



CS5120 VLSI System Design, Spring 2025

Deep Neural Networks (DNNs)

Part 1

黃稚存

Chih-Tsun Huang

cthuang@cs.nthu.edu.tw



國立清華大學
NATIONAL TSING HUA UNIVERSITY

資訊工程學系
Computer Science

Lecture 02



聲明

- ◎ 本課程之內容 (包括但不限於教材、影片、圖片、檔案資料等)，僅供修課學生個人合理使用，非經授課教師同意，不得以任何形式轉載、重製、散布、公開播送、出版或發行本影片內容 (例如將課程內容放置公開平台上，如 Facebook, Instagram, YouTube, Twitter, Google Drive, Dropbox 等等)。如有侵權行為，需自負法律責任。



Outline

- ⊙ Machine learning basics
- ⊙ Deep learning basics
- ⊙ On training deep neural networks
 - ◆ Review of DNN
 - ◆ Empirical loss functions
 - ◆ Backpropagation
 - ◆ Learning rate
 - ◆ Gradient decent
 - ◆ Overfitting and regularization
- ⊙ Tinker with a Neural Network
- ⊙ Deep learning frameworks
- ⊙ Training vs. inference



Machine Learning Basics

- » Machine Learning: Infrastructure for Everything
- » Reference: <https://www.coursera.org/learn/machine-learning>
Prof. Andrew Ng, Stanford

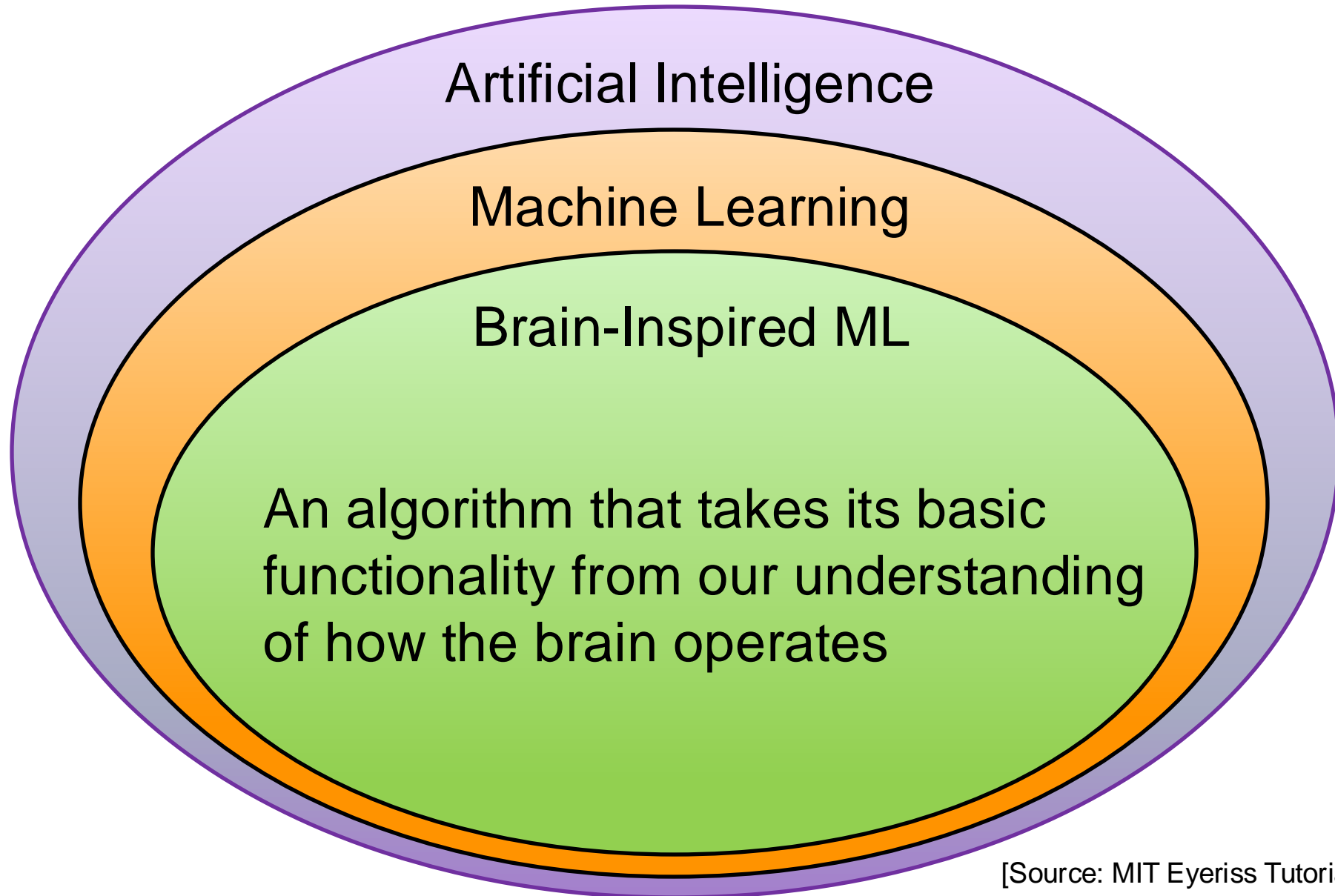


What is Machine Learning (ML)

- ⊙ “The field of study that gives computers the ability to learn without being explicitly programmed.”
Arthur Samuel (1959)
- ⊙ Artificial Intelligence (AI)?
- ⊙ Deep Learning (DL)
- ⊙ Convolutional Neural Network (CNN)?
- ⊙ Deep Neural Network (DNN)?
- ⊙ Neuromorphic Computing (NC)?



General Definitions





General Definitions (Cont)

⊙ Artificial Intelligence

- ◆ “The science and engineering of creating intelligent machines” - John McCarthy, 1956

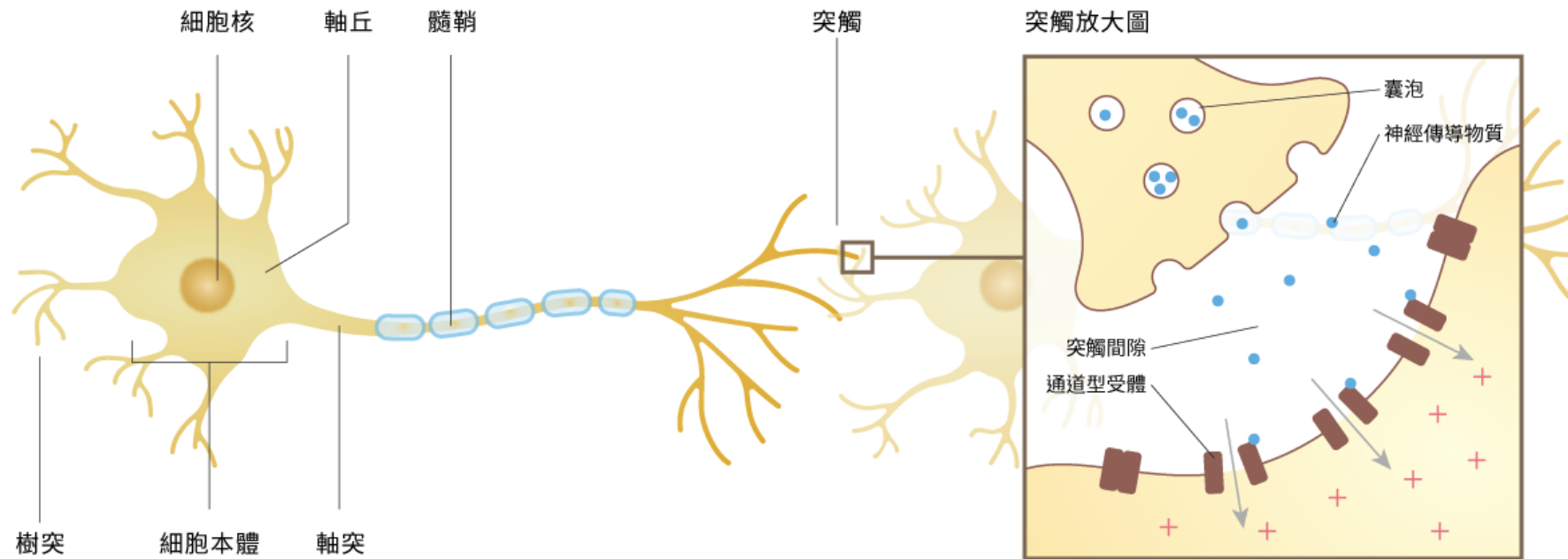
⊙ Machine Learning (ML)

- ◆ “Field of study that gives computers the ability to learn without being explicitly programmed” - Arthur Samuel, 1959

⊙ Brain-Inspired ML

- ◆ An algorithm that takes its basic functionality from our understanding of how the brain operates

How Does the Brain Work?

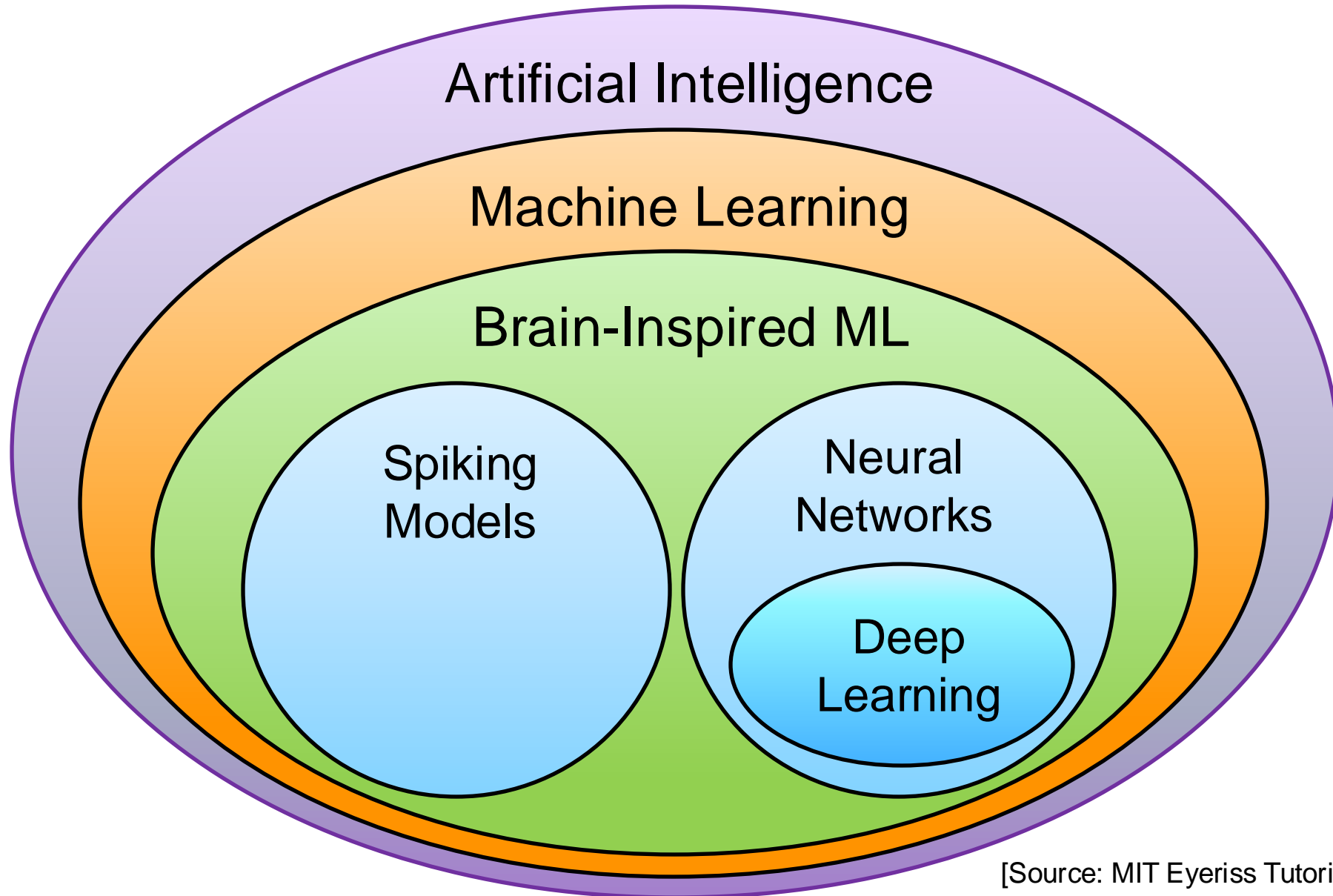


<<https://www.narlabs.org.tw/>>

- **Neuron**: basic computational unit of the brain
 - ◆ 86B neurons in the brain
 - ◆ Connected with $10^{14} - 10^{15}$ **synapses (突觸)**
- Neurons receive input signal from **dendrites (樹突)** and produce output signal along **axon (軸突)**, which interact with the dendrites of other neurons via **synaptic weights**
 - ◆ Synaptic weights: learnable & control influence strength



General Definitions





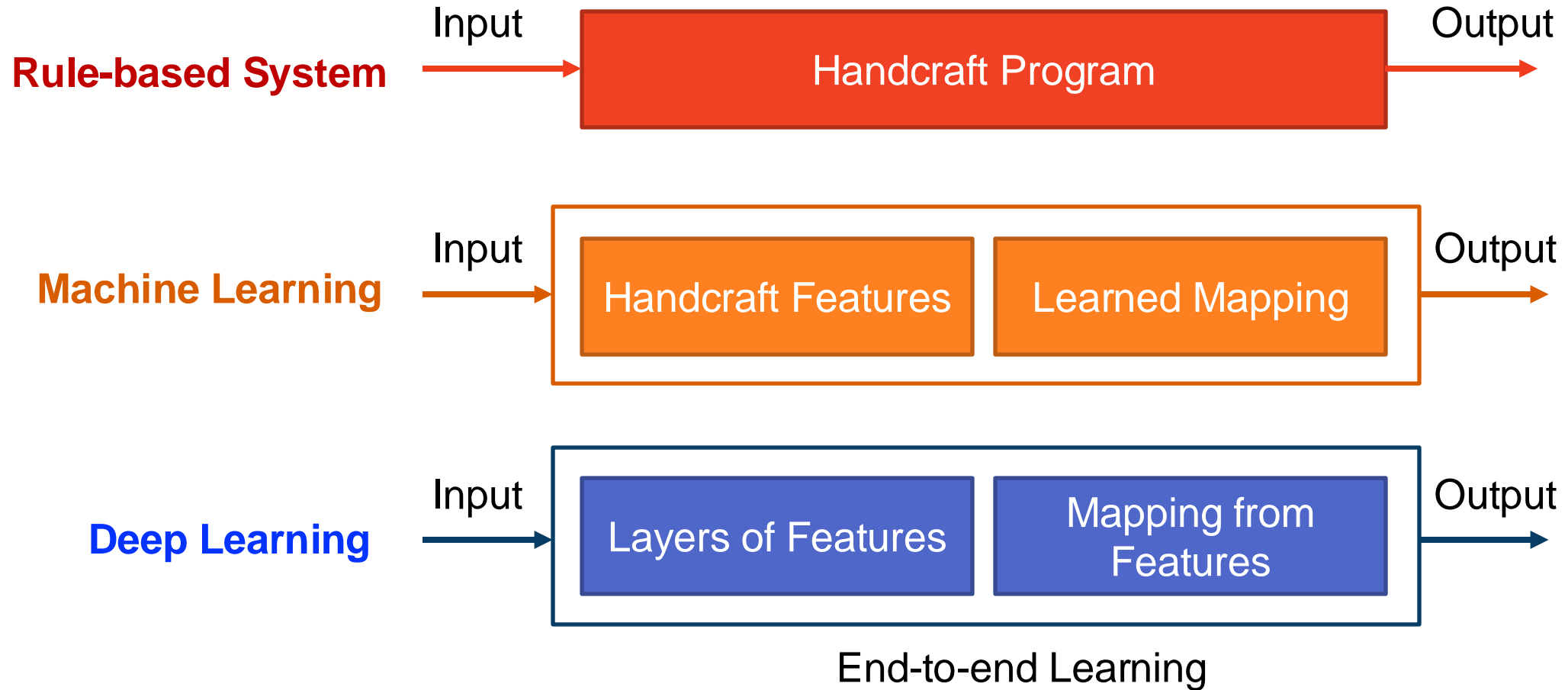
Deep Learning (DL)

◎ “Seek to exploit the unknown structure in the input distribution in order to discover good representations, often at multiple levels.”

- Yoshua Bengio, 2012



Variations of Artificial Intelligence





Is The Winter Over?

◎ Two major **AI Winters** (1974-80, 1987-93) and some smaller periods

- ◆ https://en.wikipedia.org/wiki/AI_winter

◎ It is!

- ◆ AI has surpassed human performance

- Chess (1997)
- Trivia (2011)
- Atari games (2013)
- Image recognition (2015)
- Speech recognition (2015)
- Go (2016)

◎ Maybe?!

- ◆ We need explainable AI (or general AI)



©. Artistic illustration of nuclear winter. (Elena Naylor/Shutterstock)



Machine Learning Algorithms

- ⊙ A computer program is said to **learn** from **experience (E)** with respect to some **task (T)** and some **performance measure (P)**, if its performance on T , as measured by P , improves with experience E .
 - ◆ Tom Mitchell, 1998
- ⊙ E.g., spam classification
 - ◆ Task (T): Predict emails as spam or not spam
 - ◆ Experience (E): Observe users label emails as spam or not
 - ◆ Performance (P): # of emails that are correctly predicted



Machine Learning Algorithms

- ⊙ Suppose a tumor classification program watches which tumor is being marked as “benign” or “malignant” by doctors, and it learns how to better predict benign/malignant tumors.
- ⊙ What is the task T in this setting?
 - A. The number (or fraction) of tumors correctly classified as benign/malignant
 - B. Watching doctors label tumors as benign or malignant
 - C. Classifying tumors as benign or malignant
 - D. Non of the above. This is not a machine learning problem.



