

Deep Learning and its Applications:- EGE480

Introduction to Deep Learning

Lecturer: Associate Prof. Dimitrios Rafailidis

<https://blog.boardinfinity.com/deep-learning-dive-into-the-world-of-machine-learning/>

Course structure

- Εφαρμογή τεχνικών deep learning σε πραγματικά δεδομένα της εταιρείας Hive Streaming!
- Δεν υπάρχει γραπτή εξέταση στο μάθημα, ο βαθμός θα βασίζεται σε τελική εργασία. Η εργασία θα είναι ατομική ή ομαδική με 2 φοιτητές. Η εργασία θα αποτελείται από δύο μέρη. Το πρώτο μέρος αφορά την ανάλυση πραγματικών δεδομένων και το δεύτερο μέρος την μοντελοποίηση και εφαρμογή τεχνικών deep learning. Η προθεσμία της ενδιάμεσης υποβολής του πρώτου μέρους είναι 15 Απριλίου. Η τελική υποβολή και των δύο μερών είναι 17 Ιουνίου στο draf@uth.gr και stefanos.antaris@hivestreaming.com.
- Tutorials and group meetings with S. Antaris - HiveStreaming online.
Links will be announced soon
- Lectures (Αμφ.2-101)
- Tutorials and group meetings (online)

Course schedule

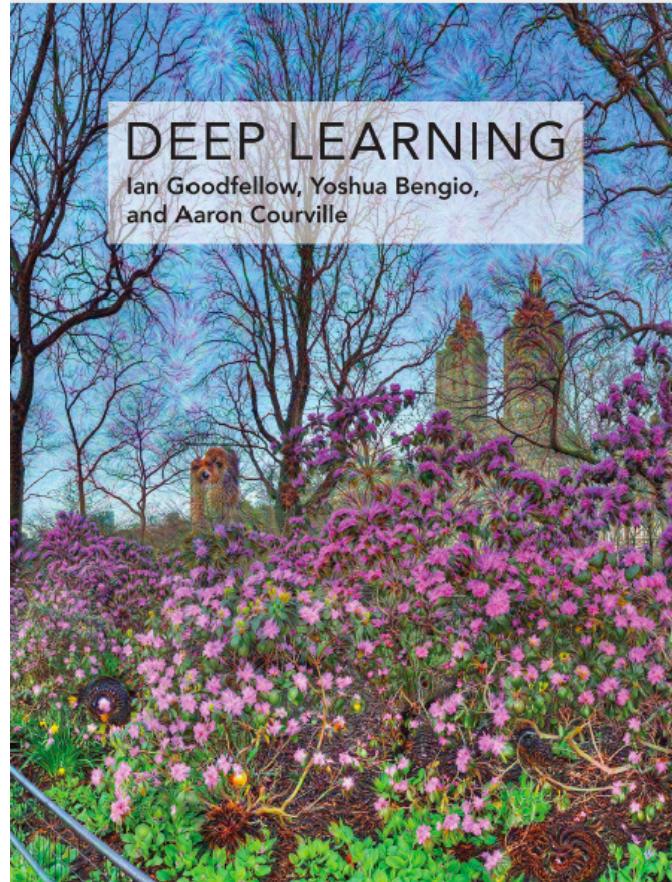
- 3 March: Introduction - Αμφ.2-101
- 4 March: Tutorial I – online
- 10 March: Assignment description – online
- 11 March: Tutorial II – online
- 17 March: Group meeting I – online
- 18 March: Group meeting I – online
- 24-25 March: -
- 31 March: Deep sequential learning - Αμφ.2-101
- 1 April: Tutorial III – online
- 7 April: Group meeting II – online
- 8 April: Group meeting II – online
- 14 April: CNNs - Αμφ.2-101
- 15 April: Deep generative models - Αμφ.2-101
- 5 May: Tutorial IV – online
- 6 May: Introduction to reinforcement learning - Αμφ.2-101
- 12 May: Group meeting III – online
- 13 May: Group meeting III – online
- 19 May: Deep reinforcement learning - Αμφ.2-101
- 20 May: -
- 26 May: Group meeting IV – online
- 27 May: Group meeting IV – online

Course material

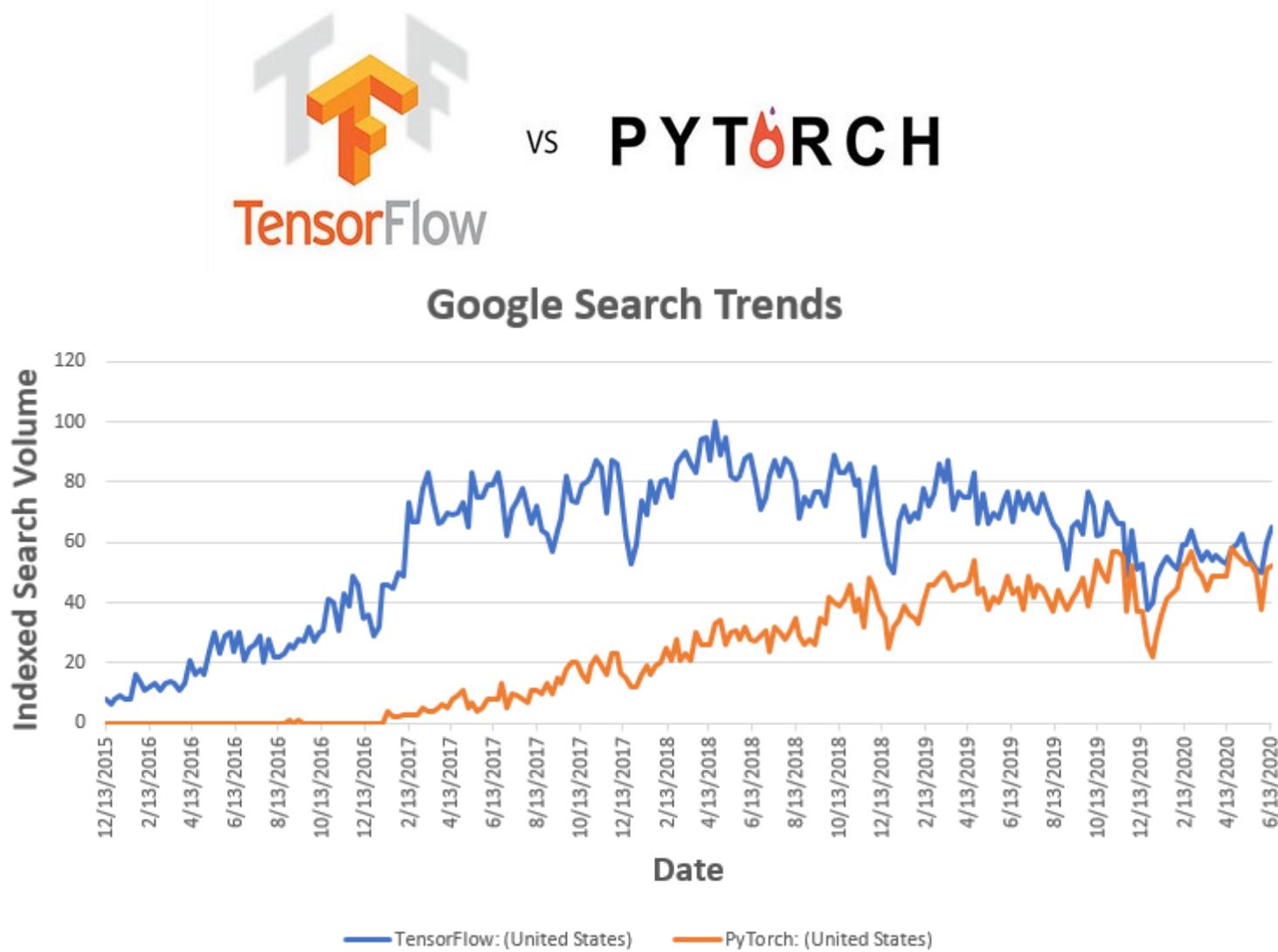
This course follows:

- MIT 6.S191: Introduction to Deep Learning by Alexander Amini and Ava Soleimany
- COMP3547: Reinforcement Learning by Chris G. Willcocks

Suggested book:
Goodfellow et al. 2016



Pytorch tutorials



What is deep learning?

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



MACHINE LEARNING

Ability to learn without explicitly being programmed



DEEP LEARNING

Extract patterns from data using neural networks



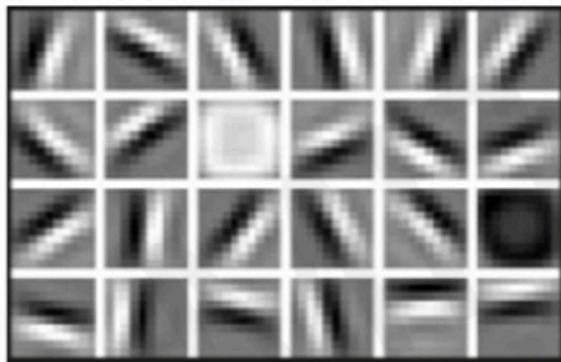
Why deep learning and Why now?

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



Lines & Edges

Mid Level Features



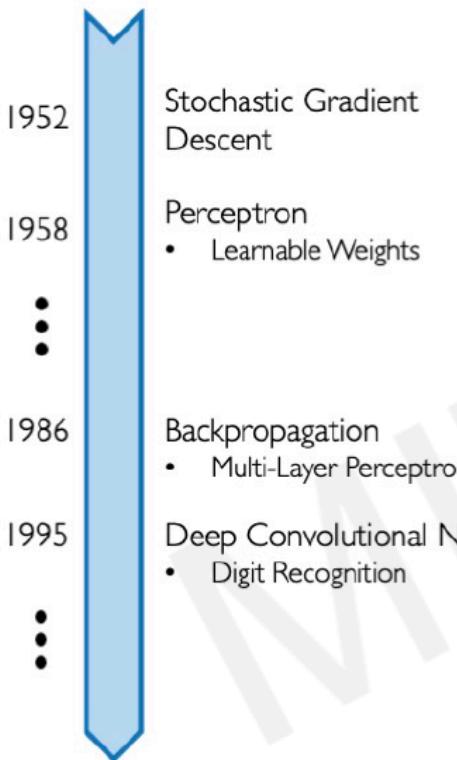
Eyes & Nose & Ears

High Level Features



Facial Structure

Why now?



Neural Networks date back decades, so why the resurgence?

I. Big Data

- Larger Datasets
- Easier Collection & Storage



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



3. Software

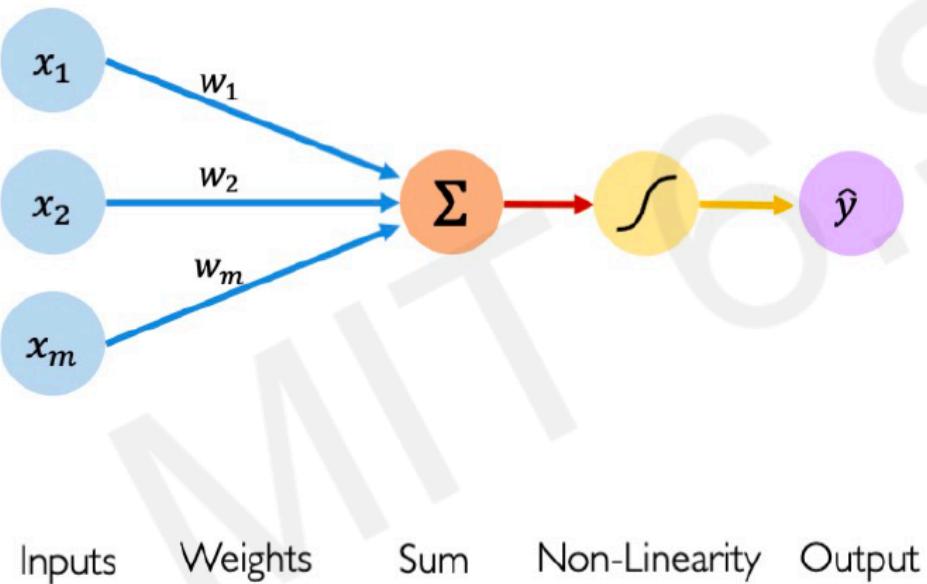
- Improved Techniques
- New Models
- Toolboxes



The Perceptron

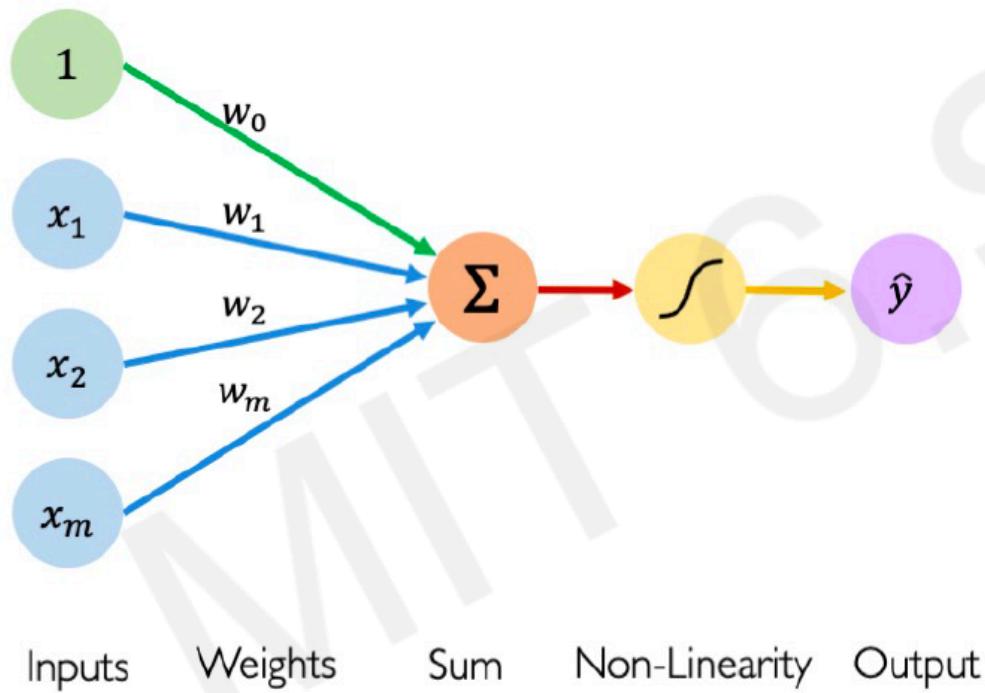
The structural building block of deep learning

The Perceptron: Forward Propagation



The diagram illustrates the computation of the output \hat{y} from inputs x_i and weights w_i through a non-linear activation function g . The inputs x_i and weights w_i are combined via a linear combination of inputs, represented by a red arrow pointing down. This result is then passed through a non-linear activation function g , represented by a yellow arrow pointing up, to produce the final output \hat{y} .

The Perceptron: Forward Propagation



Linear combination of inputs

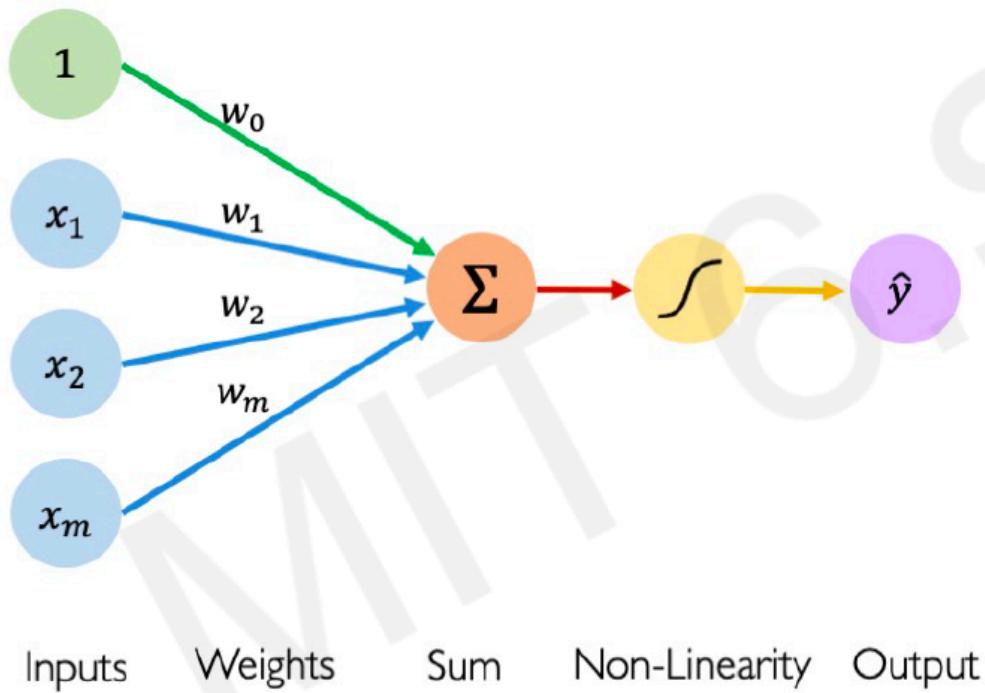
$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

