# Today: Outline

- **Neural networks:** artificial neuron, MLP, sigmoid units; neuroscience inspiration, output vs hidden layers; linear vs nonlinear networks;

- **Feed-forward networks**

- Reminders: First lab session tomorrow (Fri Sept 18)

# Fei-Fei Li



- Professor, Computer Science, Stanford University

- Co-Director of Stanford's Human-Centered AI Institute

- Previously Vice President at Google and Chief Scientist of AI/ML at Google Cloud

- Co-founder and chairperson of the national non-profit AI4ALL

- Online Deep Learning Course

- **"First, we teach them see, then they help us to see better."**

# Image Captioning



*A young boy holding a baseball bat*



*A man riding a horse next to a building*

# Introduction to Neural Networks

Motivation

# Recall: Logistic Regression

$$0 \leq h_\theta(x) \leq 1$$

$$h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

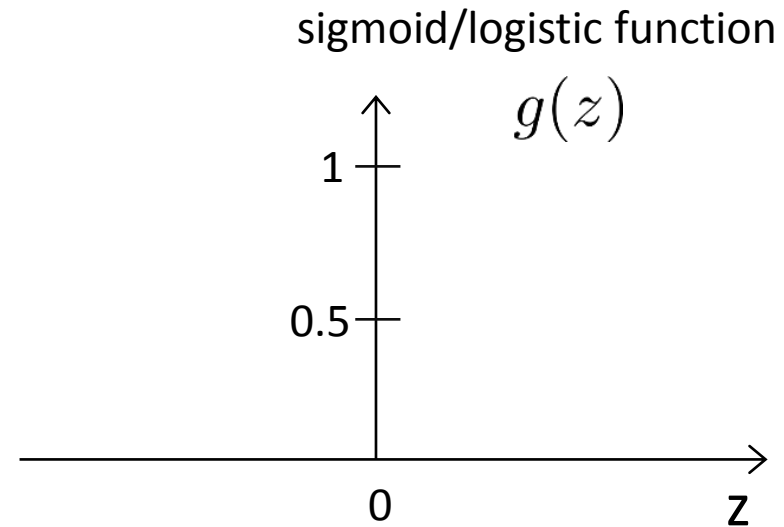$$g(z) = \frac{1}{1 + e^{-z}}$$

Output is probability of label 1 given input

$$p(y = 1|x) = \frac{1}{1 + e^{-\theta^T x}}$$

sigmoid/logistic function



$g(z)$

predict "$y = 1$" if $h_\theta(x) \geq 0.5$
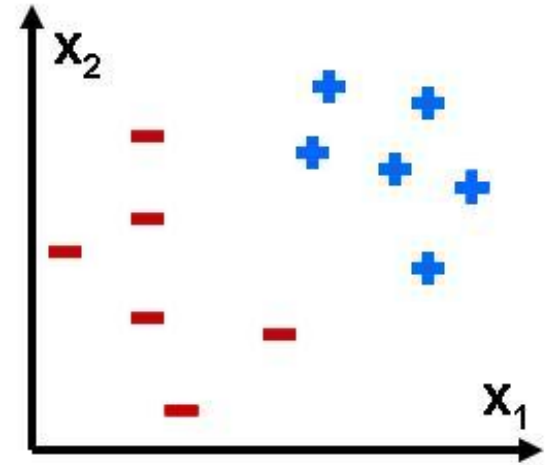
predict "$y = 0$" if $h_\theta(x) < 0.5$

# Recall: Logistic Regression Cost

Logistic Regression Hypothesis:

$$h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$
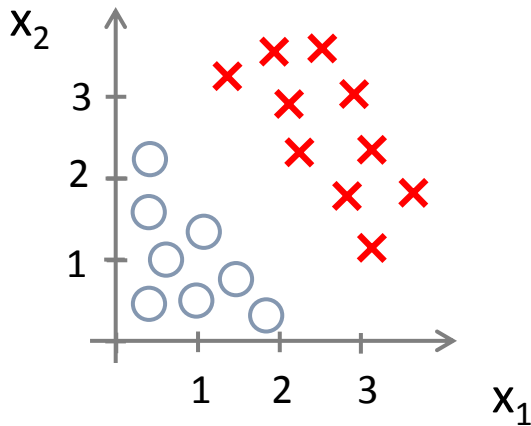
$\theta$: parameters

$D = \{x^i, y^i\}$: data

Logistic Regression Cost Function:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} [\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))]$$
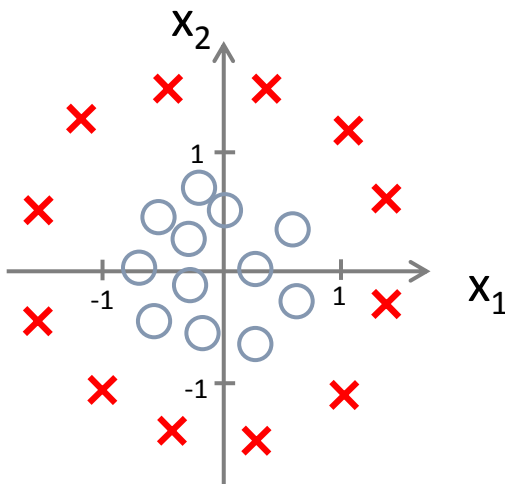
Goal: minimize cost $\quad \min_\theta J(\theta)$

# Decision boundary



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

## Non-linear decision boundaries



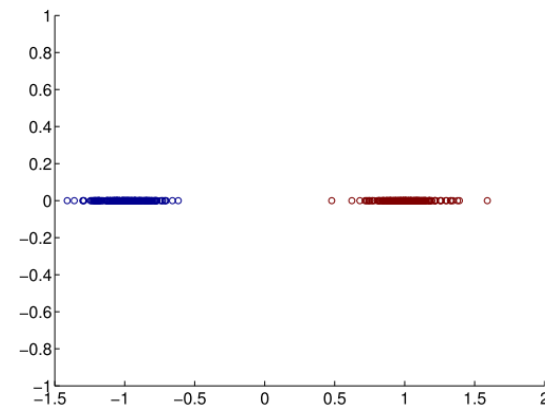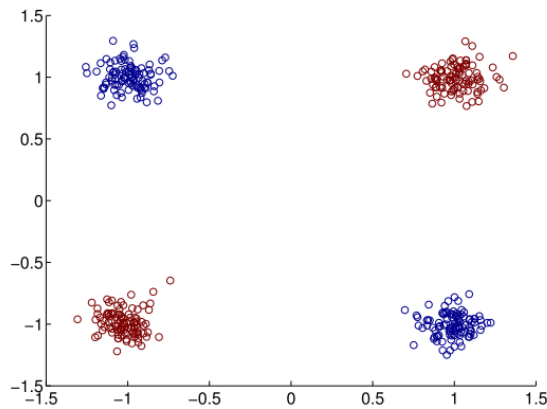Replace features with non-linear functions
e.g. log, cosine, or polynomial

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

# Nonlinear basis functions

**Transform the input/feature**

$$\phi(x) : x \in R^2 \ \rightarrow \ z \ = \ x_1 \cdot x_2$$

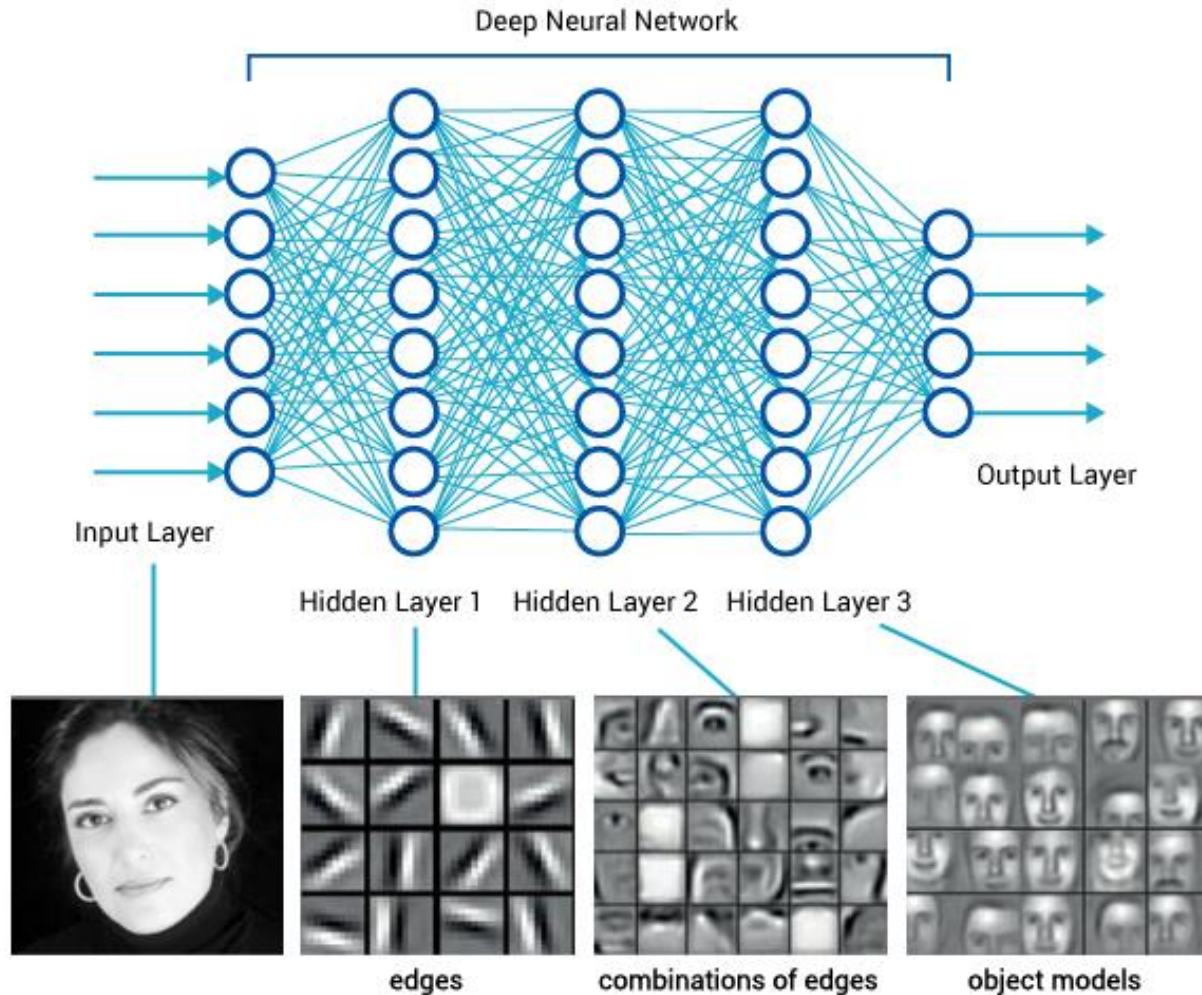**Transformed training data: linearly separable!**