Building an ML Algorithm

- ◎ Nearly all ML algorithms can be described as particular instances of a simple recipe:
 - ◆ Experience (E) \rightarrow Dataset with ground truth (labels)
 - ◆ Performance Measure (P) \rightarrow Cost (loss) function
 - ◆ Task (T) \rightarrow Model + optimization method
- ◎ Use this recipe to see the different algorithms:
 - ◆ As part of a taxonomy of methods for doing related tasks that work for similar reasons
 - ◆ Rather than as a long list of algorithms that each have separate justifications

Example: Linear Regression

Dataset:

- ◆ m training examples
- ◆ $(x, y) = (\text{size}, \text{price})$

Model:

- ◆ $h(x) = w_0 + w_1 * x$

Cost function:

- ◆ Measures the performance of an ML model for given data

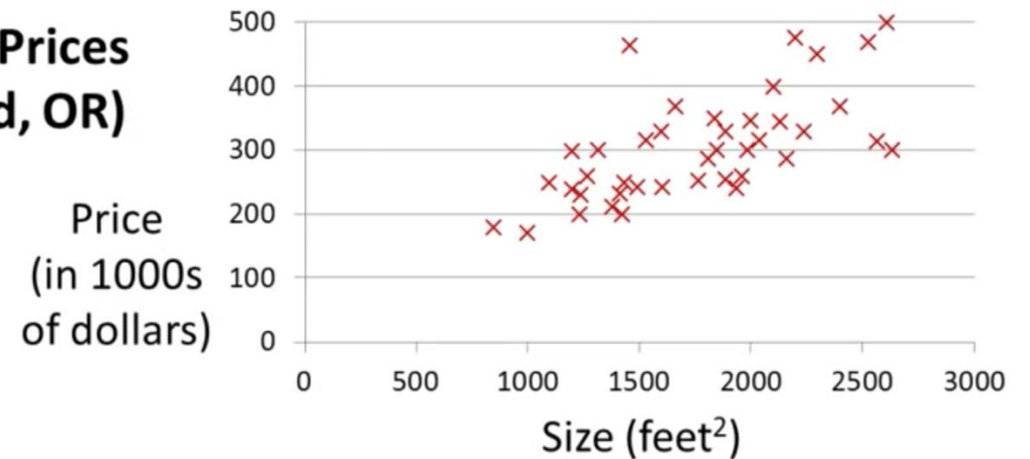
- ◆ Mean Squared Error:

$$MSE = \frac{1}{m} \sum_{i=0}^m (h(x_i) - y_i)^2$$

Optimization method:

- ◆ Solve for where its gradient is 0.
- ◆ Gradient descent

**Housing Prices
(Portland, OR)**



Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

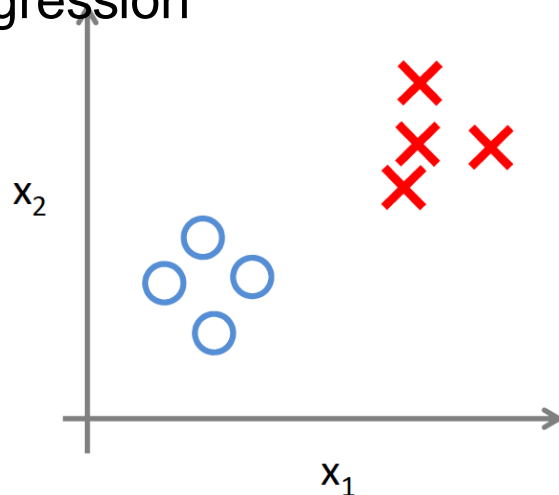
[Source: Machine Learning, Andrew Ng]



Supervised vs. Unsupervised Learning

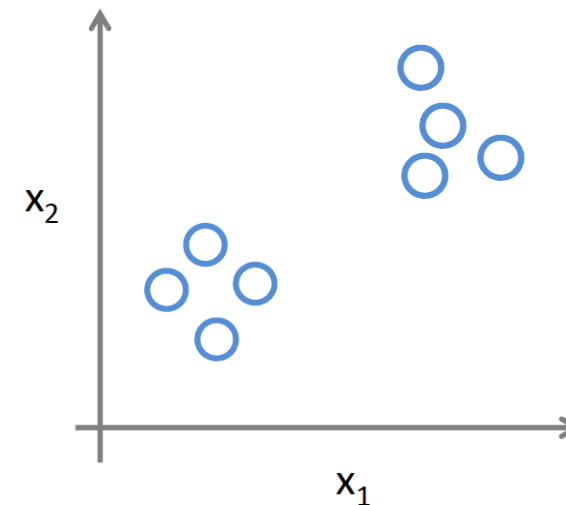
Supervised Learning

- Given right answers to learning
 - Dataset: a collection of many examples with **labels** (or **ground truth**)
- Predict proper output value
- Applications
 - Classification
 - Regression



Unsupervised Learning

- No right answers with the data
 - Without label (ground truth)
- Applications
 - Clustering
 - Anomaly detection



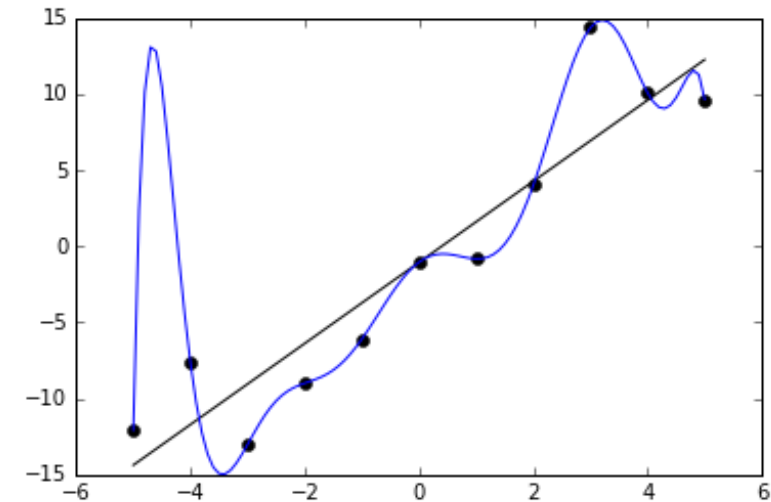
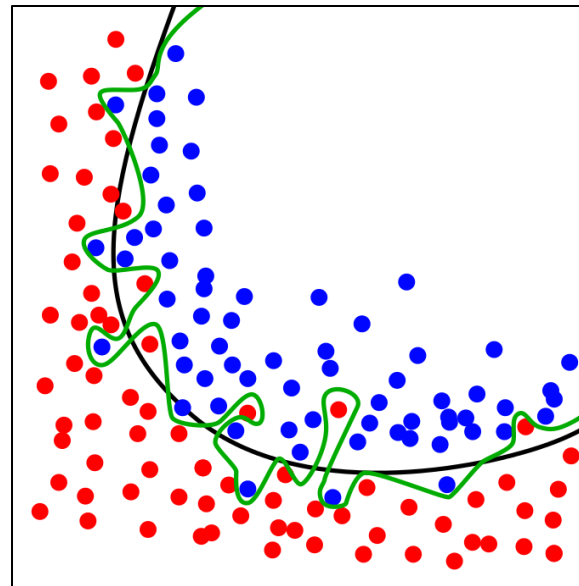
Subtlety of Supervised Learning

Hyper parameters

- ◆ Model type
- ◆ Architecture
- ◆ Training parameters (learning rate, regularization, etc.)
- ◆ Model parameters

Data selection, data size

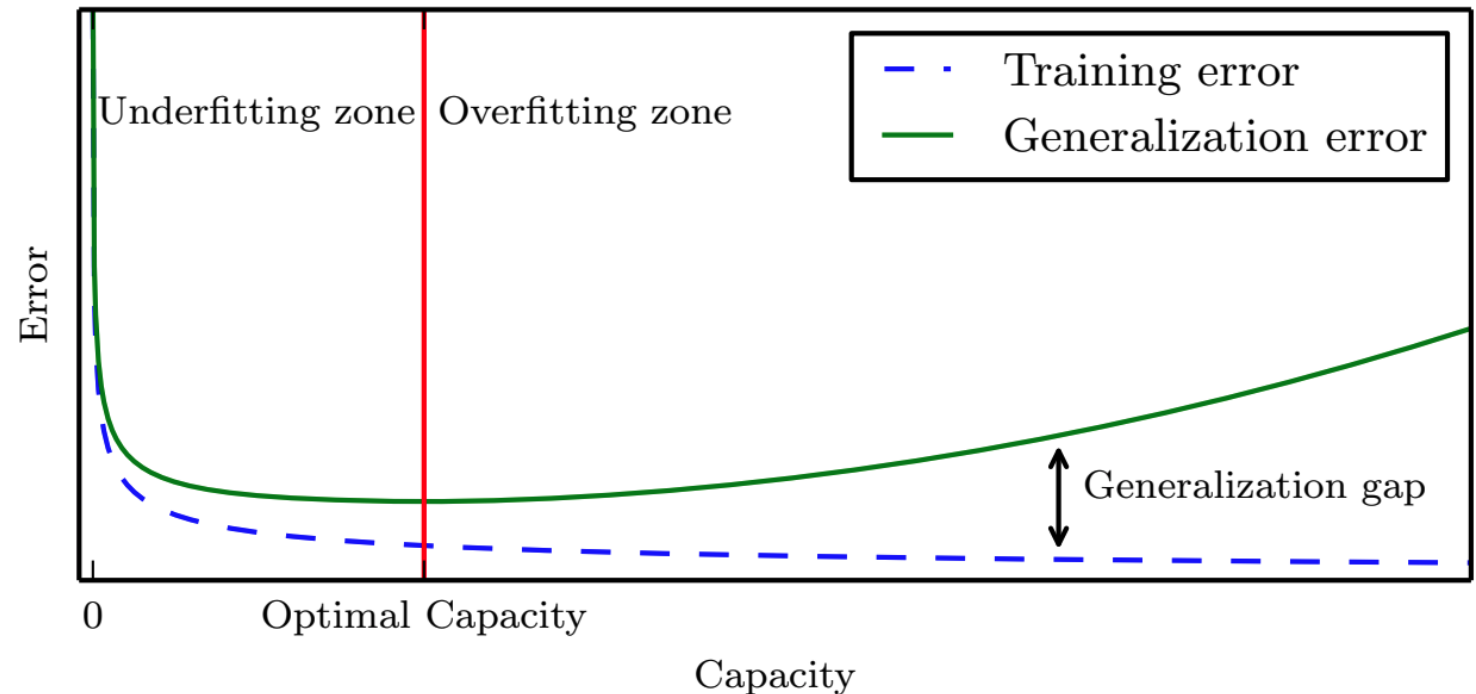
Overfitting/underfitting



<https://en.wikipedia.org/wiki/Overfitting>

Generalization of Supervised Learning

- Training with dataset
- Generalization
 - Algorithm performs well on new, unseen inputs
 - Instead of those in the dataset
- Partition the dataset into
 - Training set: where the model was trained
 - Test set: where the trained model was tested
- Overfitting
 - Large gap between training and test errors



Use Regularization to Prevent Overfitting

⊙ Regularization: modifications to reduce the generalization error but not the training errors

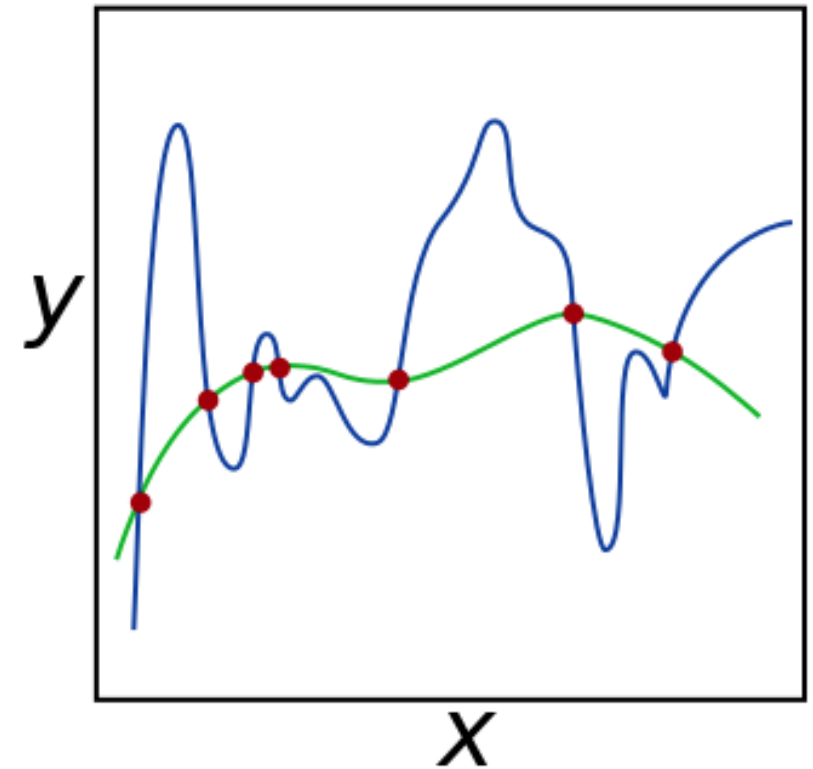
- ◆ Weight decay (L2/L1 regularization)

- Expressing preferences of smaller weights

- $CF = MSE_{train} + \lambda \sum_i w_i^2$ (L2)
(CF: cost function)

- $CF = MSE_{train} + \lambda \sum_i |w_i|$ (L1)

↑
Hyperparameter

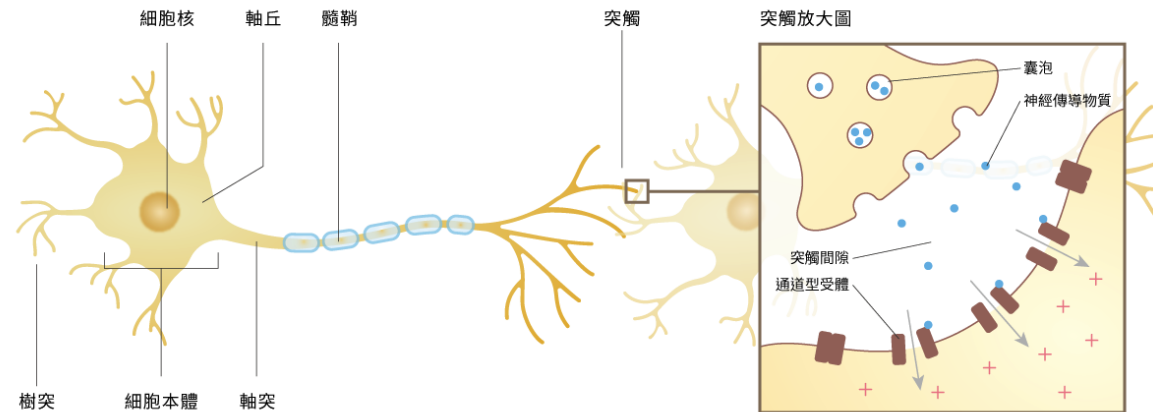
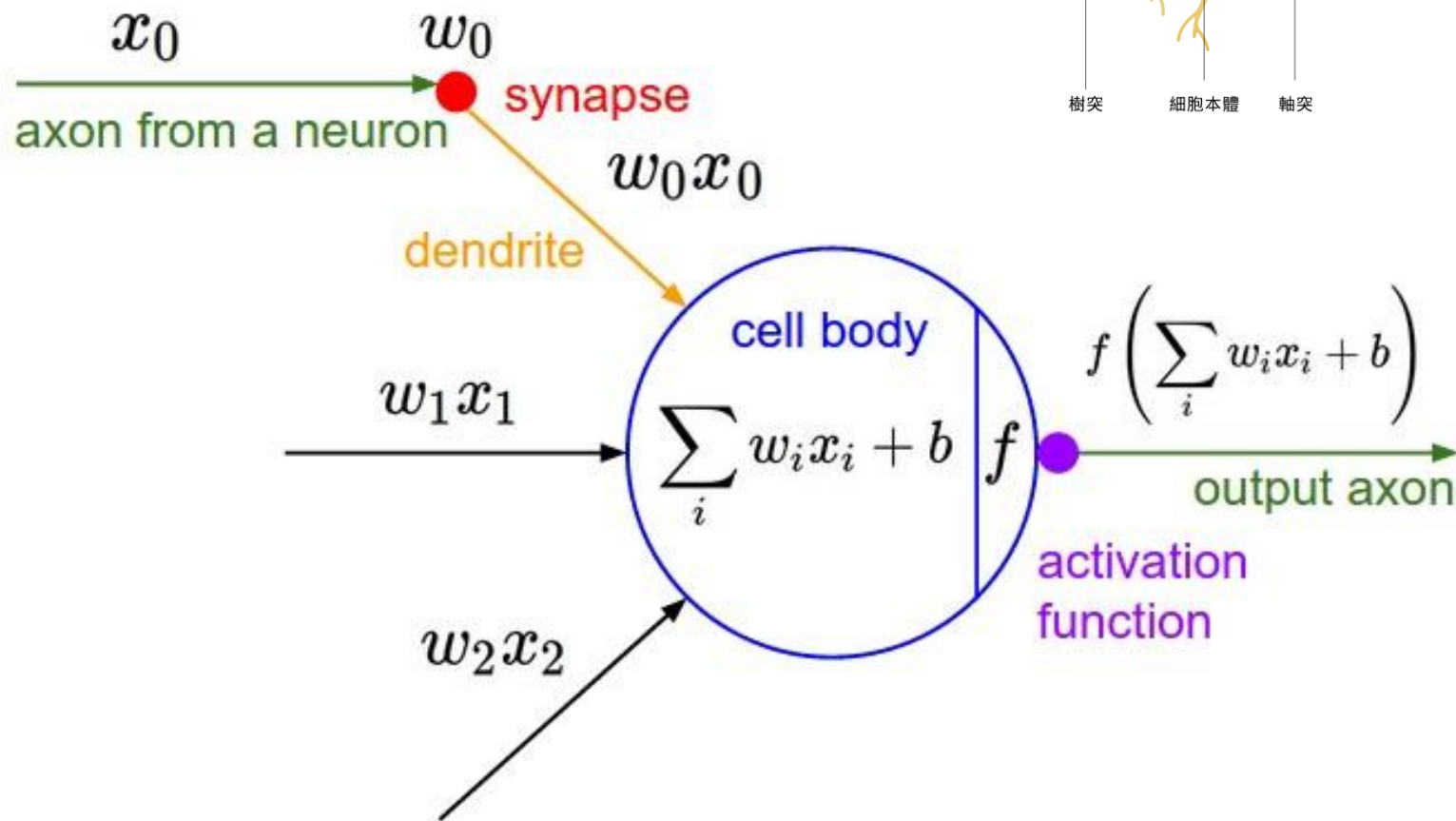