Output

Non-linear activation function

Bias

The Perceptron: Forward Propagation

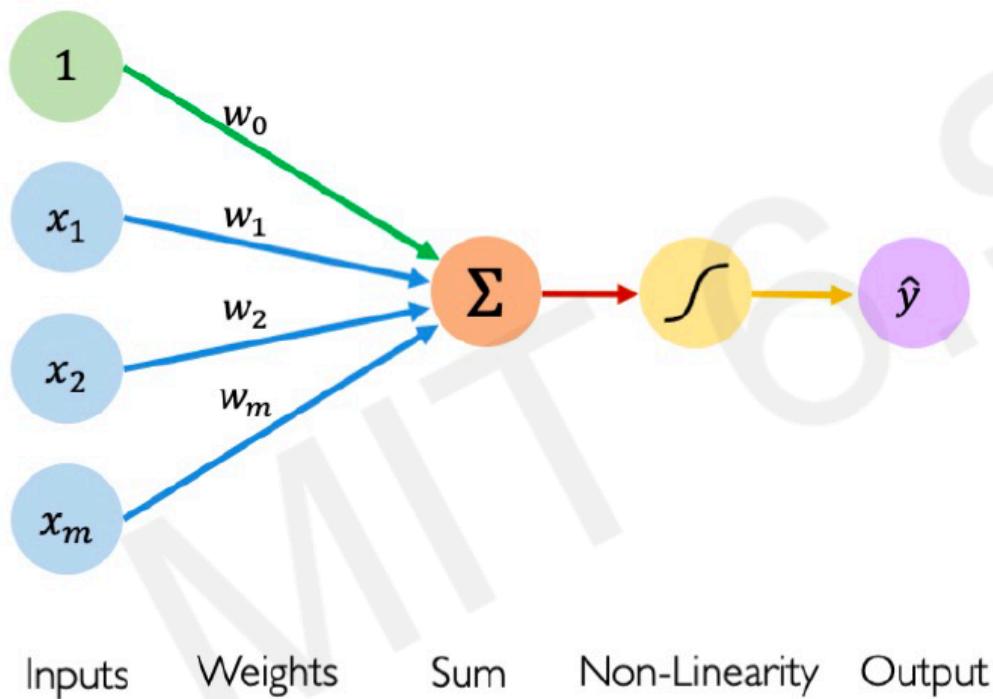


$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

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

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

The Perceptron: Forward Propagation

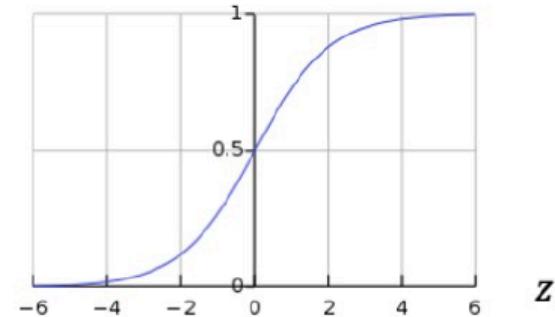


Activation Functions

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

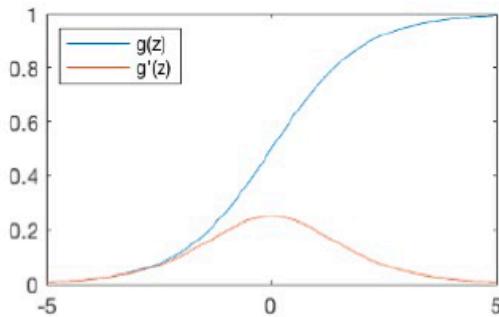
- Example: sigmoid function

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



Common activation functions

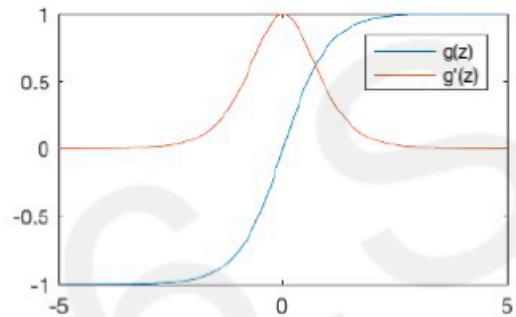
Sigmoid Function



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

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

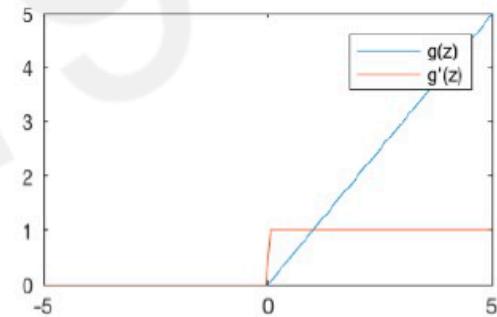
Hyperbolic Tangent



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

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

Rectified Linear Unit (ReLU)

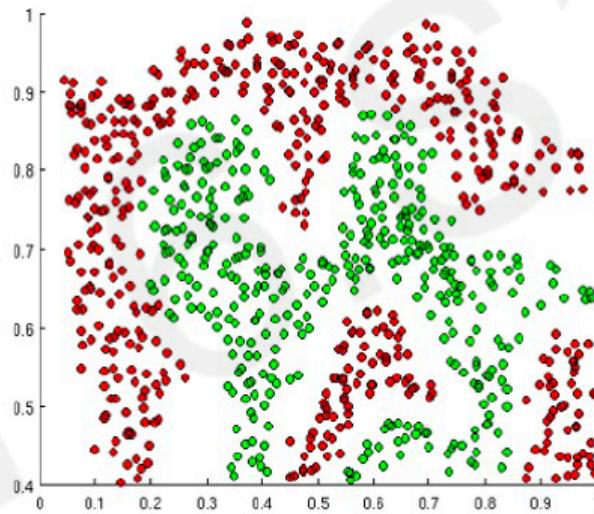


$$g(z) = \max(0, z)$$

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

Importance of activation functions

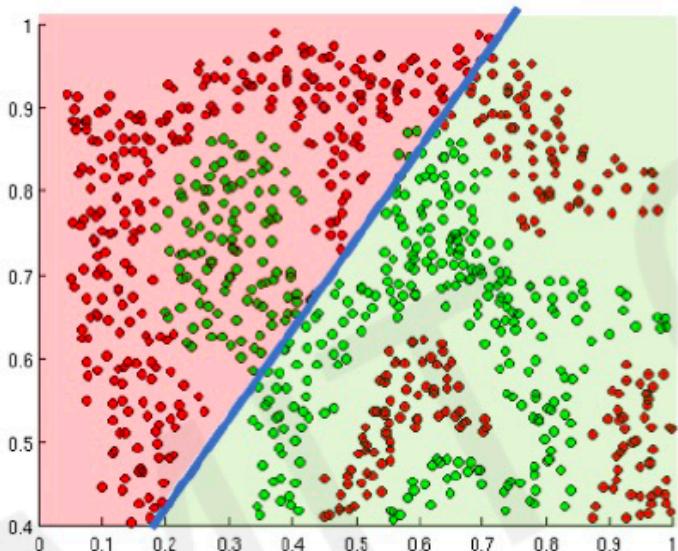
The purpose of activation functions is to **introduce non-linearities** into the network



What if we wanted to build a neural network to
distinguish green vs red points?

Importance of activation functions

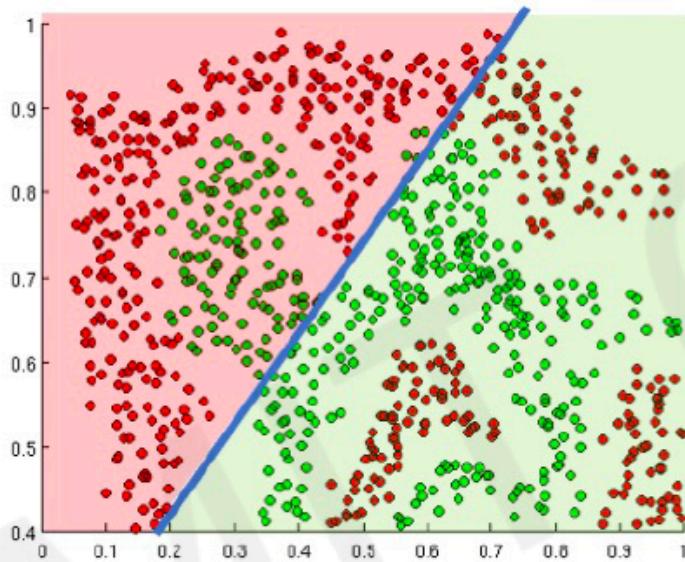
The purpose of activation functions is to **introduce non-linearities** into the network



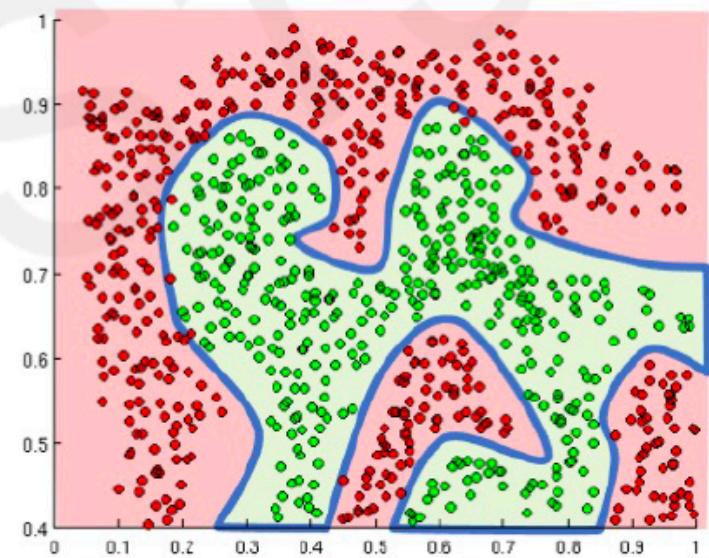
Linear activation functions produce linear decisions no matter the network size

Importance of activation functions

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

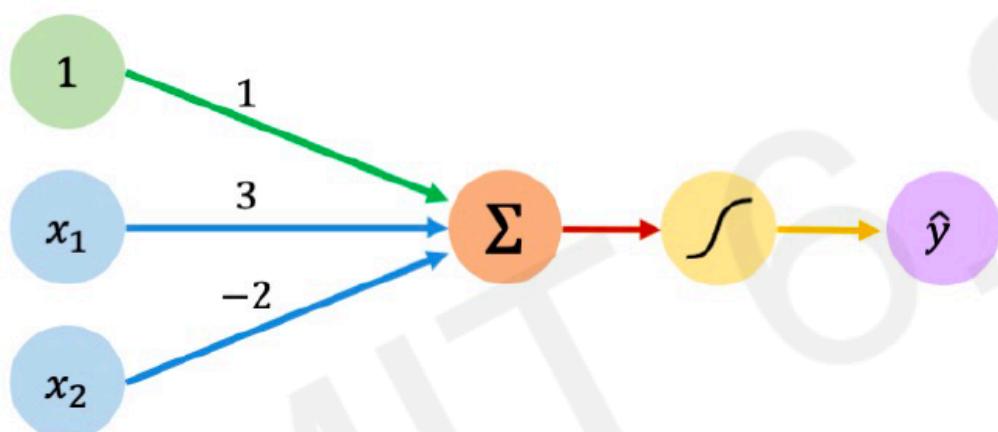


Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

The perceptron: example

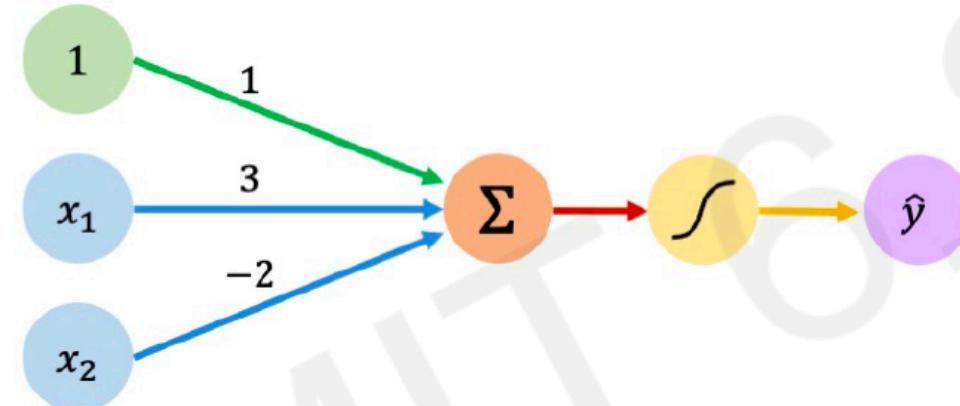


We have: $w_0 = 1$ and $\mathbf{w} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

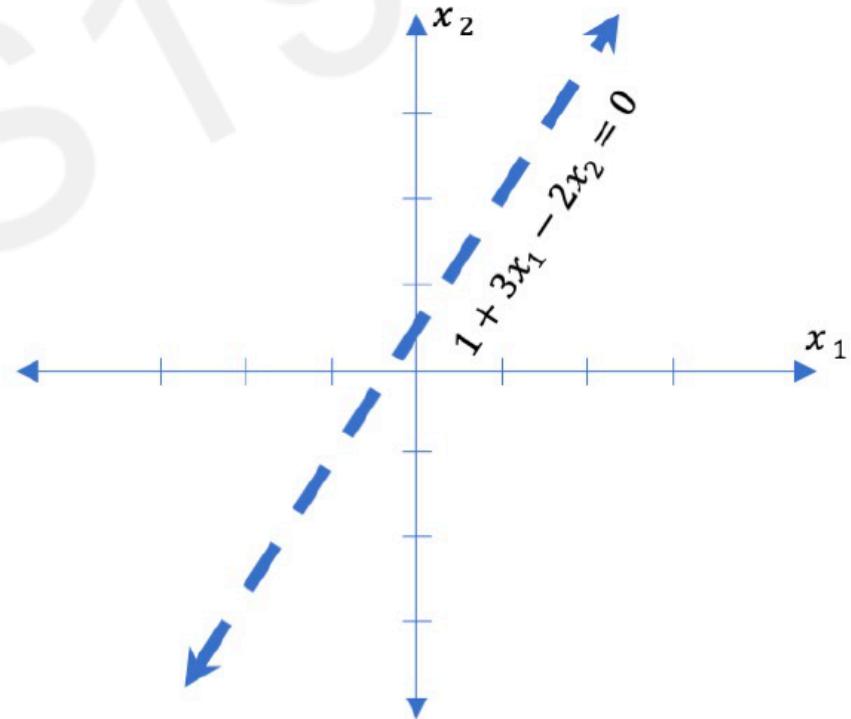
$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{w}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

