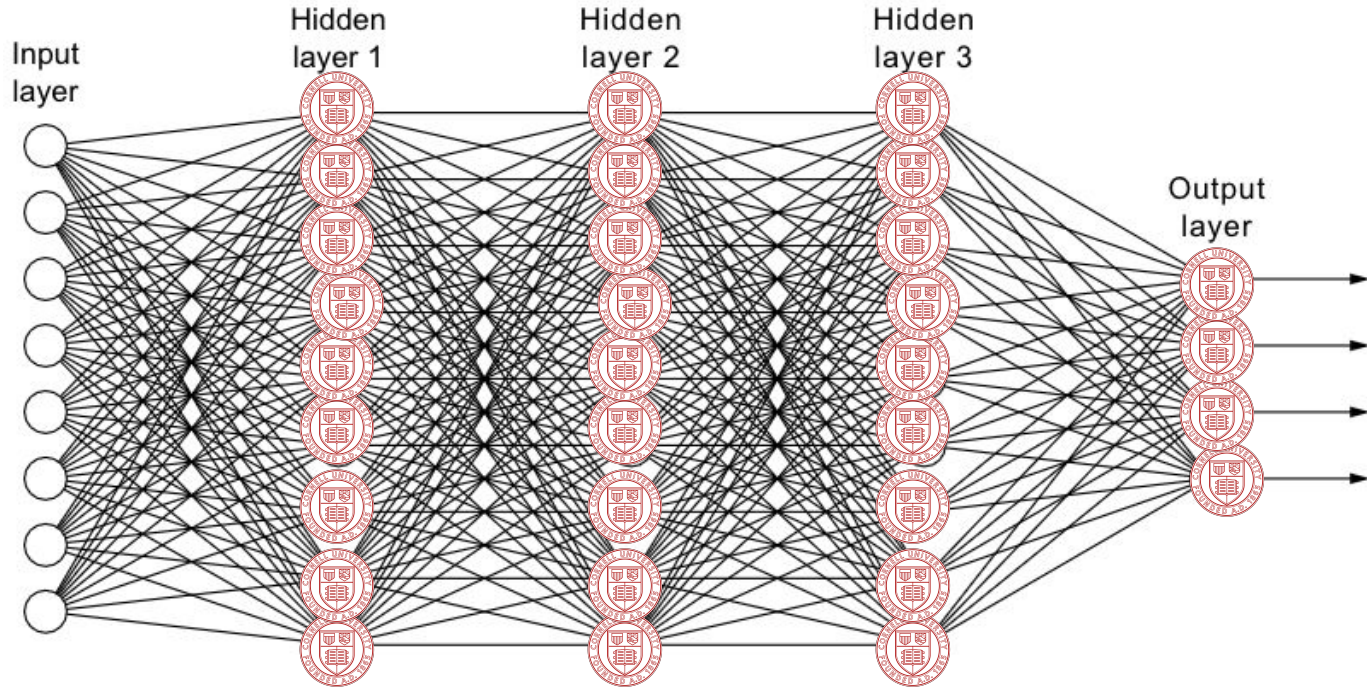
Slides at the end of this lecture cover neural network training /without/ automatic differentiation. The algebra is horrible.

# Neural Net Architectures

# Multilayer Perceptron (MLP)

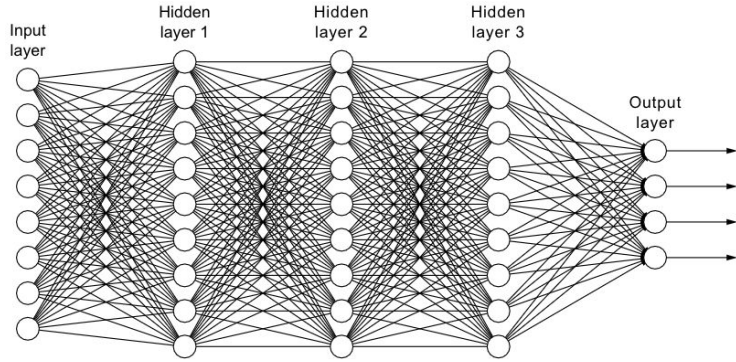


# Multilayer Perceptron (MLP)





# Multilayer Perceptron (MLP)



## Advantages:

- Doesn't assume much about the inputs
- Could permute the inputs and the learner would be oblivious

## Disadvantages:

- Doesn't assume much about the inputs
- Could permute the inputs and the learner would be oblivious

# Specialized Neural Architectures:

Images  
Text

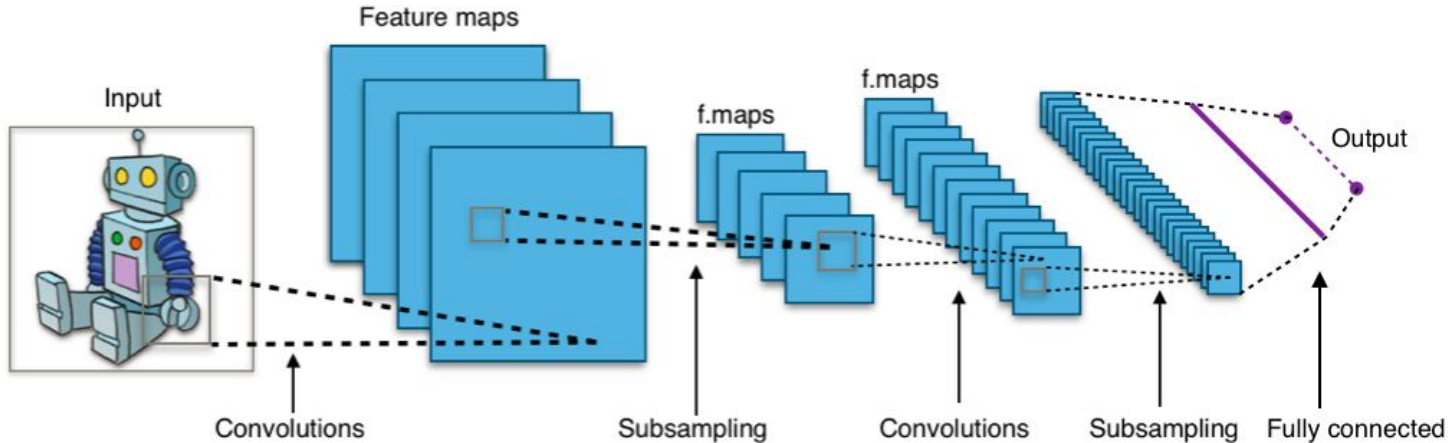
## Specialized Neural Architectures:

Images (2D arrays)  
Text (1D sequences)

