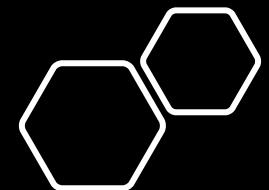


AMMI Review sessions

Deep Learning (2)
Deep Learning Basics

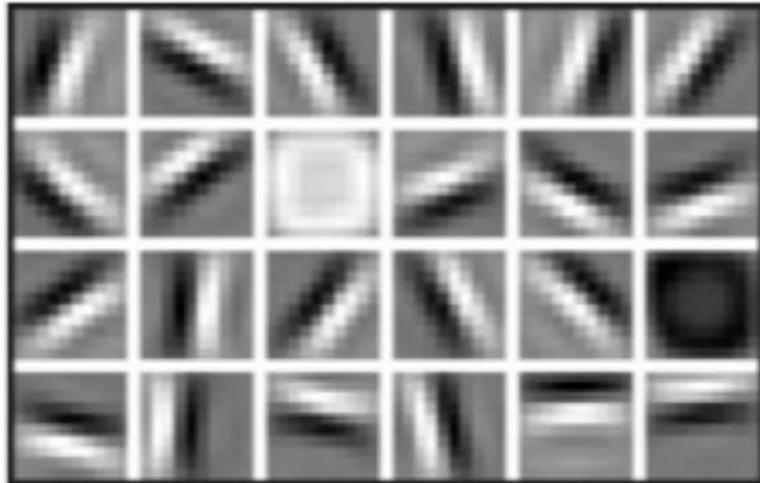


Why Deep Learning?

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



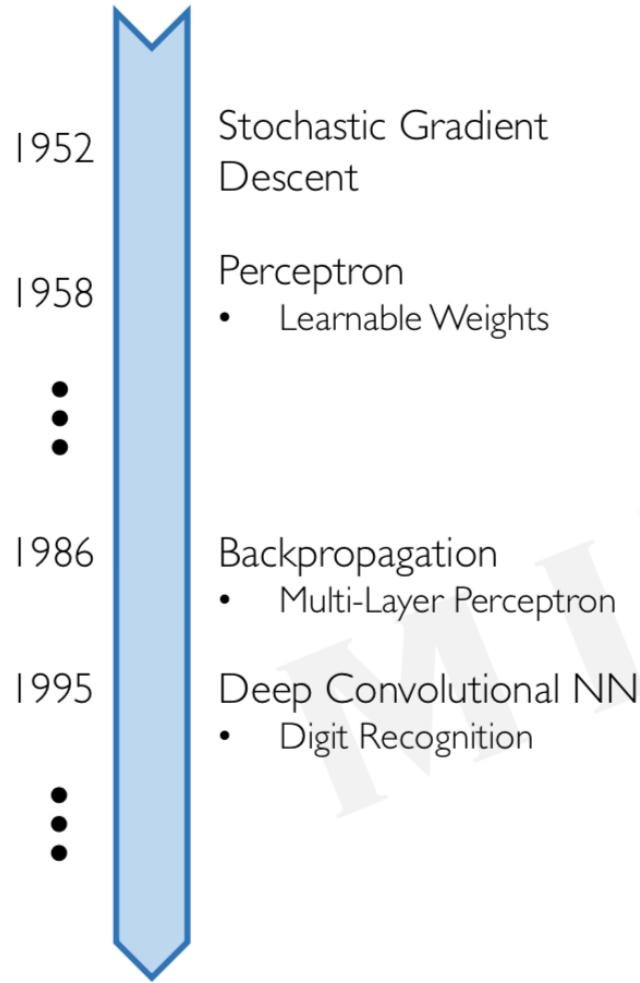
Eyes & Nose & Ears

High Level Features



Facial Structure

Why Now?



Neural Networks date back decades, so why the resurgence?