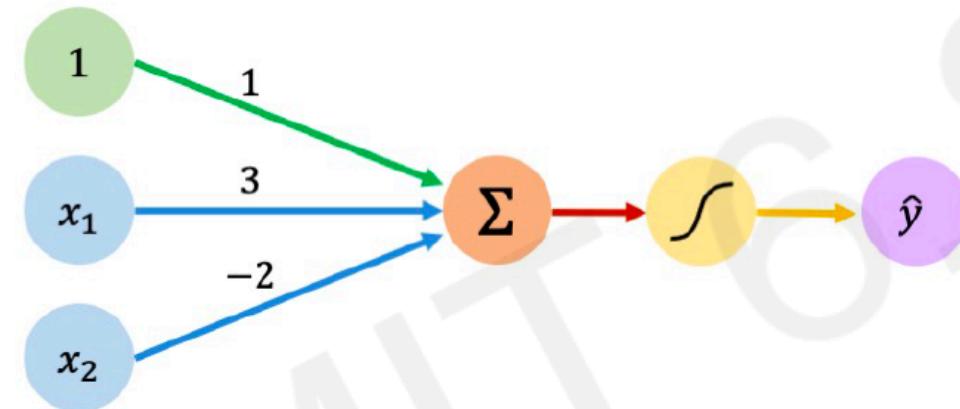
The perceptron: example



$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

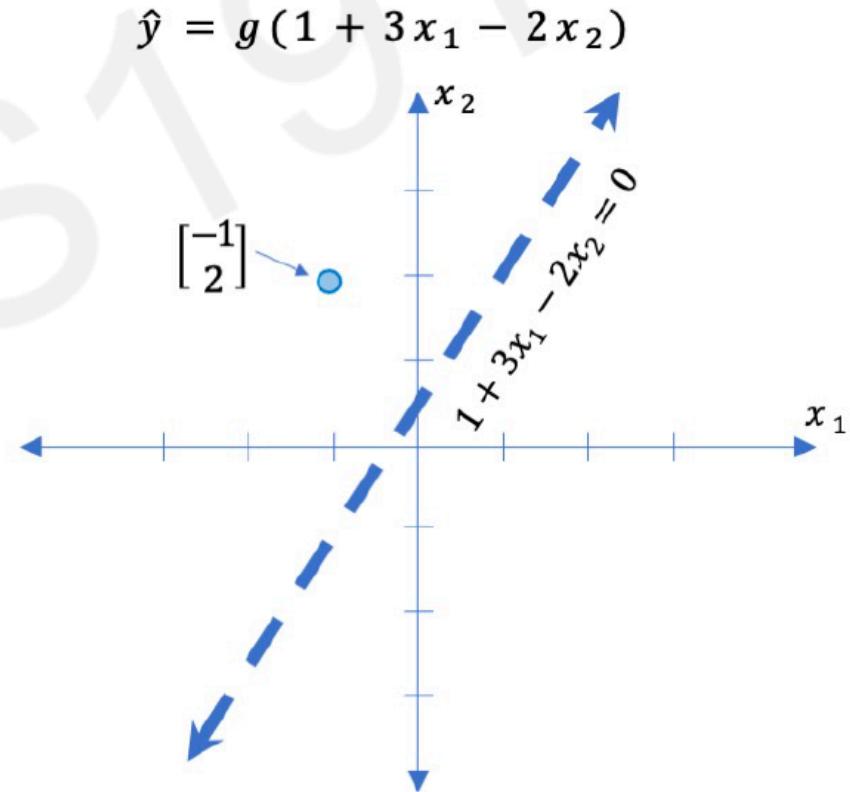


The perceptron: example

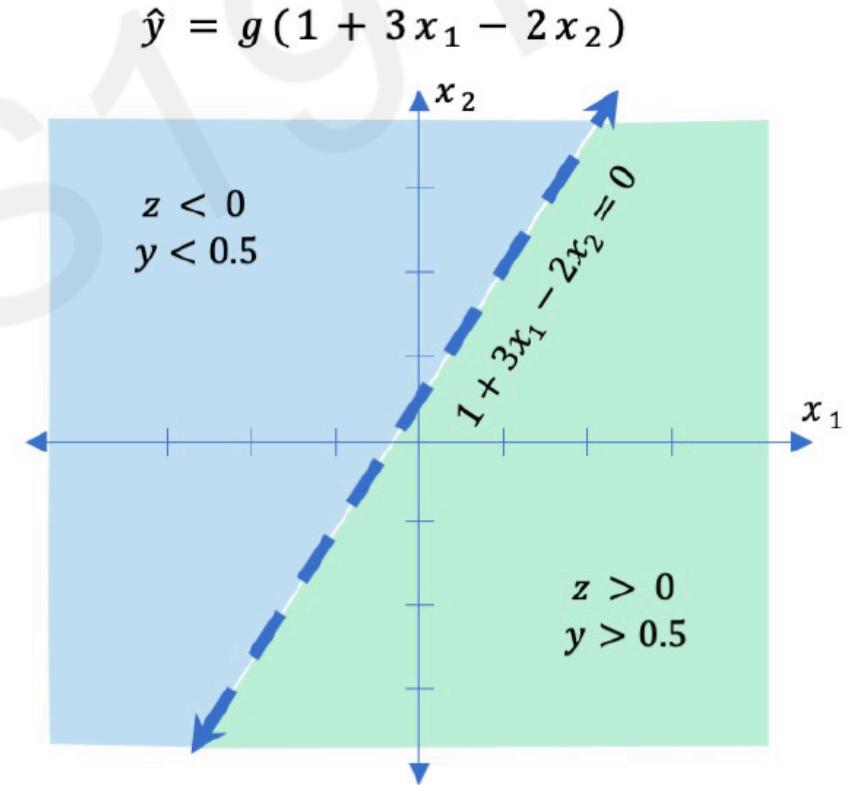
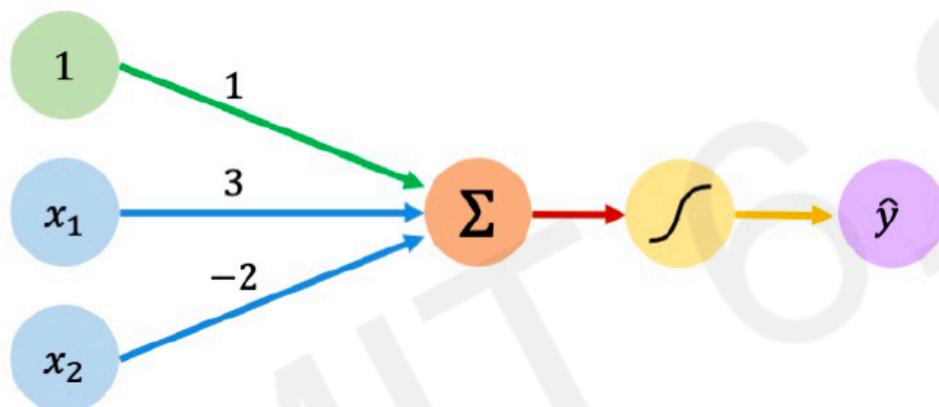


Assume we have input: $\mathbf{x} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$



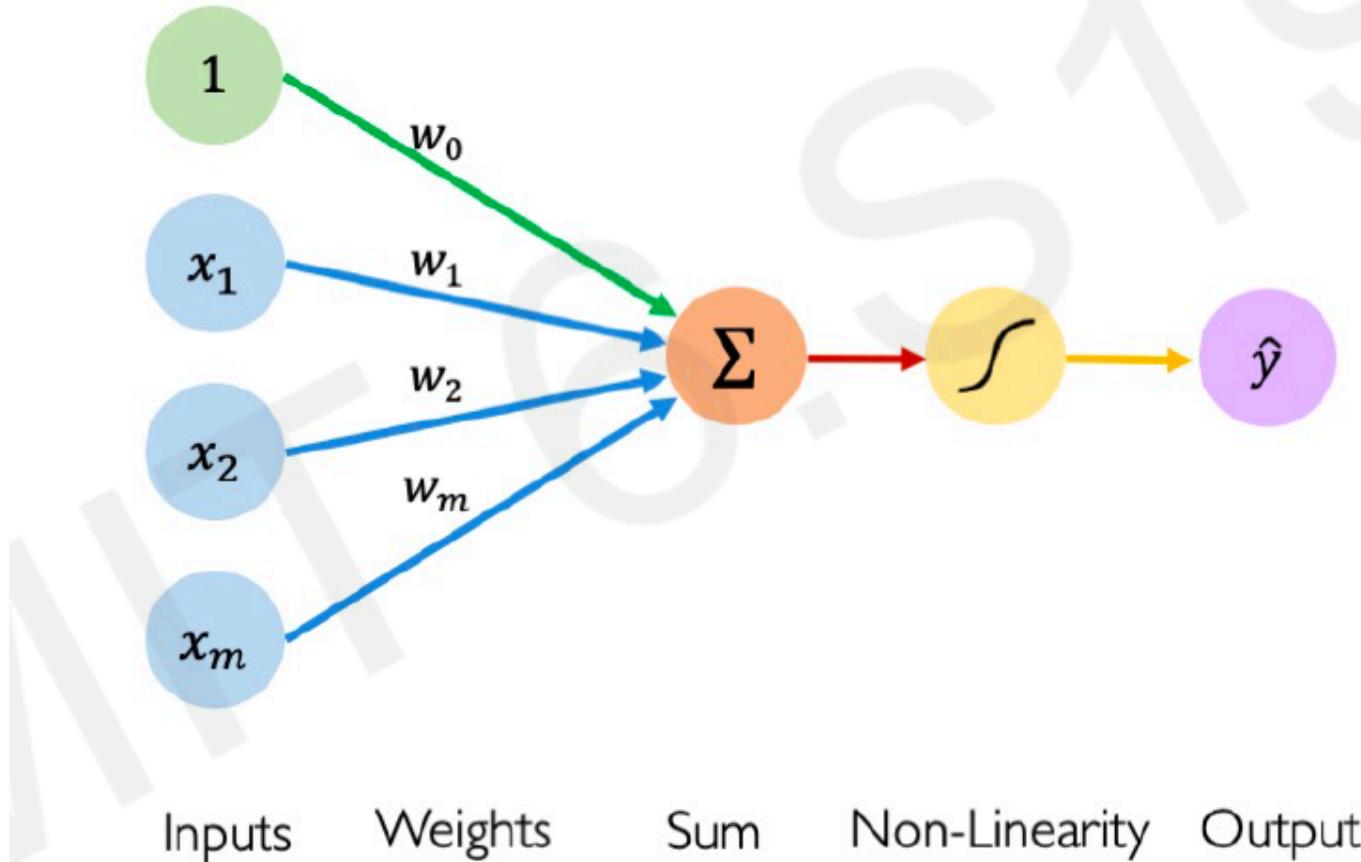
The perceptron: example



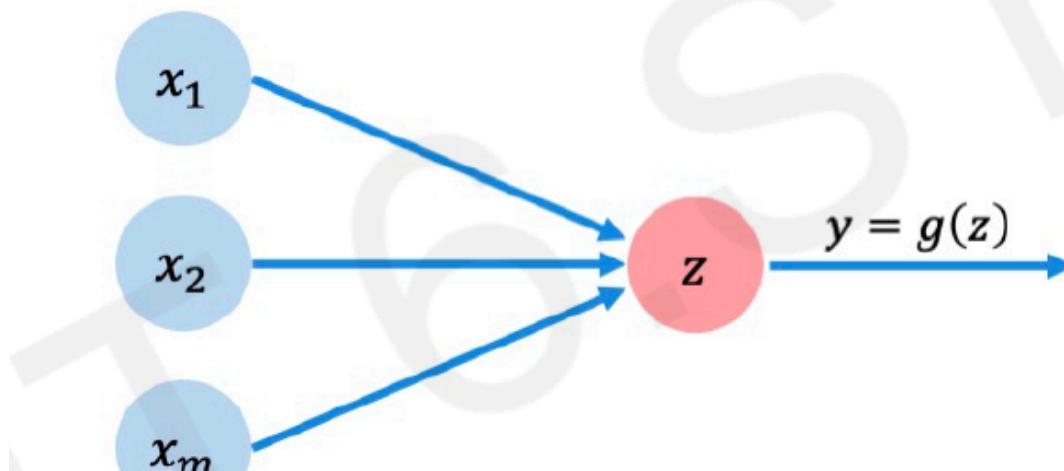
Building Neural Networks with Perceptrons

The perceptron: simplified

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



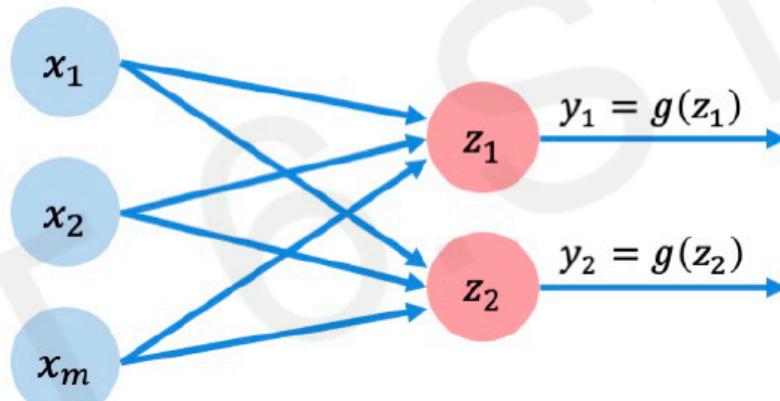
The perceptron: simplified



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

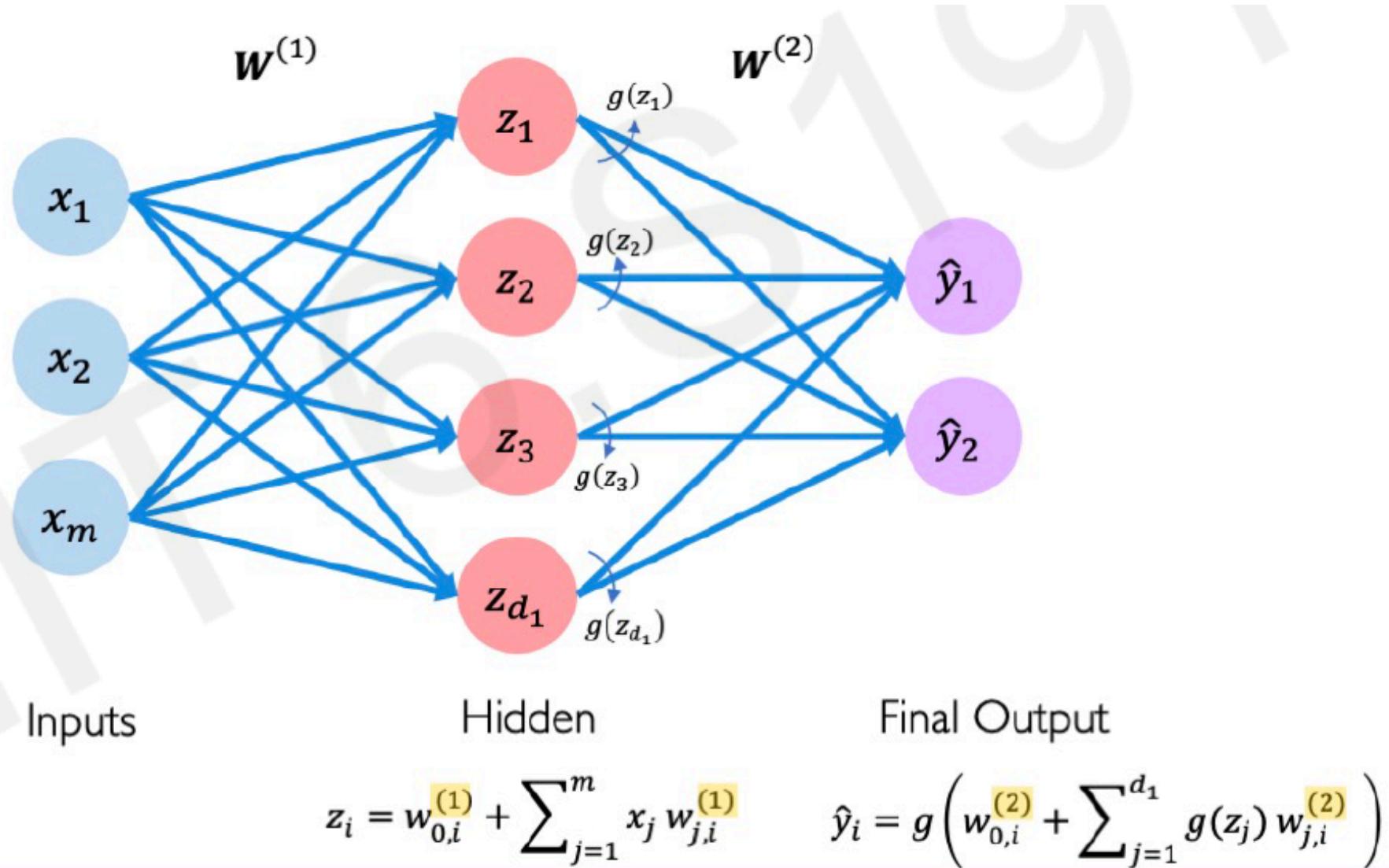
Multi output perceptron

Because all inputs are densely connected to all outputs, these layers are called **Dense** layers