Deep Learning Basics



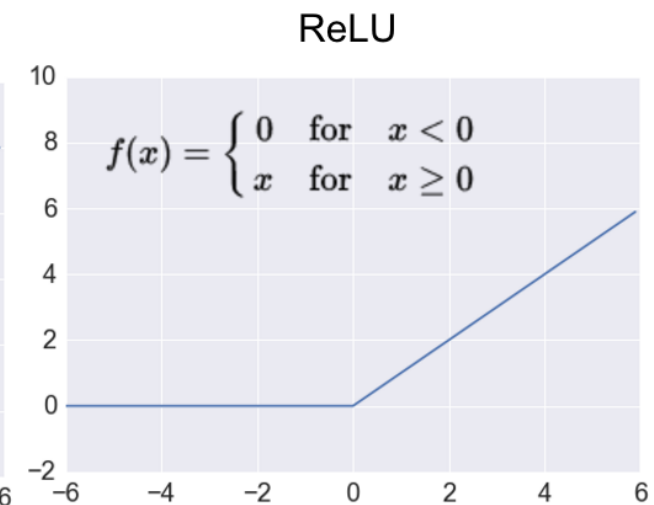
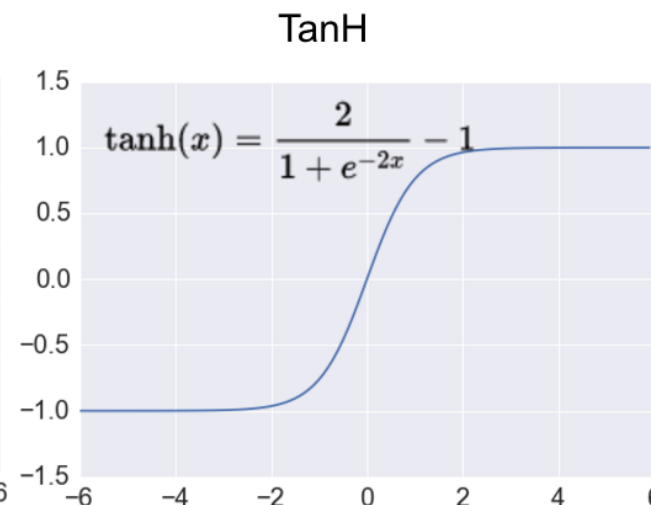
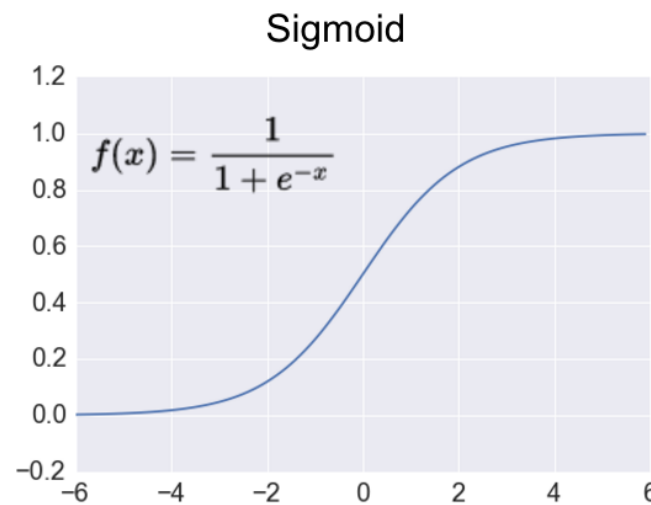
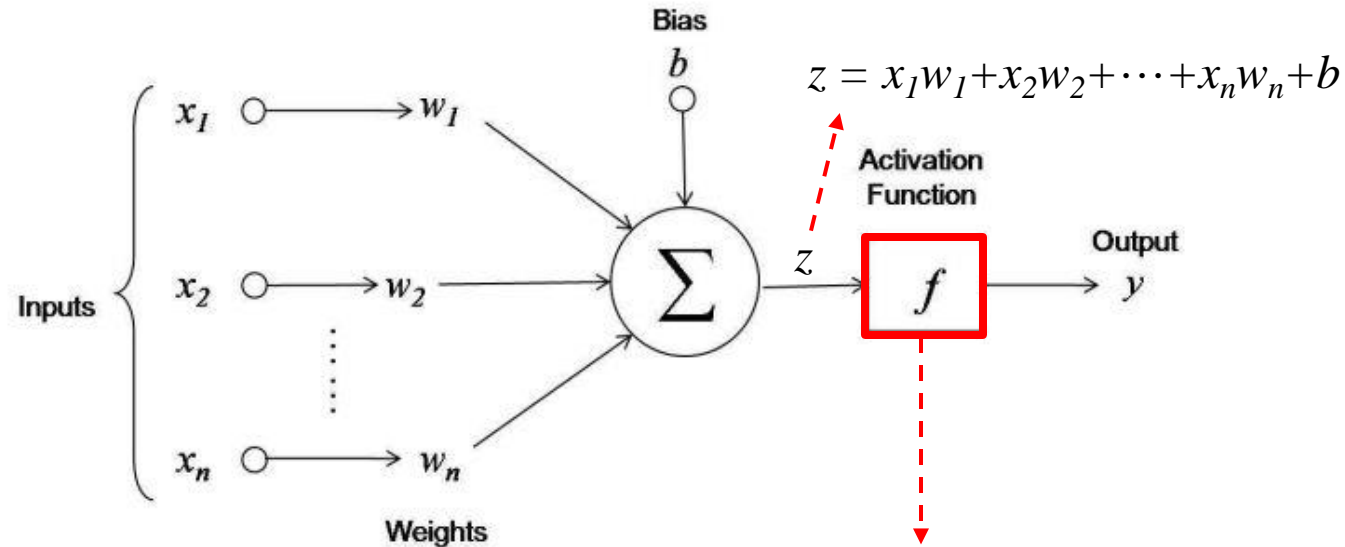
Neural Networks: Weighted Sum



<<https://www.narlabs.org.tw/>>

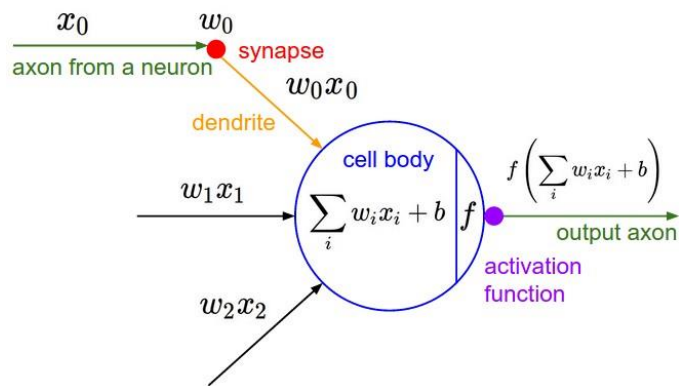
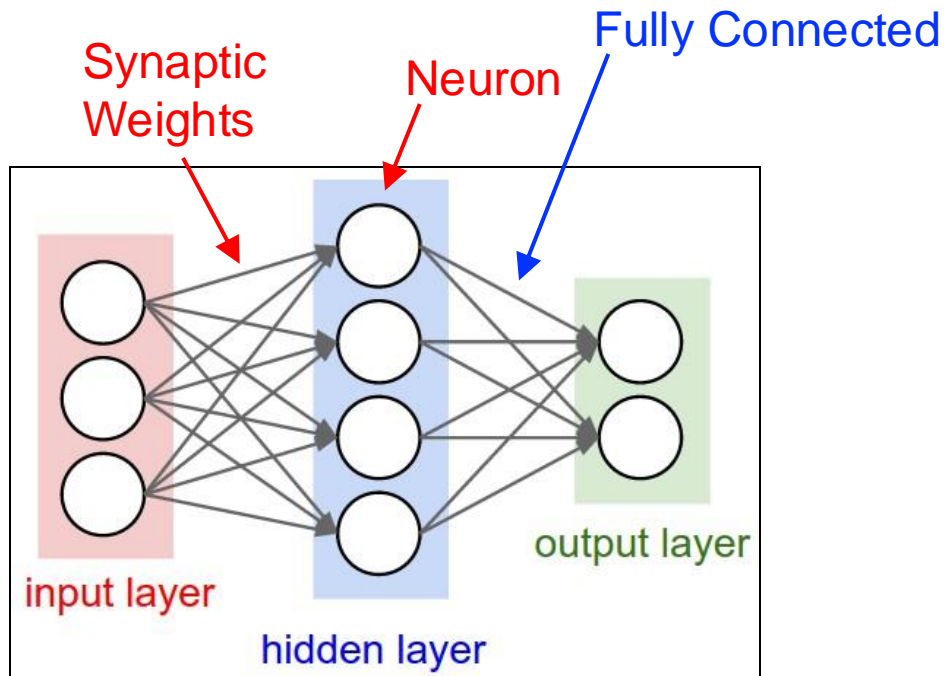
[Source: Stanford CS231n]

Neural Networks: Non-Linear Activation Function

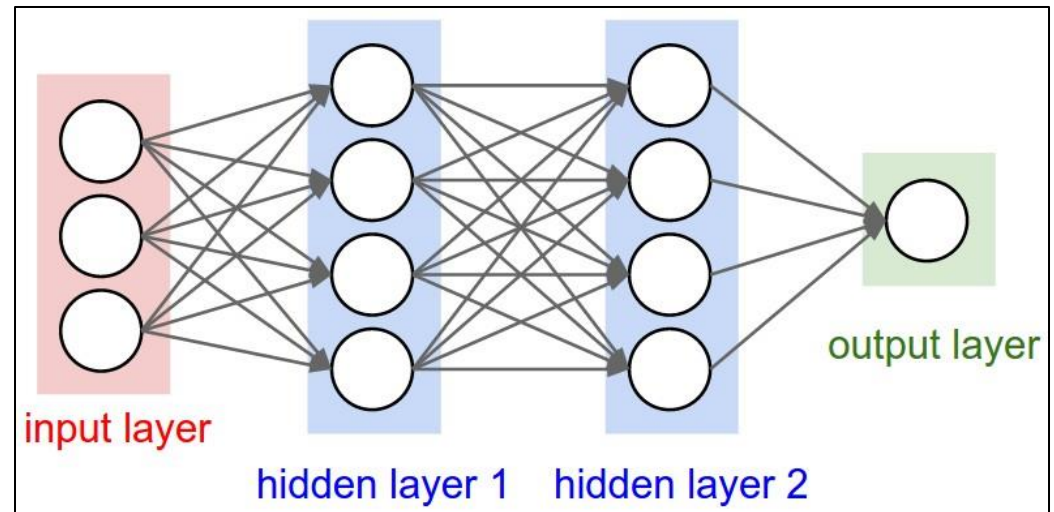




Neural Networks: Layer Architecture



Deep Learning
↓
“Deep” Neural Network (DNN)



(Multilayer Perceptron, MLP)

[Source: Stanford CS231n]



Deep Learning

- ⦿ A subset of machine learning algorithms
- ⦿ Artificial neural networks with **more than one hidden layers**
- ⦿ Using a cascade of many layers of nonlinear processing units for feature extraction and transformation
- ⦿ Each successive layer uses the output from the previous layer as input
- ⦿ Using the **learned weights** instead of hand-coded features



The Neural Network Zoo

A mostly complete chart of

Neural Networks

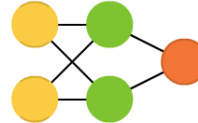
©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

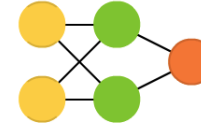
Perceptron (P)



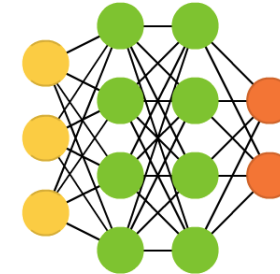
Feed Forward (FF)



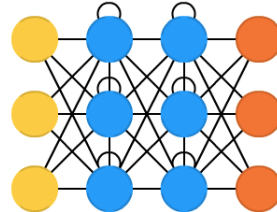
Radial Basis Network (RBF)



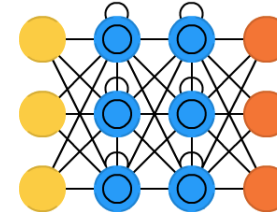
Deep Feed Forward (DFF)



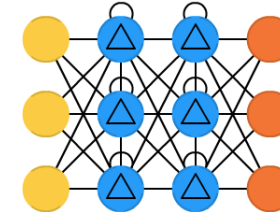
Recurrent Neural Network (RNN)



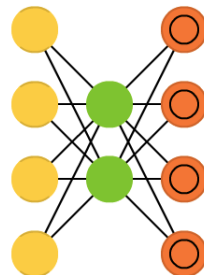
Long / Short Term Memory (LSTM)



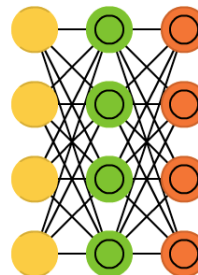
Gated Recurrent Unit (GRU)



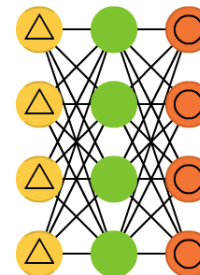
Auto Encoder (AE)



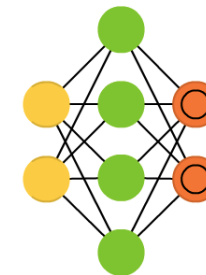
Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)





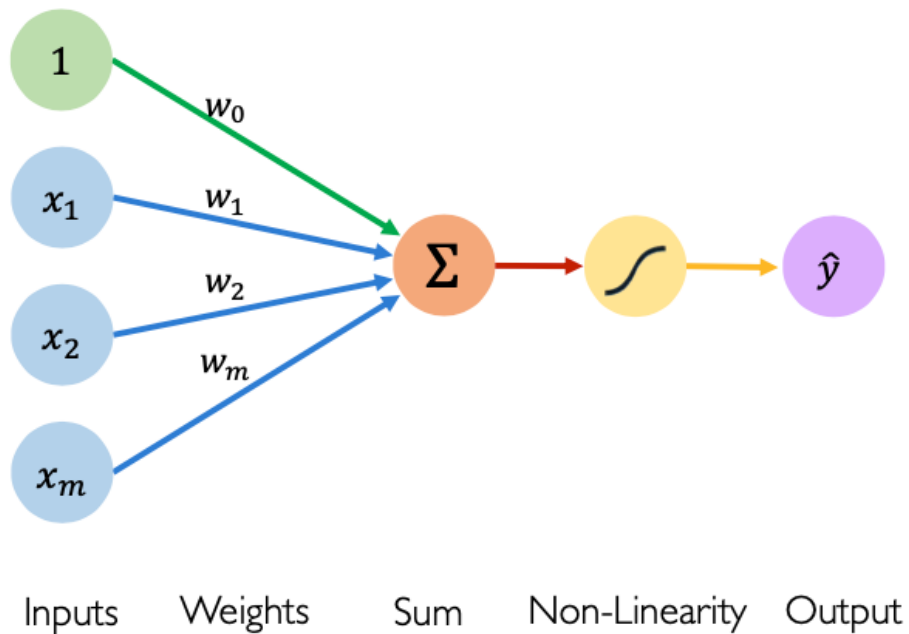
On Training Deep Neural Networks

- » Prof. Alexander Amini,
Lecture 1: Introduction to Deep Learning, MIT 6.S191
- » <http://introtodeeplearning.com>
https://youtu.be/ErnWZxJovaM?si=_qZNfceTdMLInBkI



The Perceptron: Forward Propagation

(Single neuron)



Output

Linear combination of inputs

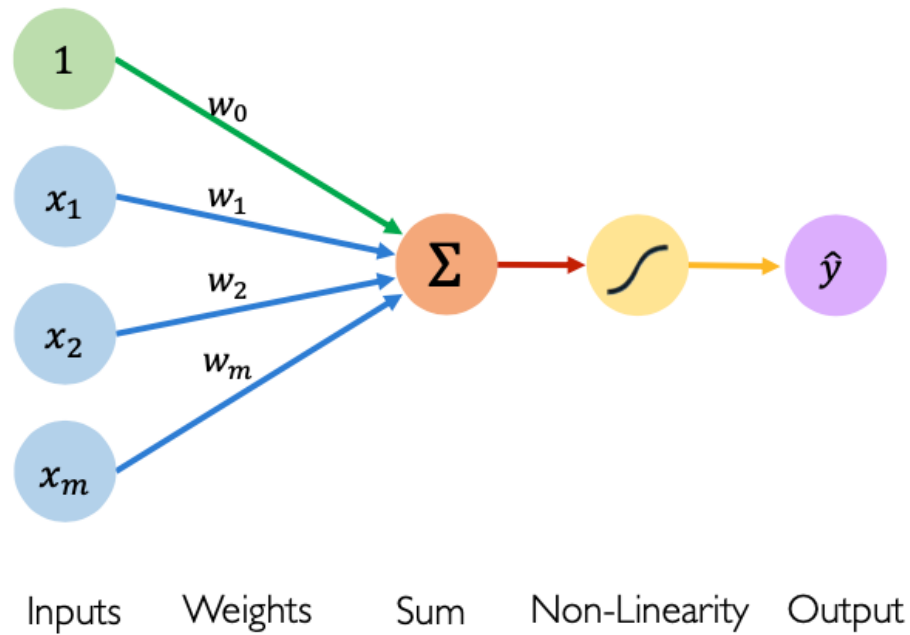
$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

Non-linear activation function

Bias



The Perceptron: Forward Propagation



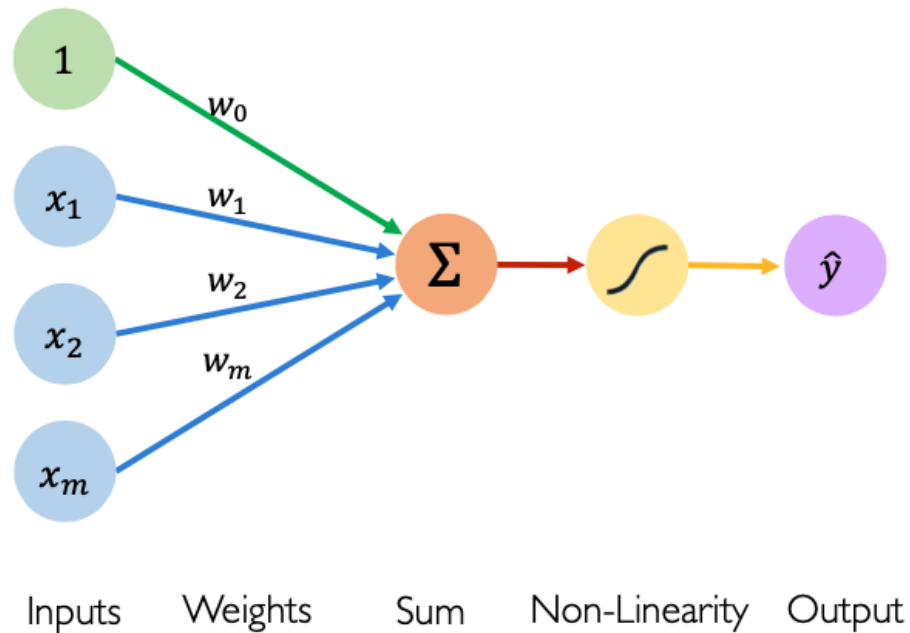
$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

$$\hat{y} = g (w_0 + \mathbf{X}^T \mathbf{W})$$

$$\text{where: } \mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } \mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$