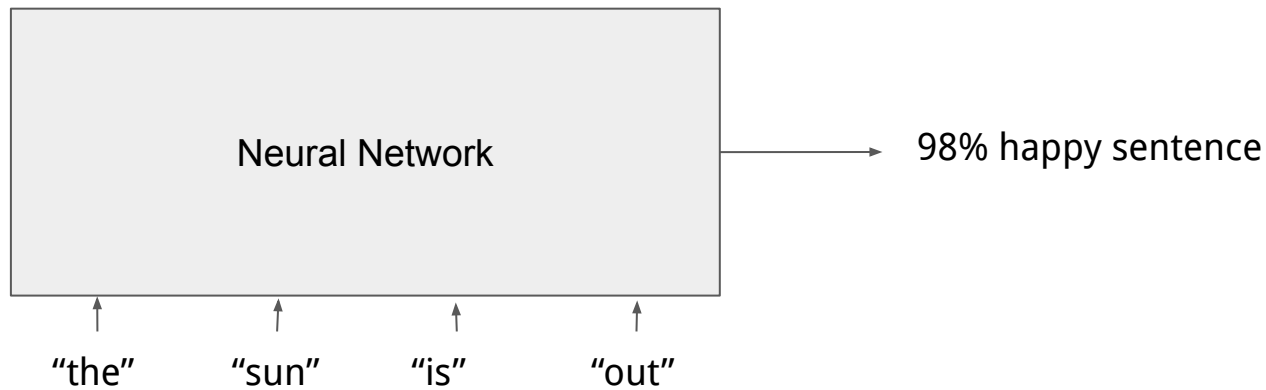
# Images (2D array)

## Convolutional neural networks

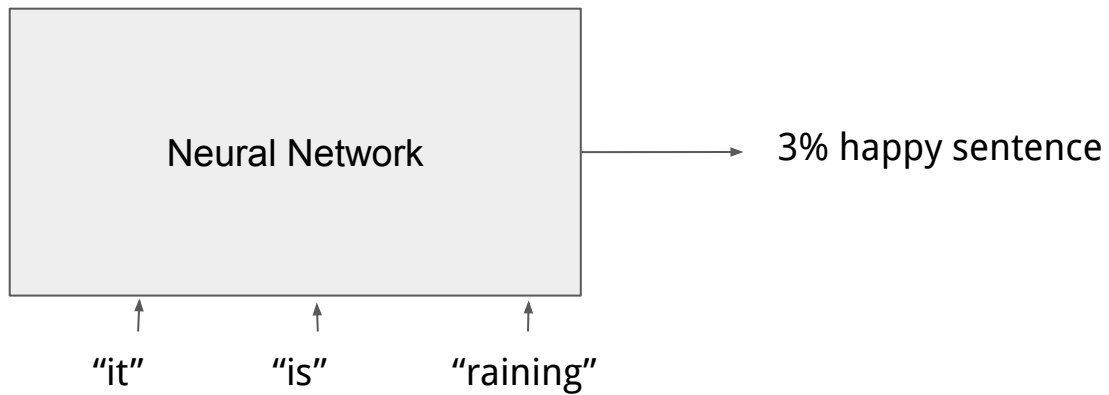
- Neuron that takes 5x5 grid from an image as input (“receptive field”)
- Make copies of it for every 5x5 window in image
- Assign them all the same weights
- Pooling layer: Is something present in *any* of the windows



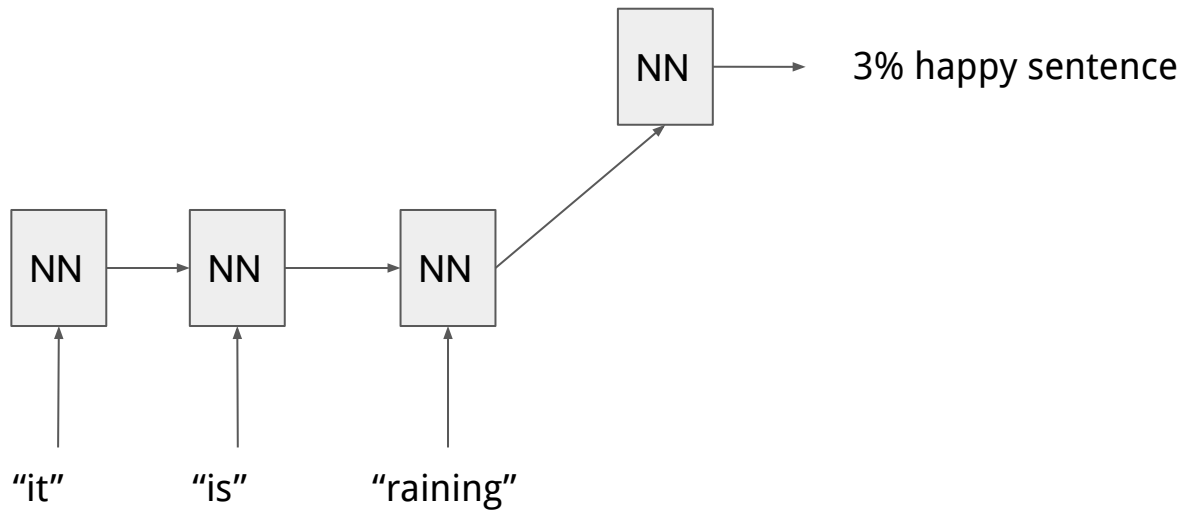
# Text (1D Sequences)



# Text (1D Sequences)



# Text (1D Sequences)



This is a neural network with parameters  $(w_1, w_2)$

$$\begin{aligned} \text{output}_2 &= \text{neural\_network}_2(\text{output}_1, w_2) \\ \text{output}_1 &= \text{neural\_network}_1(\text{input}, w_1) \end{aligned}$$

Real-valued, can be a single number, a vector of numbers, a 2D array of numbers, ...