I. Big Data

- Larger Datasets
- Easier Collection & Storage

IMAGENET



WIKIPEDIA
The Free Encyclopedia



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable

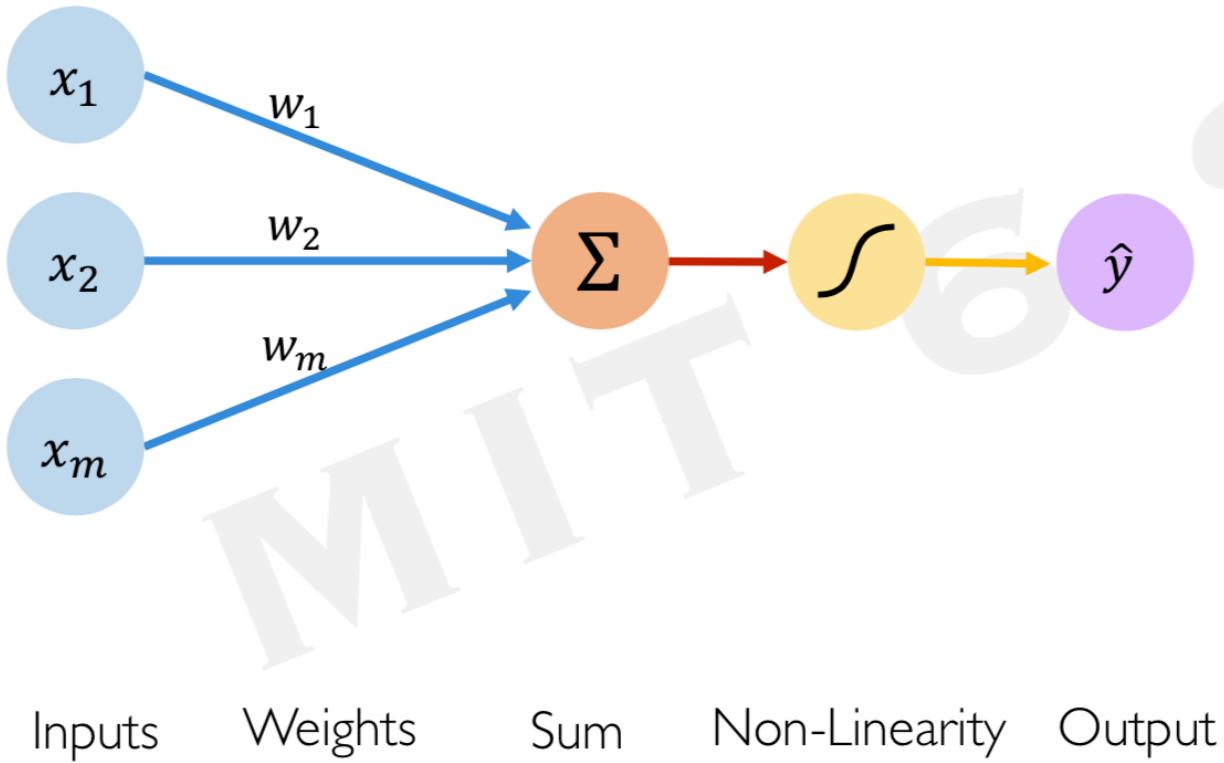


3. Software

- Improved Techniques
- New Models
- Toolboxes



The Perceptron: Forward Propagation