The Perceptron: Forward Propagation

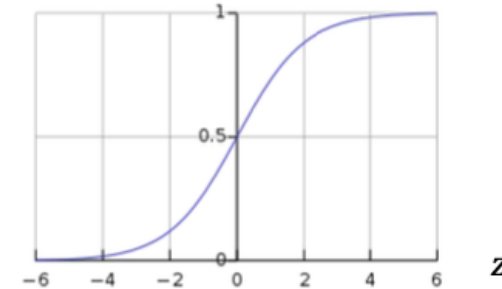


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

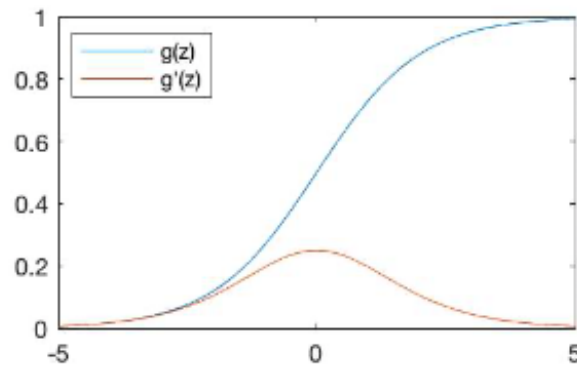
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$





Common Activation Functions

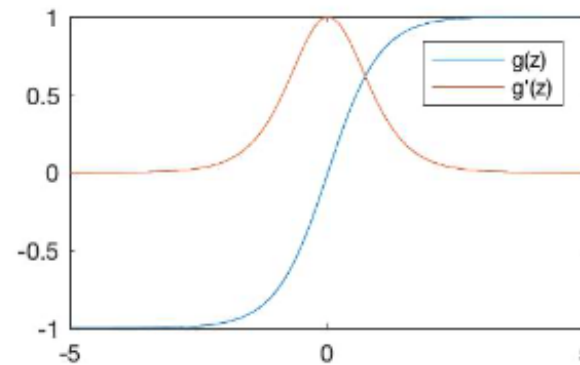
Sigmoid Function



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

$$g'(z) = g(z)(1 - g(z))$$

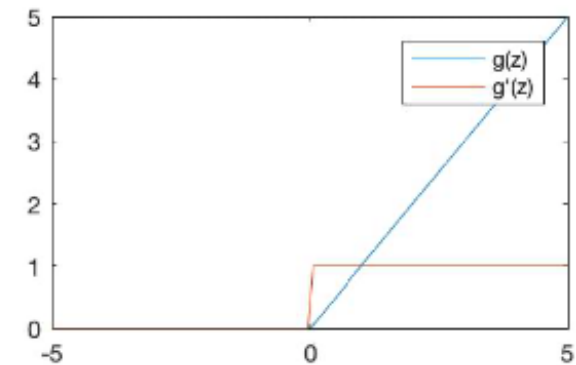
Hyperbolic Tangent



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

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



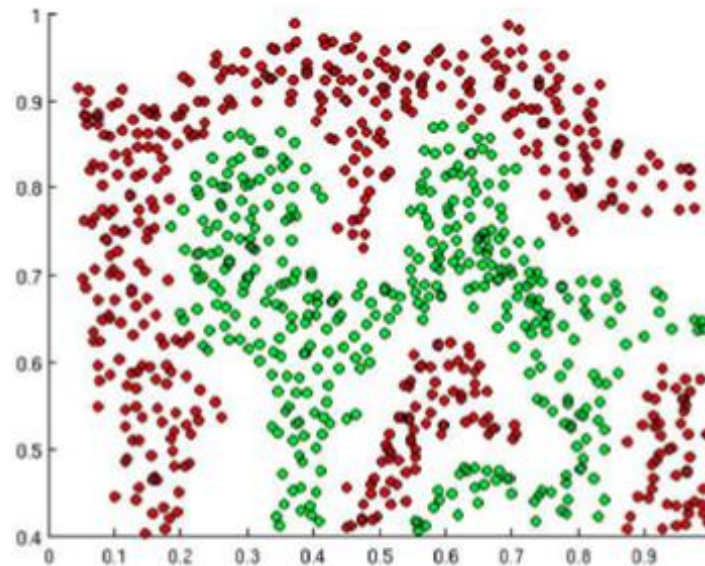
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$



Importance of Activation Functions

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

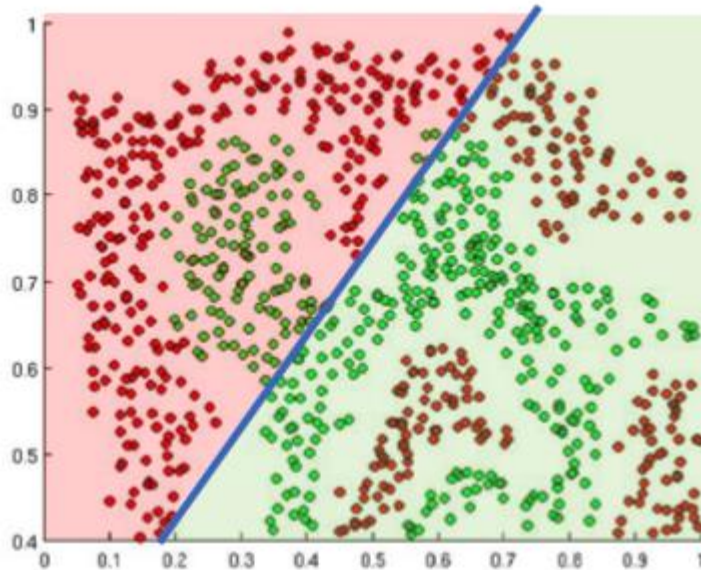


What if we wanted to build a Neural Network to distinguish green vs red points?

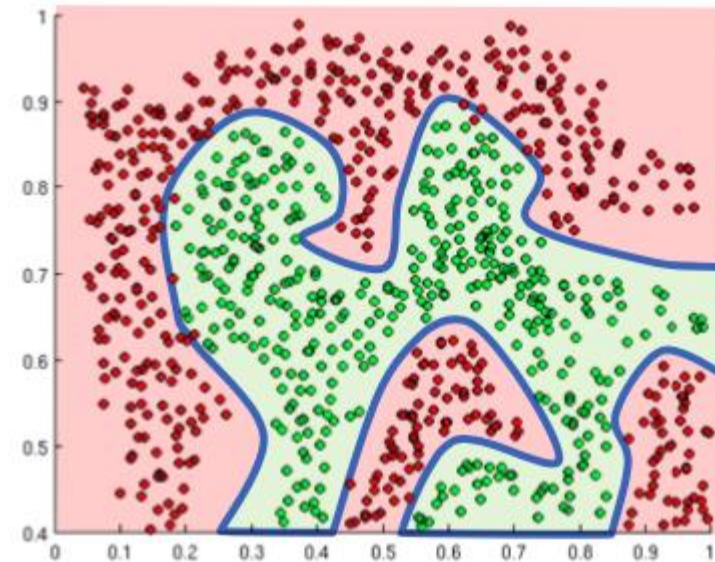


Importance of Activation Functions

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



Linear Activation functions produce linear decisions no matter the network size

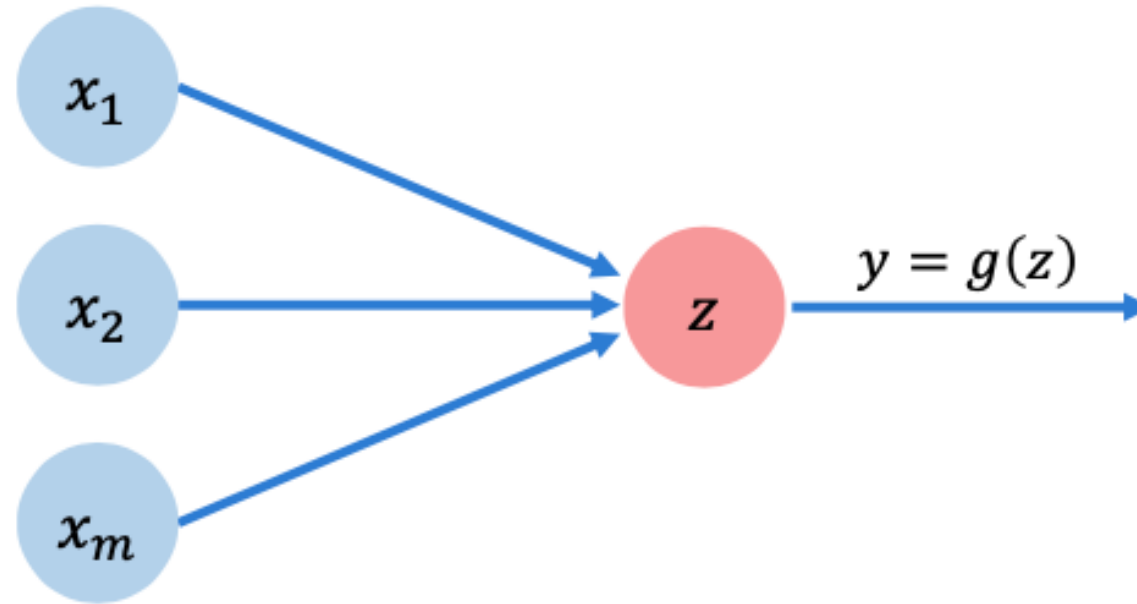


Non-linearities allow us to approximate arbitrarily complex functions



The Perceptron: Simplified

(Bias removed for simplicity)

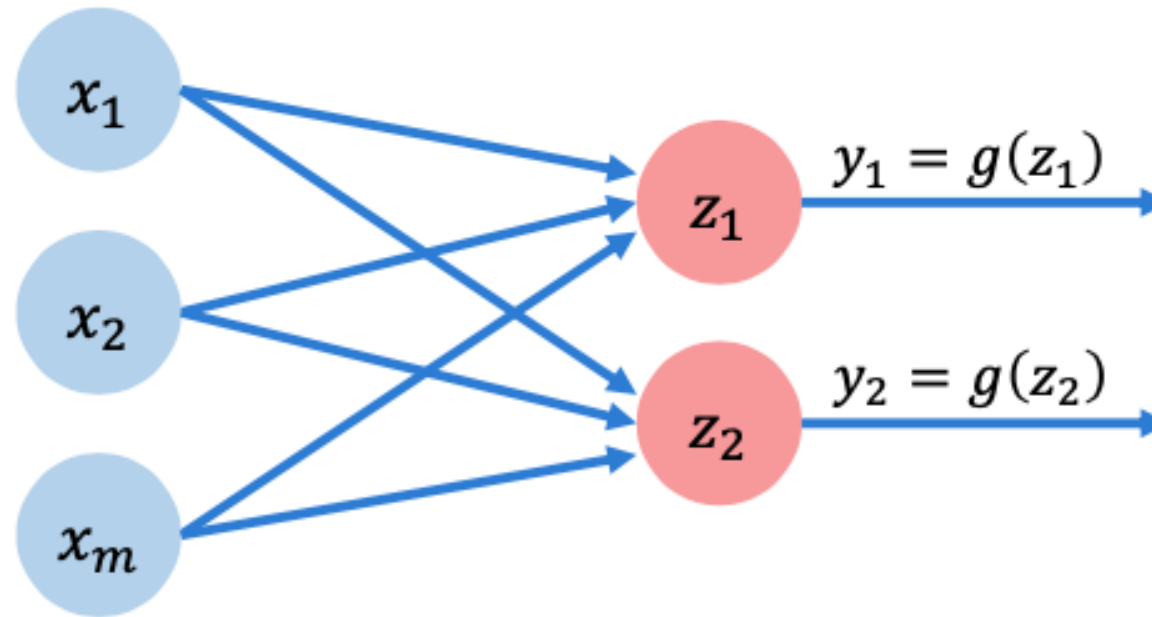


$$z = w_0 + \sum_{j=1}^m x_j w_j$$



Multi-Output Perceptron

All inputs are densely connected to all outputs: **dense layer** or **fully connected layer**

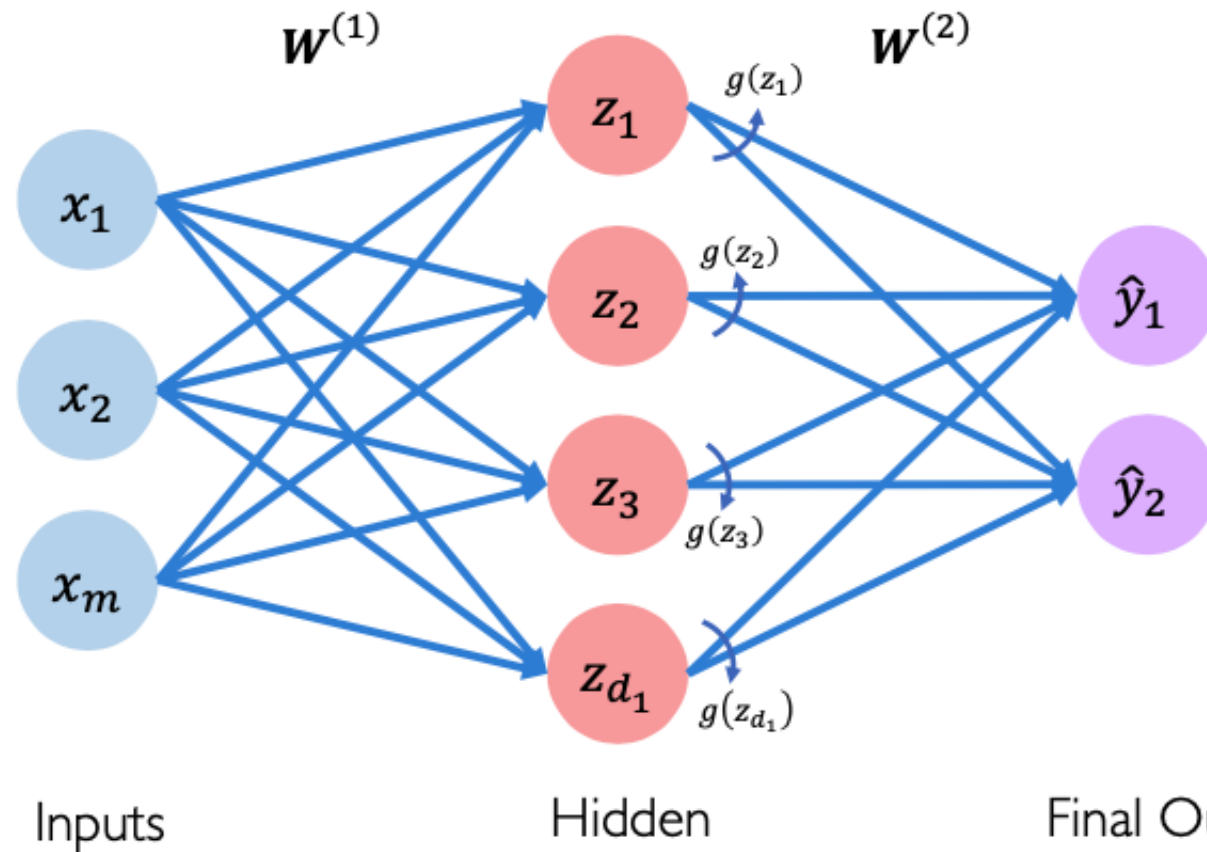


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$



Single-Layer Neural Network

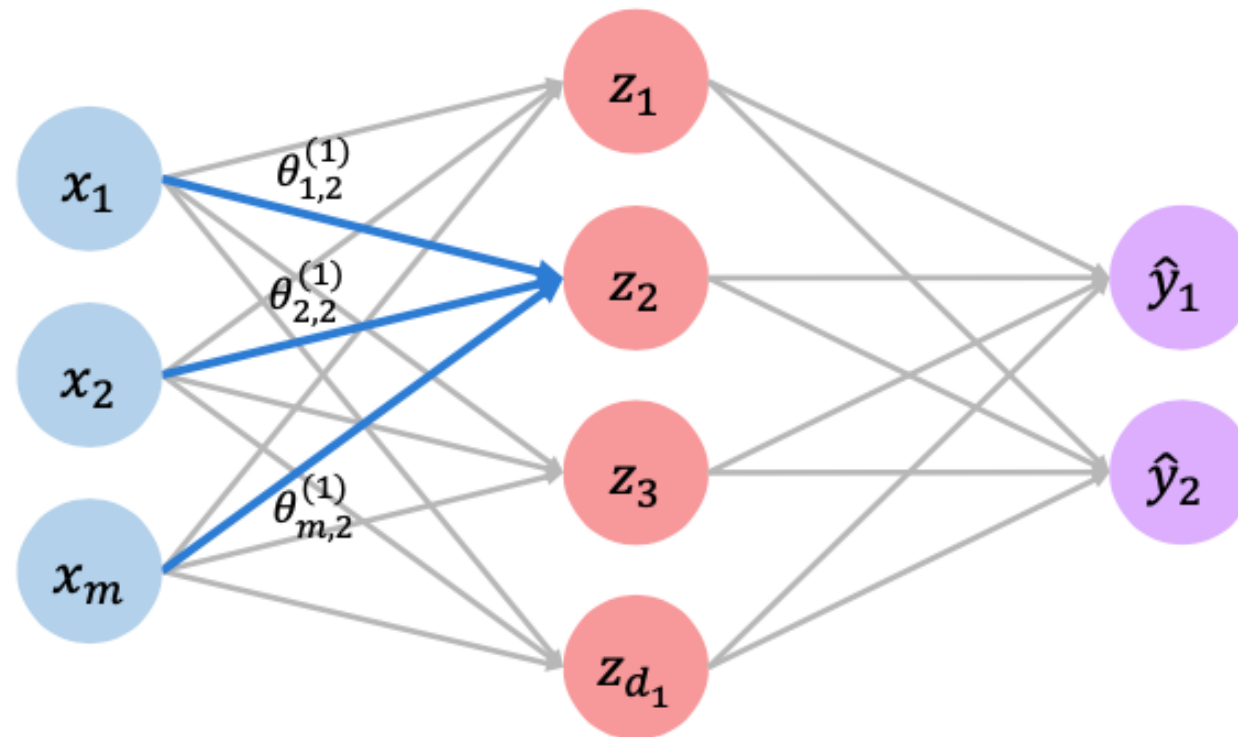
(Single hidden layer)



$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g \left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)} \right)$$



Single-Layer Neural Network



$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$



Example

Will I pass this class?