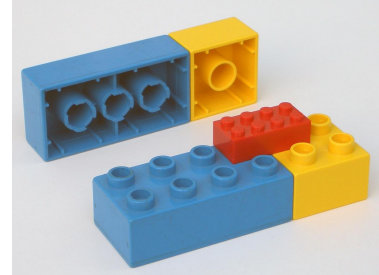
Network input

"Parameters" -- a (typically huge) vector of real numbers  
Weights and biases

**Differentiable:**

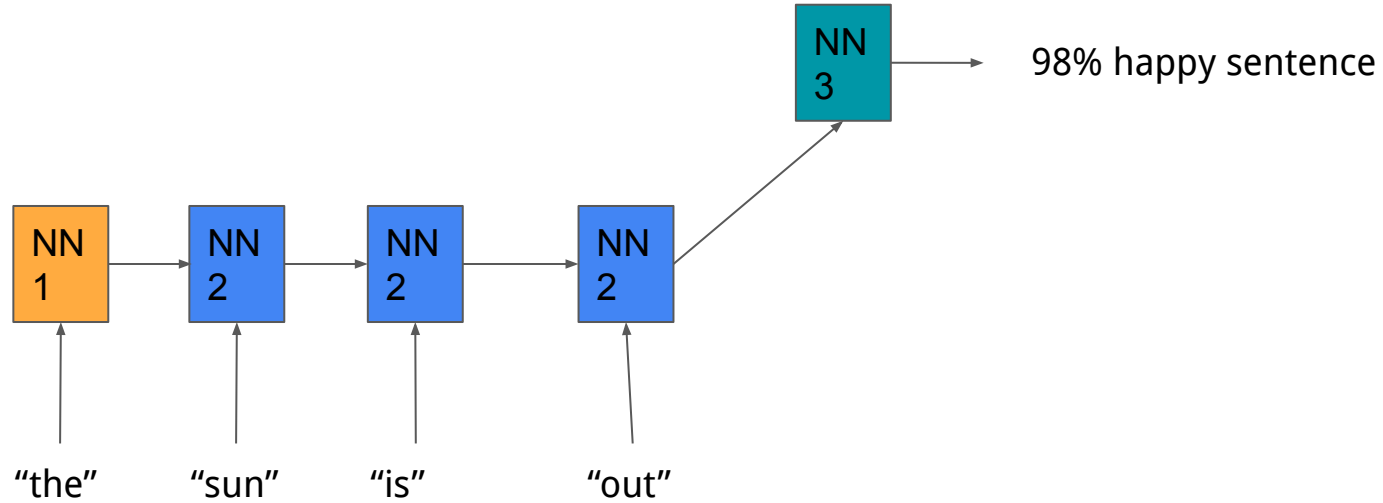
Can calculate:

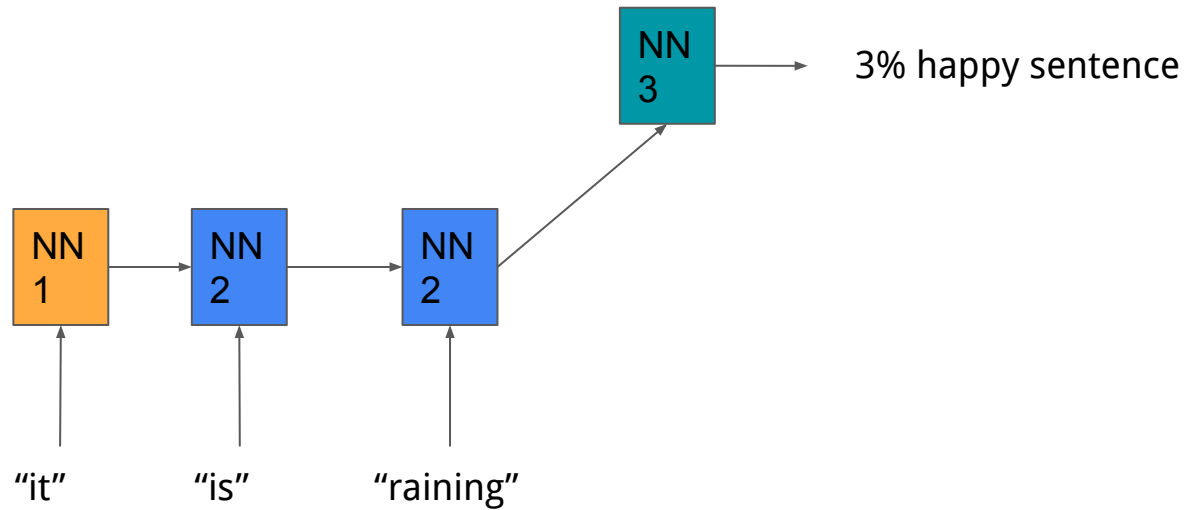
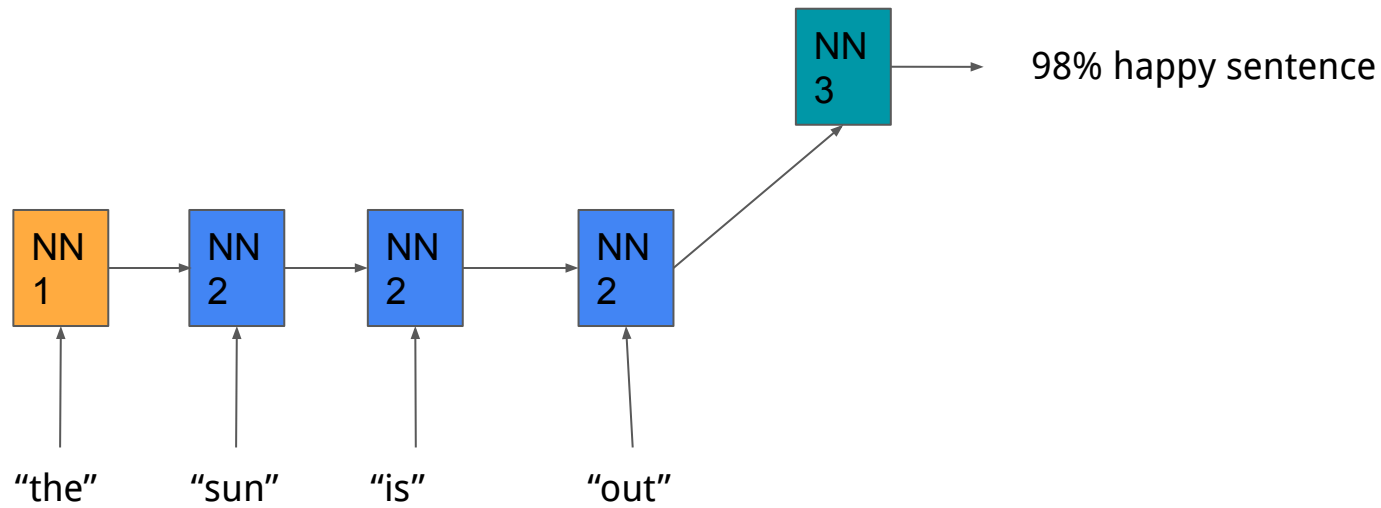
$d/dw \text{ neural\_network}(\text{input}, w)$





# Text (1D Sequences)

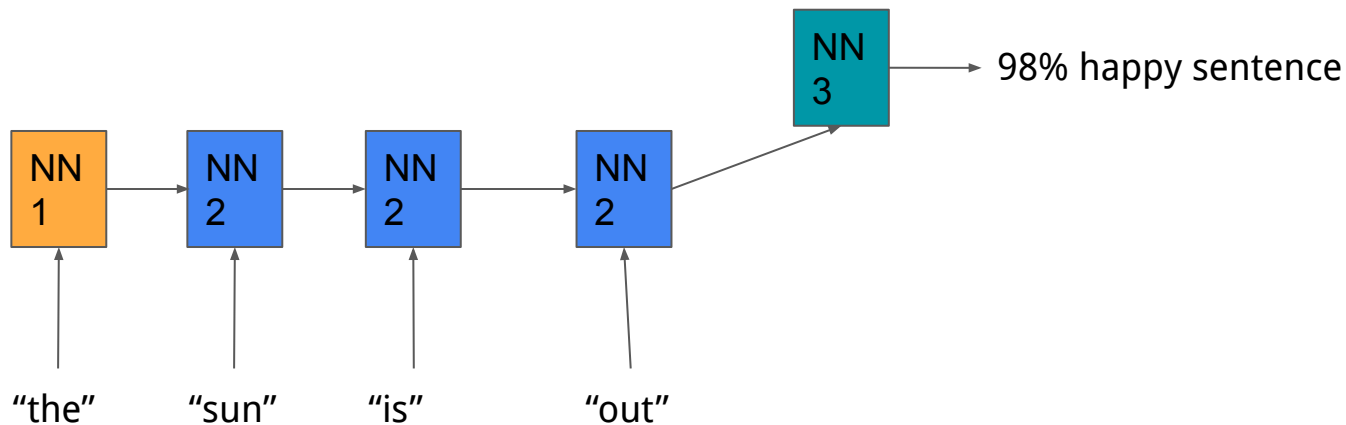




# Recurrent Neural Network

Further details:

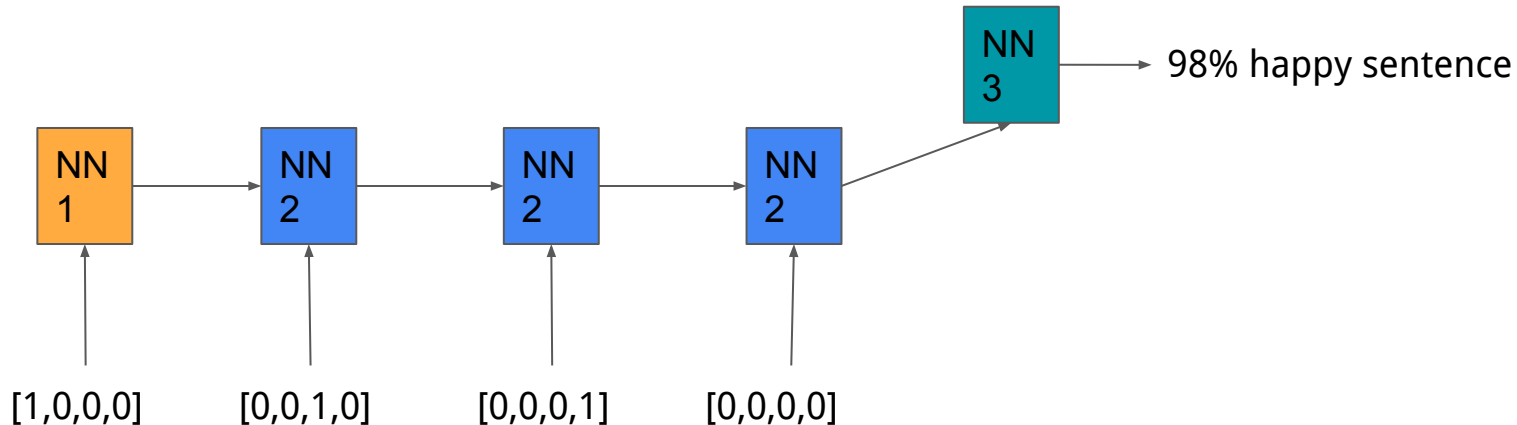
1. Each word is turned into a one-hot vector
2. Usually NN1 is just NN2 but with a dummy input fed in on the left
3. Usually NN3 is just a logistic regressor



# Recurrent Neural Network

Further details:

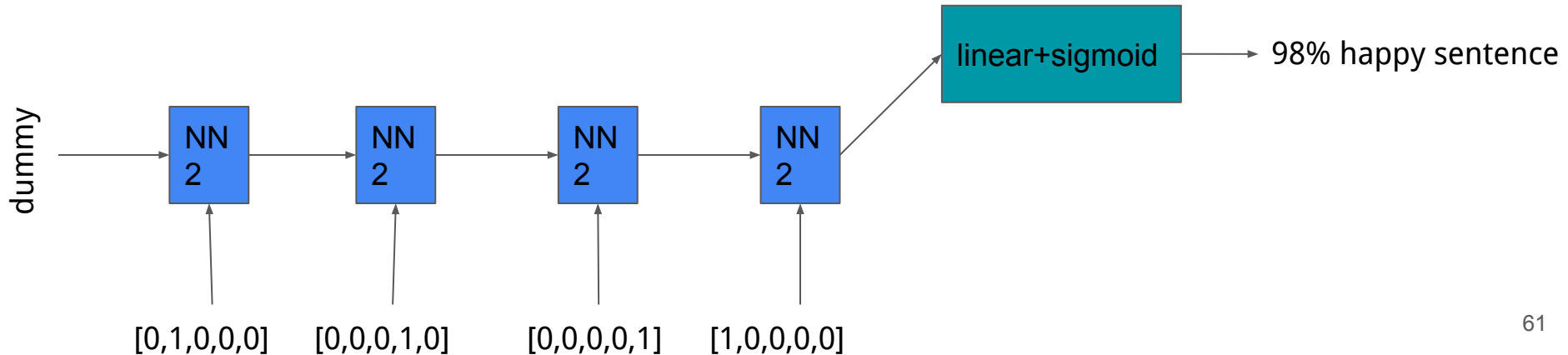
1. Each word is turned into a one-hot vector
2. Usually NN1 is just NN2 but with a dummy input fed in on the left
3. Usually NN3 is just a logistic regressor



# Recurrent Neural Network

Further details:

1. Each word is turned into a one-hot vector
2. Usually NN1 is just NN2 but with a dummy input fed in on the left
3. Usually NN3 is just a logistic regressor



# Summary, architecture

## Convolutional Neural Network:

- Runs the same neural net computation over sliding windows of the image
- Downsamples the input
- Repeatedly alternates sliding windows and down sampling

## Recurrent Neural Network:

- Runs the same neural net computation over each element of a sequence
- Feeds the output of earlier inputs into later input computations

# Bonus Slides on backpropagation

Automatic differentiation makes it so that you don't have to do the computations on the following slides.

None of this is on any of your homework or exams.

It should make you appreciate automatic differentiation, however.