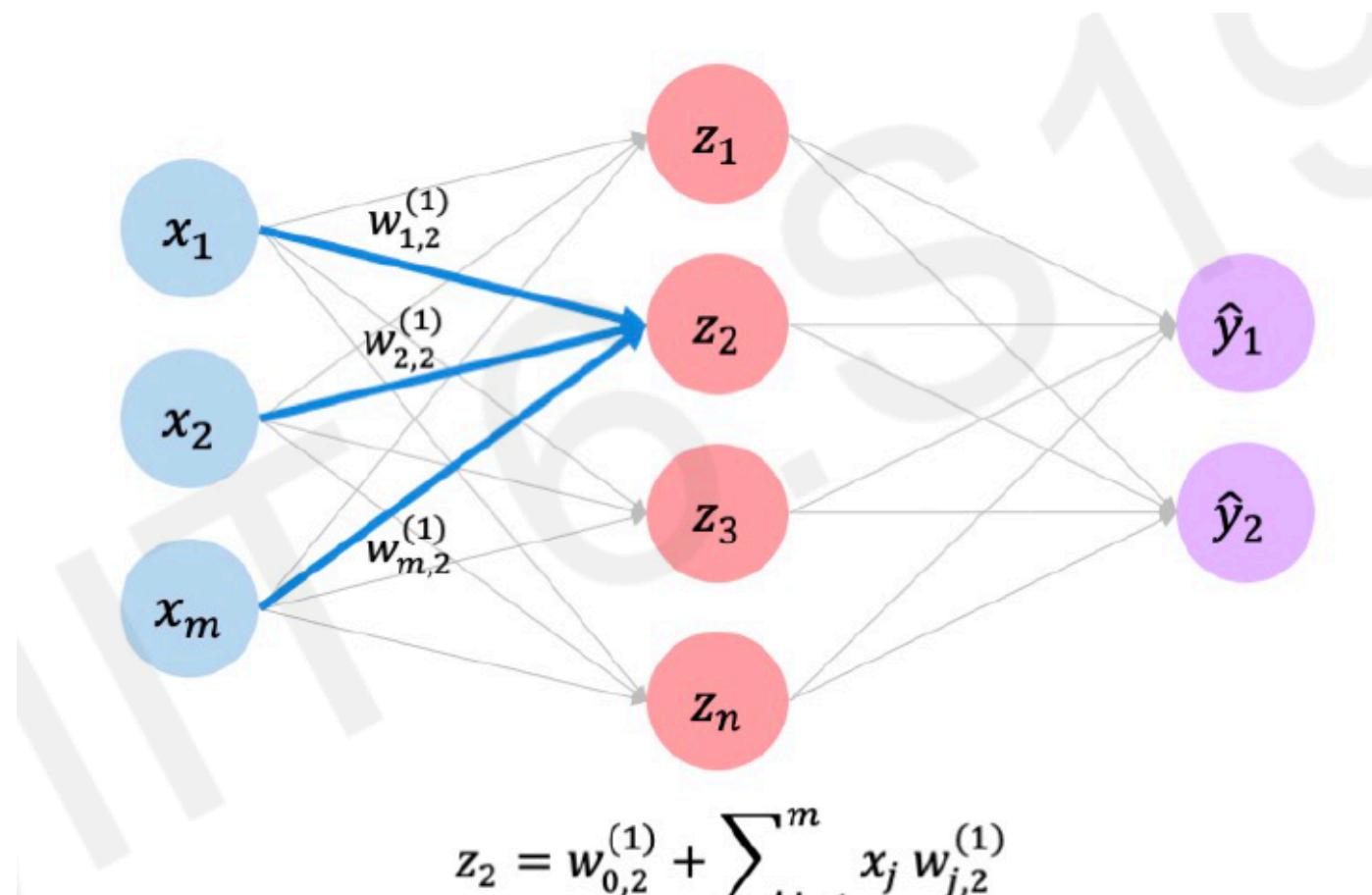


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

Single layer neural network

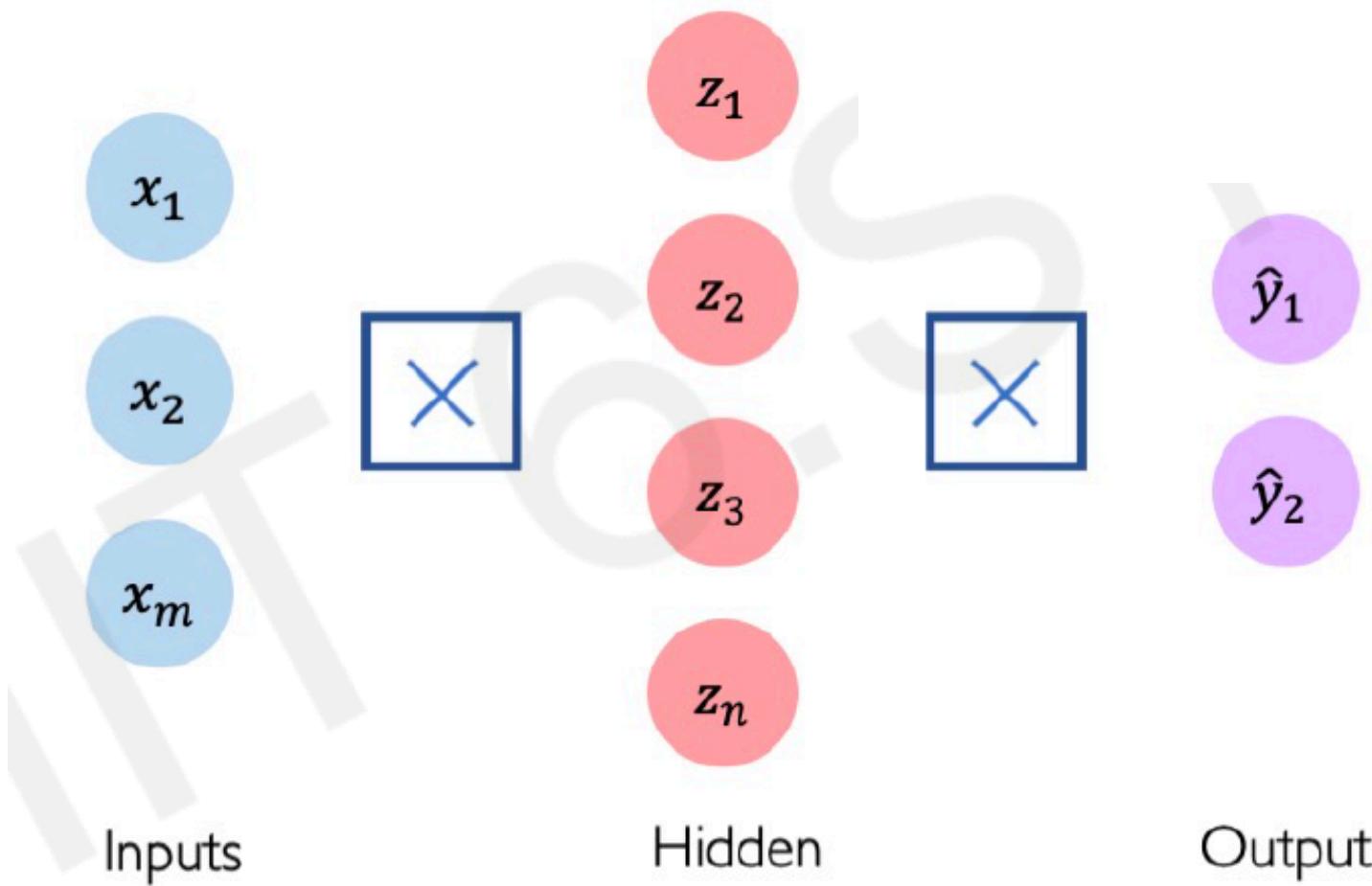


Single layer neural network

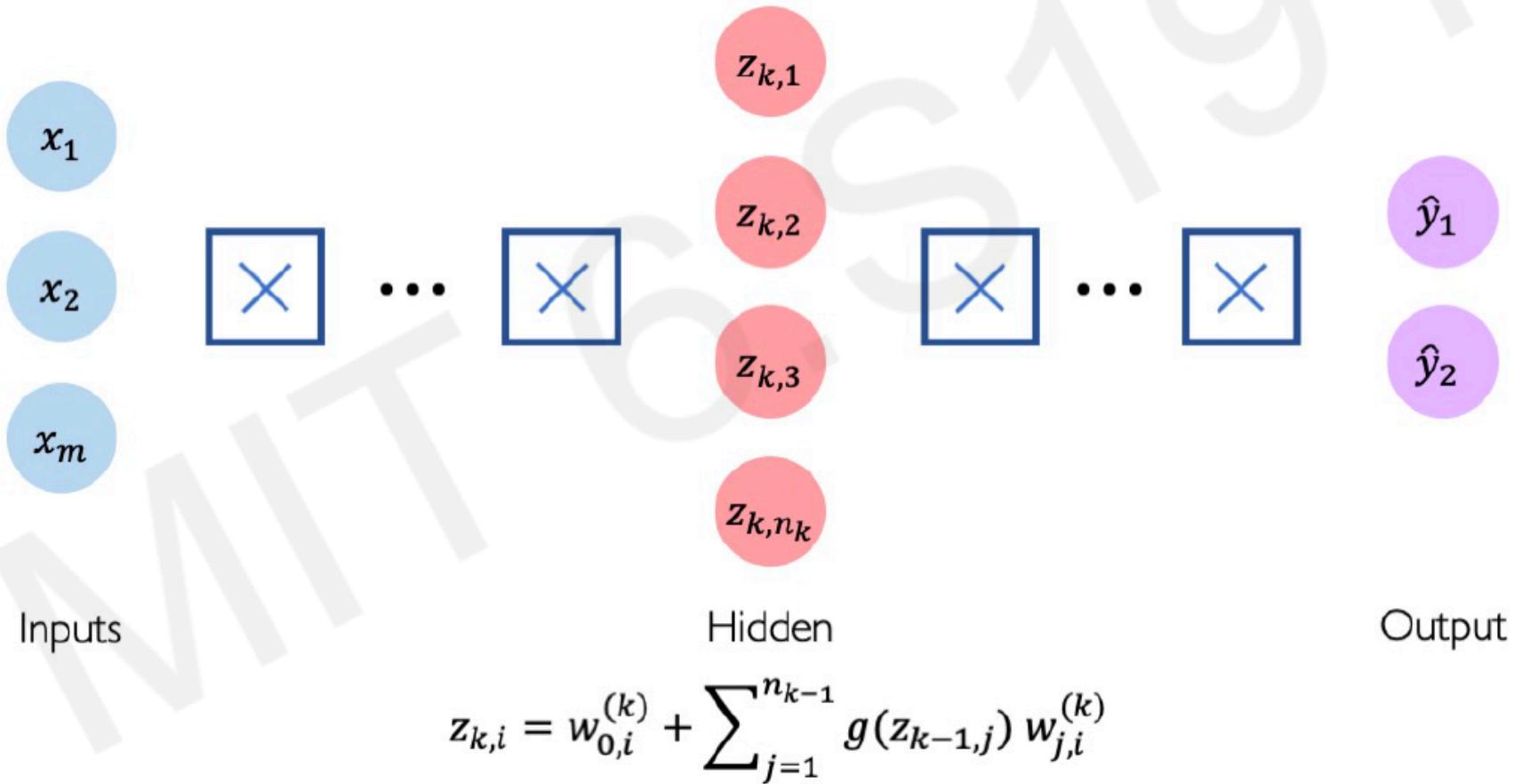


$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$

Multi output perceptron



Deep neural network



Applying Neural Networks

Example problem

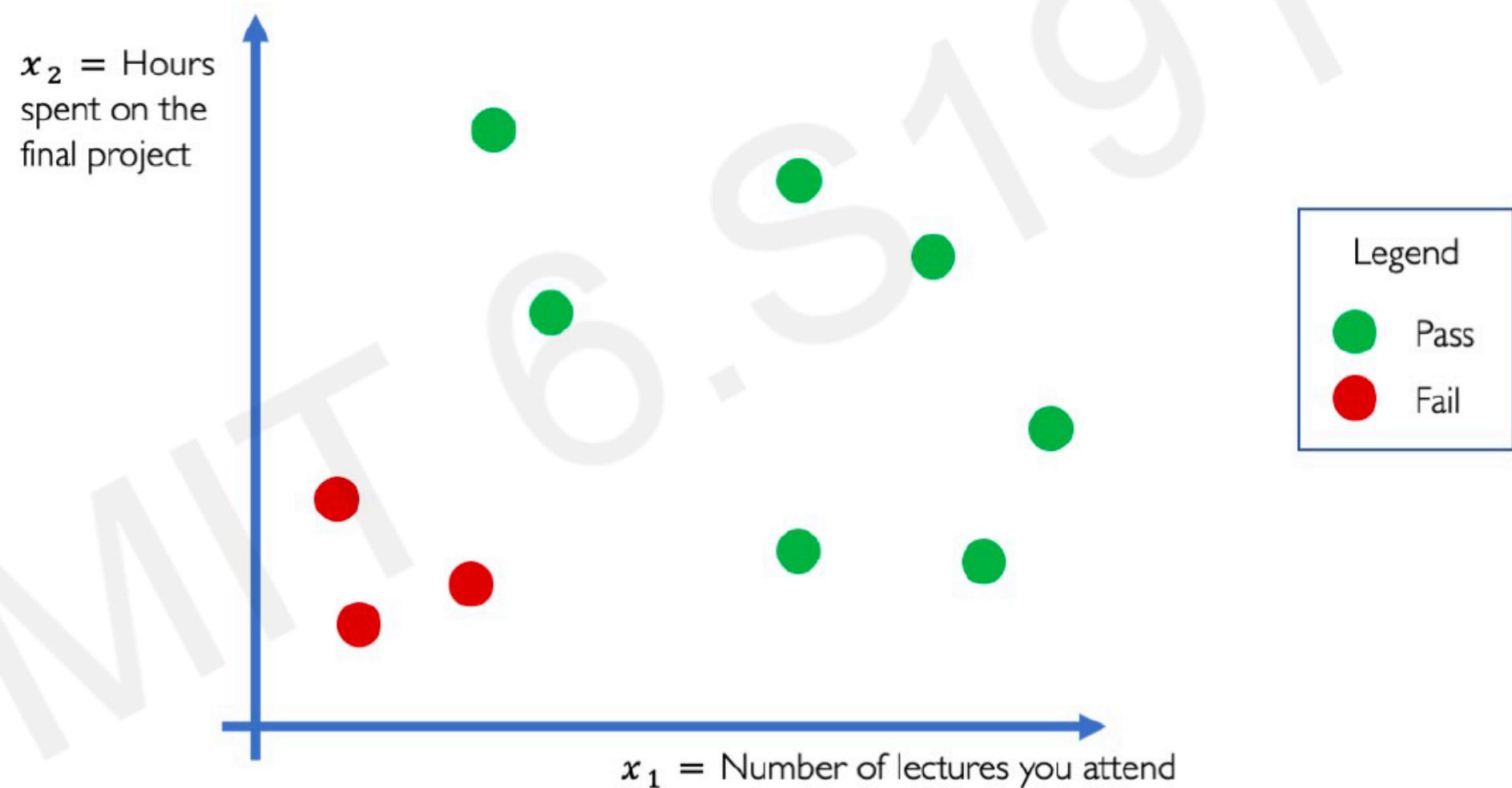
Will I pass this class?

Let's start with a simple two feature model

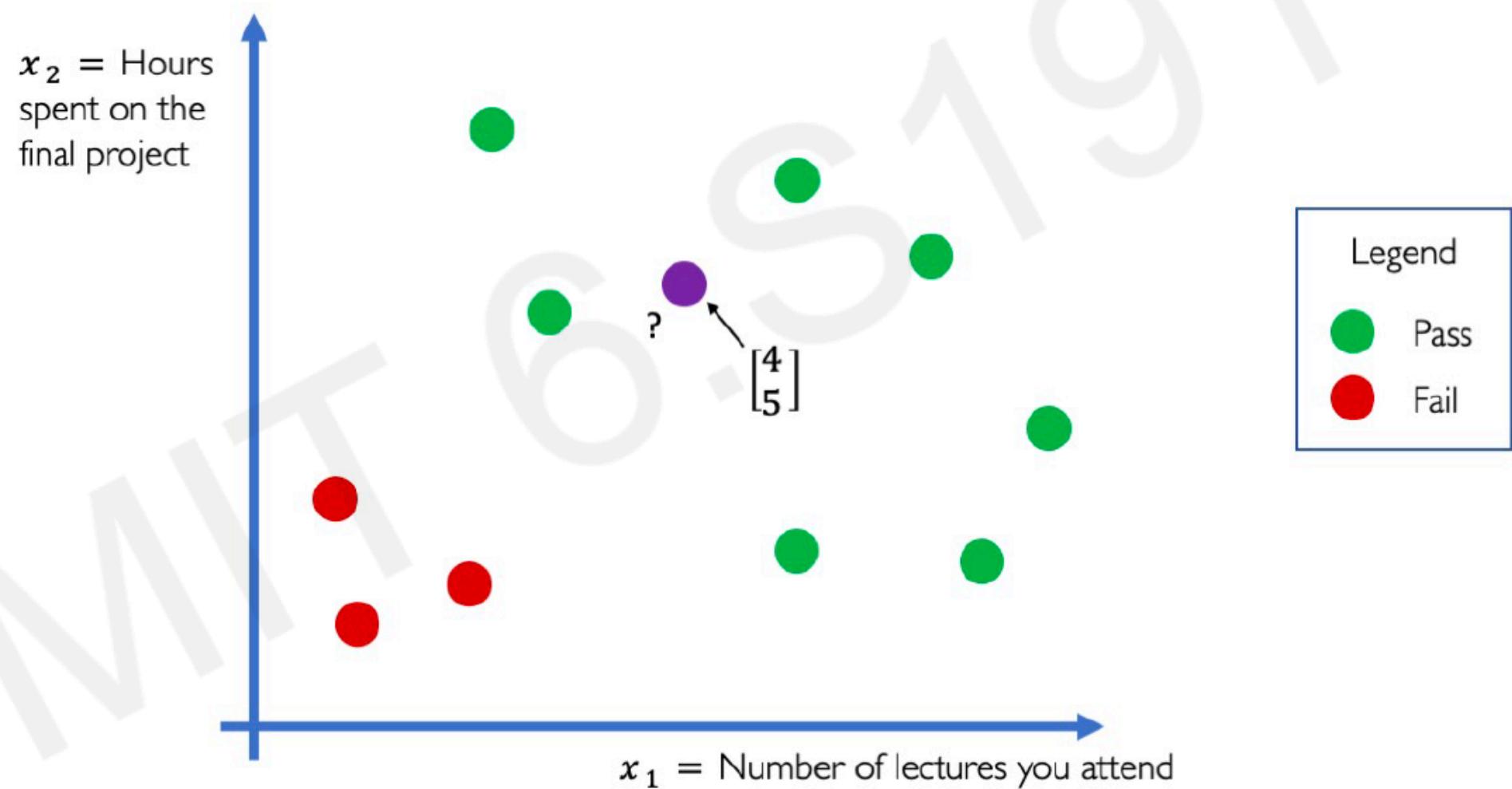
x_1 = Number of lectures you attend

x_2 = Hours spent on the final project

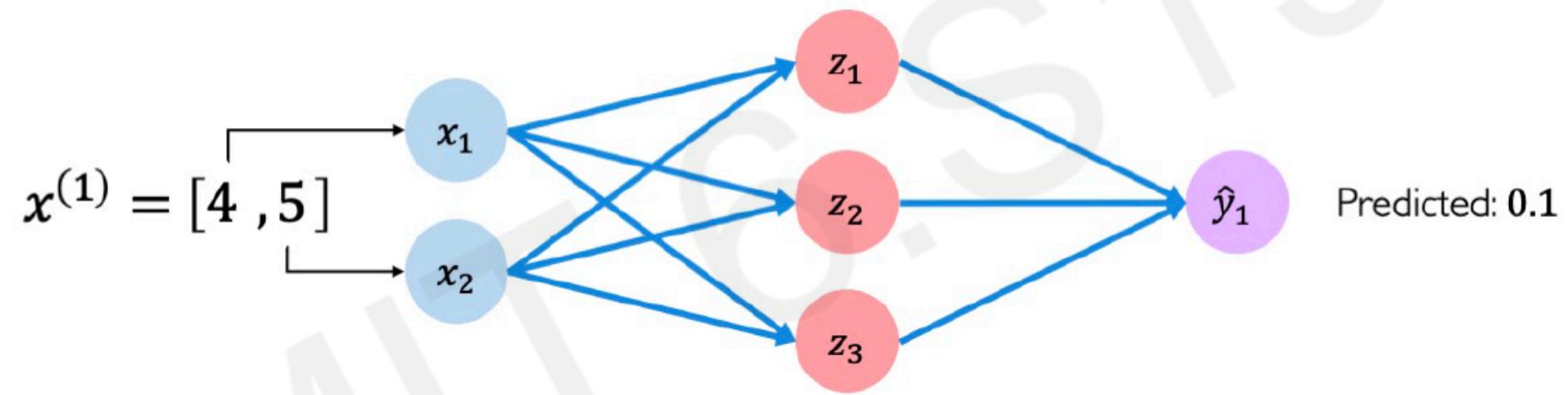
Example problem: will I pass this class?