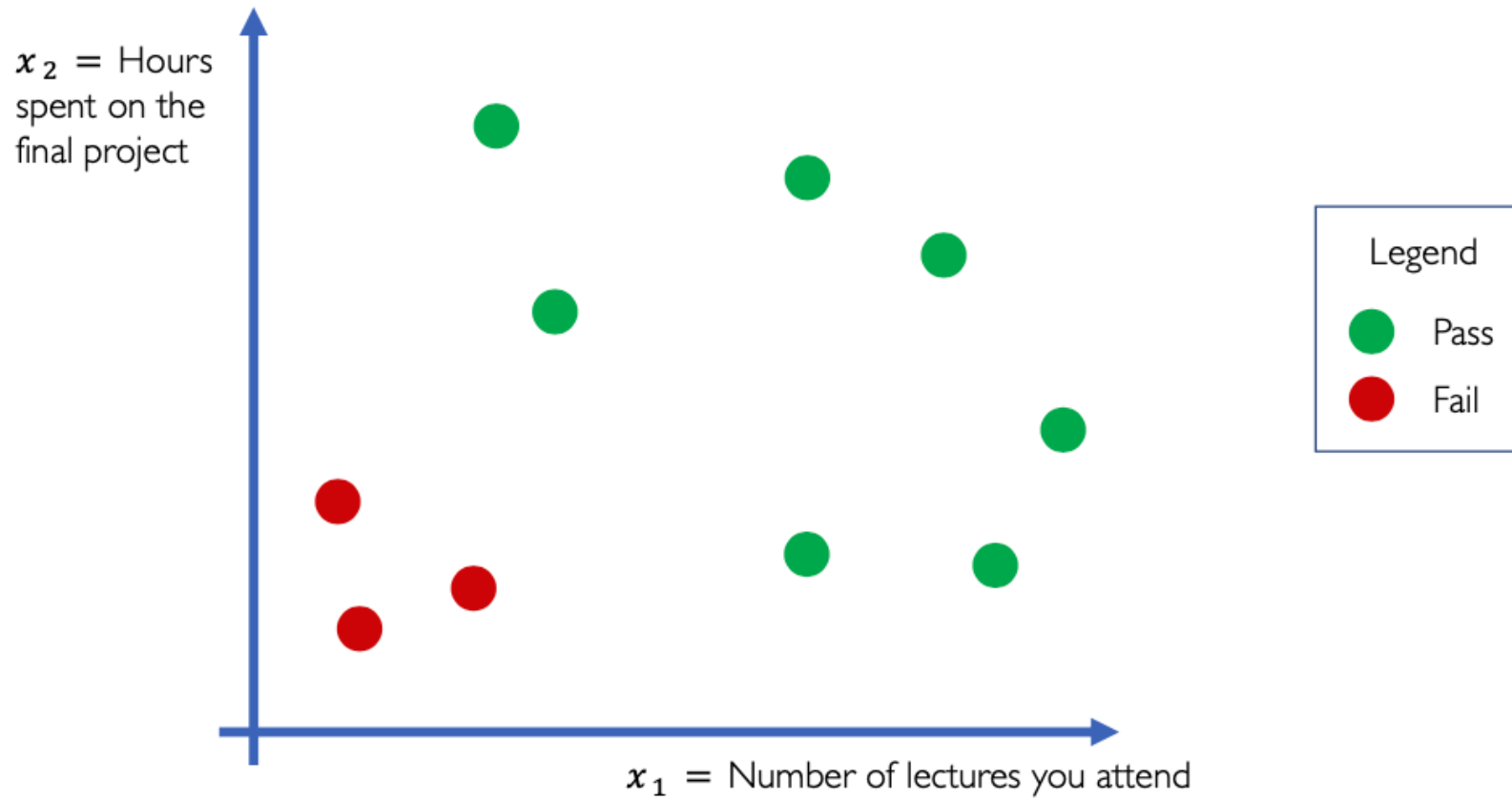
Let's start with a simple two feature model

x_1 = Number of lectures you attend

x_2 = Hours spent on the final project

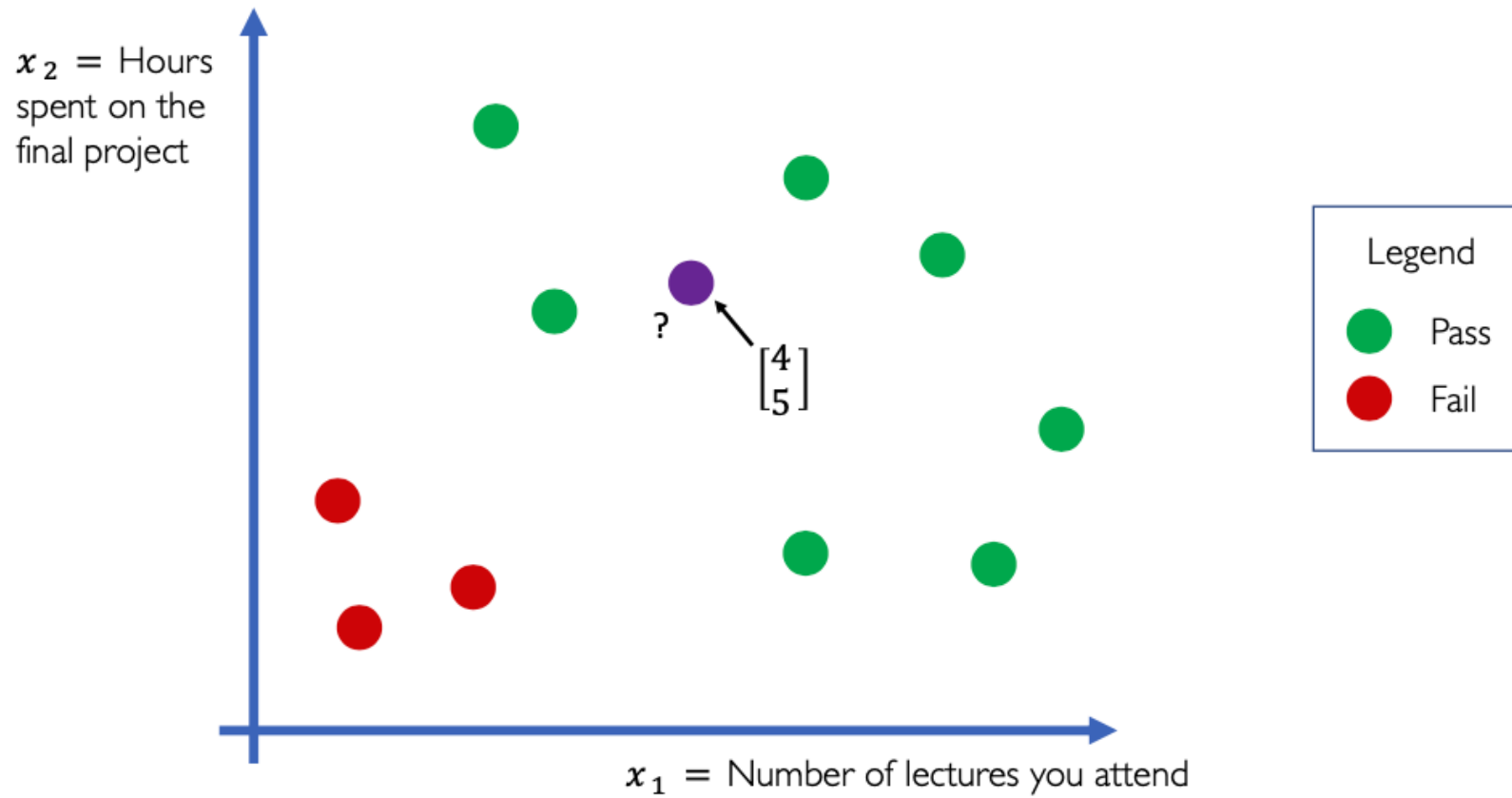


Example Problem: Will I pass this class?





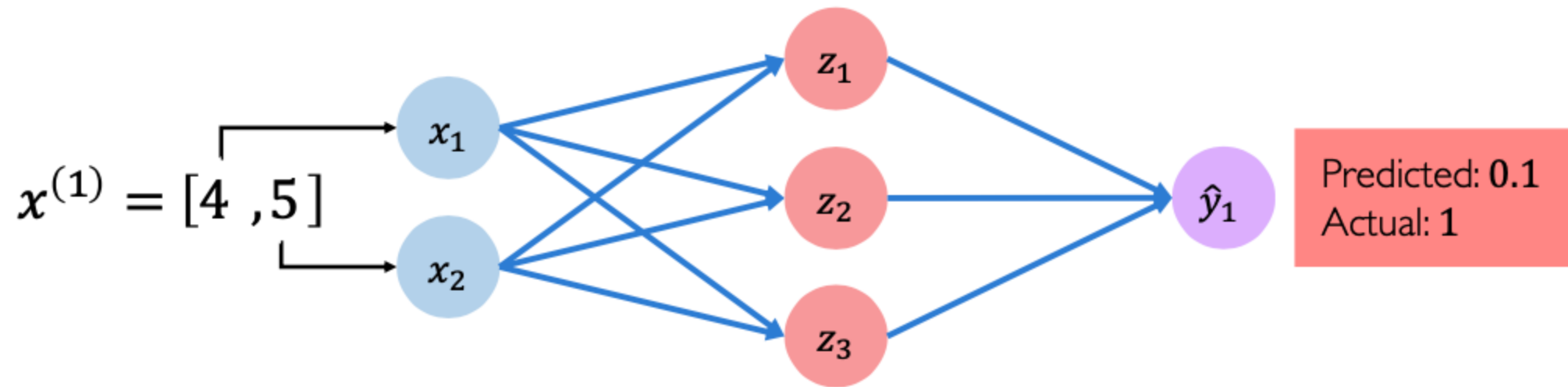
Example Problem: Will I pass this class?





Quantifying Loss

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

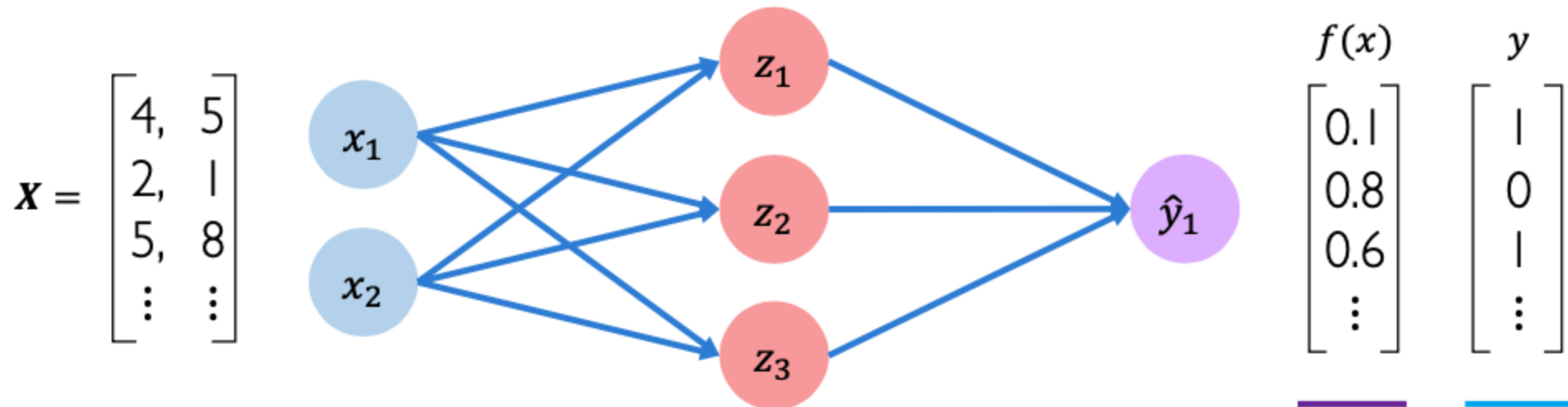


$$\mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$



Empirical Loss

The **empirical loss** measures the total loss over our entire dataset



Also known as:

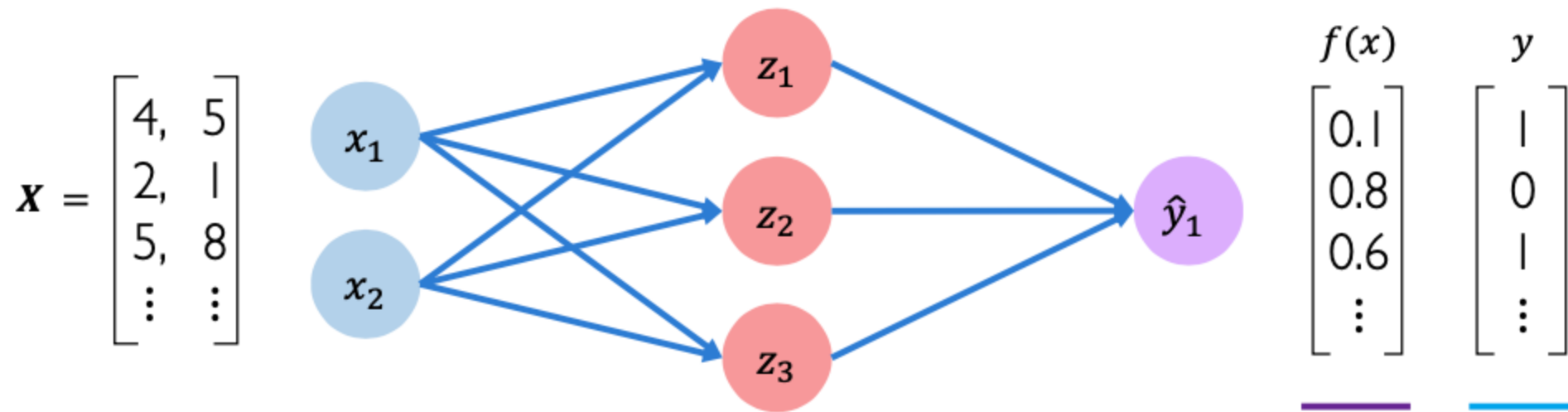
- Objective function
- Cost function
- Empirical Risk

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$



Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1

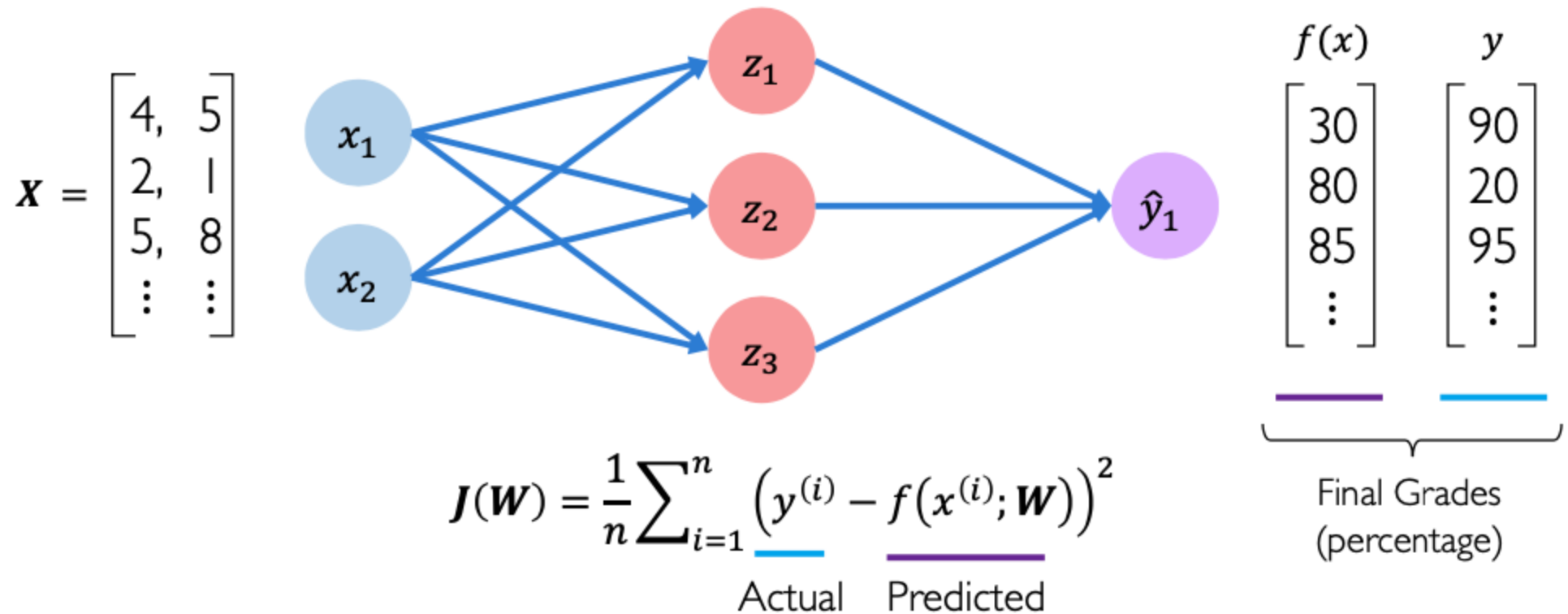


$$J(W) = - \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)}}_{\text{Actual}} \log \left(\underbrace{f(x^{(i)}; W)}_{\text{Predicted}} \right) + (1 - \underbrace{y^{(i)}}_{\text{Actual}}) \log \left(1 - \underbrace{f(x^{(i)}; W)}_{\text{Predicted}} \right)$$



Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers





Loss Optimization

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

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

Remember:

$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

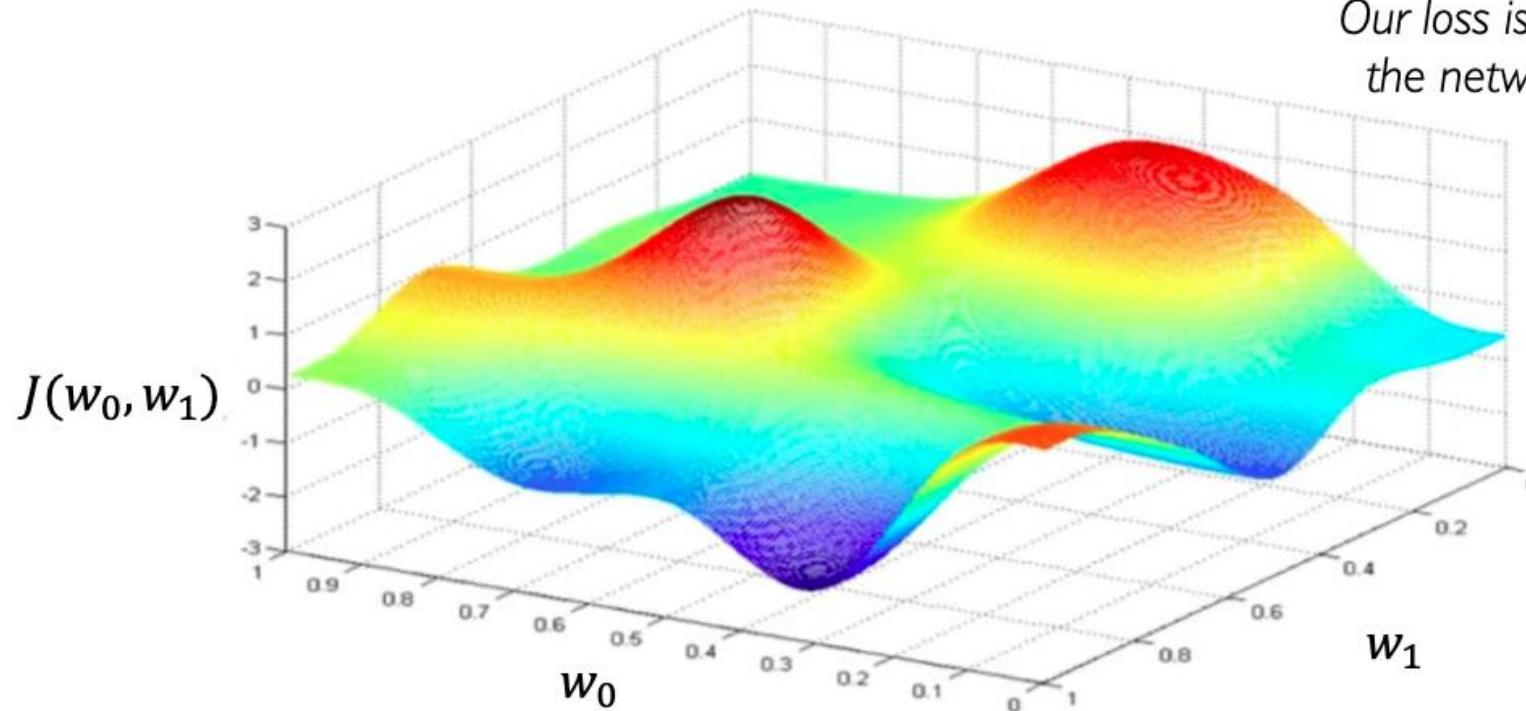


Loss Optimization

To minimize $J(W)$

$$W^* = \operatorname{argmin}_W J(W)$$

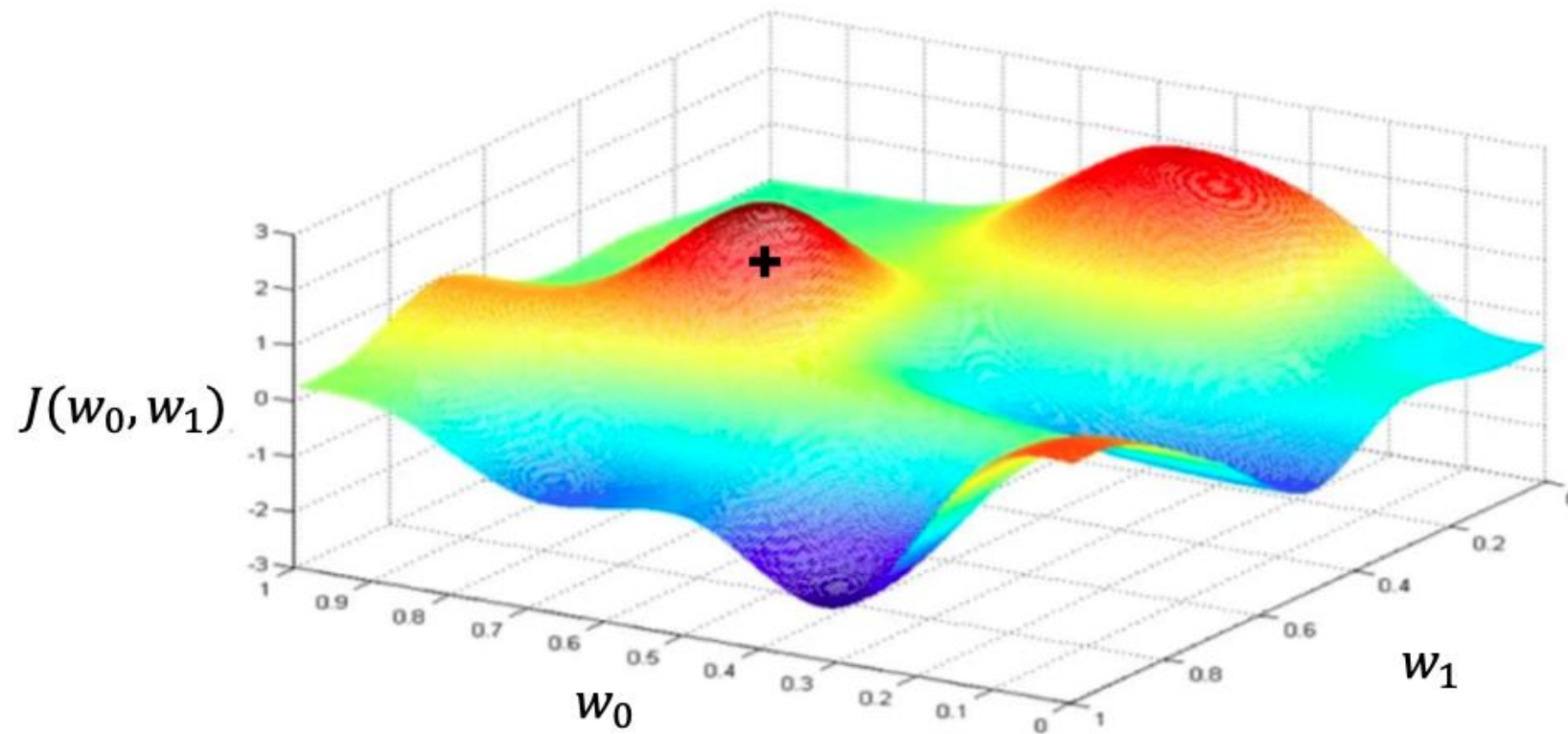
Remember:
*Our loss is a function of
the network weights!*





Loss Optimization

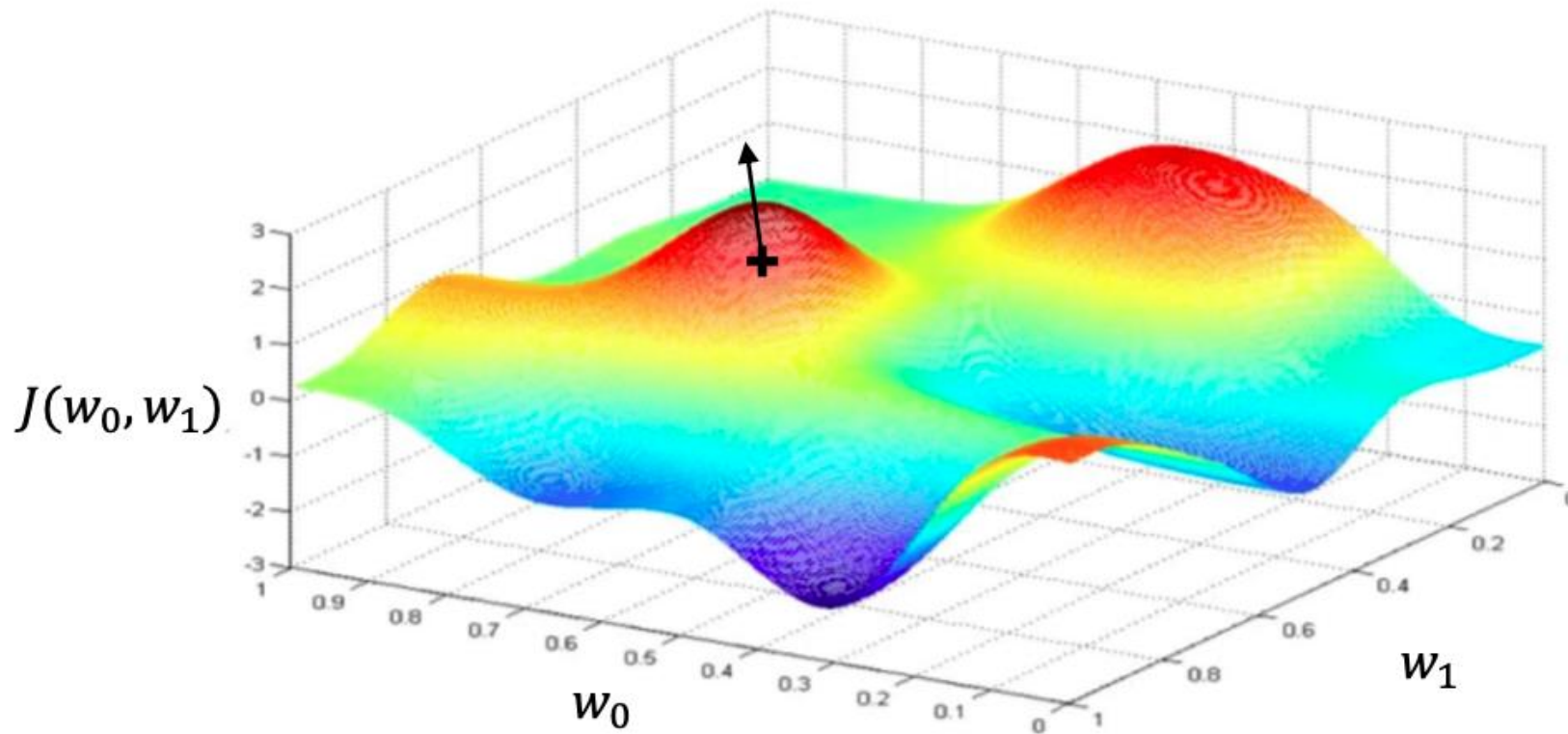
Randomly pick an initial (w_0, w_1)





Loss Optimization

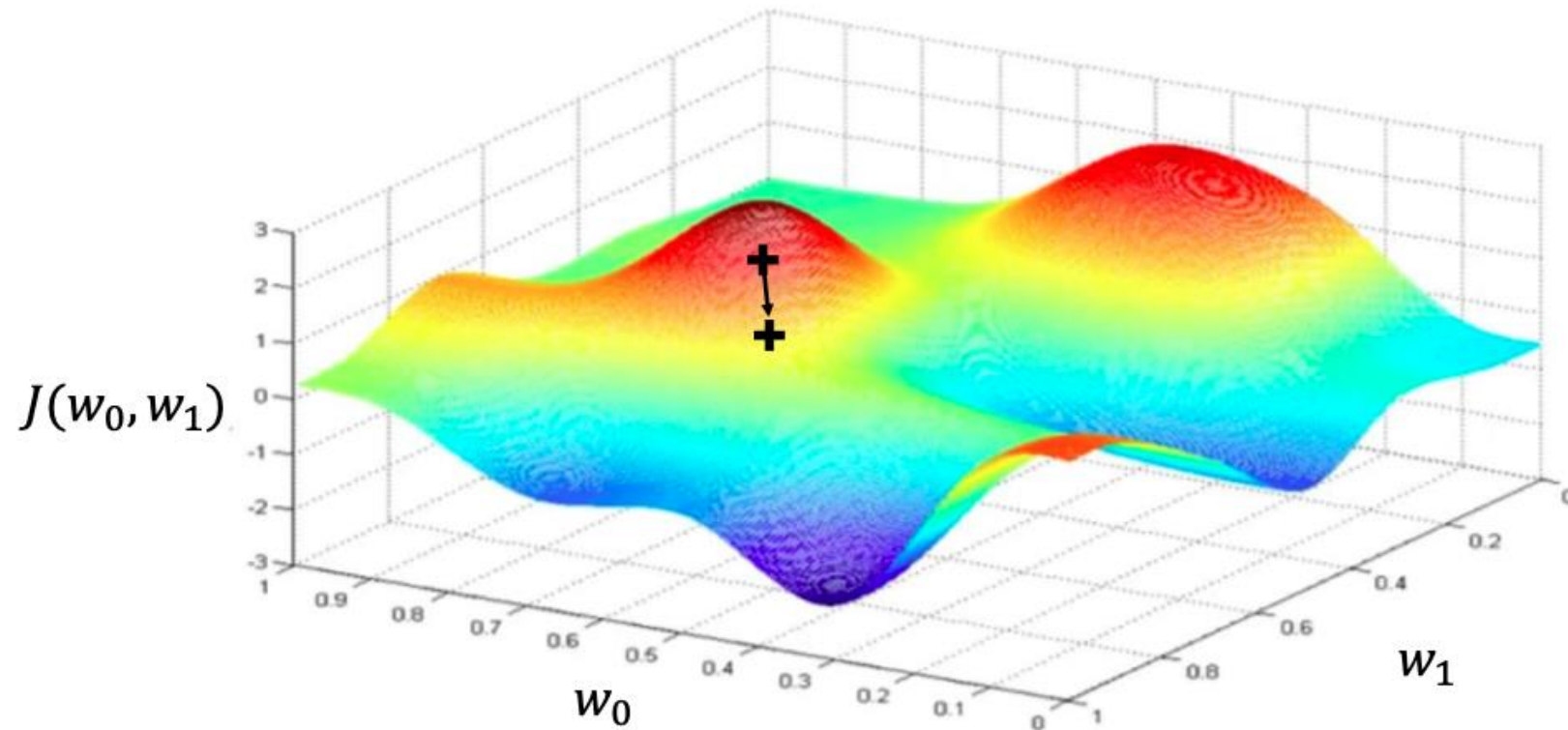
Compute gradient, $\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$





Loss Optimization

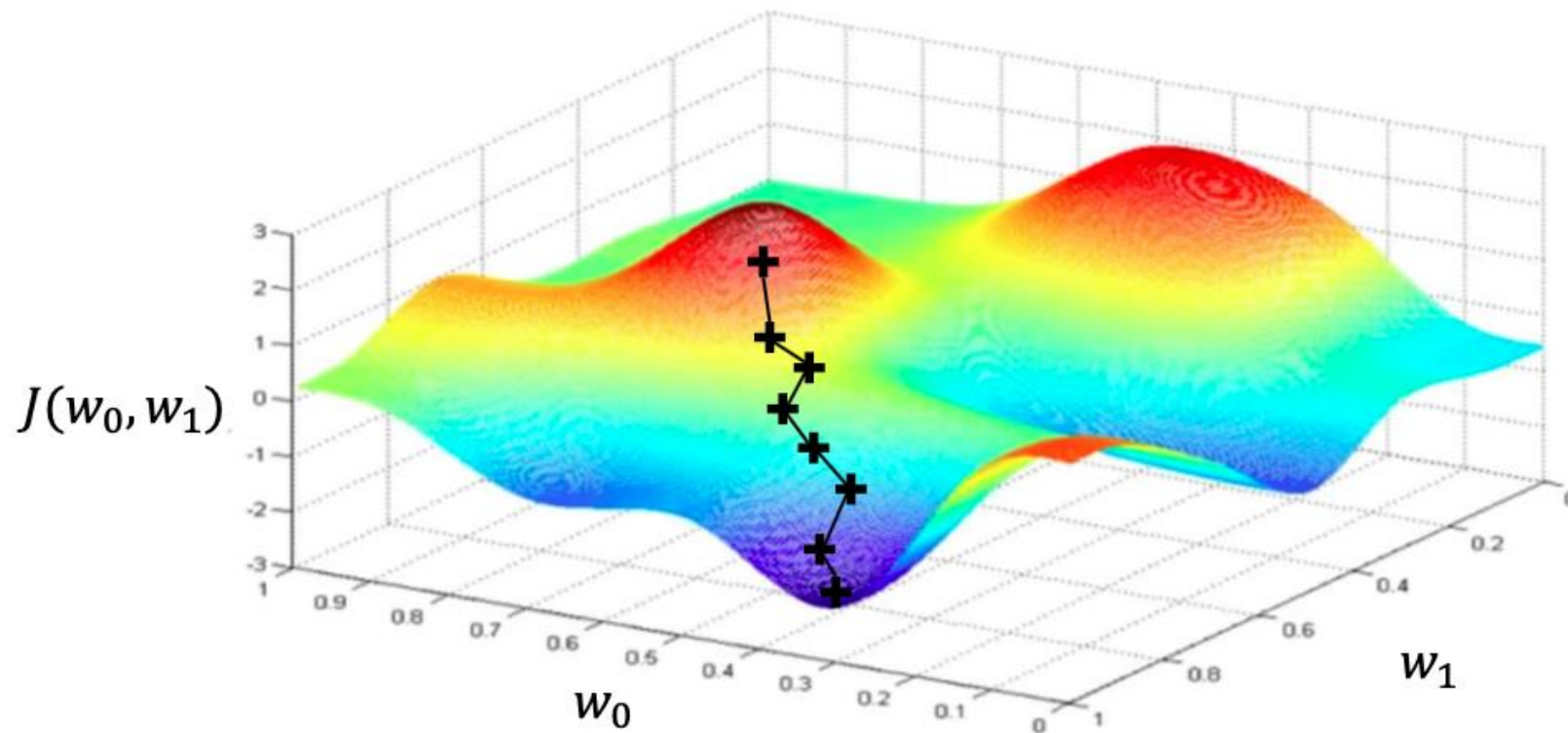
Take small step in opposite direction of gradient





Loss Optimization


Repeat until convergence





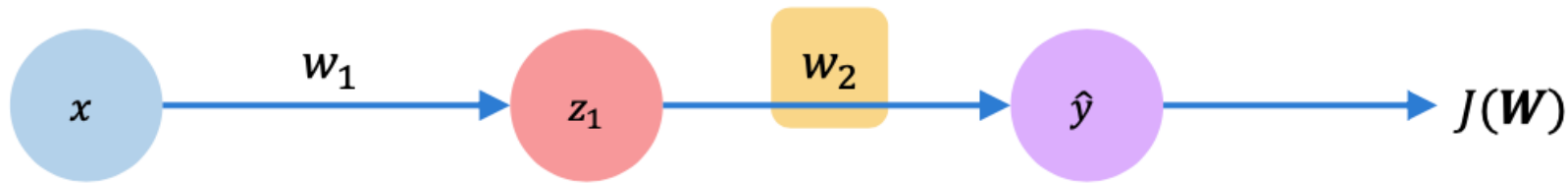
Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
 Learning rate
5. Return weights