## Chain Rule

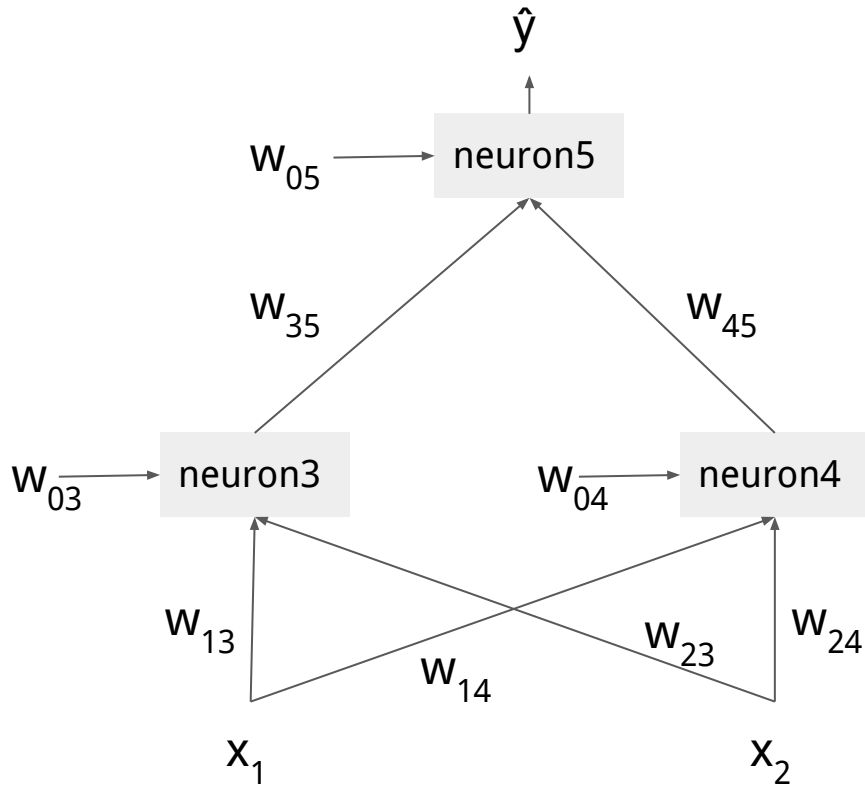
$$\frac{d}{dx} f(g(x)) = f'(g(x)) \frac{d}{dx} g(x)$$

$$z = g(x)$$

$$y = f(z)$$

$$\frac{d}{dx} y = \frac{d}{dx} f(g(x)) = \frac{df(z)}{dz} \frac{dz}{dx} = \frac{df(z)}{dz} \frac{dg(z)}{dx}$$