MIT S19.1

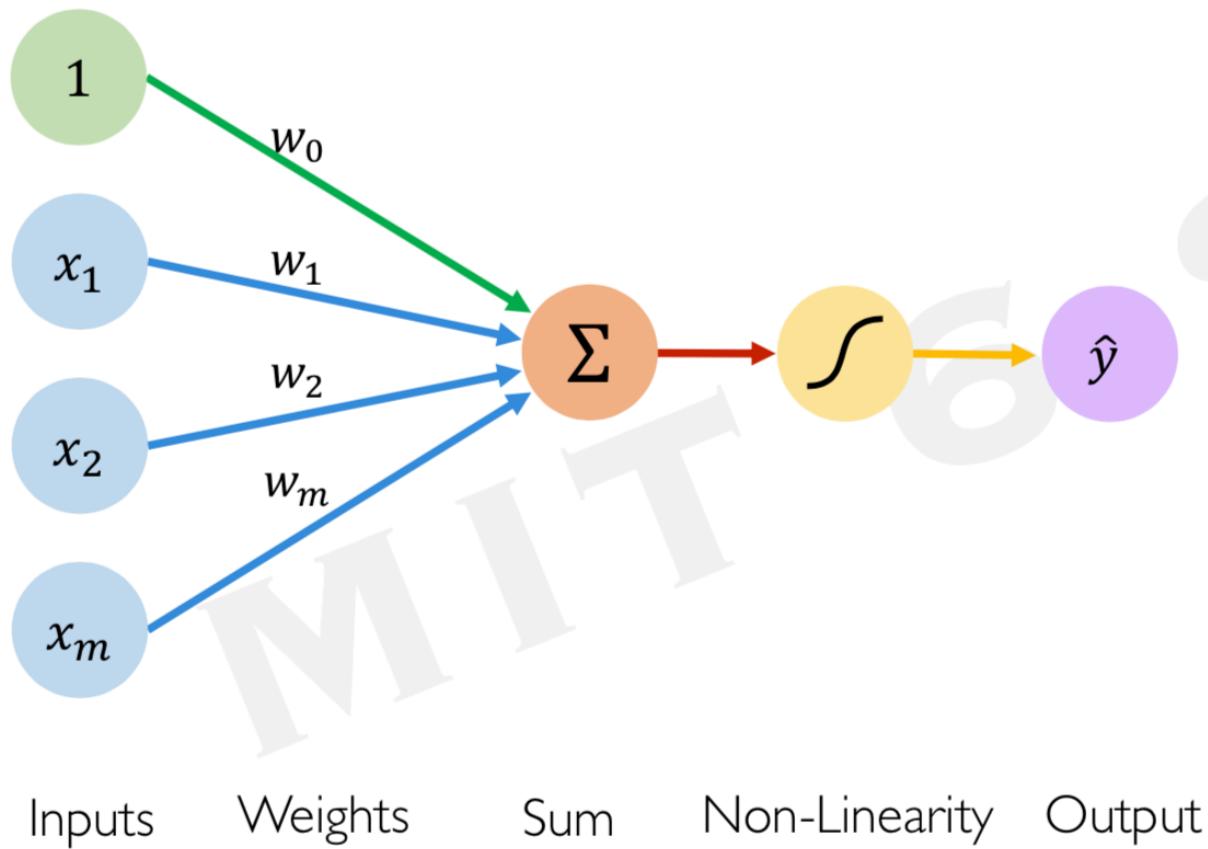
Linear combination of inputs

Output

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

Non-linear activation function

The Perceptron: Forward Propagation



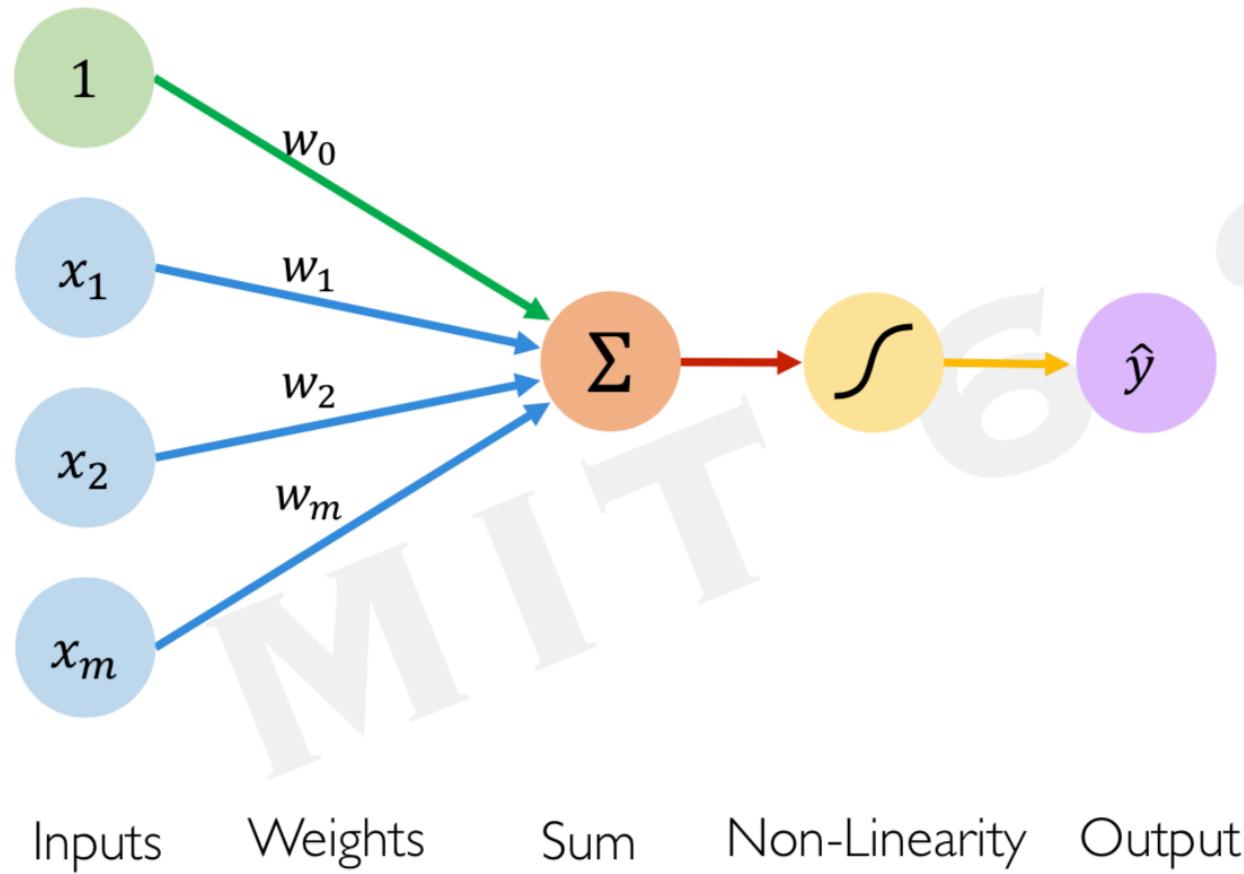
MIT S19

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

Annotations for the equation:

- Output:** Points to the output variable \hat{y} .
- Linear combination of inputs:** Points to the term $\sum_{i=1}^m x_i w_i$.
- Bias:** Points to the term w_0 .
- Non-linear activation function:** Points to the function g .