Computing Gradients: Backpropagation



How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$

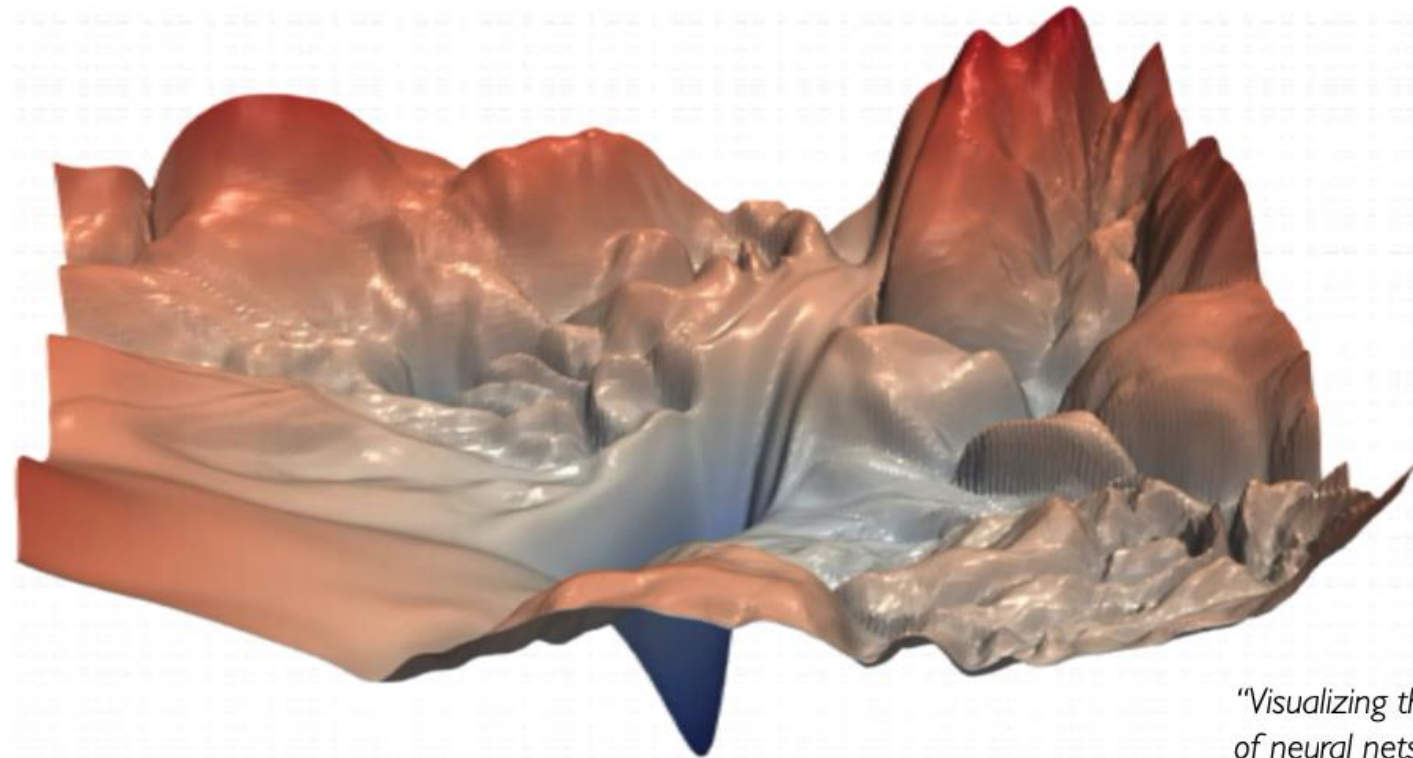
 ← Chain rule!

$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

Repeat this for **every weight in the network** using gradients from later layers



Training Neural Networks is Difficult



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



Loss Functions Can Be Difficult to Optimize

Remember:

Optimization through gradient descent

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

How can we set the
learning rate?

Hyperparameter!

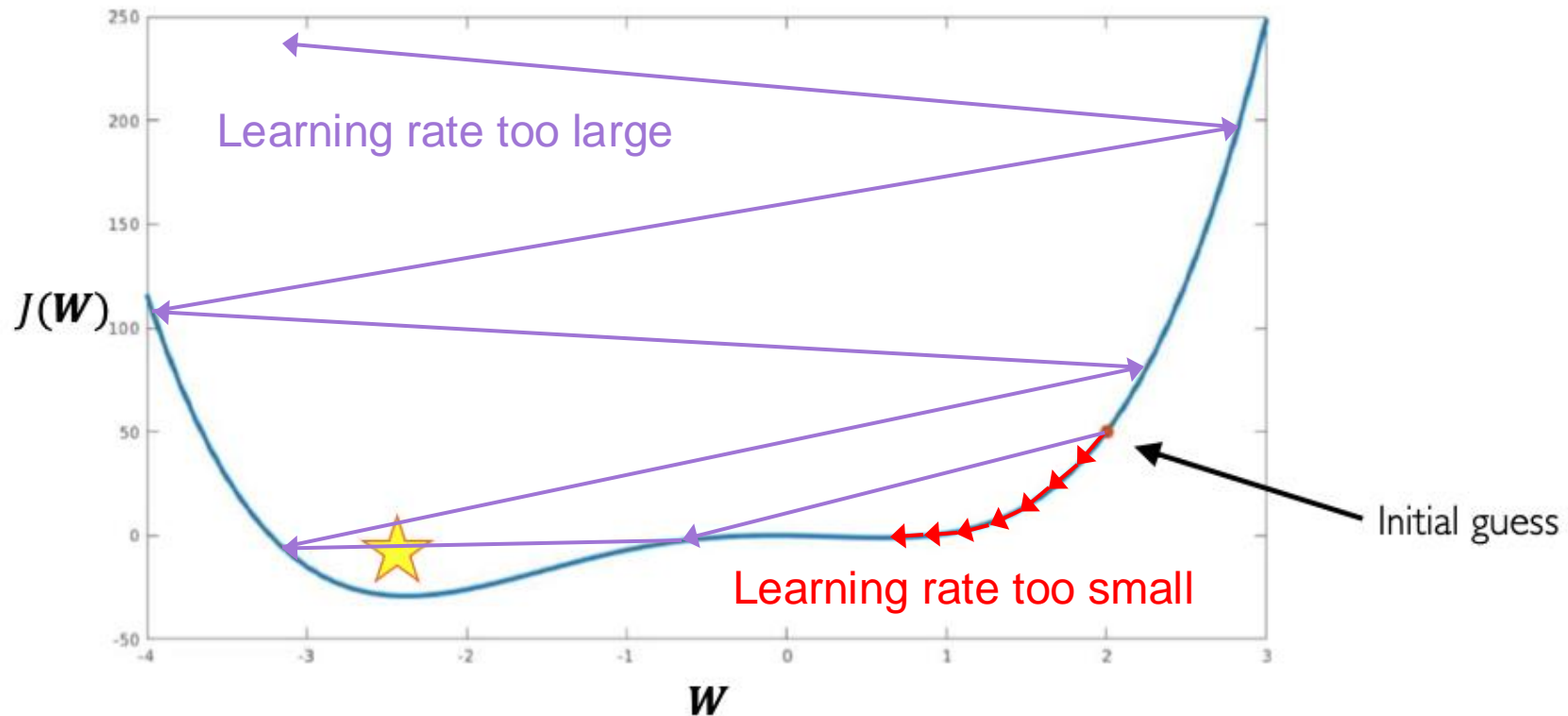


Setting the Learning Rate

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

Large learning rates overshoot, become unstable and diverge

Stable learning rates converge smoothly and avoid local minima





How to Choose A Good Learning Rate?

- ⦿ Try lots of different learning rates and see what works
"just right"
- ⦿ Use an adaptive learning rate that "adapts" to the landscape
 - ◆ Learning rates can be larger or smaller dynamically, depending on
 - ▣ How large gradient is
 - ▣ How fast learning is happening
 - ▣ Size of particular weights
 - ▣ Etc.

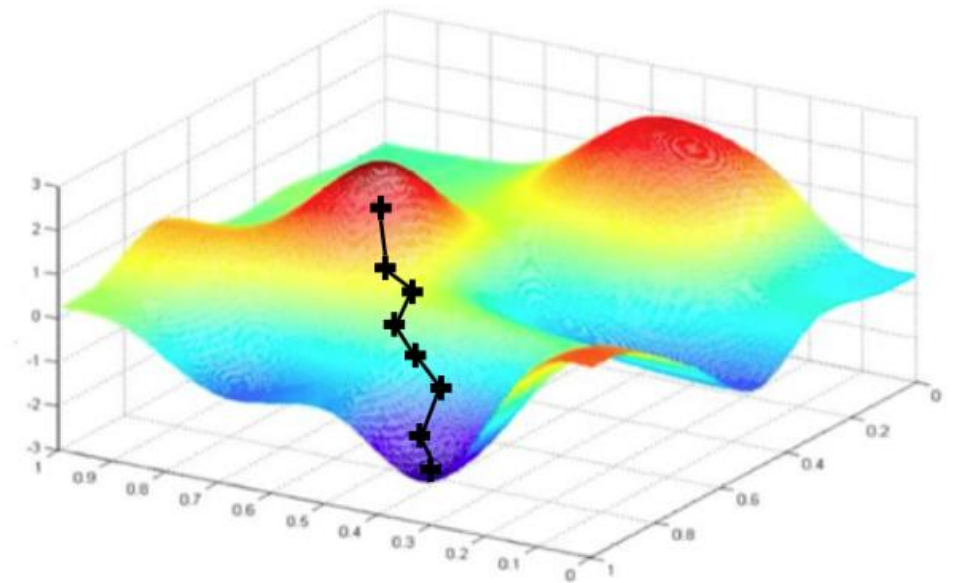


Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Can be very
computational to
compute!



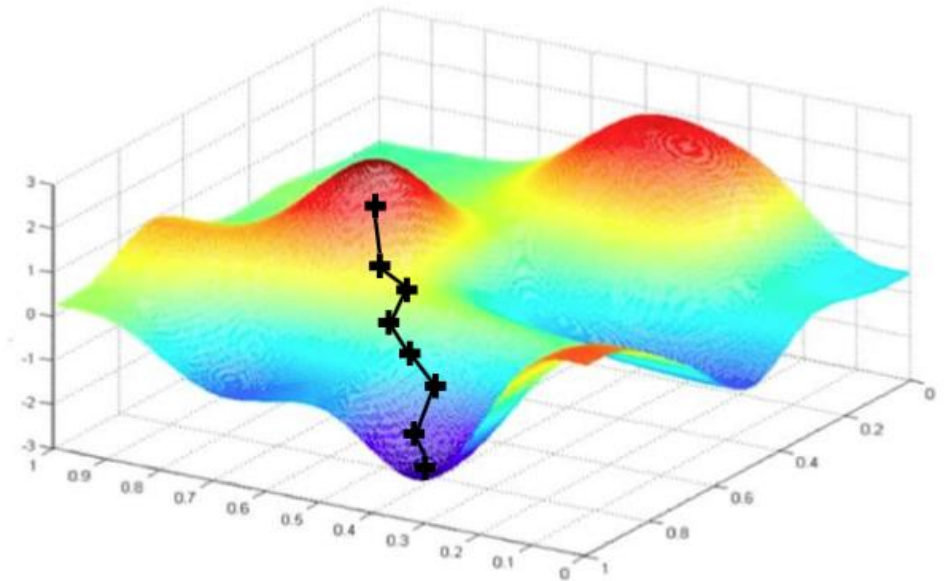


Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(W)}{\partial W}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(W)}{\partial W}$
6. Return weights

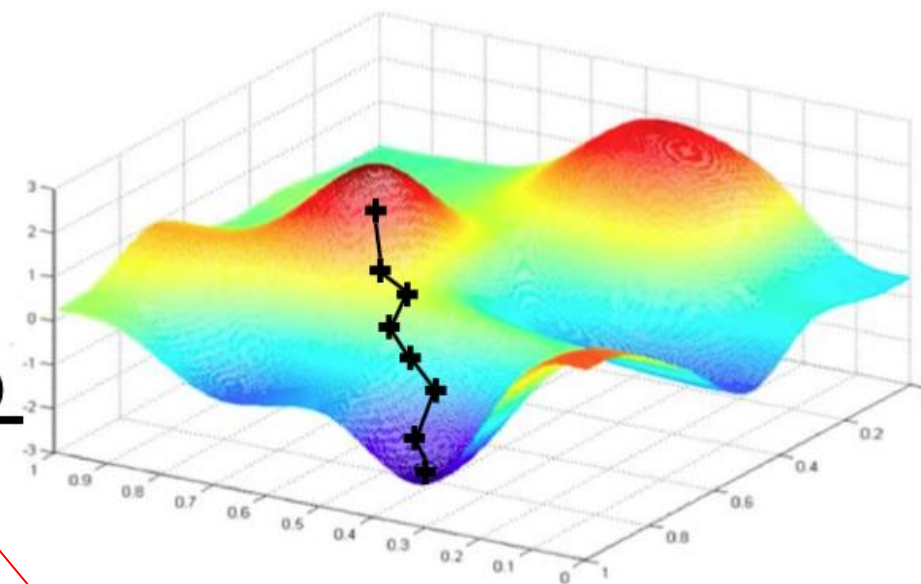
Easy to compute but
very noisy
(stochastic)!



Stochastic Gradient Descent (with Mini-Batches)

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



Fast to compute
A much better estimation of the true gradient!

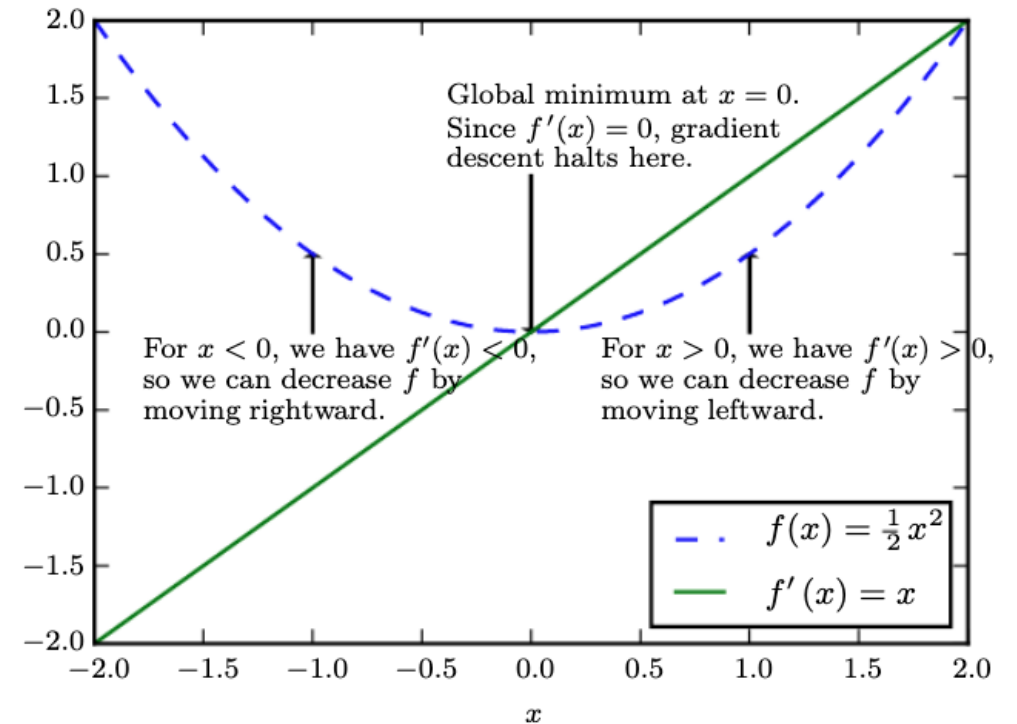


Training with Mini-batches

- ⊙ More accurate estimation of gradient
 - ◆ Smoother convergence
 - ◆ Allows for larger learning rates
- ⊙ Mini-batches lead to fast training
 - ◆ Can parallelize computation
 - ◆ Achieve significant speed increases on GPU

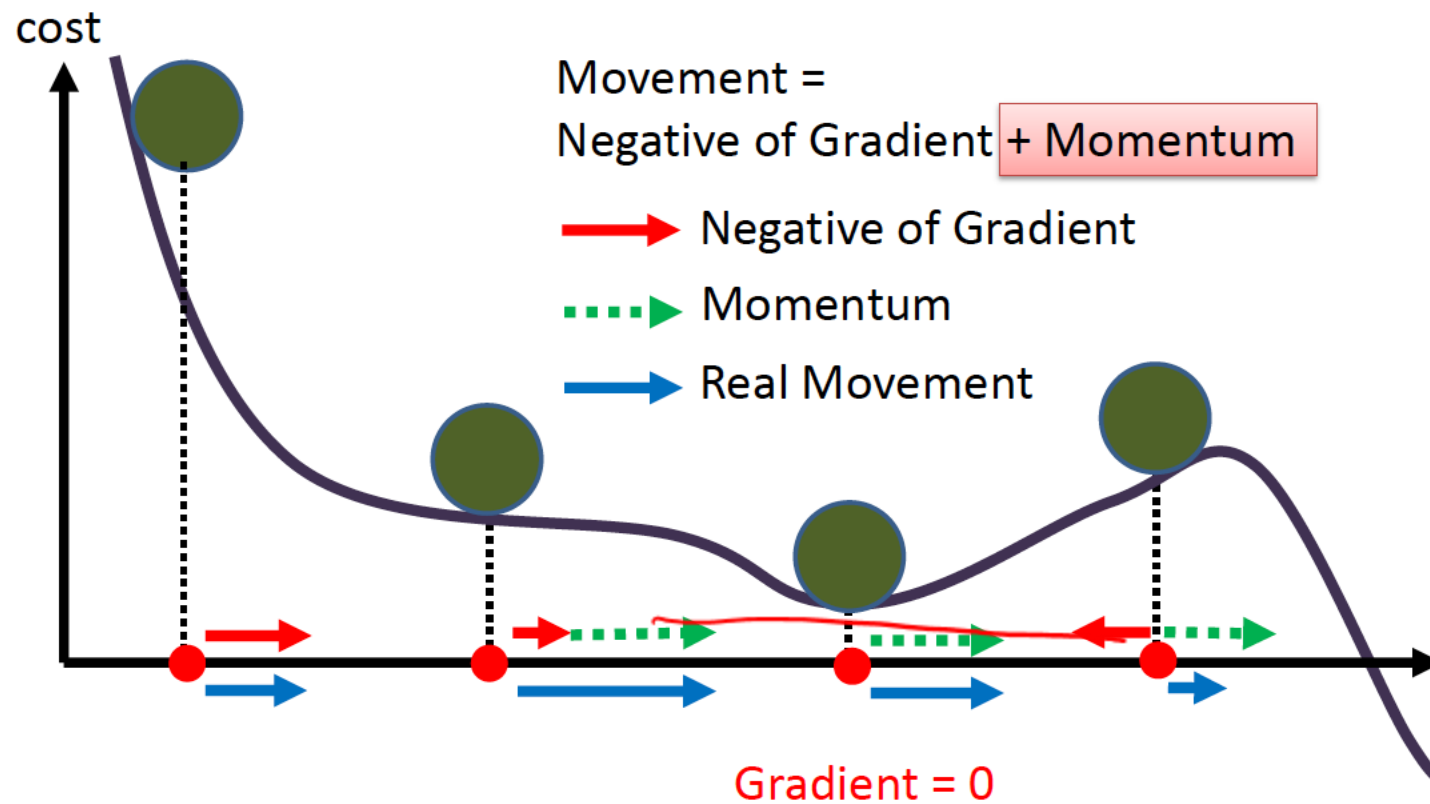
Optimization

- Method to minimize the cost function by updating weights
- Gradient descent
 - Iteratively moving in the direction of steepest descent as defined by the negative of the gradient
- Stochastic gradient descent (SGD)
 - To handle large training sets
 - Only run a subset of the training sets (i.e., batch/minibatch) for each update
 - Easier to converge