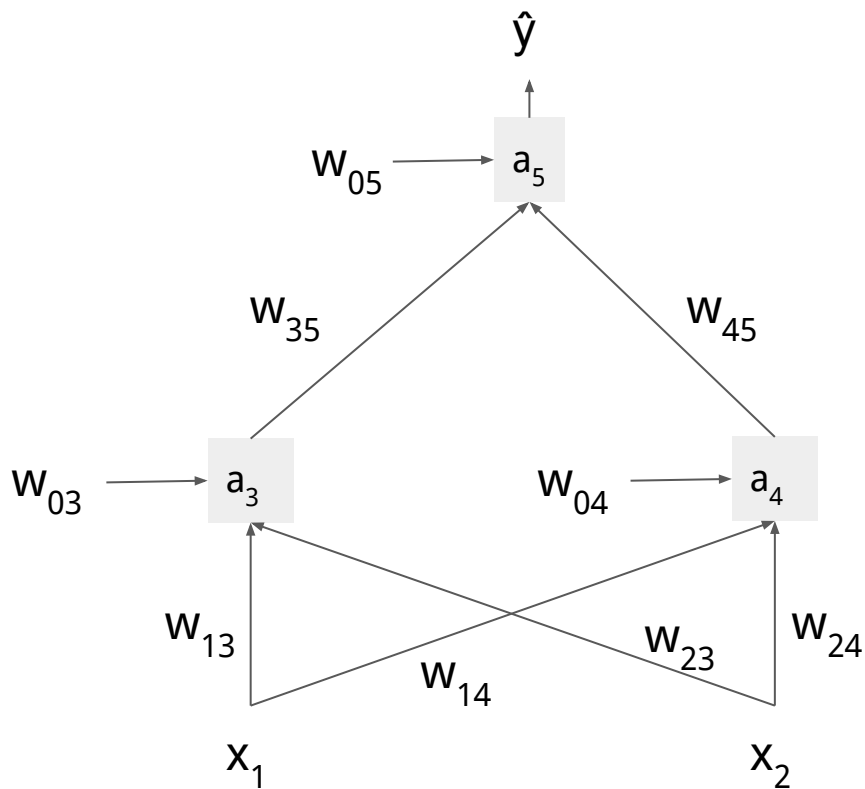


$$\text{Loss}(w) = (y - \hat{y})^2$$

Target output, from  
training example

Predicted output,  
function of  $w$  and  
training input  $x$

Bias for neuron  $i$ :  $w_{0i}$   
Weight from  $i \rightarrow j$ :  $w_{ij}$

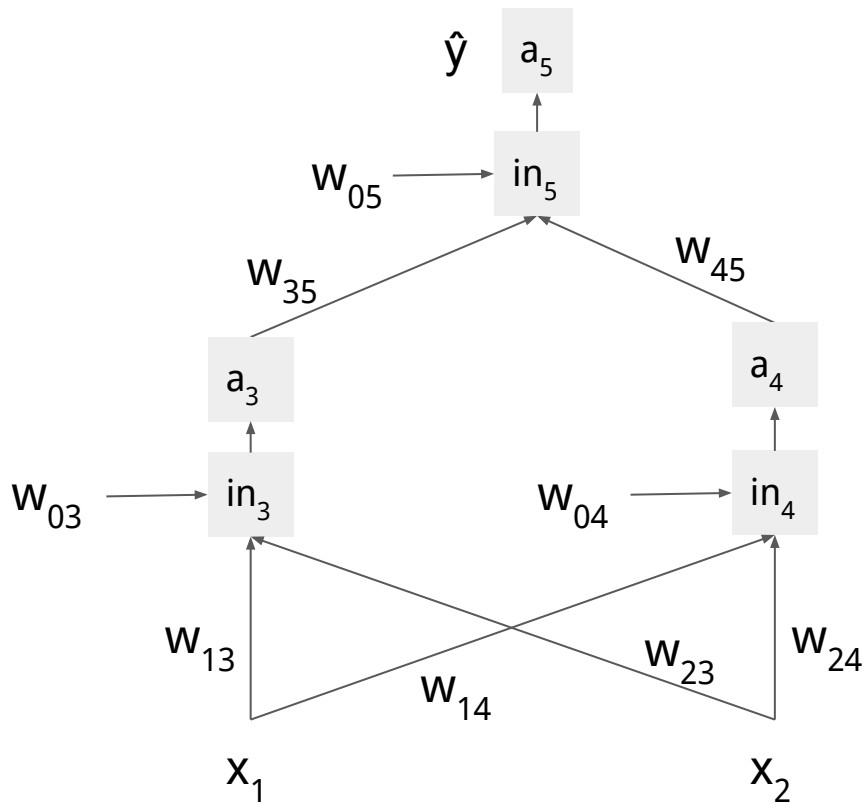


$$\text{Loss}(w) = (y - \hat{y})^2$$

Target output, from  
training example

Predicted output,  
function of  $w$  and  
training input  $x$

Bias for neuron  $i$ :  $w_{0i}$   
Weight from  $i \rightarrow j$ :  $w_{ij}$



$$\text{Loss}(w) = (y - \hat{y})^2$$

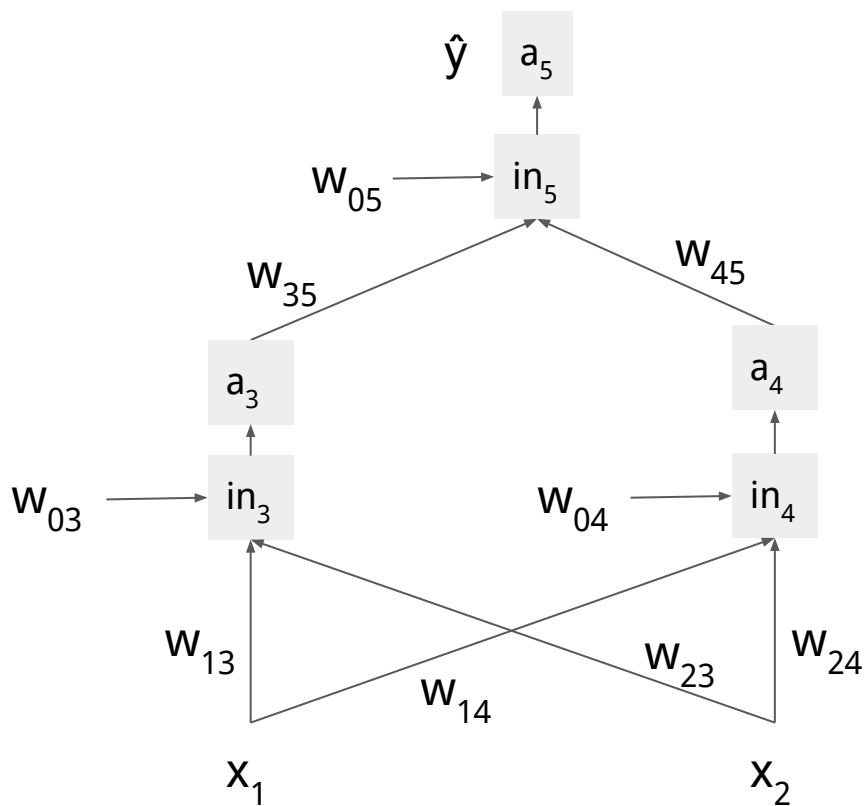
Target output, from training example

Predicted output, function of w and training input x

Activation output:  $a_i = g(\text{in}_i)$

Linear input:  $\text{in}_i = w_{0i} + \sum_{j \rightarrow i} w_{ji} a_j$

Bias for neuron i:  $w_{0i}$   
Weight from i  $\rightarrow$  j:  $w_{ij}$



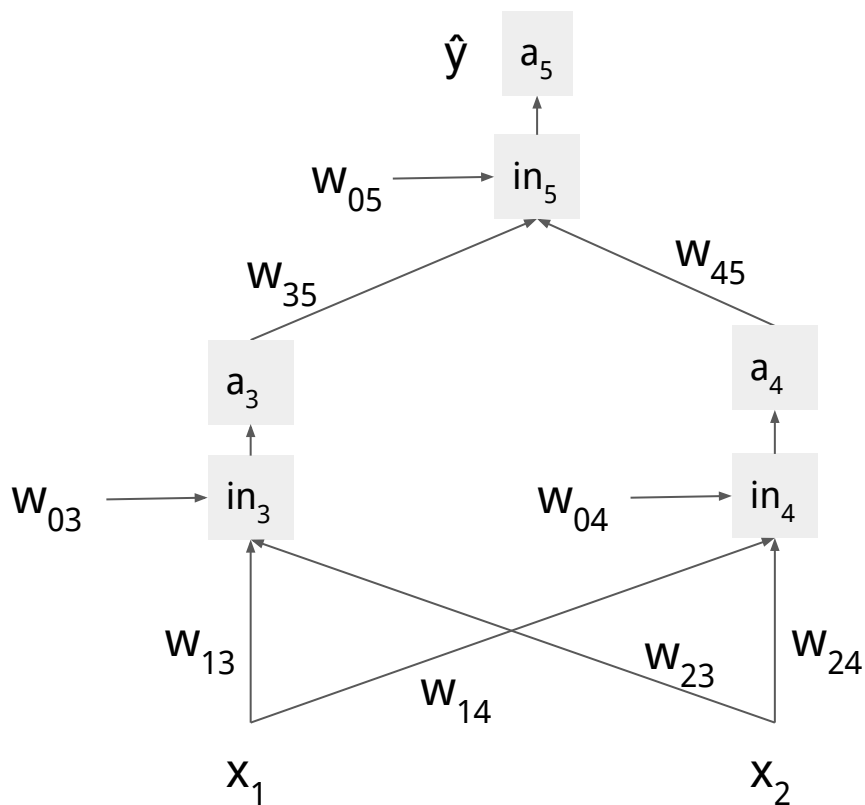
$$\text{Loss}(w) = (y - \hat{y})^2$$

$$\begin{aligned}
 \frac{d\text{Loss}}{dw_{05}} &= \frac{d\text{Loss}}{da_5} \frac{da_5}{dw_{05}} \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} \frac{din_5}{dw_{05}} \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} \frac{d}{dw_{05}} (w_{05} + w_{35}a_3 + w_{45}a_4) \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} (1 + 0 + 0) \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} \\
 &= \frac{d\text{Loss}}{da_5} g'(in_5) = 2(y - a_5)g'(in_5)
 \end{aligned}$$