# Limitations of linear models

- Logistic regression and other linear models cannot handle nonlinear decision boundaries
  - Must use non-linear feature transformations
  - Up to designer to specify which one

- Can we instead learn the transformation?
  - Yes, this is what neural networks do!

- A Neural network chains together many layers of "neurons" such as logistic units (logistic regression functions)
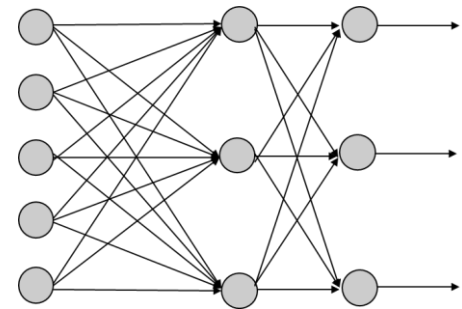
# Neural Networks learn features
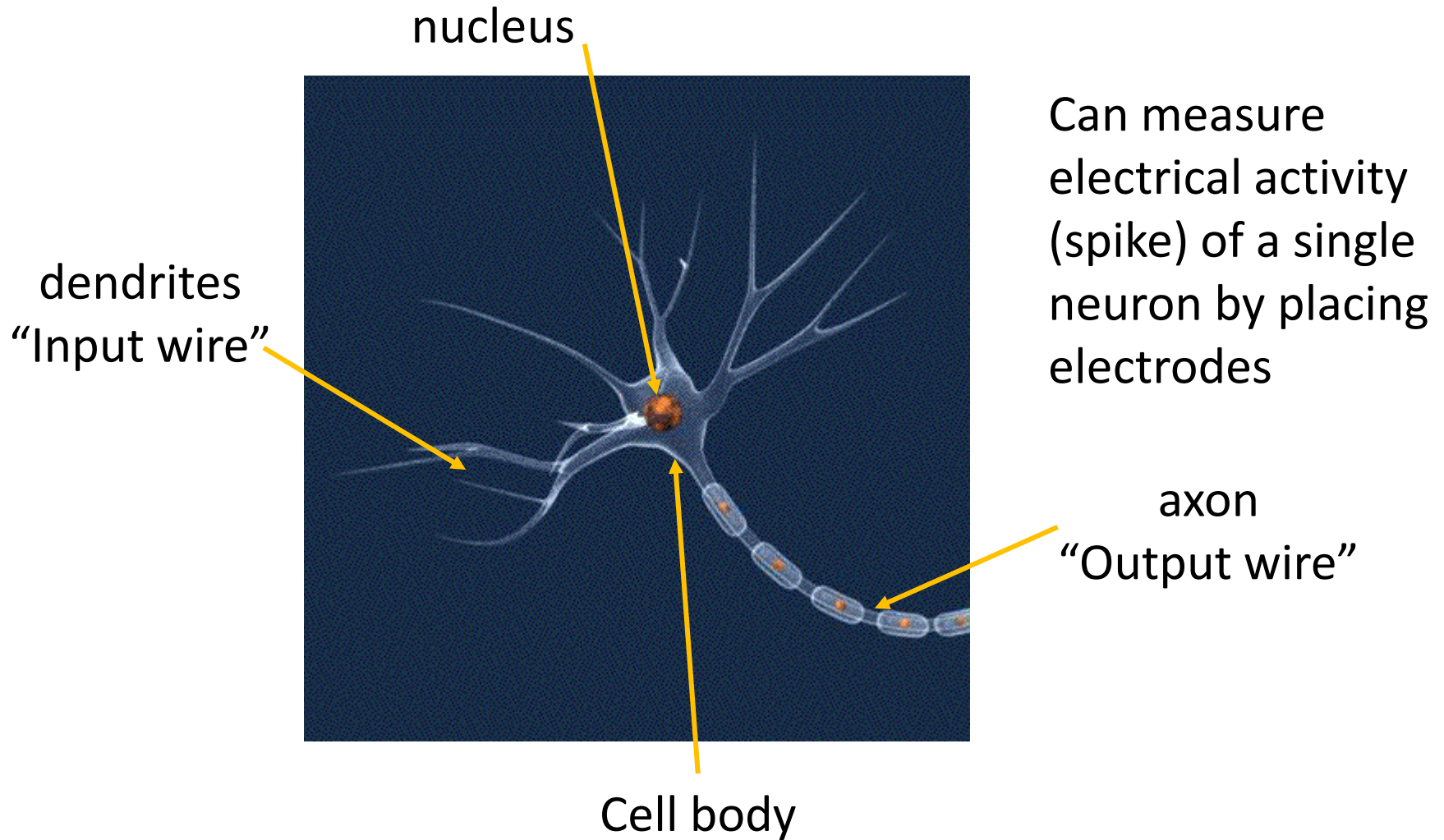
# Neurons in the Brain
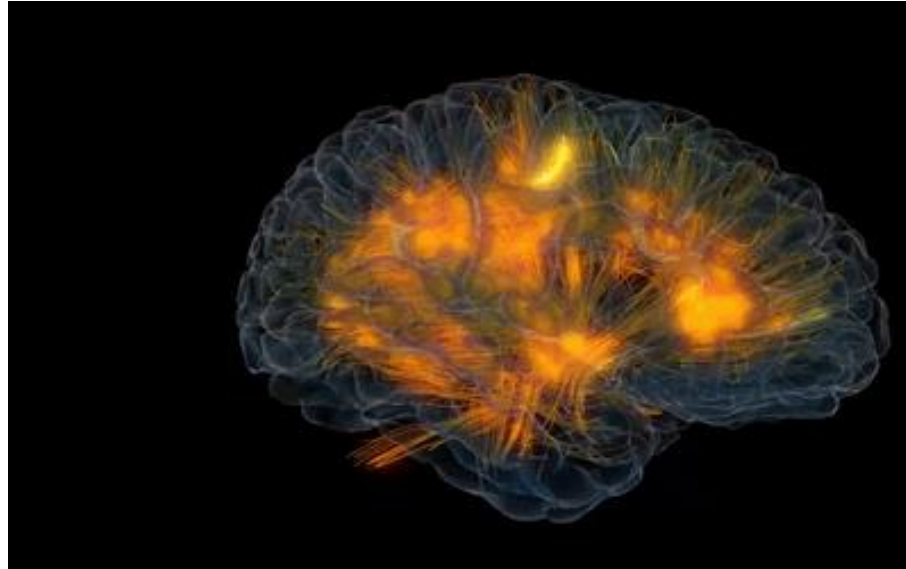


Inspired "Artificial Neural Networks"



Neurons are cells that process chemical and electrical signals and transmit these signals to neurons and other types of cells
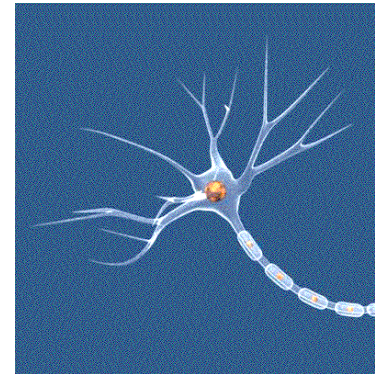
# Neuron in the brain



nucleus

dendrites
"Input wire"

Cell body

axon
"Output wire"