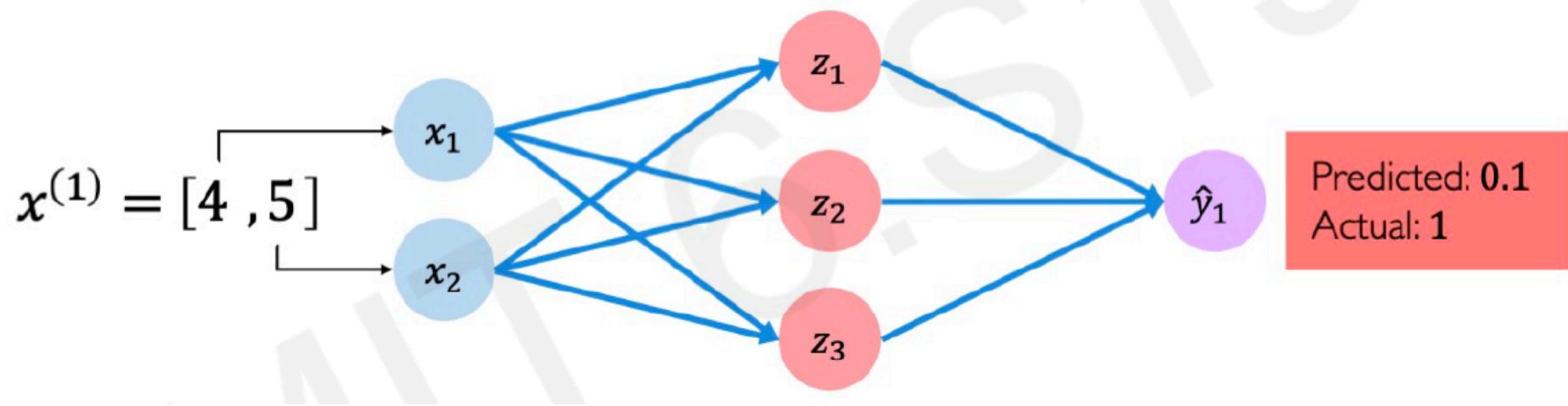
Example problem: will I pass this class?



Example problem: will I pass this class?

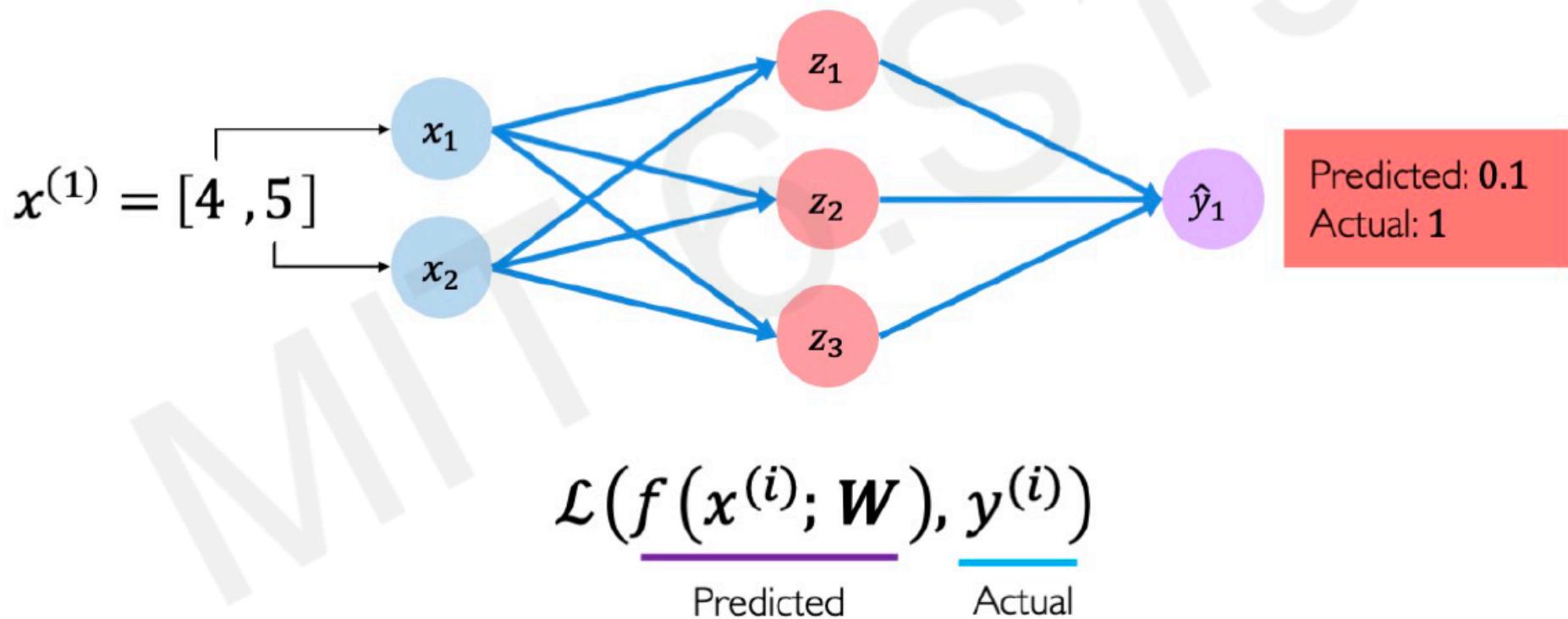


Example problem: will I pass this class?



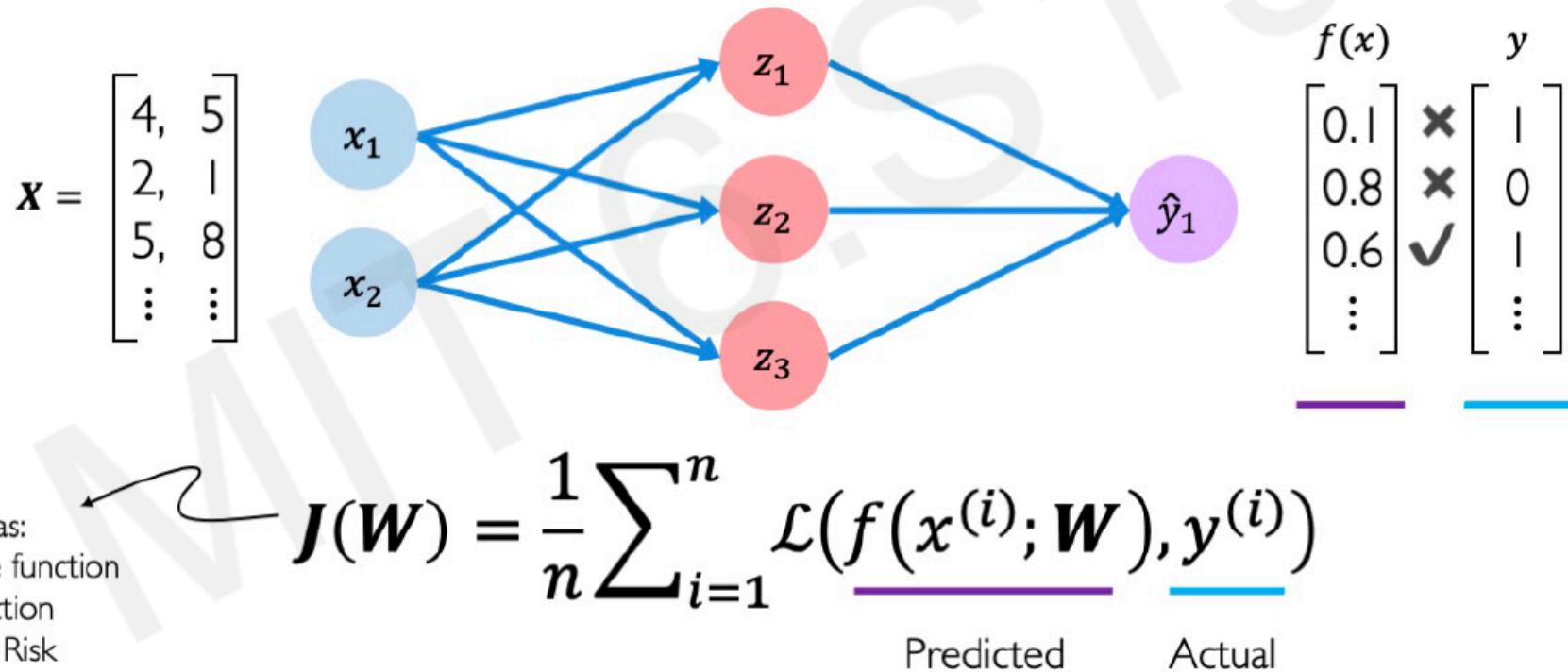
Quantifying loss

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



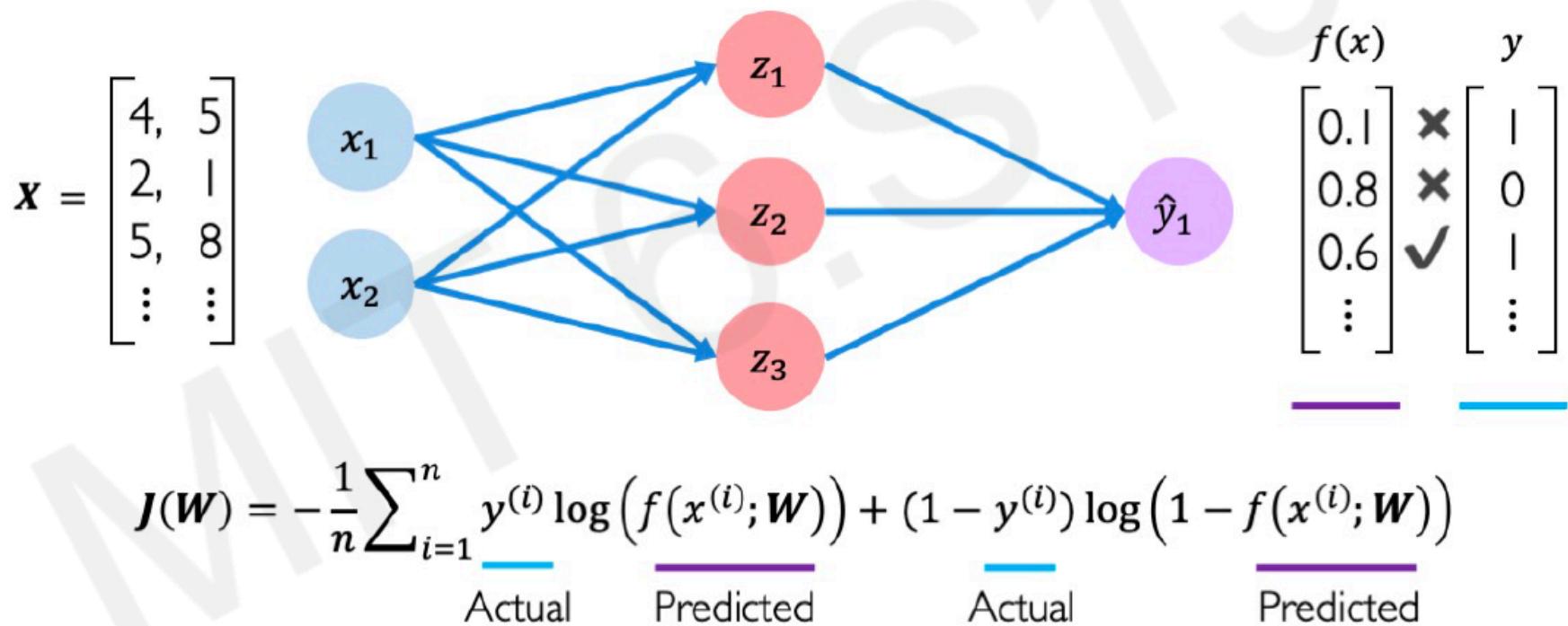
Empirical loss

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



Binary cross entropy loss

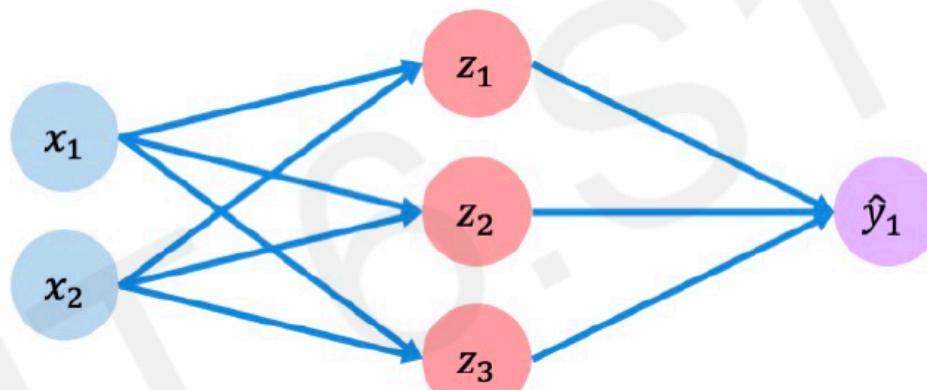
Cross entropy loss can be used with models that output a probability between 0 and 1



Mean squared error loss

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

$$\mathbf{x} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \frac{\underline{(y^{(i)} - f(x^{(i)}; \mathbf{W}))^2}}{\text{Actual} \quad \text{Predicted}}$$

$f(x)$	y
30	✗ 90
80	✗ 20
85	✓ 95
\vdots	\vdots

Final Grades
(percentage)

Training Neural Networks

Loss optimization

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

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

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

Loss optimization

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

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

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



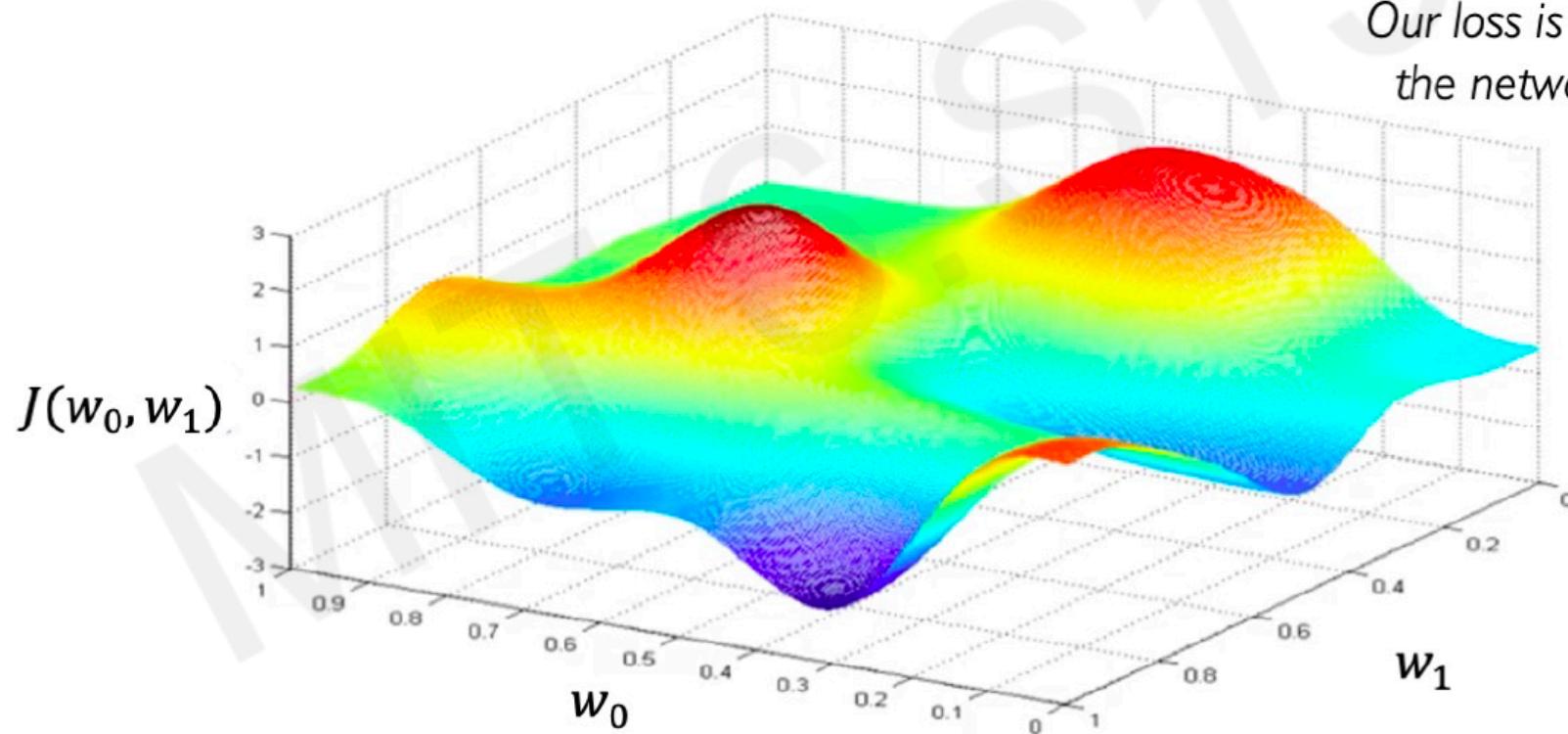
Remember:

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

Loss optimization

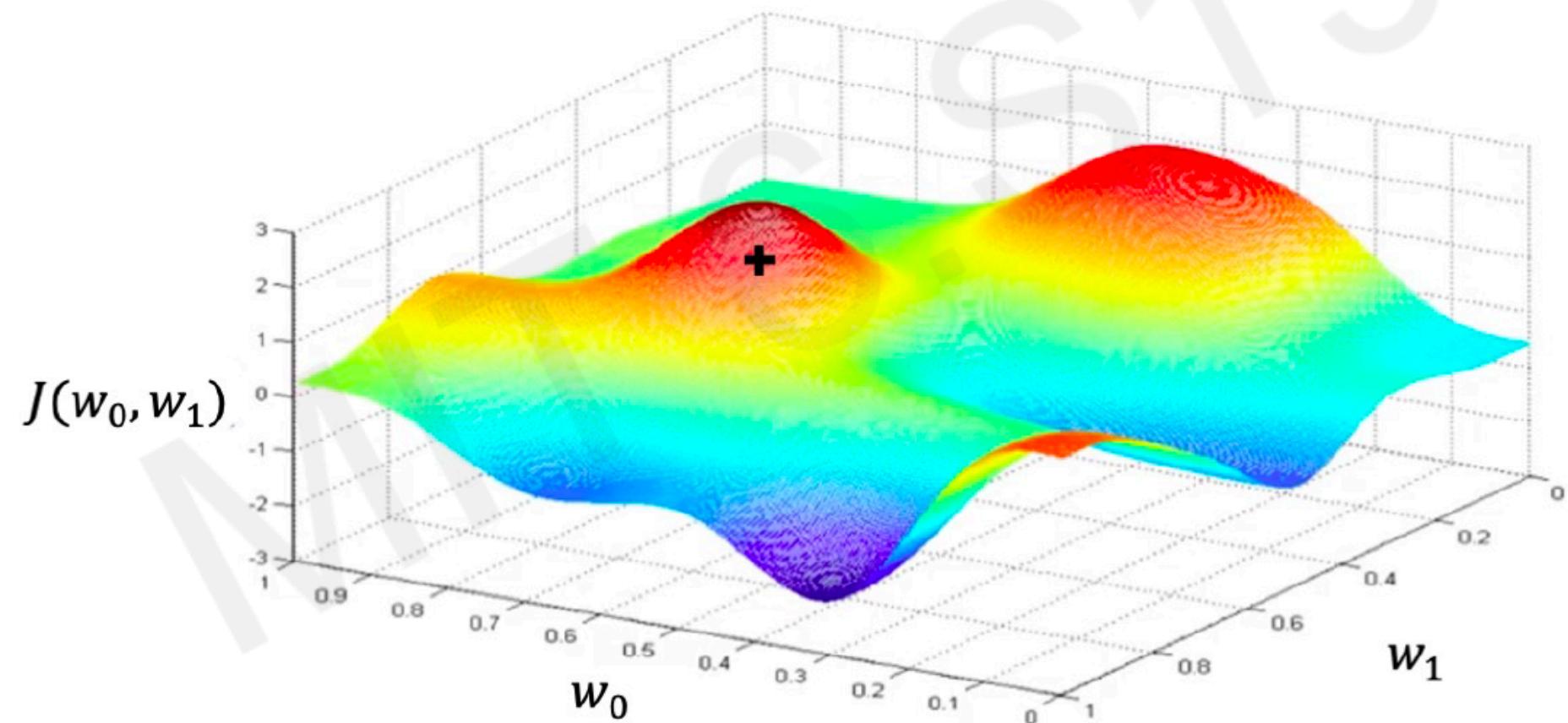
$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$

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



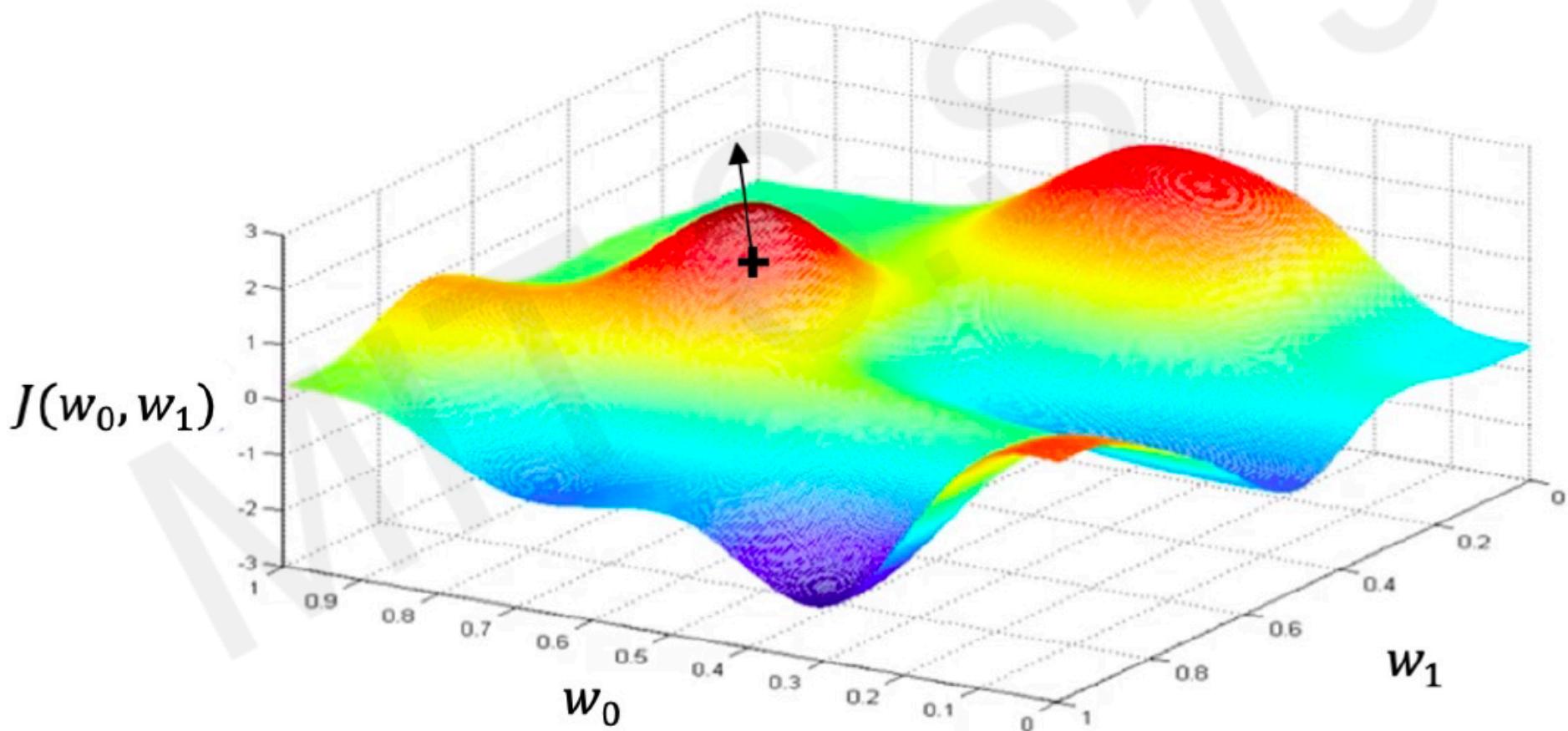
Loss optimization

Randomly pick an initial (w_0, w_1)



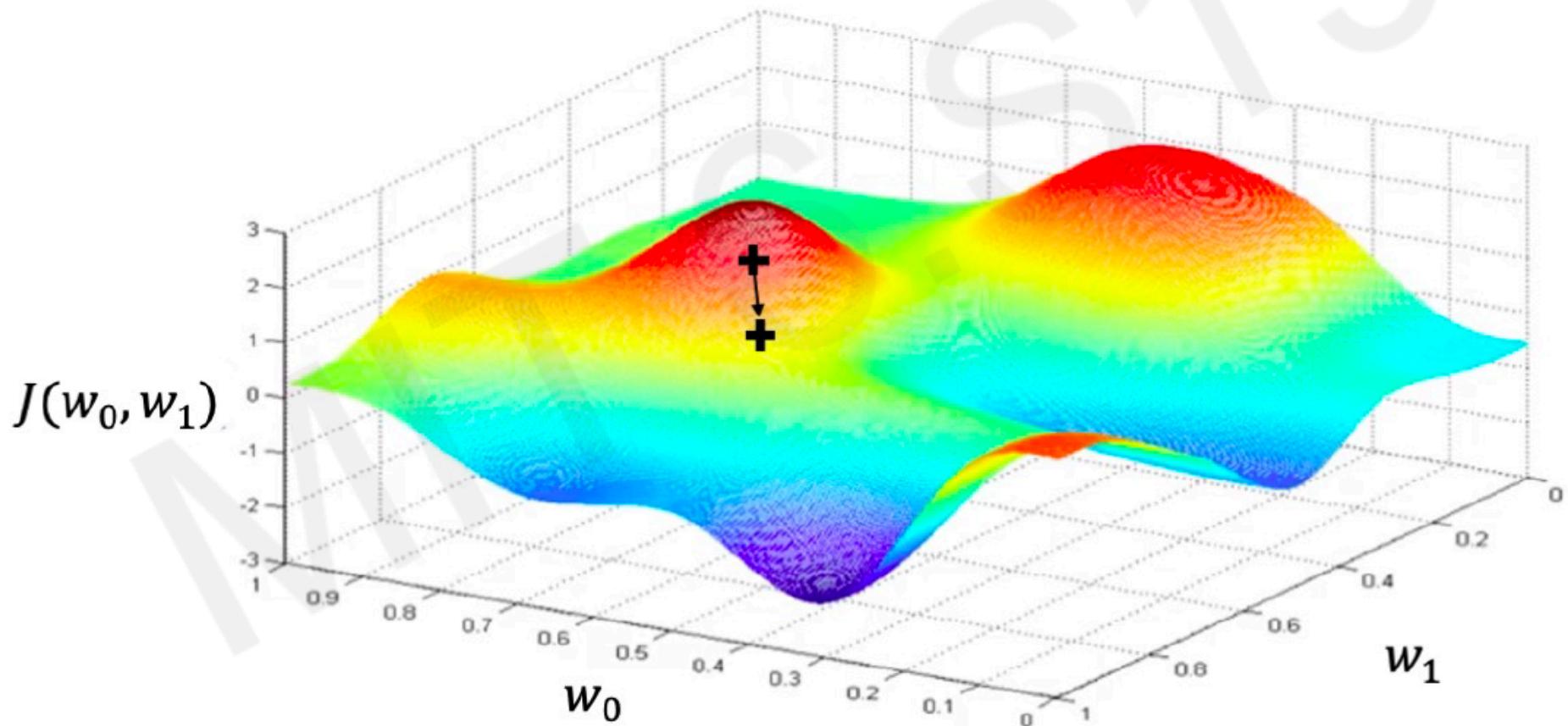
Loss optimization

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



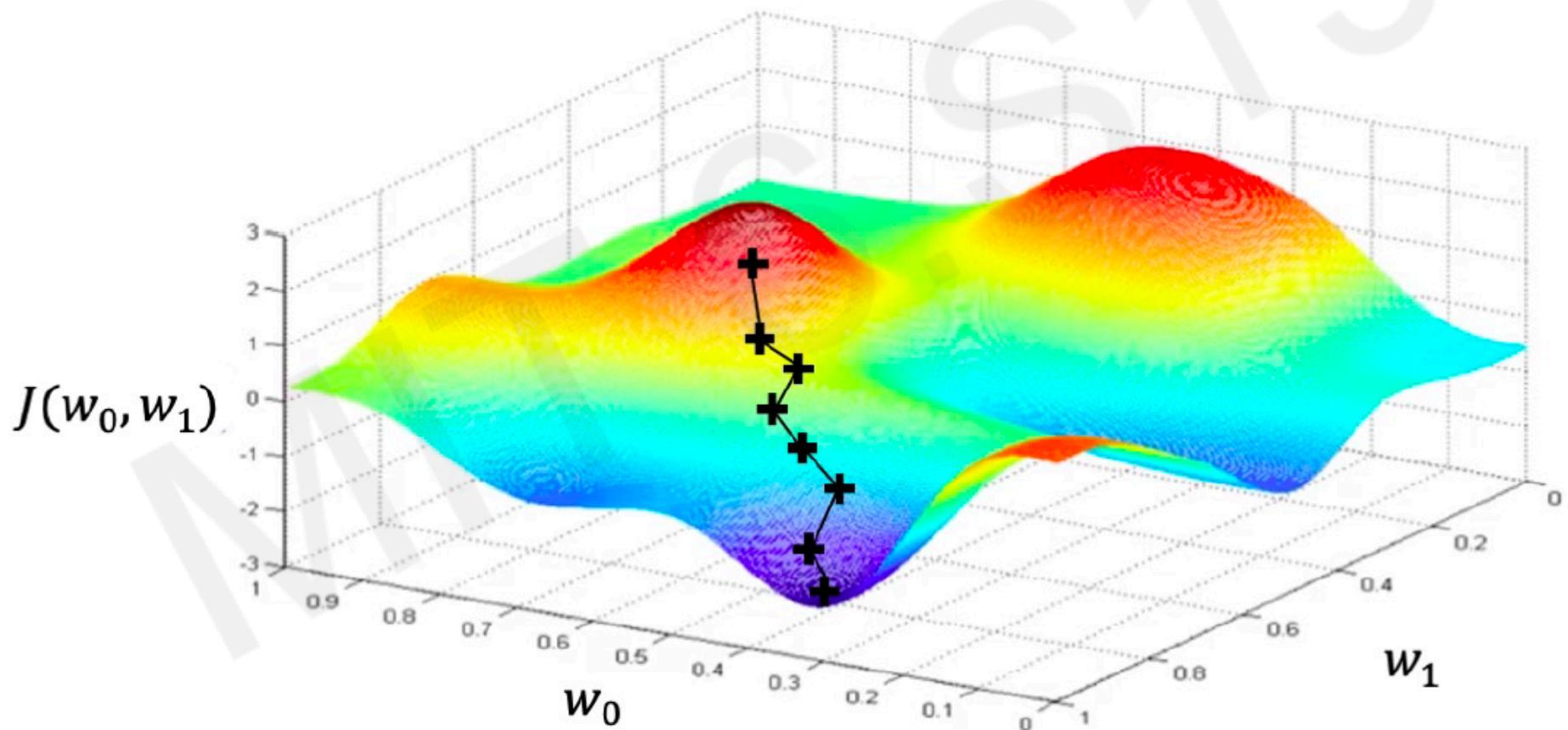
Loss optimization

Take small step in opposite direction of gradient



Gradient descent

Repeat until convergence



Gradient descent

Algorithm

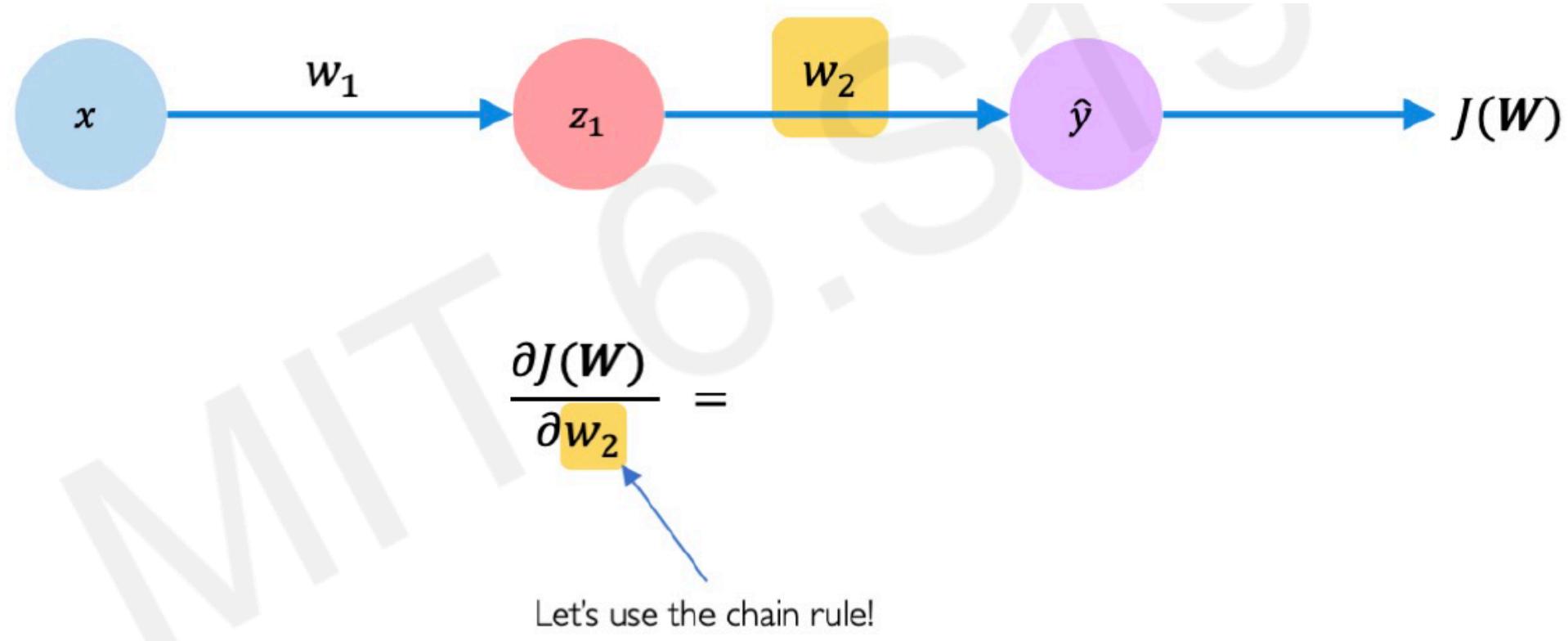
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Computing gradients: backpropagation