The Perceptron: Forward Propagation

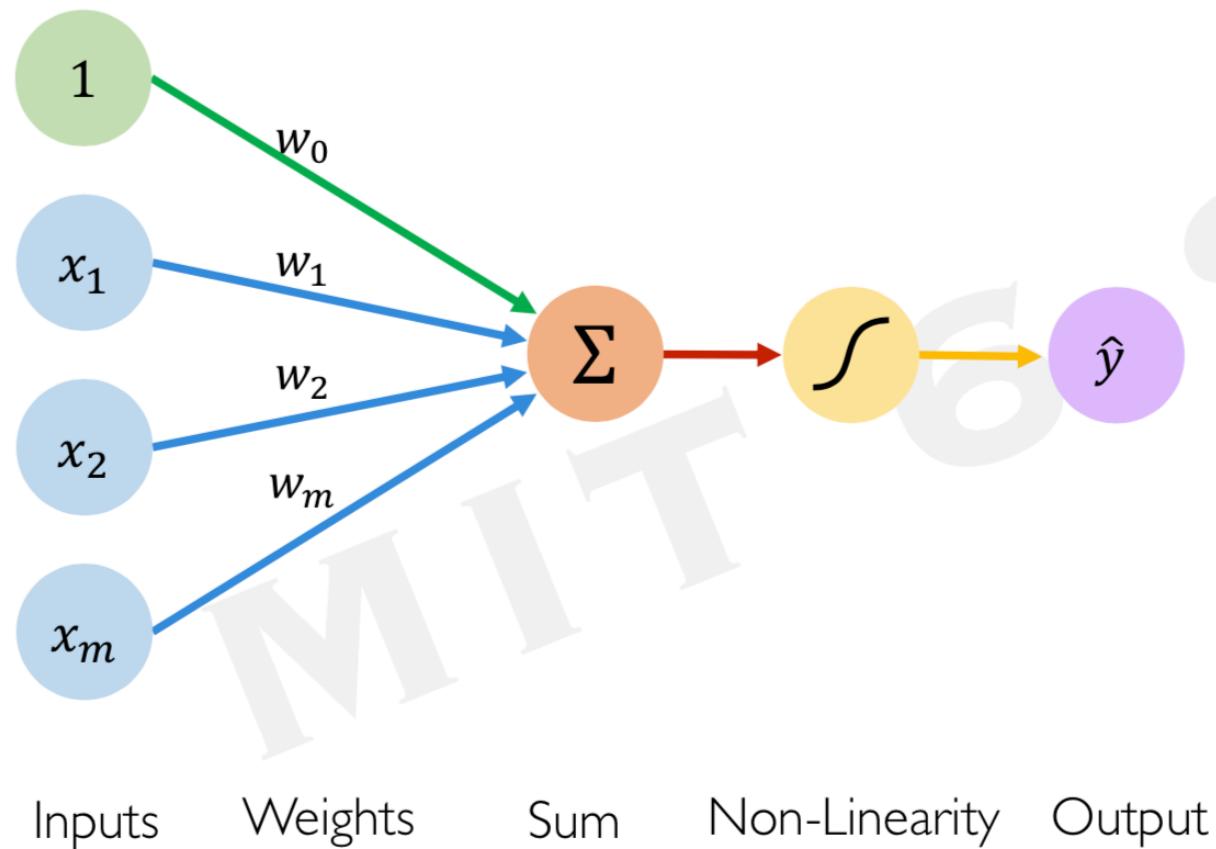


$$\hat{y} = g\left(w_0 + \sum_{i=1}^m x_i w_i\right)$$

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

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

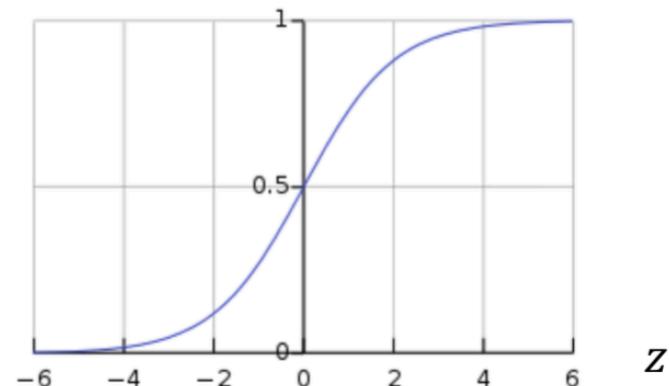


Activation Functions

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

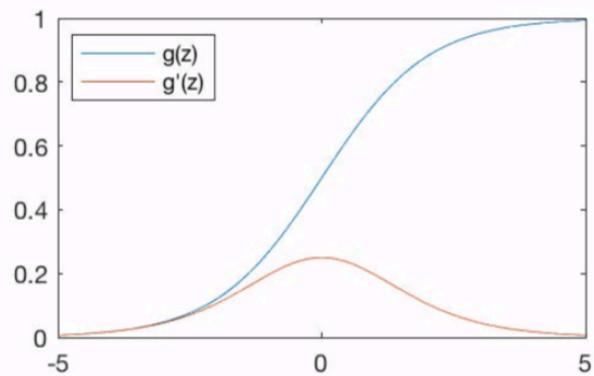
- Example: sigmoid function