Can measure electrical activity (spike) of a single neuron by placing electrodes
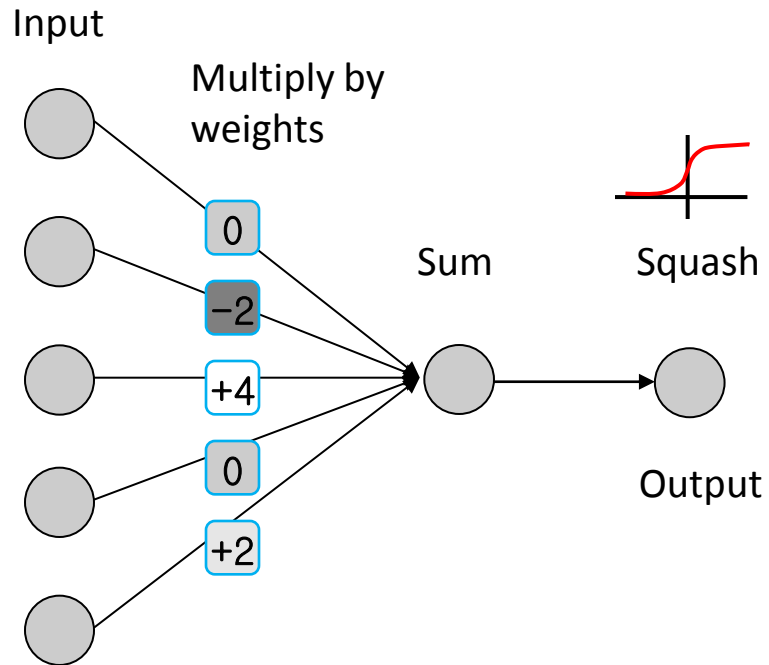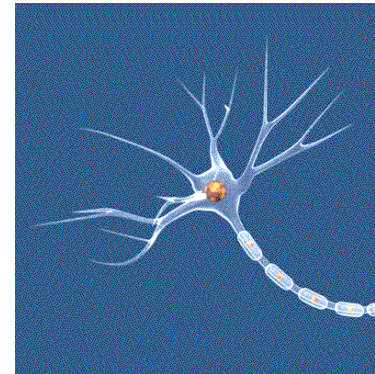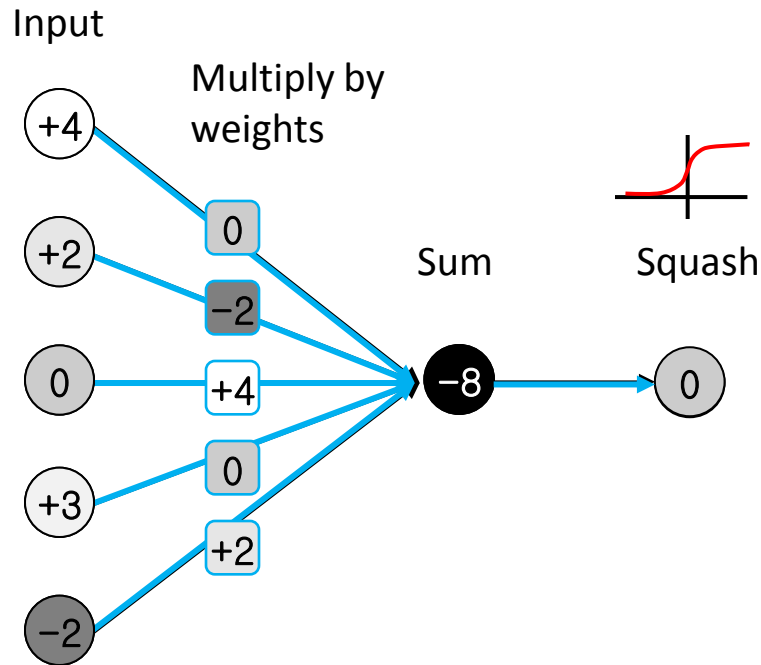
# Neural network in the brain



- **Micro networks**: several connected neurons perform sophisticated tasks: mediate reflexes, process sensory information, generate locomotion and mediate learning and memory.

- **Macro networks**: perform higher brain functions such as object recognition and cognition.

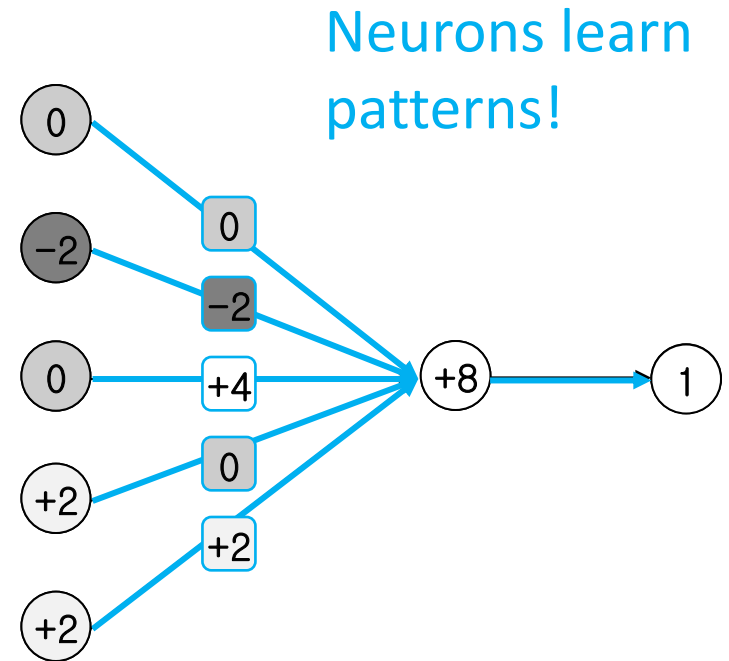# Logistic Unit as Artificial Neuron

Input

Multiply by weights
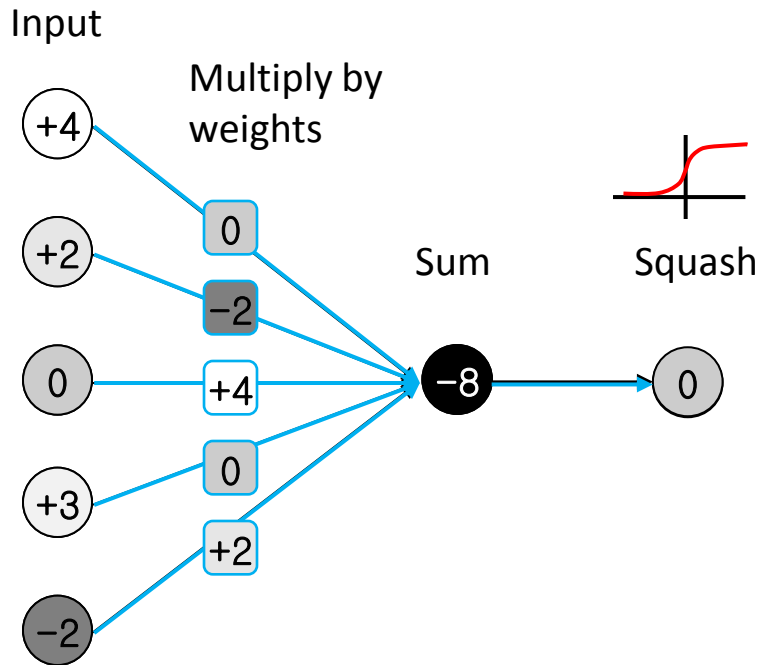
0

−2

+4

0

+2

Sum

Squash

Output

# Logistic Unit as Artificial Neuron

Input

Multiply by weights

+4

+2

0

+3

−2

0

−2

+4

0

+2

Sum

−8

Squash

0
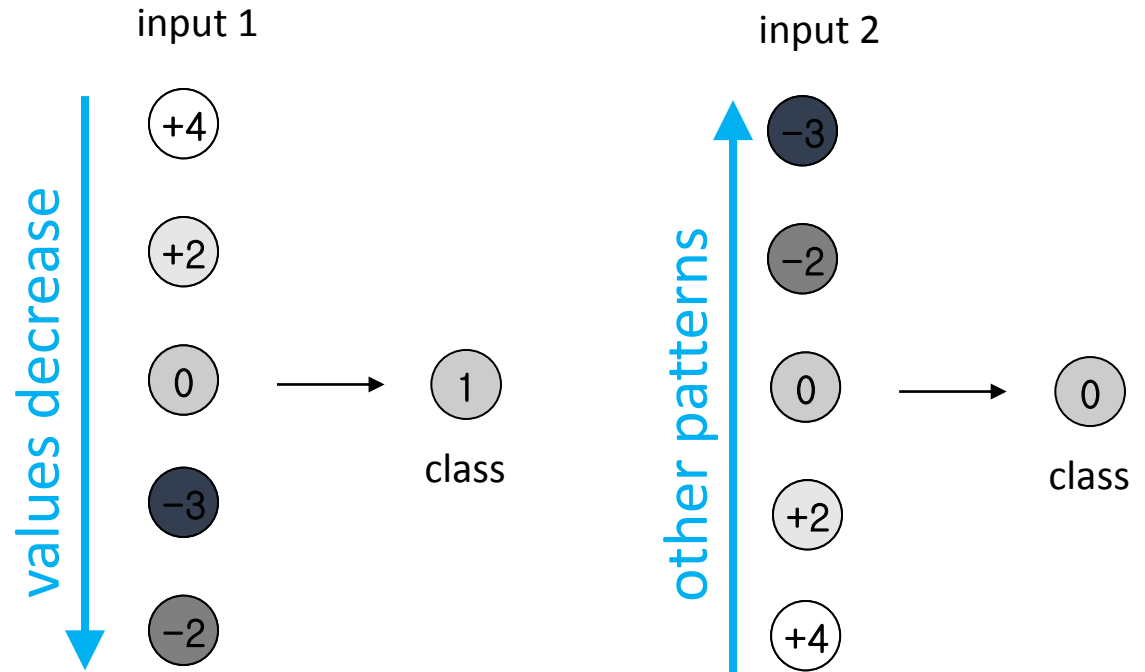
$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

# Logistic Unit as Artificial Neuron

Input

Multiply by weights

Sum

Squash

Neurons learn patterns!