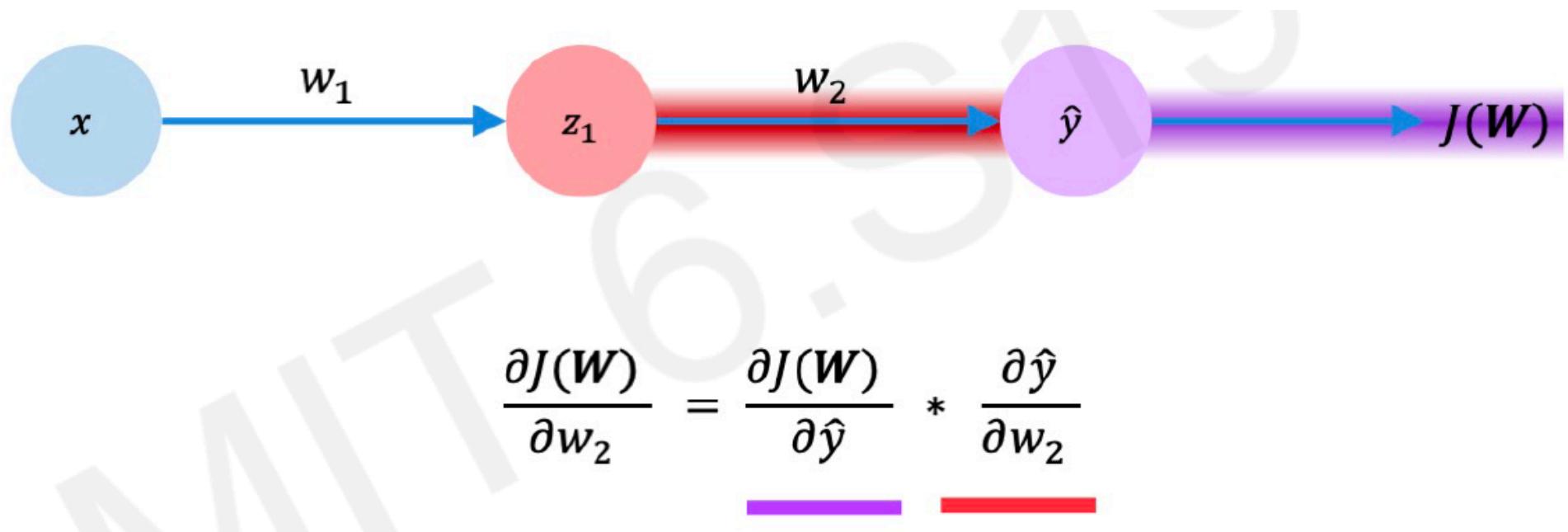


How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

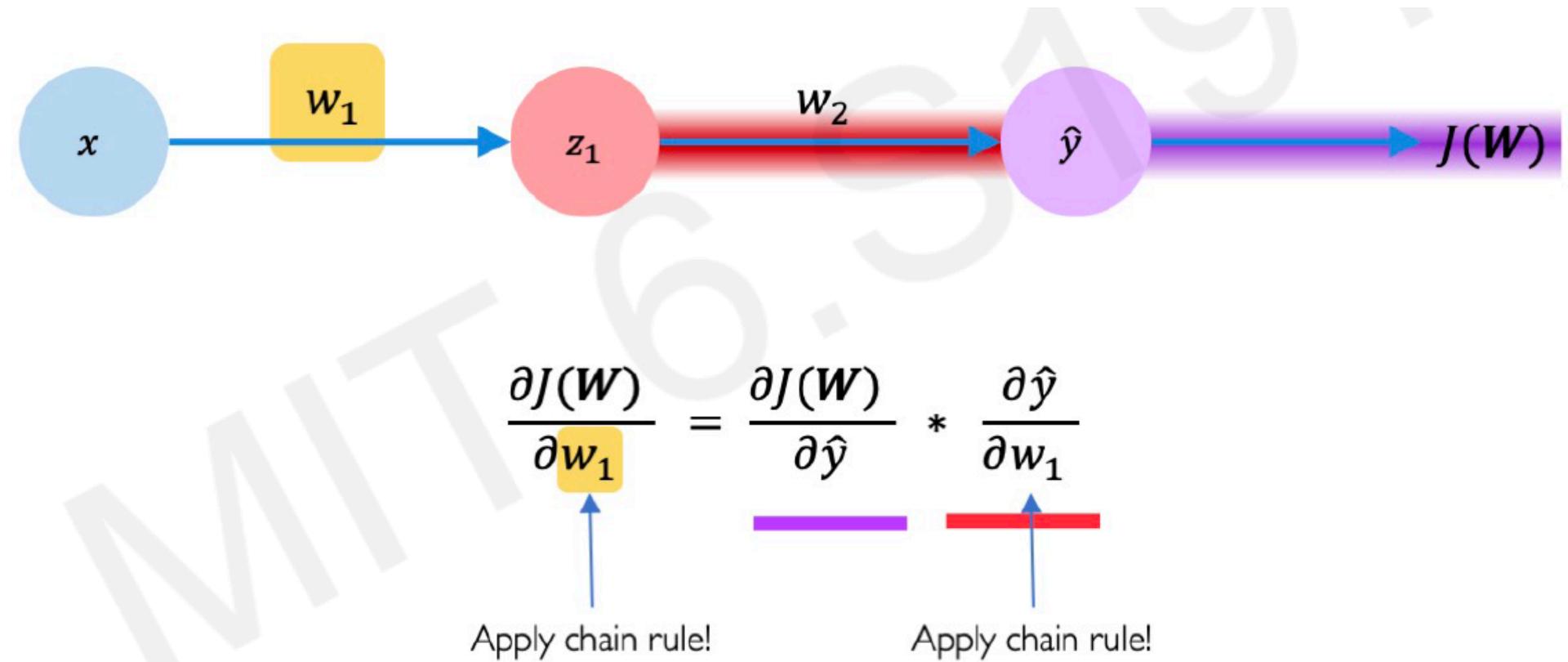
Computing gradients: backpropagation



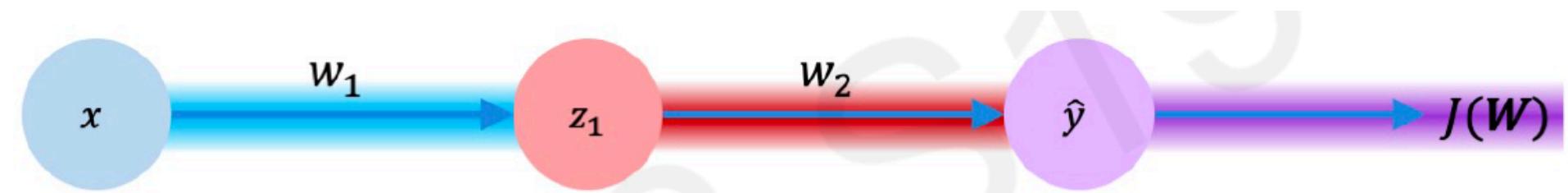
Computing gradients: backpropagation



Computing gradients: backpropagation

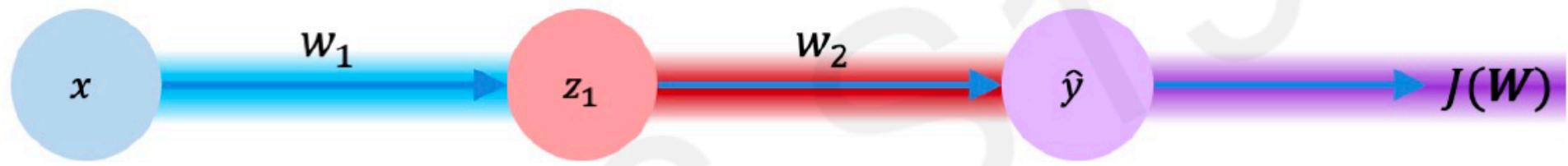


Computing gradients: backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple bar}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red bar}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{cyan bar}}$$

Computing gradients: backpropagation

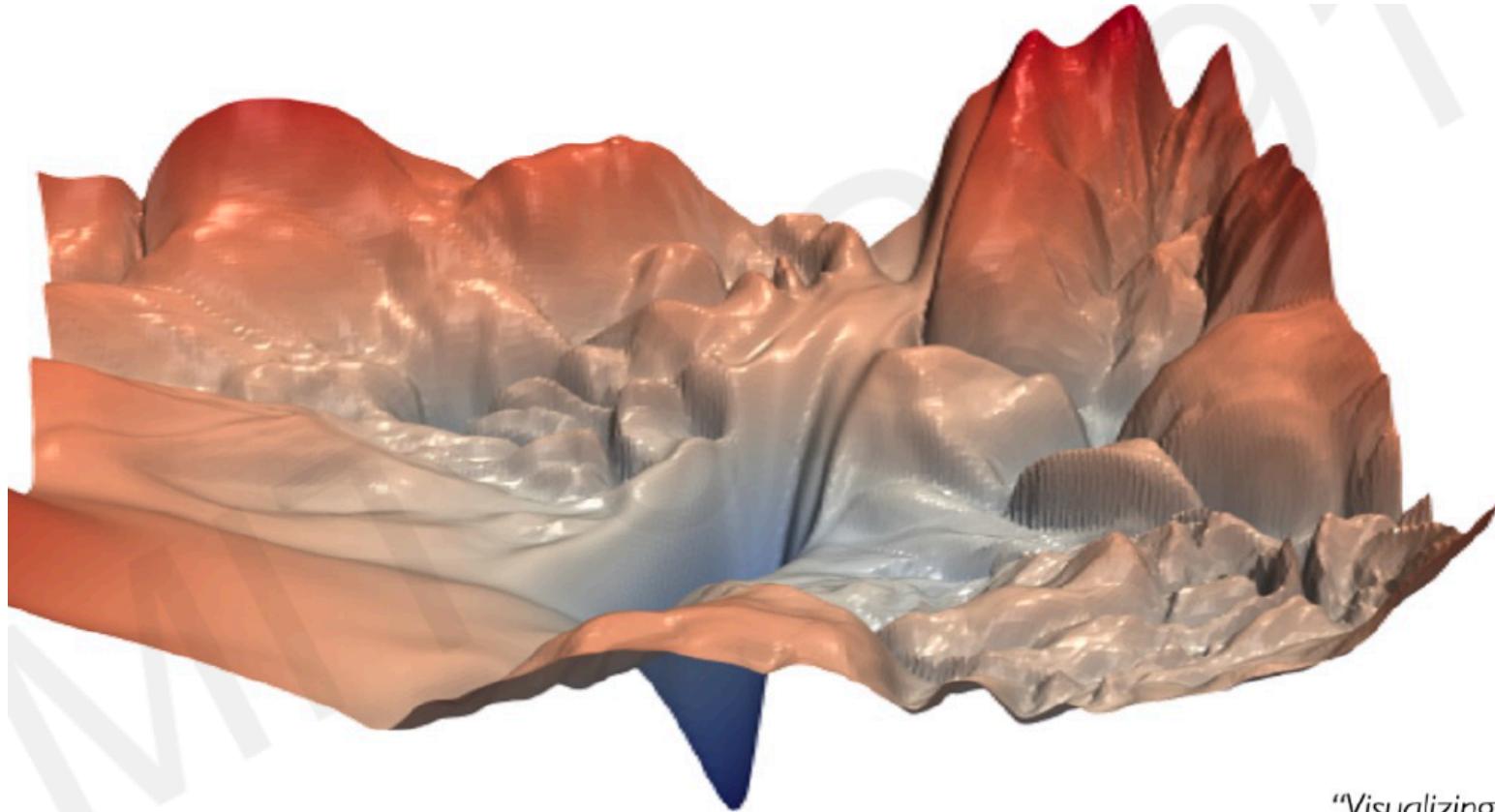


$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple bar}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red bar}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue bar}}$$

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

Neural Networks in Practice: Optimization

Training neural networks is difficult



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

Loss function can be difficult to optimize

Remember:

Optimization through gradient descent

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

Loss function can be difficult to optimize

Remember:

Optimization through gradient descent

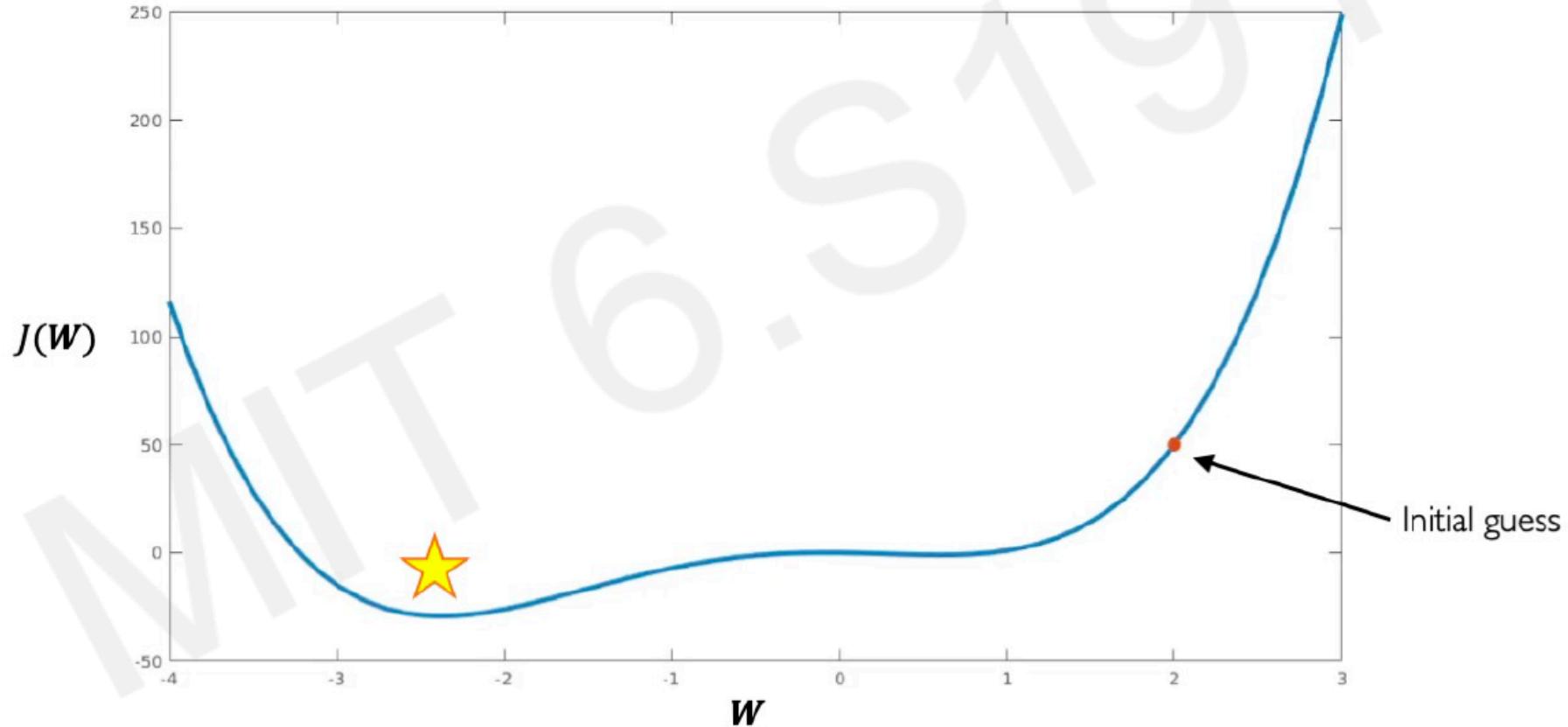
$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$



How can we set the
learning rate?