[Source: Prof. Sophia Shao, EE290-2, Berkeley]

Optimization with Momentum



SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:  
    dx = compute_gradient(x)  
    x -= learning_rate * dx
```

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

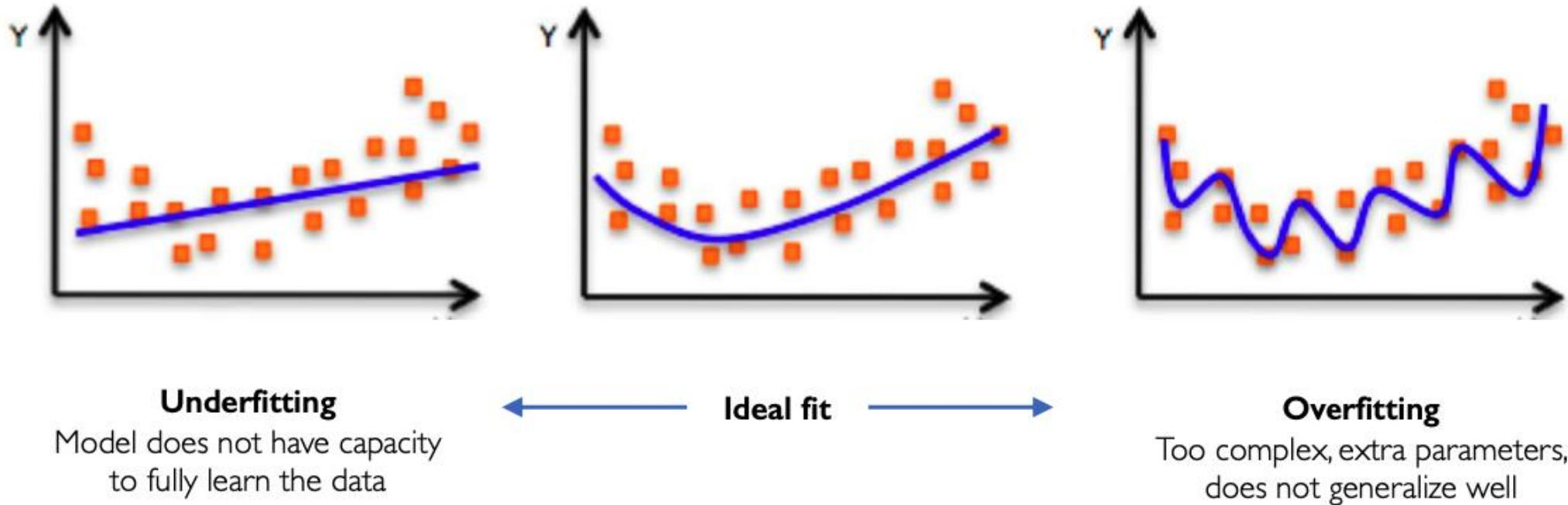
$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0  
while True:  
    dx = compute_gradient(x)  
    vx = rho * vx + dx  
    x -= learning_rate * vx
```

[Source: Prof. Sophia Shao, EE290-2, Berkeley]



The Problem of Overfitting



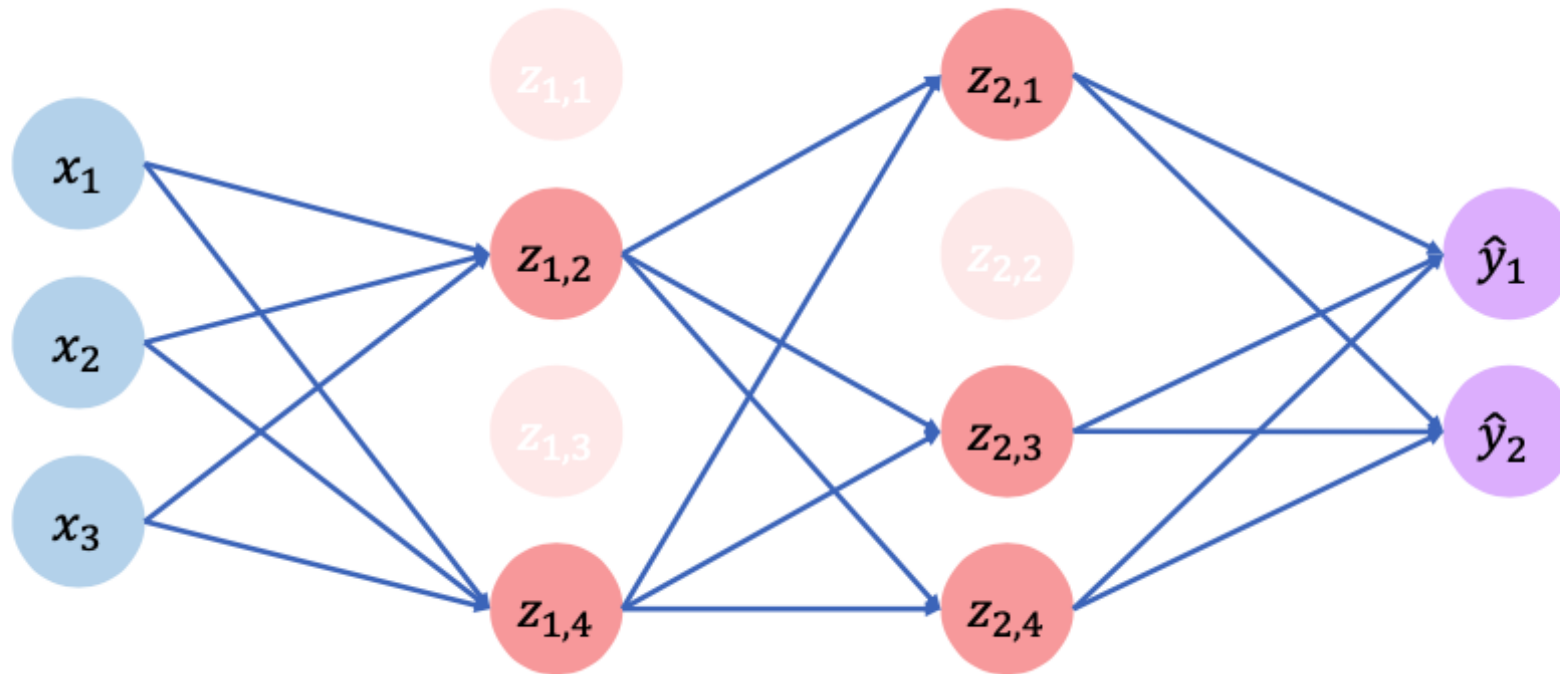
Regularization:

- ◆ Technique in the optimization loop to **discourage** complex models
- ◆ Improve generalization of our model on **unseen data**



Regularization: Dropout

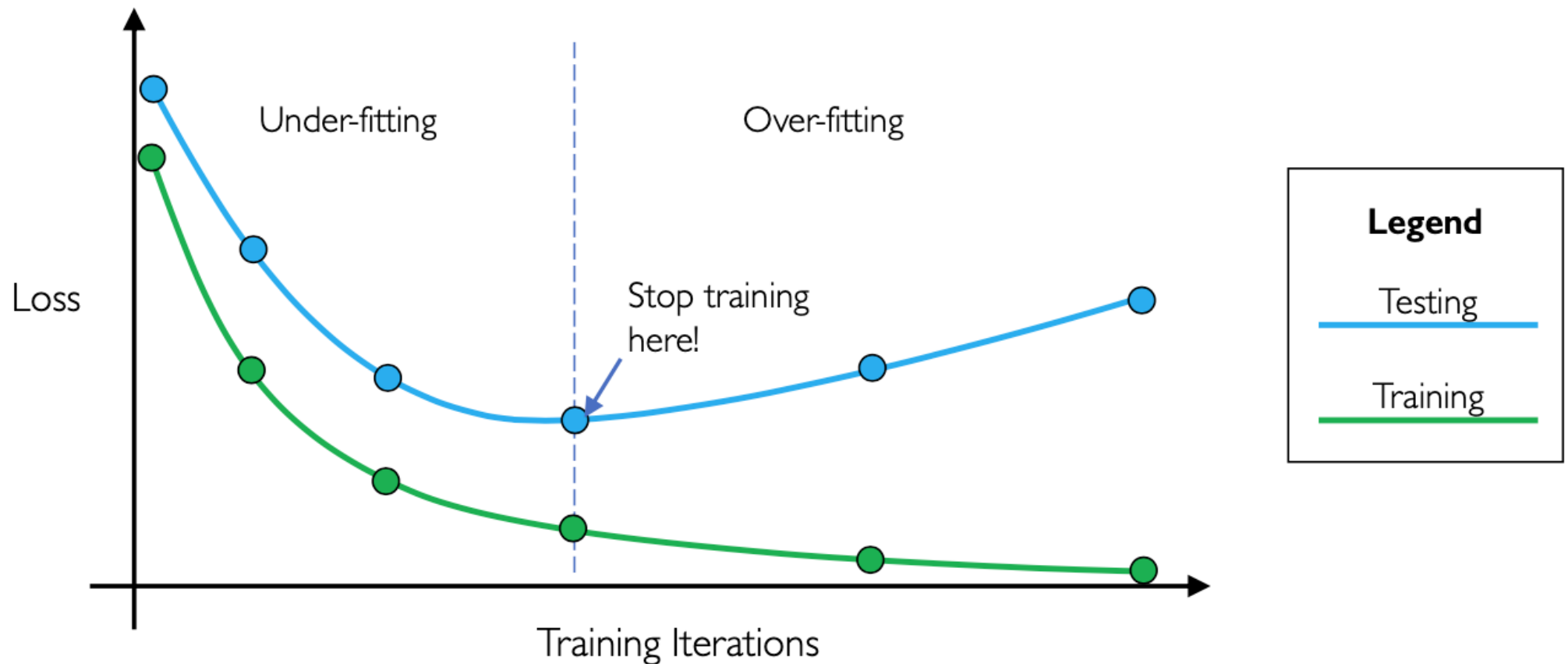
- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node
- (Also faster to train!)





Regularization: Early Stopping

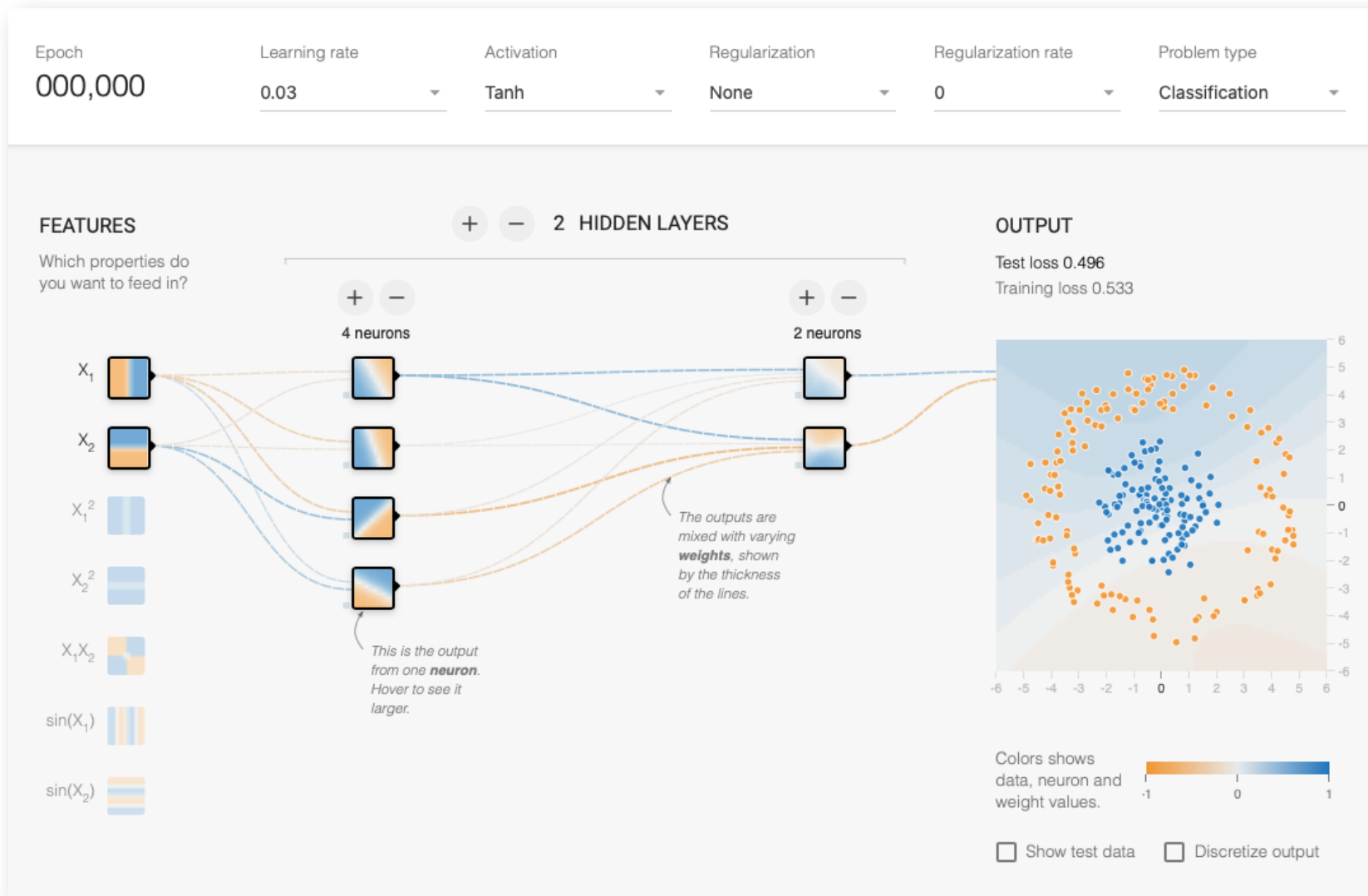
- Stop training before we have a chance to overfit





Tinker with a Neural Network

🔗 <http://playground.tensorflow.org>





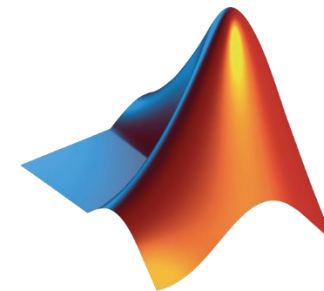
Deep Learning Frameworks

Caffe



theano

PYTORCH



Ref: <https://developer.nvidia.com/deep-learning-frameworks>

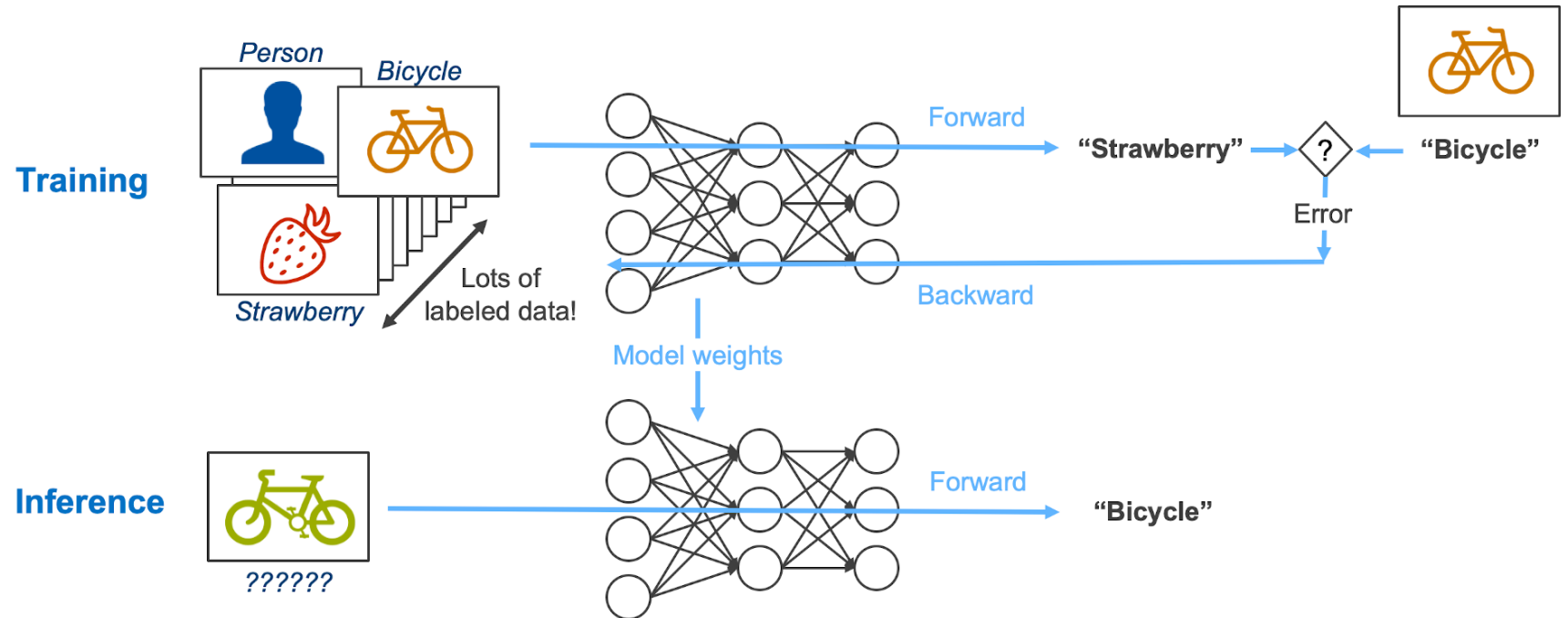
Training vs. Inference

Training

- ◆ Dataset
- ◆ Cost function
- ◆ Optimization function
- ◆ Model

Inference

- ◆ Dataset or realistic data
- ◆ Model



Src: <https://www.intel.com/content/www/us/en/artificial-intelligence/posts/deep-learning-training-and-inference.html>