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



Common Activation Functions

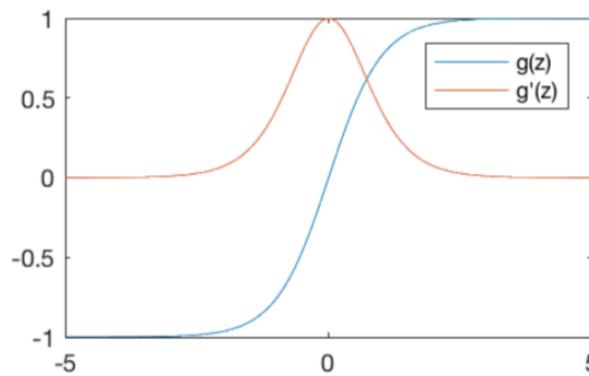
Sigmoid Function



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

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

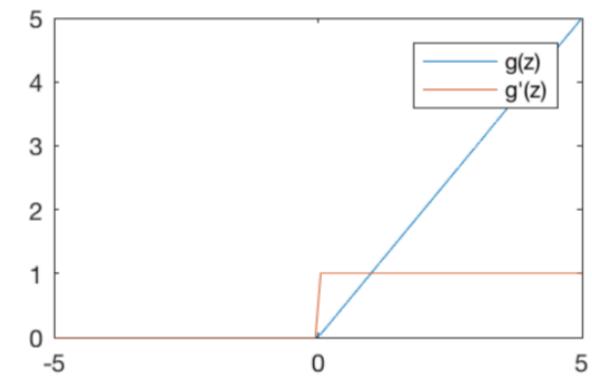
Hyperbolic Tangent



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

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

Rectified Linear Unit (ReLU)