# Artificial Neuron Learns Patterns

- Classify input into class 0 or 1

- Teach neuron to predict correct class label

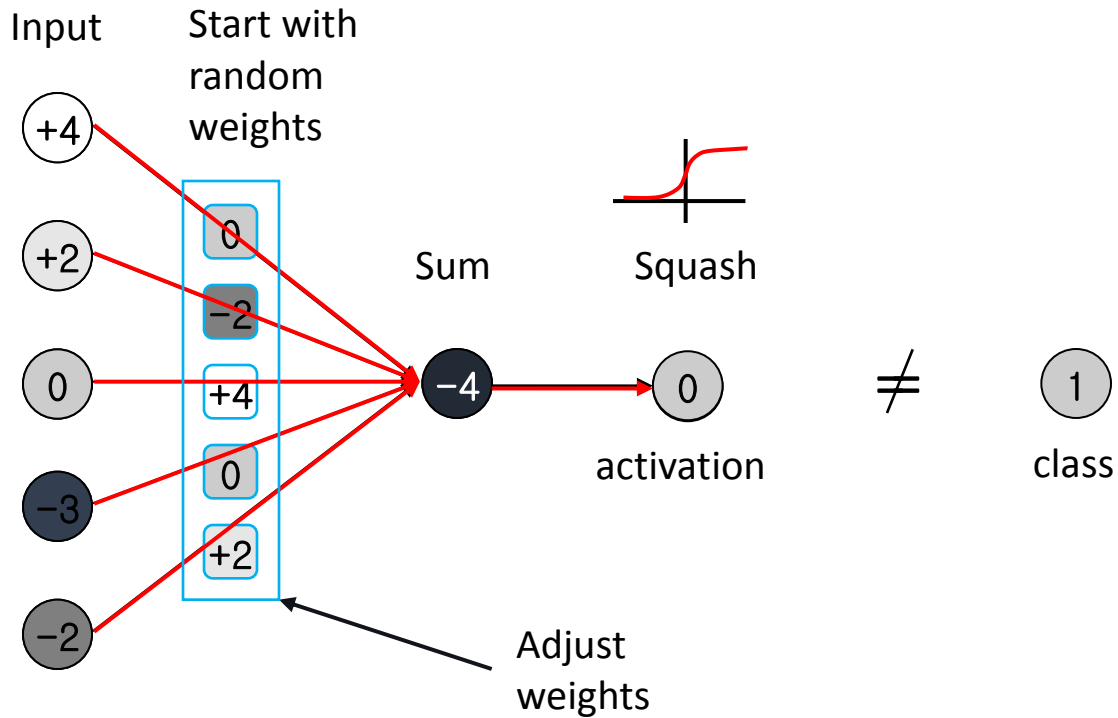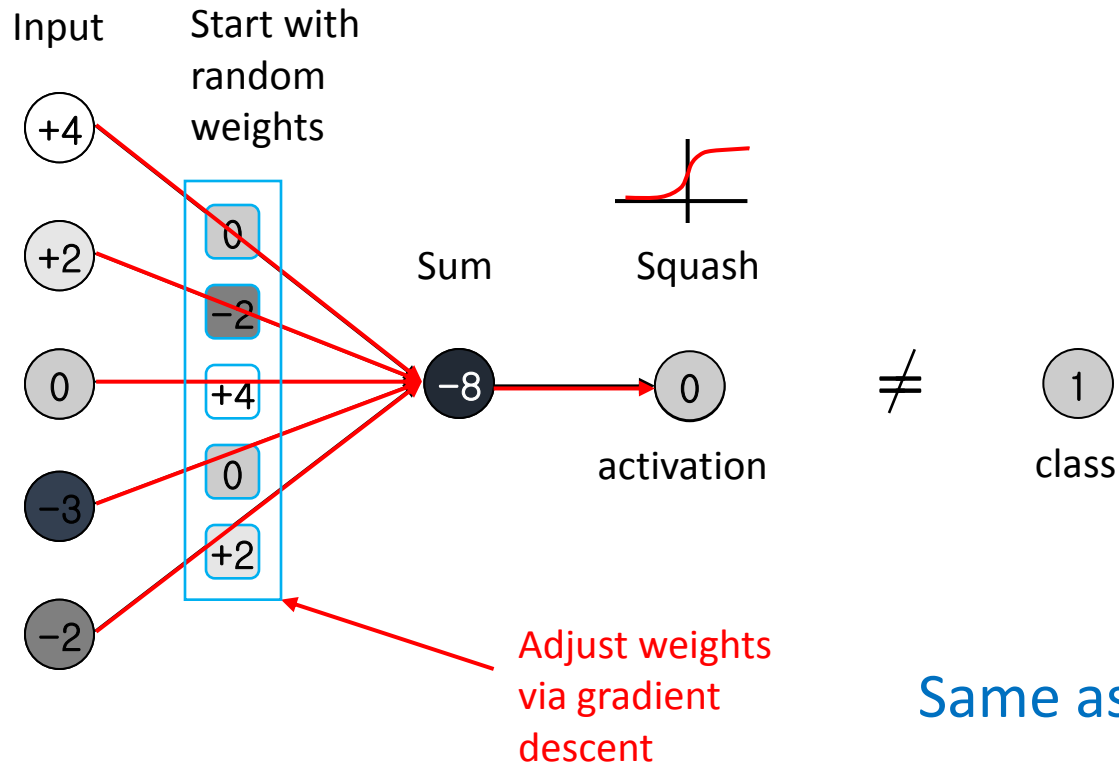- Detect presence of a simple "feature"