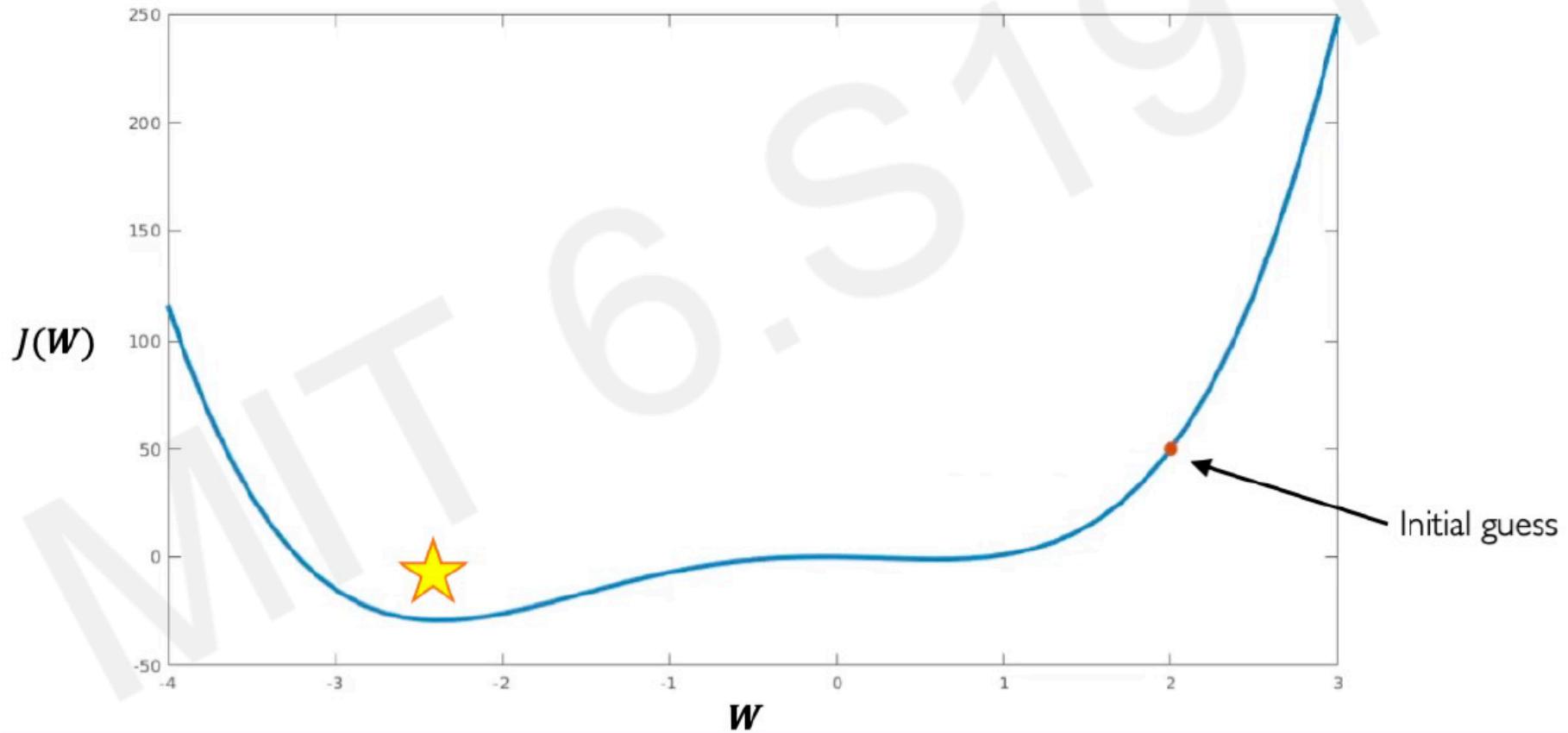
Setting the learning rate

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



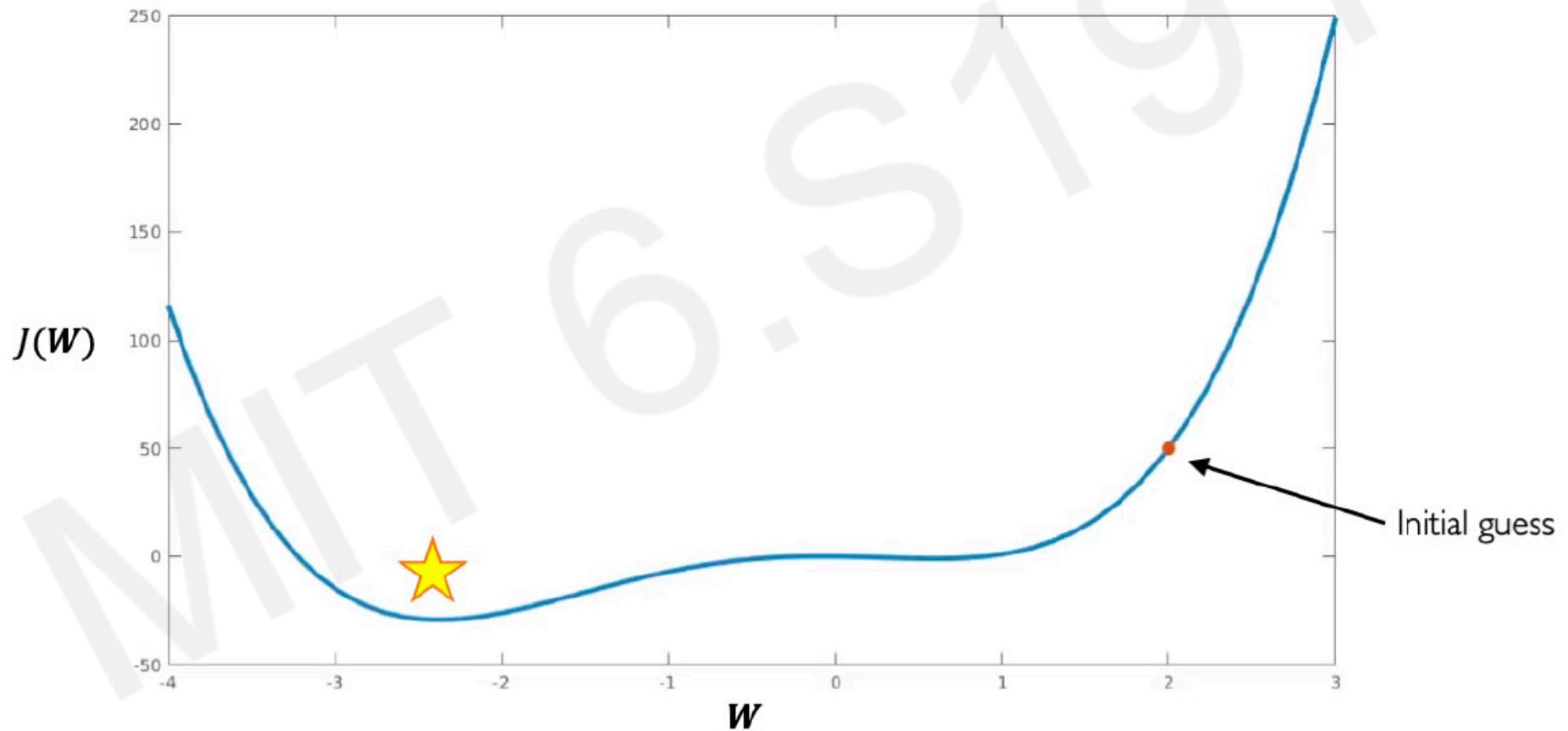
Setting the learning rate

Large learning rates overshoot, become unstable and diverge



Setting the learning rate

Stable learning rates converge smoothly and avoid local minima



How to deal with this?

Idea I:

Try lots of different learning rates and see what works “just right”

How to deal with this?

Idea 1:

Try lots of different learning rates and see what works “just right”

Idea 2:

Do something smarter!

Design an adaptive learning rate that “adapts” to the landscape

Adaptive learning rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - how large gradient is
 - how fast learning is happening
 - size of particular weights
 - etc...

Gradient descent algorithms

Algorithm

- SGD
- Adam
- Adadelta
- Adagrad
- RMSProp

Reference

Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952.

Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.

Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.

Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.

Additional details: <http://ruder.io/optimizing-gradient-descent/>

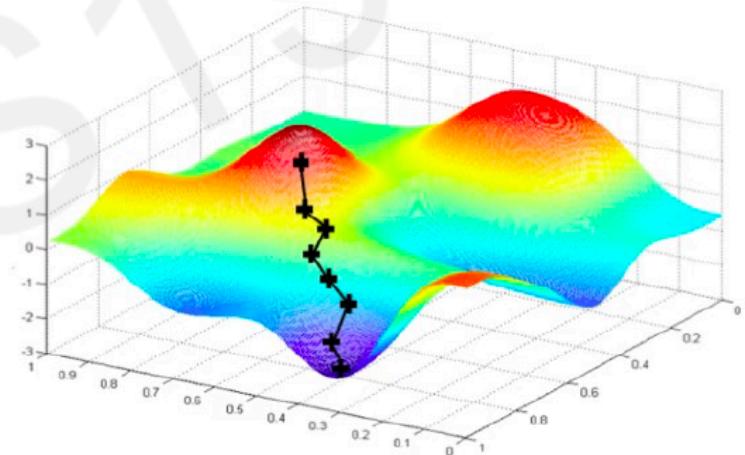
Neural Networks in Practice: Mini-batches

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
computationally
intensive to compute!

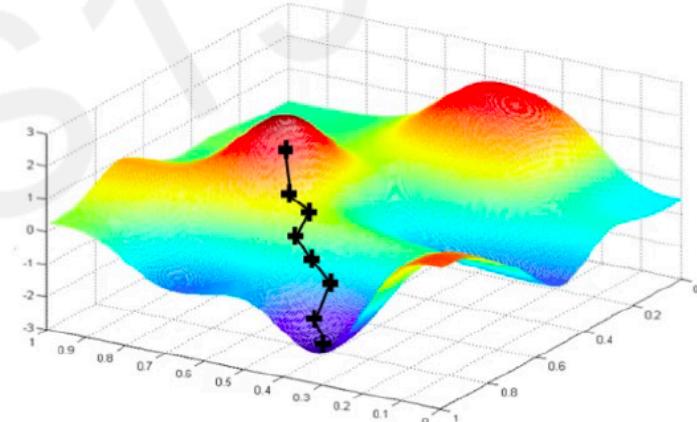


Stochastic gradient descent

Algorithm

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

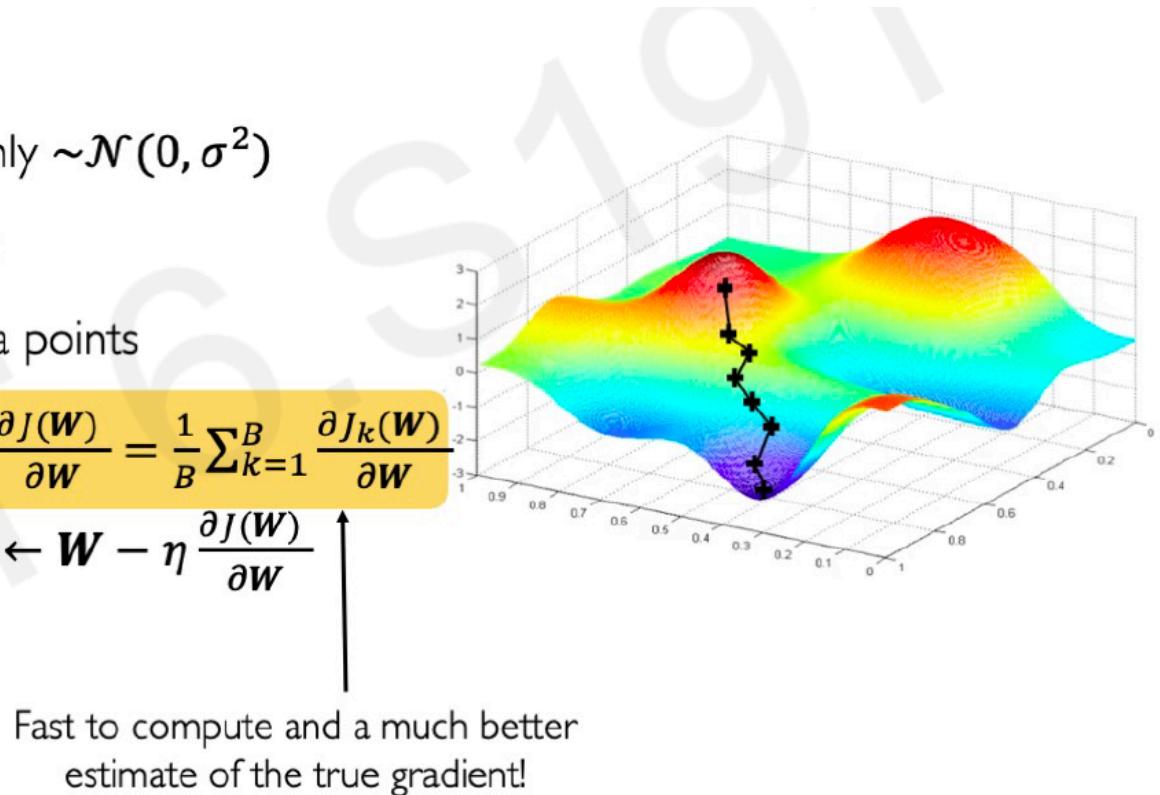
Easy to compute but
very noisy (stochastic)!



Mini-batches while training

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

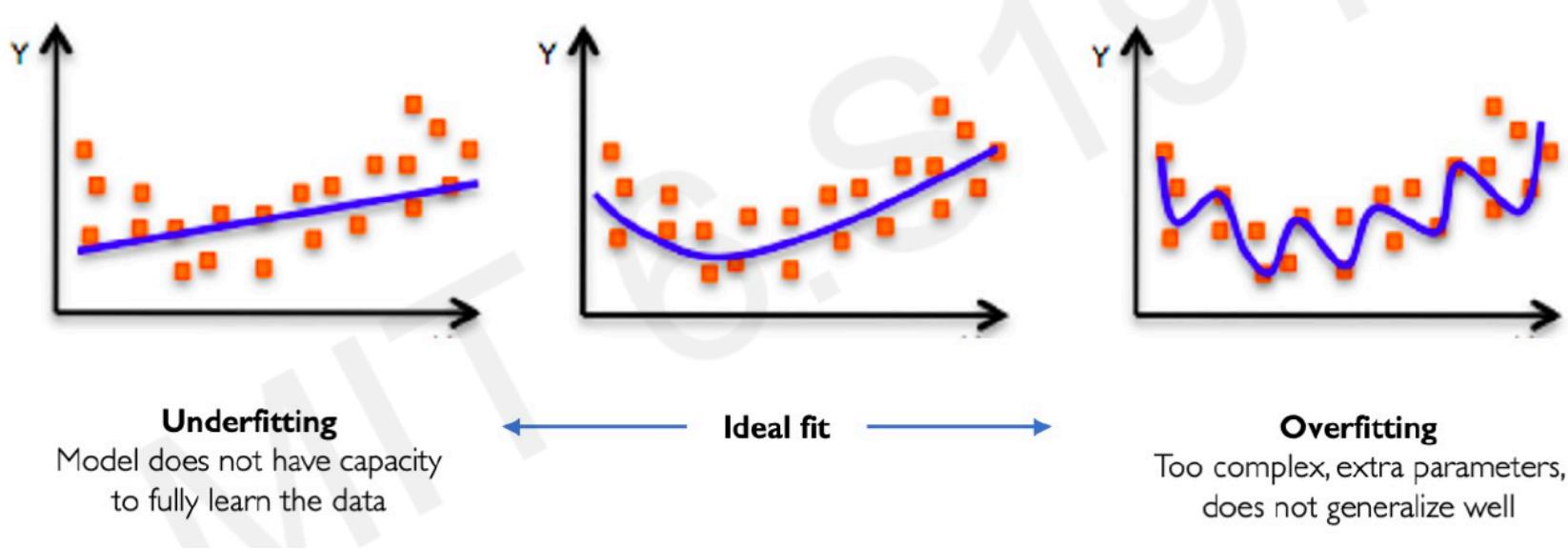


More accurate estimation of gradient

Smoother convergence
Allows for larger learning rates

Neural Networks in Practice: Overfitting

The problem of overfitting



Regularization

What is it?

Technique that constrains our optimization problem to discourage complex models

Regularization

What is it?

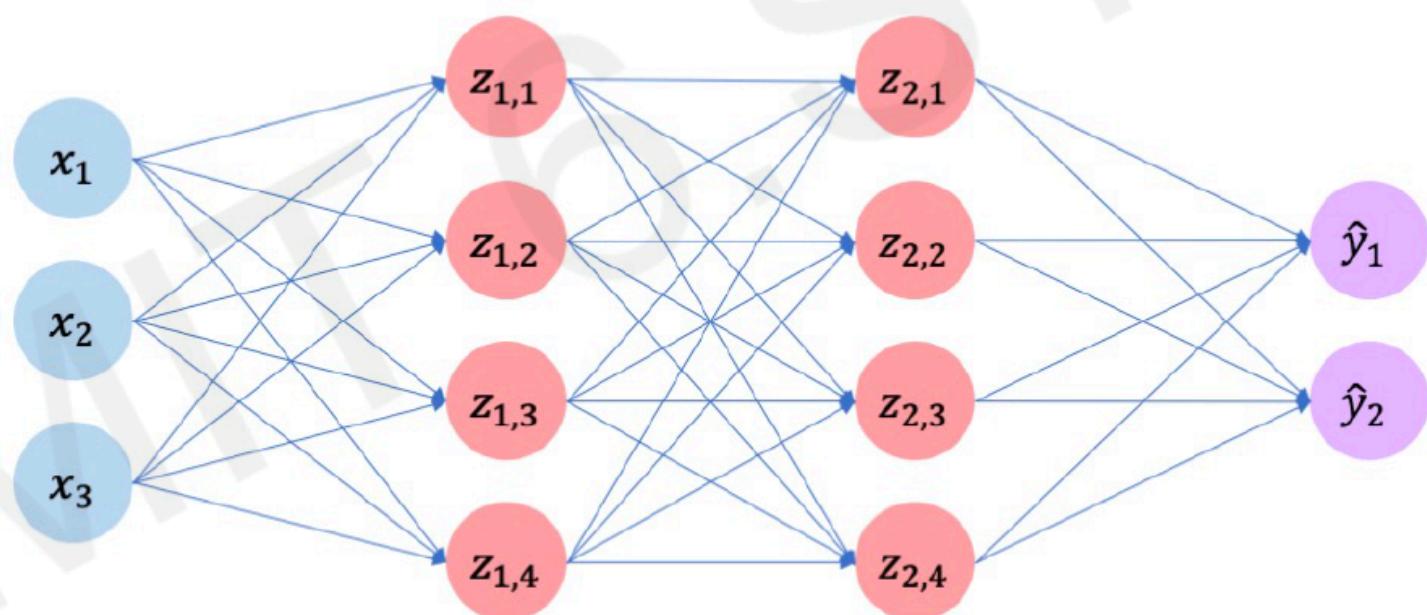
Technique that constrains our optimization problem to discourage complex models

Why do we need it?

Improve generalization of our model on unseen data

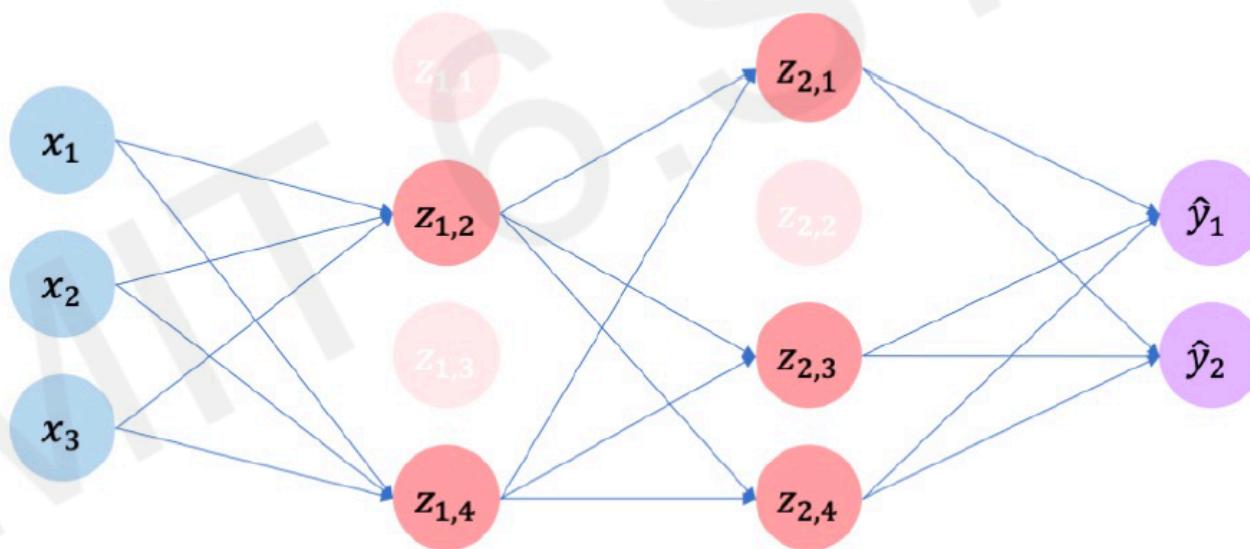
Regularization I: dropout

- During training, randomly set some activations to 0



Regularization I: 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



Regularization II: early stopping

- Stop training before we have a chance to overfit



Regularization II: early stopping

- Stop training before we have a chance to overfit