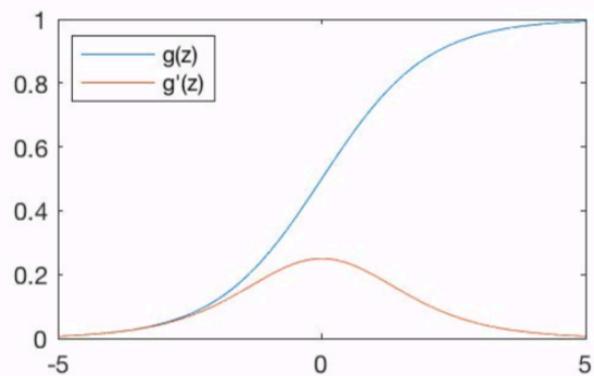


$$g(z) = \max(0, z)$$

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

Common Activation Functions

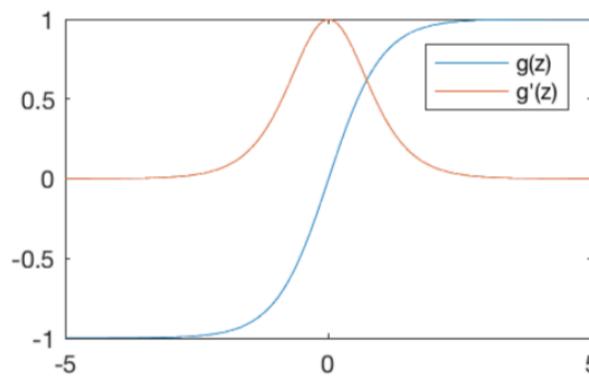
Sigmoid Function



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

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

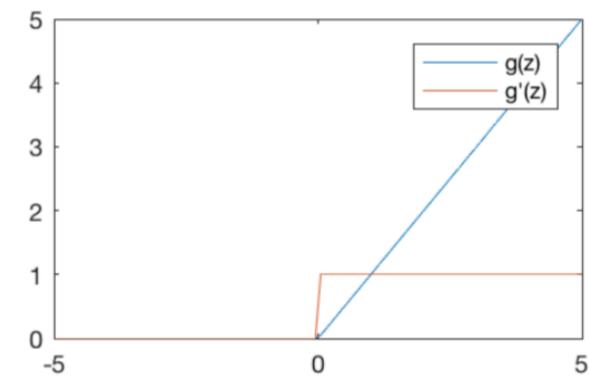
Hyperbolic Tangent



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

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

Rectified Linear Unit (ReLU)

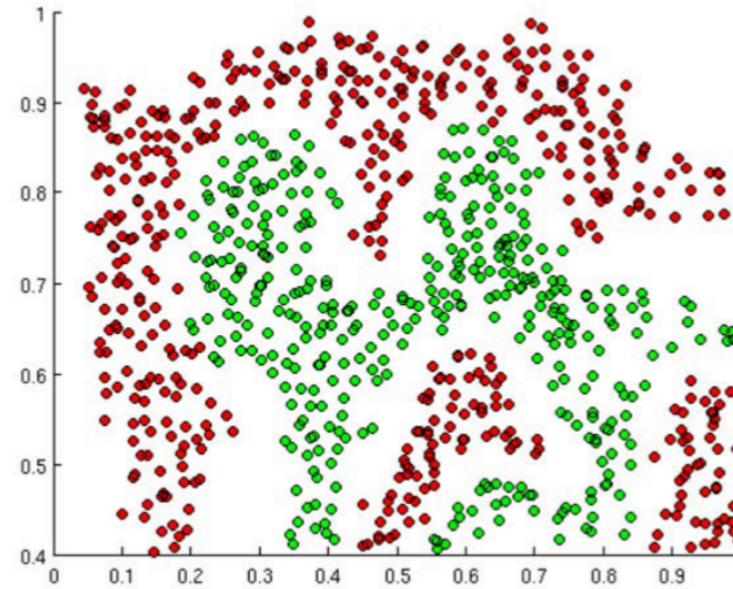


$$g(z) = \max(0, z)$$

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

Importance of Activation Functions

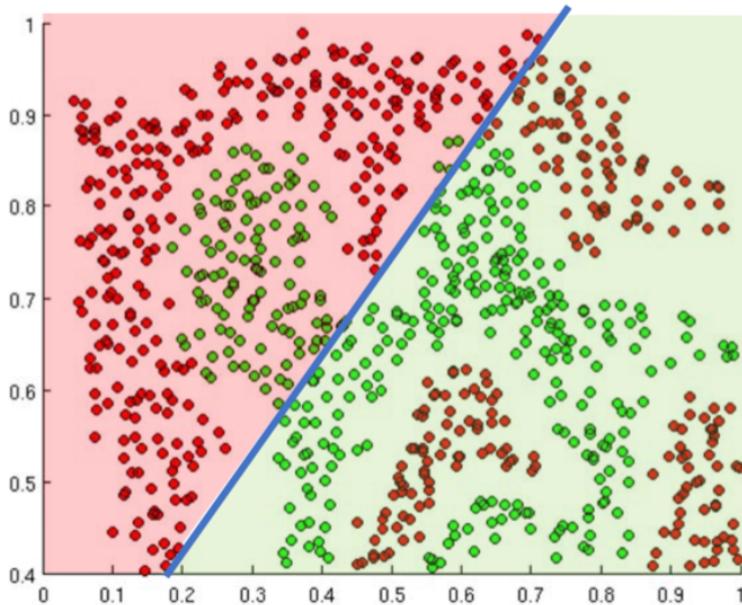
The purpose of activation functions is to *introduce non-linearities* into the network