Example

# Neural Networks: Learning

Intuition

# Artificial Neuron: Learning

Input

Start with random weights

+4

+2

0

−3

−2

0

−2

+4

0

+2

Adjust weights

Sum

−4

Squash

0

activation

≠

1

class

# Artificial Neuron: Learning

Input

Start with random weights

+4

+2

0

−3

−2

0

−2

+4

0

+2

Sum

Squash

−8

0

activation

≠

1

class

Adjust weights via gradient descent

Same as in logistic regression

# Artificial Neuron: Learning

Input

Multiply by
weights

+4

+2

0

−3

−2

+1

0

+4

0

+1

Sum

+2

Squash

1

activation

=

1

class

Adjust
weights

# Artificial Neuron: Learning

Input

Multiply by weights

Sum

Squash

+4

+2

0

−3

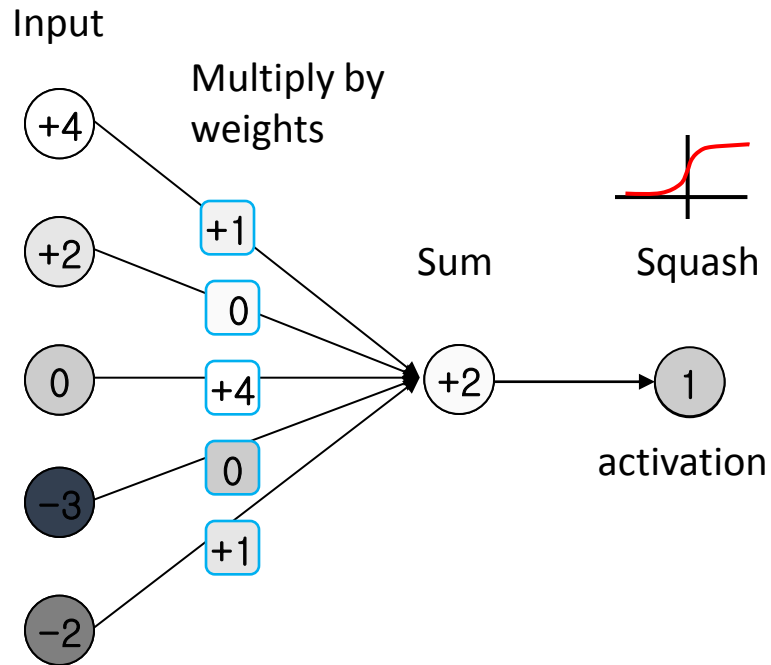−2

+1

0

+4

0

+1

+2

1

activation

= 

1

class

Adjust weights

*Forward propagation of information through a neuron*

# Neural Networks: Learning

Multi-layer network

# Artificial Neuron: simplify

Input

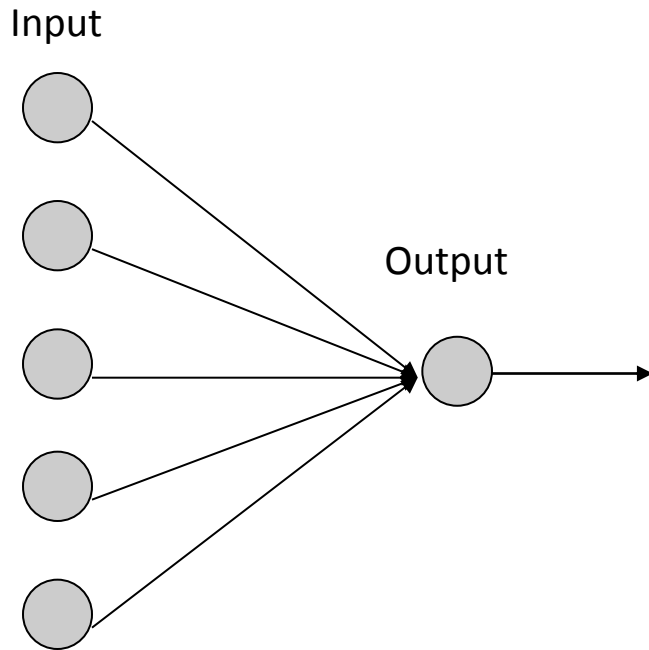Multiply by
weights

+4

+2          +1

0           0           Sum          Squash

            +4          +2    →    1

-3          0                      activation

            +1

-2

# Artificial Neuron: simplify

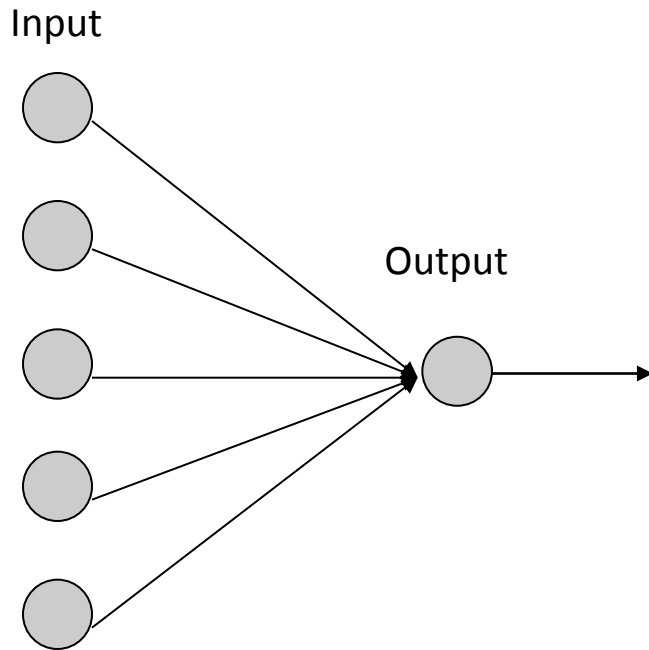Input

Weights

Output

*A single neuron is also called a perceptron*

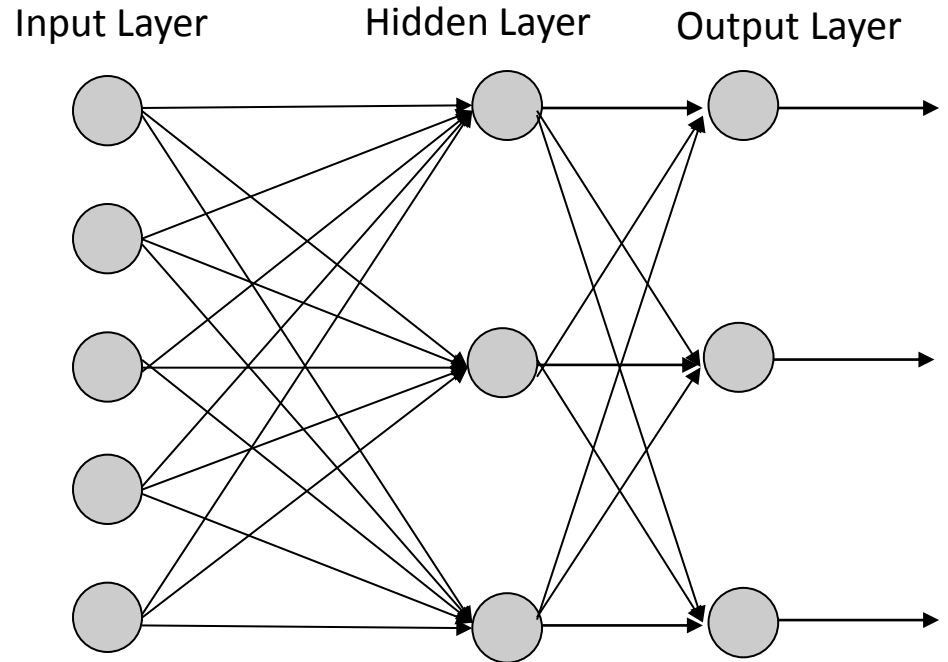# Artificial Neuron: simplify

Input

Output

# Artificial Neural Network
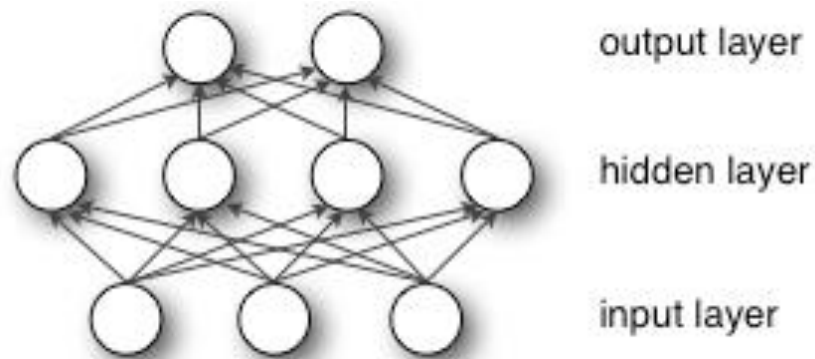


Single Neuron

Neural Network (fully connected)

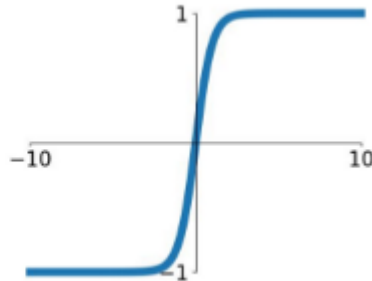Deep Network: many hidden layers

# Multi-layer perceptron (MLP)

- Just another name for a feed-forward neural network

- Logistic regression is a special case of the MLP with no hidden layer and sigmoid output.



output layer
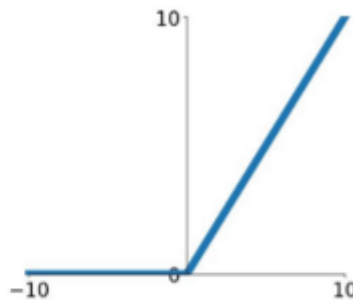
hidden layer

input layer

# Other Non-linearities

- Also called activation functions

**tanh**

$\tanh(x)$

$$tanh(x) = \frac{2}{1+e^{-2x}} - 1$$
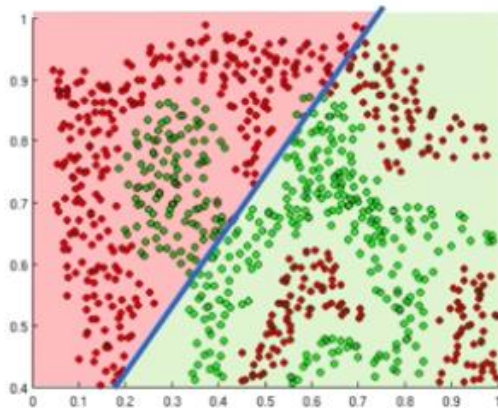
**ReLU**

$\max(0, x)$
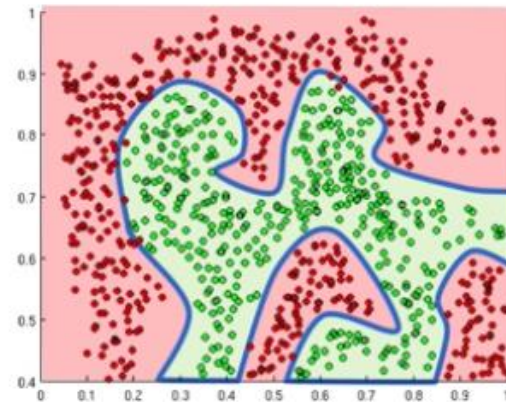
$$RELU(x) = \begin{cases} 0 \ if \ x < 0 \\ x \ if \ x >= 0 \end{cases}$$