Bias for neuron i:  $w_{0i}$   
 Weight from i  $\rightarrow$  j:  $w_{ij}$

Activation output:  $a_i = g(in_i)$

Linear input:  $in_i = w_{0i} + \sum_{j \rightarrow i} w_{ji} a_j$



$$\text{Loss}(w) = (y - \hat{y})^2$$

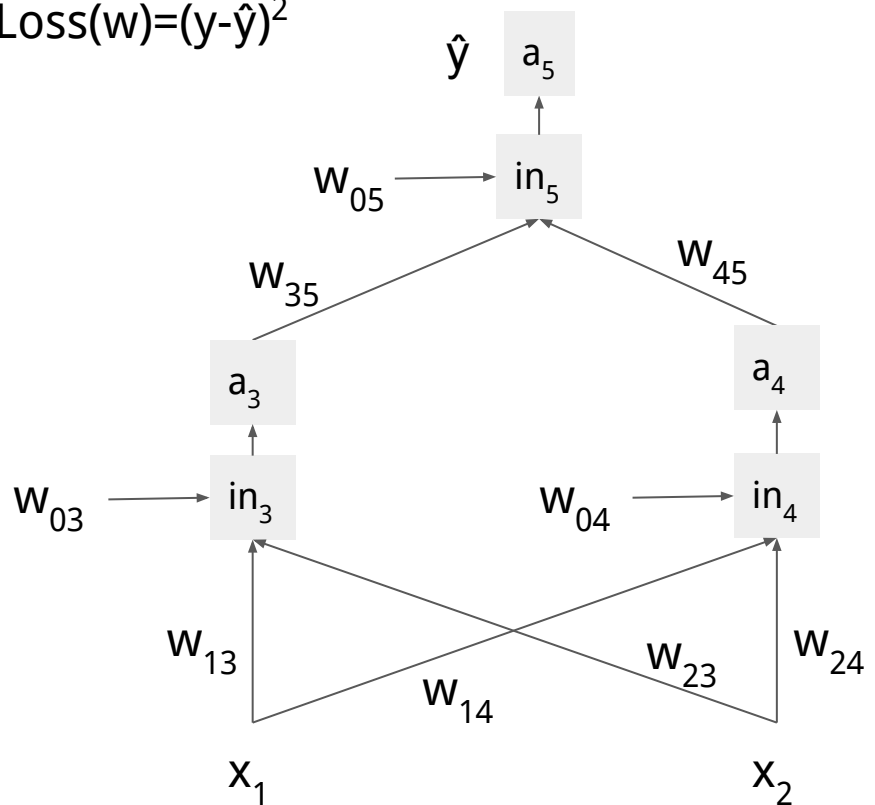
$$\begin{aligned}
 \frac{d\text{Loss}}{dw_{35}} &= \frac{d\text{Loss}}{da_5} \frac{da_5}{dw_{35}} \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} \frac{din_5}{dw_{35}} \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} \frac{d}{dw_{35}} (w_{05} + w_{35}a_3 + w_{45}a_4) \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} (0 + a_3 + 0) \\
 &= \frac{d\text{Loss}}{da_5} \frac{da_5}{din_5} a_3 \\
 &= \frac{d\text{Loss}}{da_5} g'(\text{in}_5) a_3 = 2(y - a_5) g'(\text{in}_5) a_3
 \end{aligned}$$

Bias for neuron  $i$ :  $w_{0i}$   
 Weight from  $i \rightarrow j$ :  $w_{ij}$

Activation output:  $a_i = g(\text{in}_i)$   
 Linear input:  $\text{in}_i = w_{0i} + \sum_{j \rightarrow i} w_{ji} a_j$

<many derivatives later>

$$\text{Loss}(w) = (y - \hat{y})^2$$



Activation output:  $a_i = g(in_i)$

Bias for neuron i:  $w_{0i}$

Weight from i  $\rightarrow$  j:  $w_{ij}$

$$\text{Linear input: } in_i = w_{0i} + \sum_{j \rightarrow i} w_{ji} a_j$$

$$\frac{d\text{Loss}}{din_5} = 2(y - a_5)g'(in_5)$$

$$\frac{d\text{Loss}}{dw_{05}} = \frac{d\text{Loss}}{din_5} \times 1$$

$$\frac{d\text{Loss}}{dw_{35}} = \frac{d\text{Loss}}{din_5} \times a_3$$

$$\frac{d\text{Loss}}{din_3} = \frac{d\text{Loss}}{din_5} w_{35} g'(in_3)$$

$$\frac{d\text{Loss}}{dw_{03}} = \frac{d\text{Loss}}{din_3} \times 1$$

$$\frac{d\text{Loss}}{dw_{13}} = \frac{d\text{Loss}}{din_3} \times x_1$$

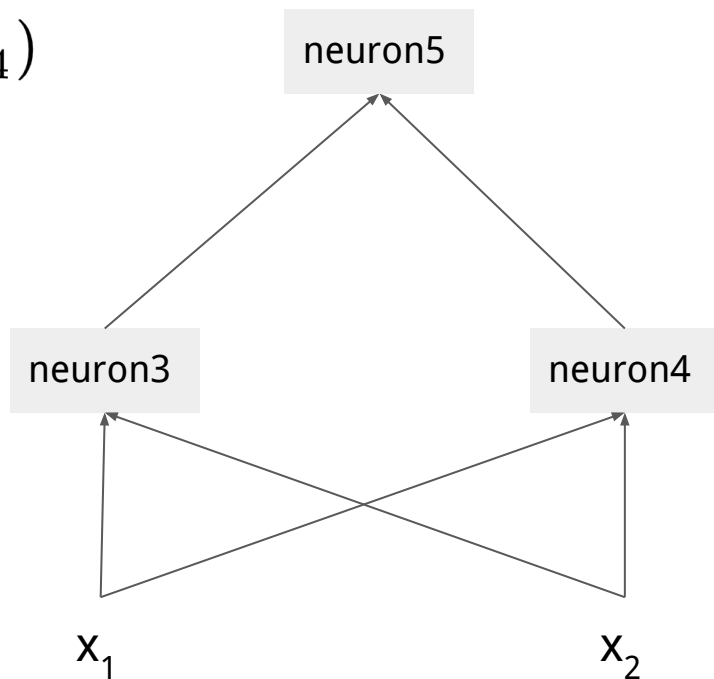
# A pattern is emerging

$$\frac{d\text{Loss}}{din_4} = \frac{d\text{Loss}}{din_5} w_{45} g'(\text{in}_4)$$

$$\frac{d\text{Loss}}{dw_{04}} = \frac{d\text{Loss}}{din_4} \times 1$$

$$\frac{d\text{Loss}}{dw_{14}} = \frac{d\text{Loss}}{din_4} \times x_1$$

$$\frac{d\text{Loss}}{dw_{24}} = \frac{d\text{Loss}}{din_4} \times x_2$$



$$\frac{d\text{Loss}}{din_3} = \frac{d\text{Loss}}{din_5} w_{35} g'(\text{in}_3)$$