# Importance of Non-linearities

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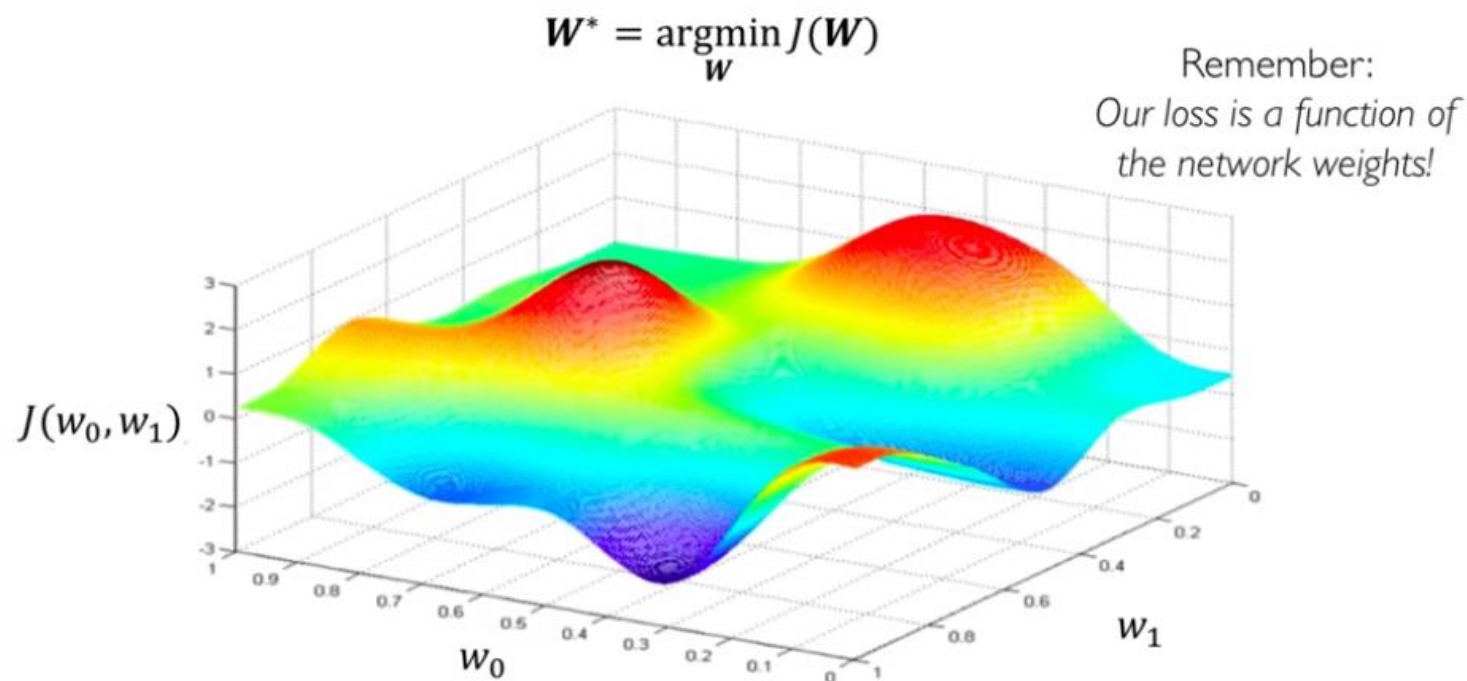


Linear activation functions produce linear
decisions no matter the network size

Non-linearities allow us to approximate
arbitrarily complex functions

# Loss Optimization

- Neural network parameters $\boldsymbol{\theta}$ are often referred to as weights $\boldsymbol{W}$.



$$W^* = \underset{W}{\mathrm{argmin}}\, J(W)$$

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

# Gradient Descent

**Algorithm**

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.      Compute gradient, $\dfrac{\partial J(W)}{\partial W}$

4.      Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

5. Return weights

# Gradient Descent

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.     Compute gradient, $\dfrac{\partial J(W)}{\partial W}$    *?*

4.     Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

5. Return weights

# Gradient Descent

**Algorithm**

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.     Compute gradient, $\dfrac{\partial J(W)}{\partial W}$     *Backpropagation*

4.     Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

5. Return weights

# Gradient Descent

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.     Compute gradient, $\dfrac{\partial J(W)}{\partial W}$

   *Not feasible to compute over all dataset*

4.     Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

5. Return weights

# Gradient Descent

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3.     Compute gradient, $\dfrac{\partial J(W)}{\partial W}$

*Compute over a mini-batch*

4.     Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$
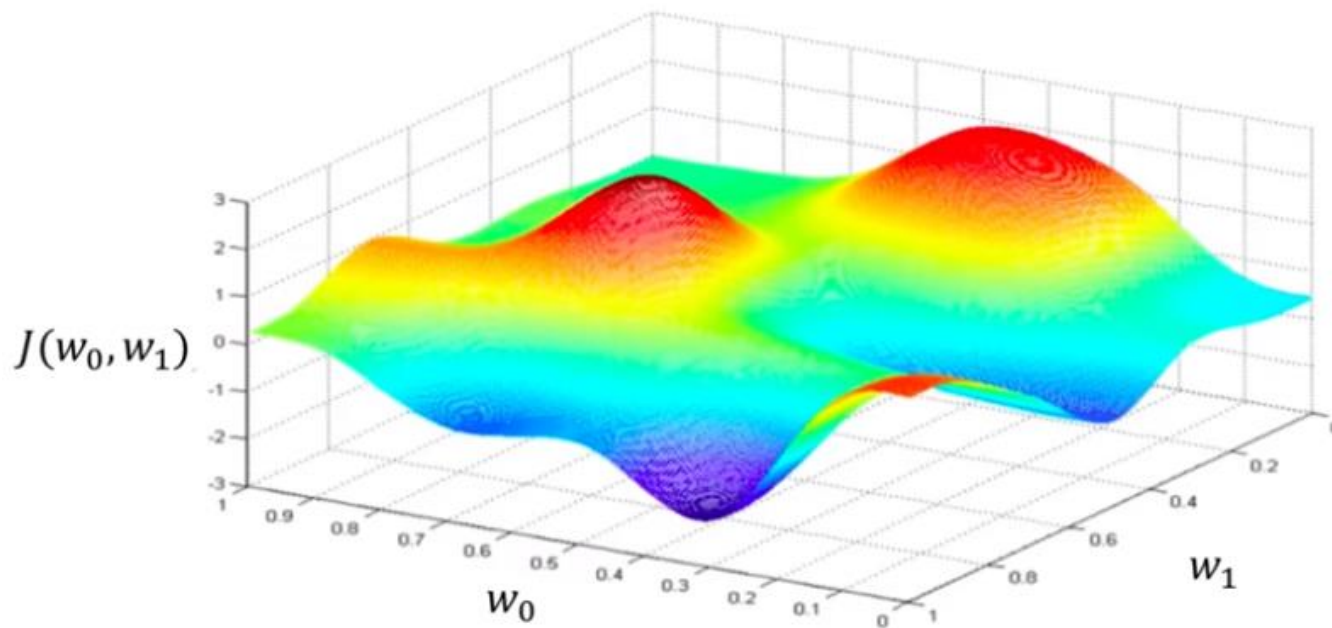
5. Return weights
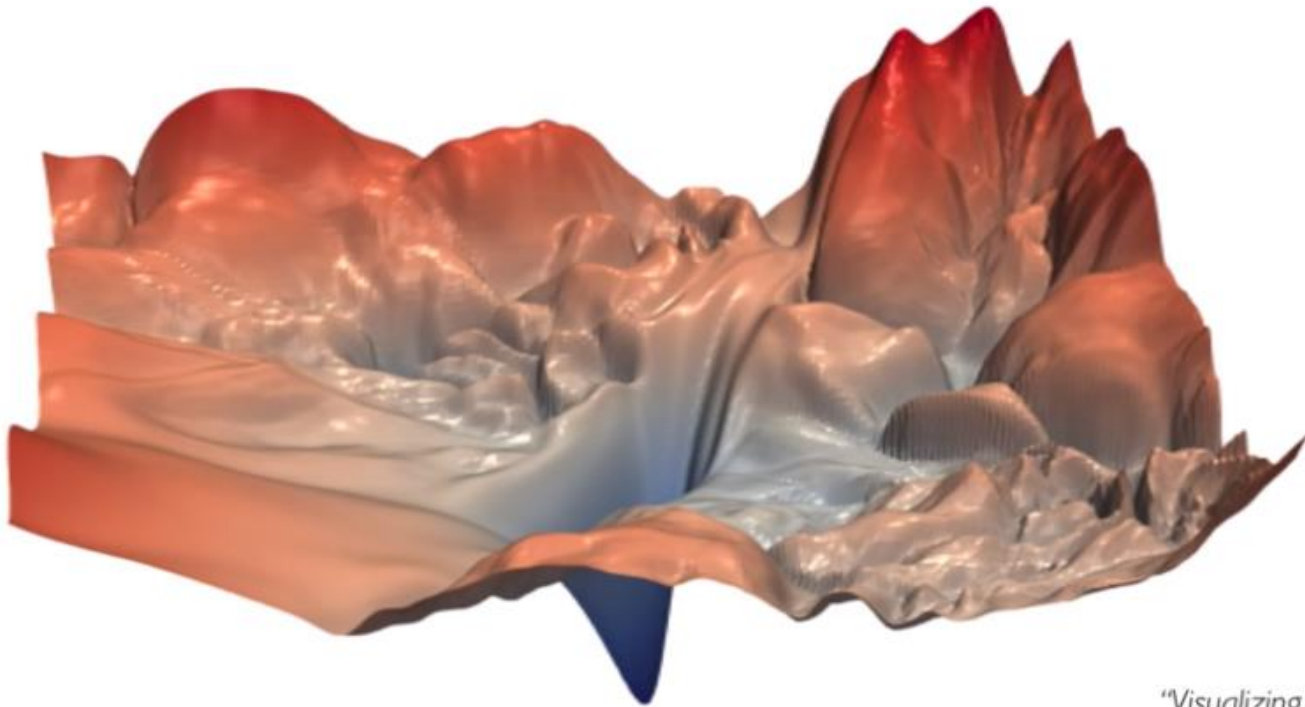
# Gradient Descent

## Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

3. Compute gradient, $\dfrac{\partial J(W)}{\partial W}$

   *Compute over a mini-batch*

4. Update weights, $W \leftarrow W - \eta \dfrac{\partial J(W)}{\partial W}$

5. Return weights

*Parallelization: Batches can be split onto multiple GPUs*

# Loss/Cost Function



$J(w_0, w_1)$

$w_0$
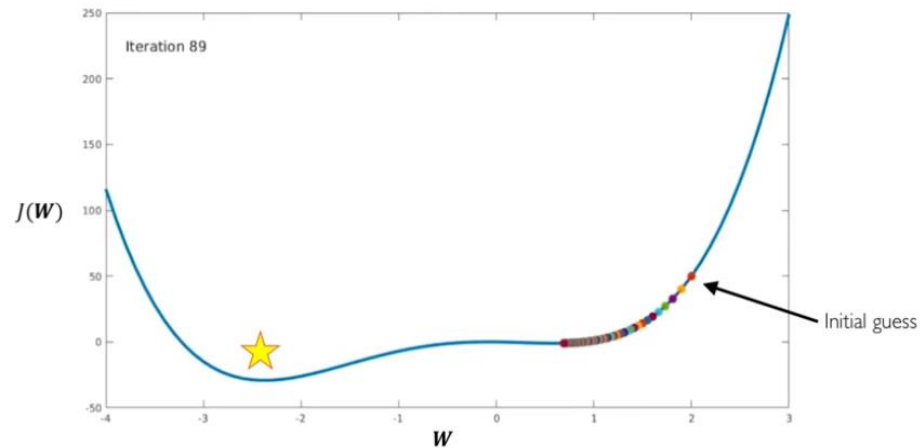
$w_1$

# Landscape Visualization



"Visualizing the loss landscape
of neural nets". Dec 2017.

# Setting the Learning Rate



**Small learning rate** converges slowly and gets stuck in false local minima

Iteration 89

Initial guess
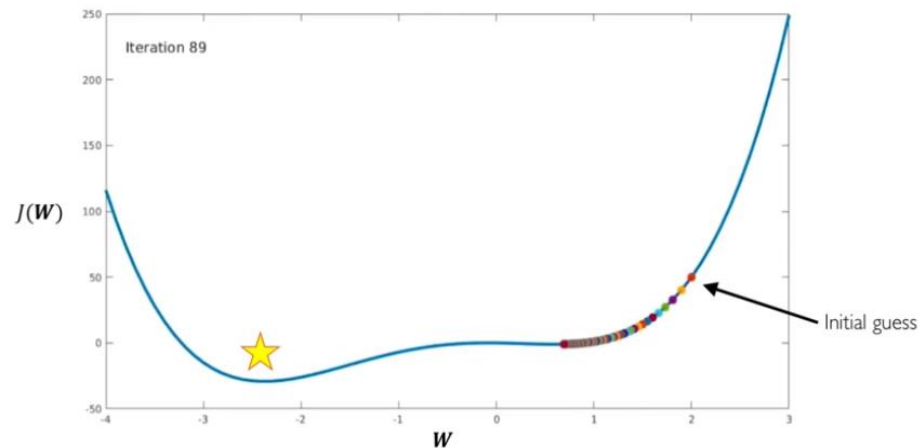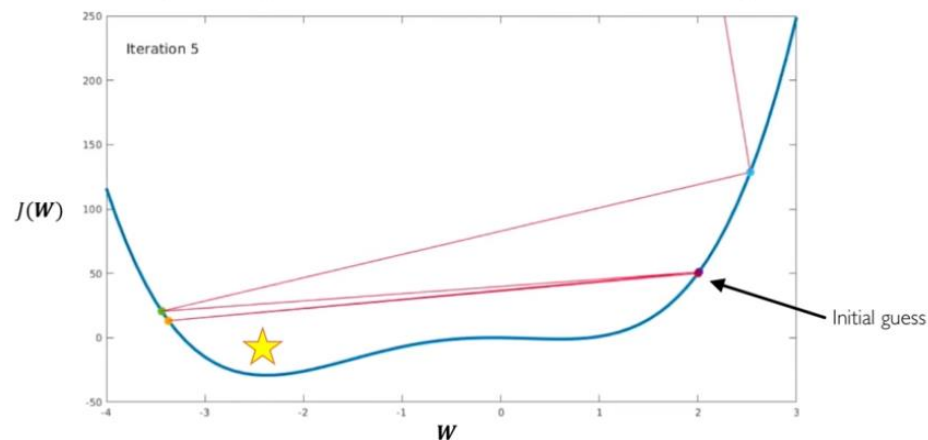
# Setting the Learning Rate

Small learning rate converges slowly and gets stuck in false local minima



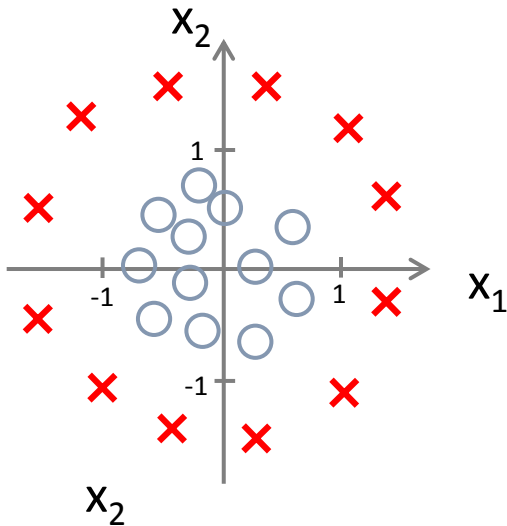Large learning rates overshoot, become unstable and diverge

# Setting the Learning Rate

- How to select the learning Rate?

  - Try several, and see which works best

  - Start with a learning rate, and change it *adaptively* as the model trains

  - Many are implemented in Neural Network Tools

# Neural Networks Learn Features

logistic regression unit == artificial neuron

chain several units together == neural network
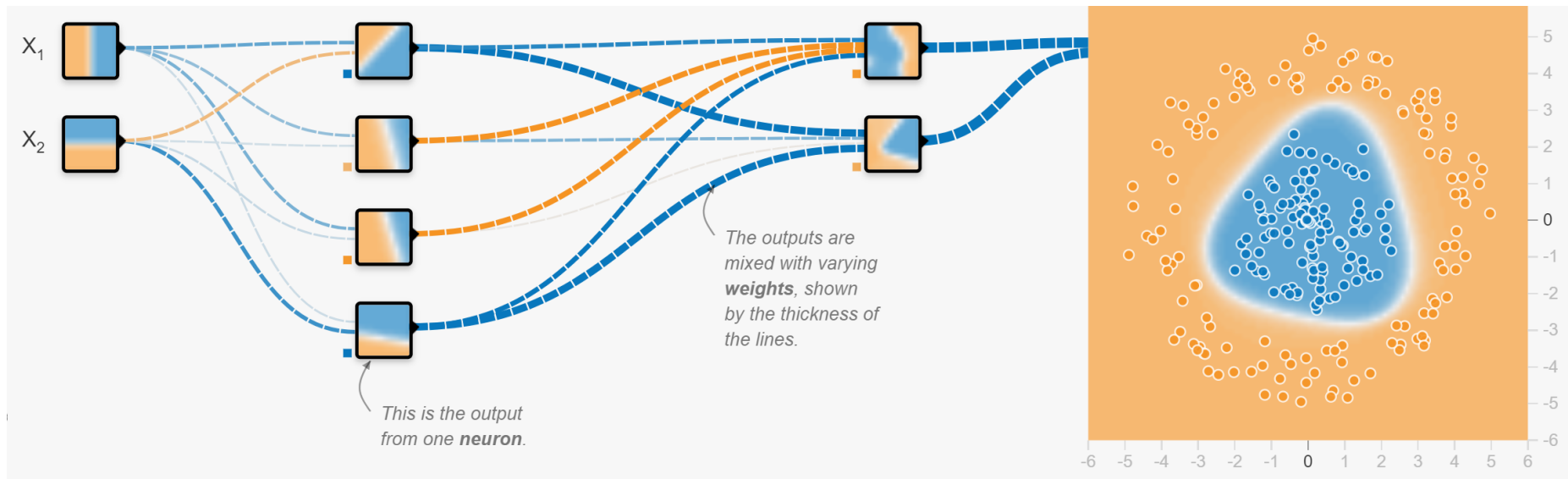
"earlier" units learn non-linear feature transformation

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

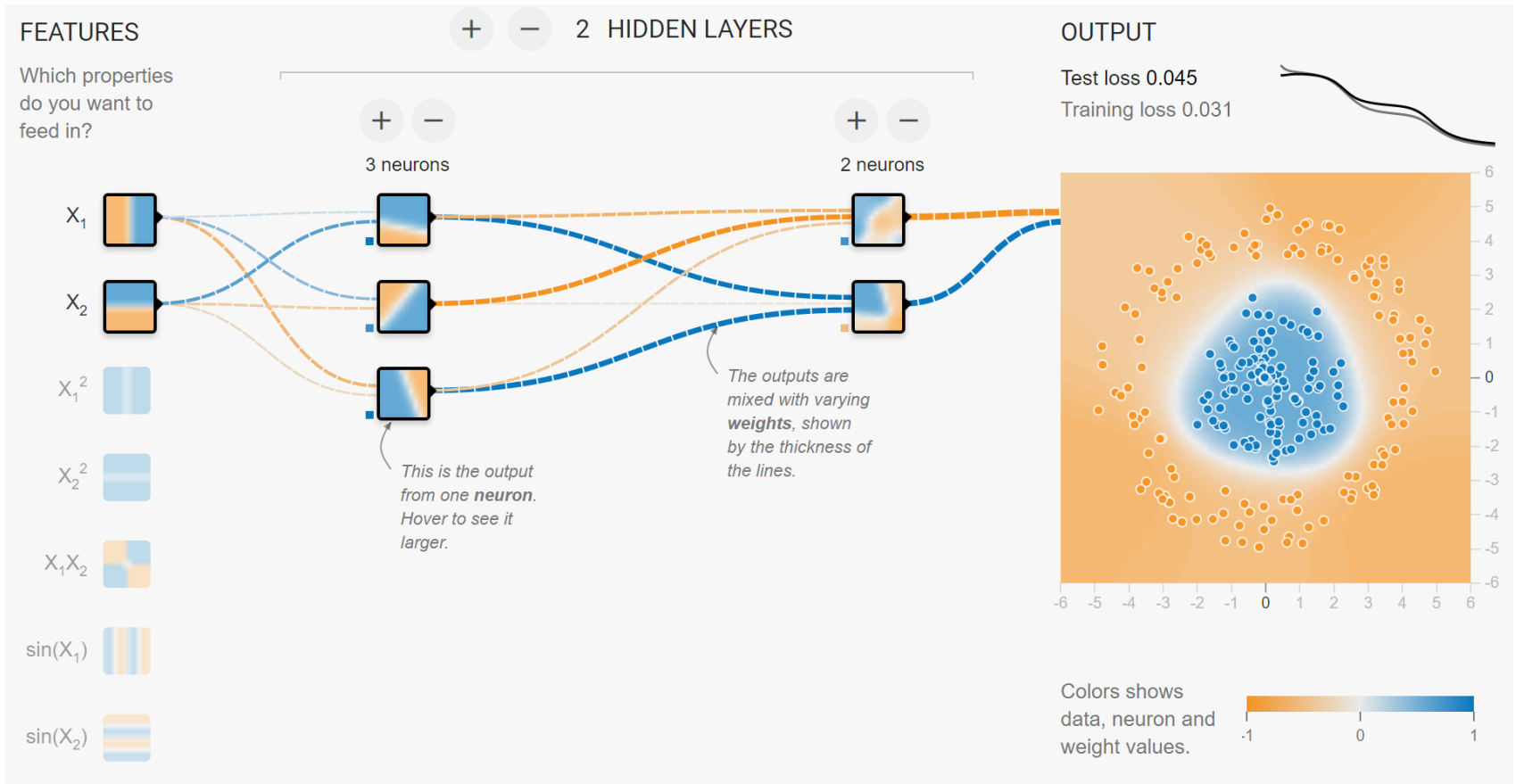simple neural network

$$h(x) = g(\theta + \theta_1 h^{(1)}(x) + \theta_2 h^{(2)}(x) + \theta_3 h^{(3)}(x))$$

# Example



The outputs are
mixed with varying
*weights*, shown
by the thickness of
the lines.

This is the output
from one *neuron*.

# Training a neural net: Demo



[Tensorflow playground](Tensorflow playground)

# Deep Learning

Architectures

# What is Deep Learning?
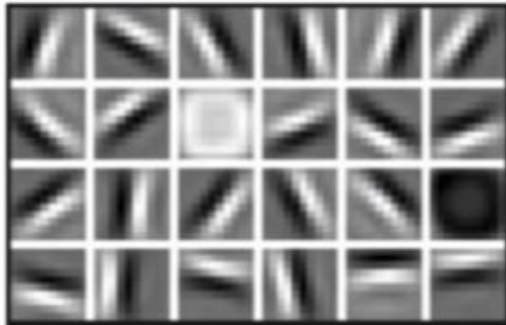
# Why Deep Learning?

Hand engineered features are time consuming, brittle, and not scalable in practice

Can we learn the **underlying features** directly from data?

| Low Level Features | Mid Level Features | High Level Features |
|---|---|---|



Lines & Edges          Eyes & Nose & Ears          Facial Structure
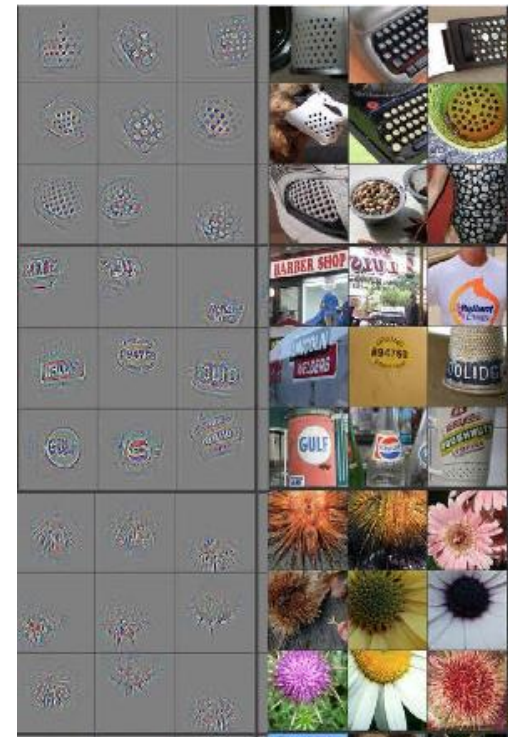
# Why Deep Learning?
## The Unreasonable Effectiveness of Deep Features
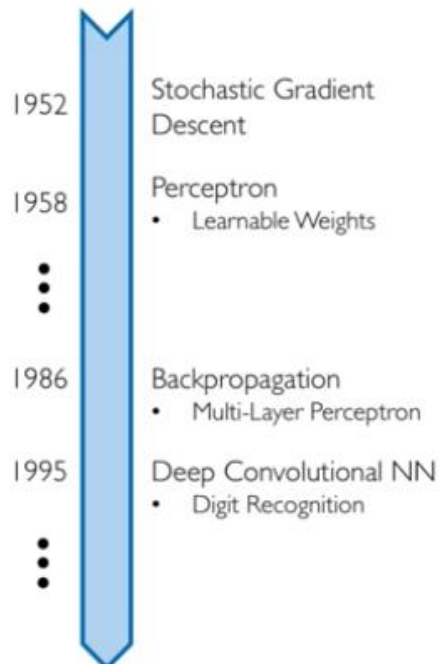


Maximal activations of pool$_5$ units   [R-CNN]

Rich visual structure of features deep in hierarchy.



conv$_5$ DeConv visualization
[Zeiler-Fergus]

# Why Now?

Neural Networks date back decades, so why the resurgence?

| | | |
|---|---|---|
| **1. Big Data** | **2. Hardware** | **3. Software** |

1952 — Stochastic Gradient Descent

1958 — Perceptron
- Learnable Weights

⋮

1986 — Backpropagation
- Multi-Layer Perceptron

1995 — Deep Convolutional NN
- Digit Recognition

⋮

**1. Big Data**
- Larger Datasets
- Easier Collection & Storage

IMAGENET

WIKIPEDIA
The Free Encyclopedia

**2. Hardware**
- Graphics Processing Units (GPUs)
- Massively Parallelizable

**3. Software**
- Improved Techniques
- New Models
- Toolboxes

TensorFlow