$$\frac{d\text{Loss}}{dw_{03}} = \frac{d\text{Loss}}{din_3} \times 1$$

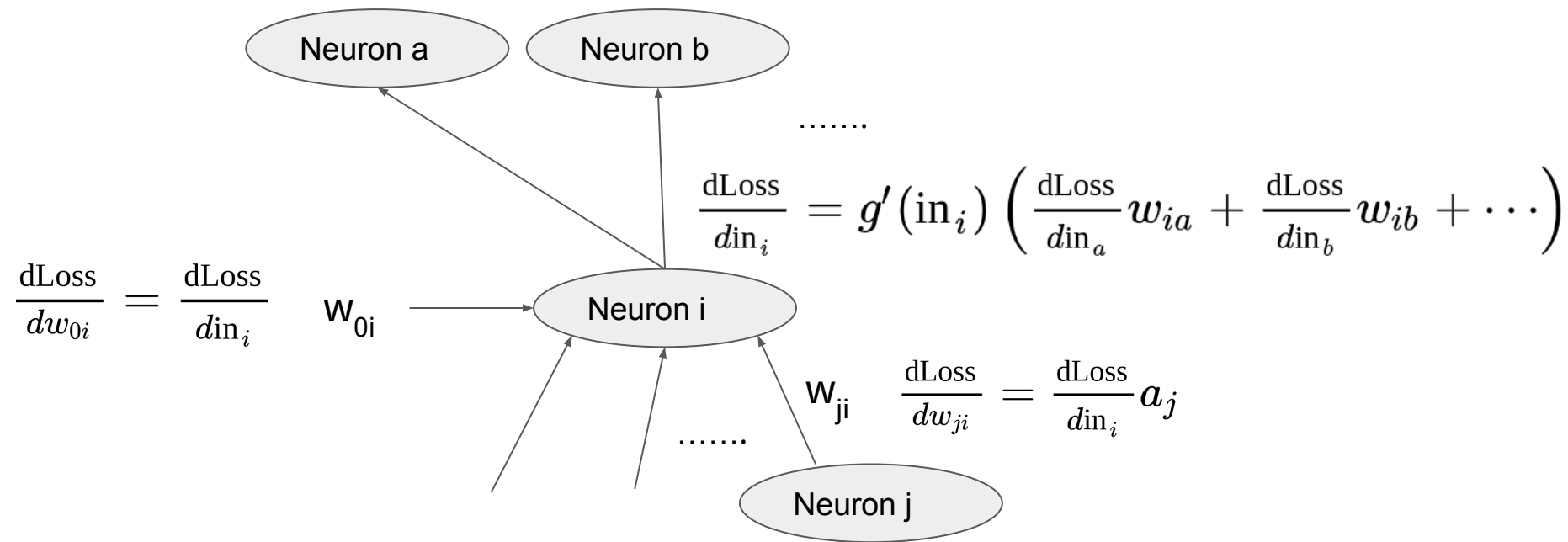
$$\frac{d\text{Loss}}{dw_{13}} = \frac{d\text{Loss}}{din_3} \times x_1$$

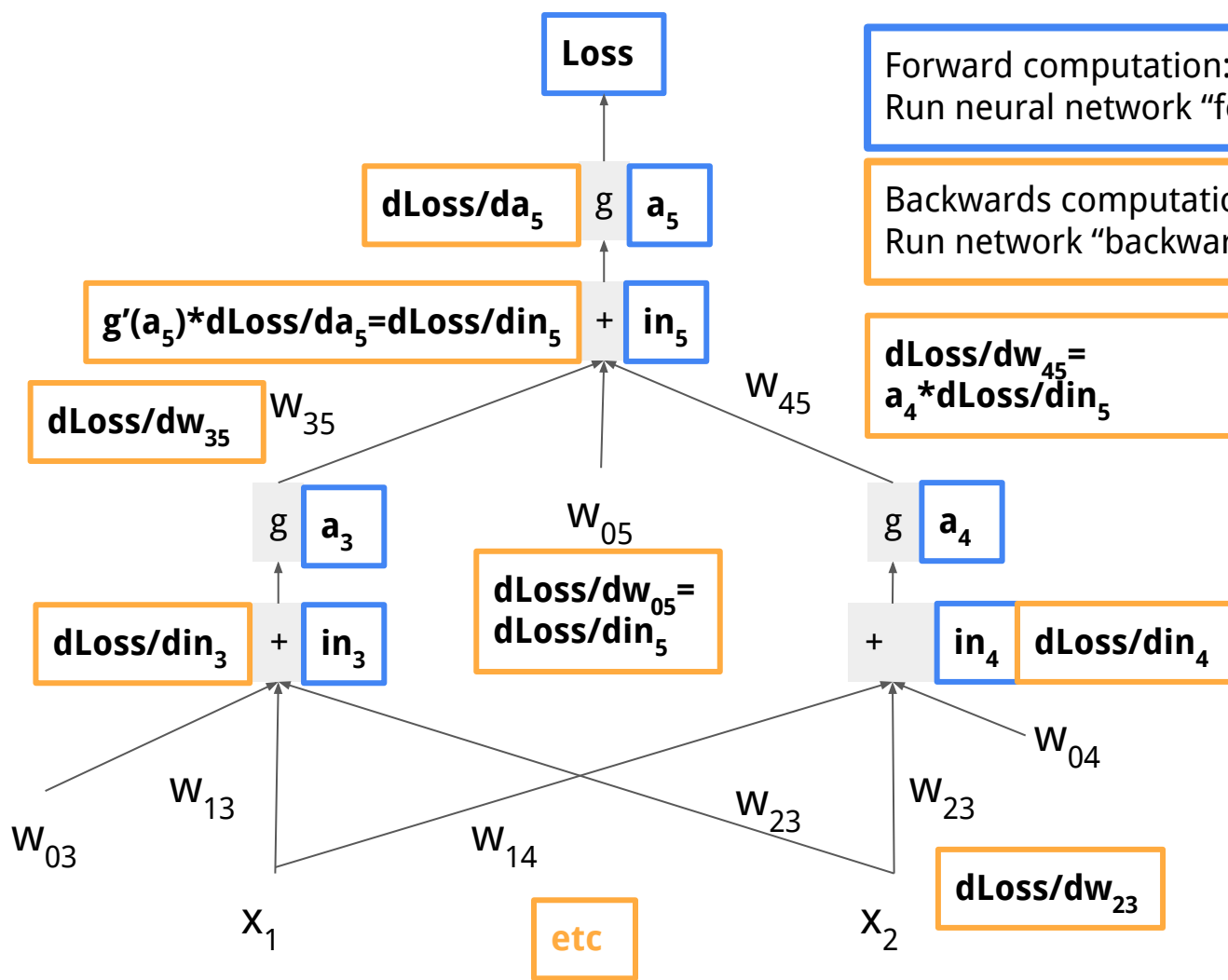
$$\frac{d\text{Loss}}{dw_{23}} = \frac{d\text{Loss}}{din_3} \times x_2$$



# Generalizing the pattern

Key quantity:  $\frac{d\text{Loss}}{din_i}$  (how much loss changes as input to neuron i changes)  
 (“perceived error” at neuron i)





Forward computation:  
Run neural network "forward" and compute loss

Backwards computation:  
Run network "backward" and compute derivatives

$$dLoss/dw_{45} = a_4 * dLoss/din_5$$

$$dLoss/dw_{05} = dLoss/din_5$$

$$dLoss/dw_{23}$$

# Backpropagation

Objective: compute  $d\text{Loss} / dw_{ij}$  for each weight/bias  $w_{ij}$

1. Run the network forward to compute each neuron's output
2. Compute the loss
3. Compute the derivative of the loss with respect to network output:

$$d\text{Loss}/dy$$

4. Compute derivative of loss with respect to each neuron's total input

$$\frac{d\text{Loss}}{din_i} = g'(\text{in}_i) \sum_{i \rightarrow j} \frac{d\text{Loss}}{din_j} w_{ij}$$

5. Compute derivative of loss with respect to weights/biases

$$\frac{d\text{Loss}}{dw_{ji}} = \frac{d\text{Loss}}{din_i} a_j$$