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

Importance of Activation Functions

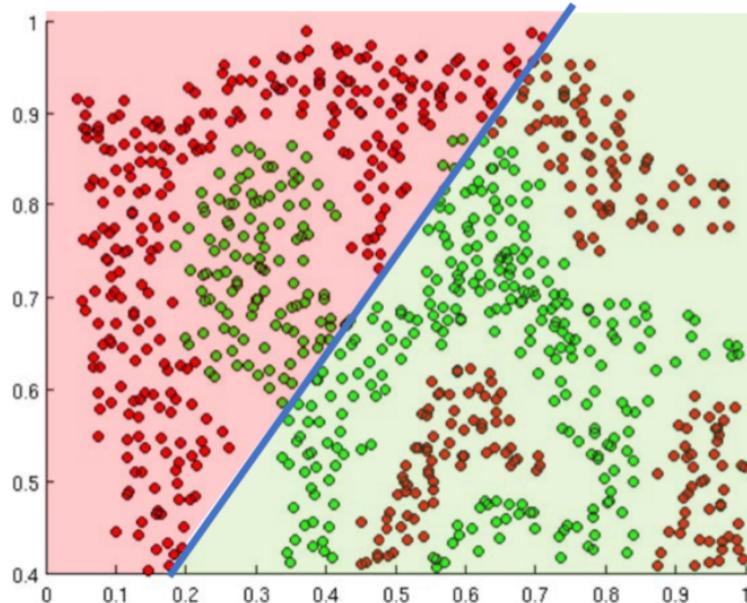
The purpose of activation functions is to **introduce non-linearities** into the network



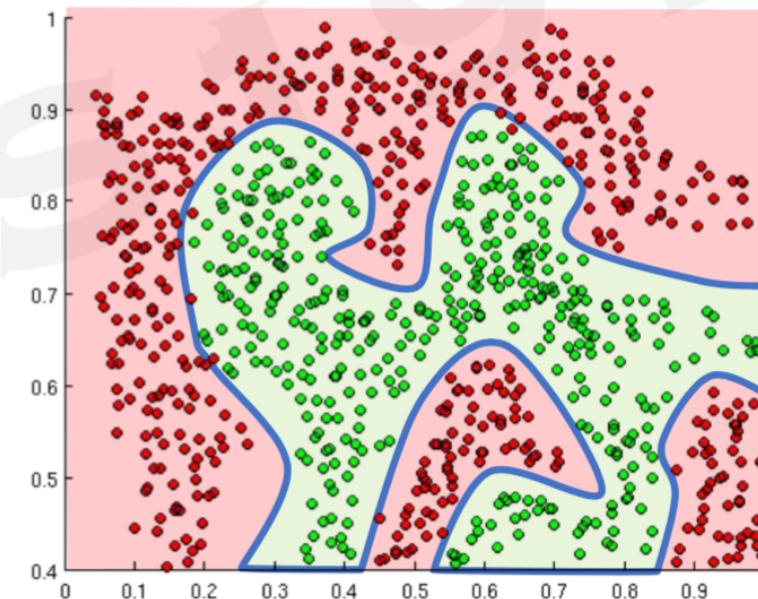
Linear activation functions produce linear decisions no matter the network size

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

Geometric Transformations of Hidden Layers

- Hidden layers are made up of a matrix, a bias term, and an activation function. Each of these components can be understood by how they geometrically effect their input
- Matrices are linear maps. Matrices “act on” an object in only two ways geometrically:
 1. Stretching, rotating, and translation.

Geometric Transformations of Hidden Layers

- Hidden layers are made up of a matrix, a bias term, and an activation function. Each of these components can be understood by how they geometrically effect their input
- Matrices are linear maps. Matrices “act on” an object in only two ways geometrically:
 1. Stretching, rotating, and translation.
 2. nonlinear functions perform “bending” squashing, or “folding” the space over onto itself

Geometric Transformations of Hidden Layers

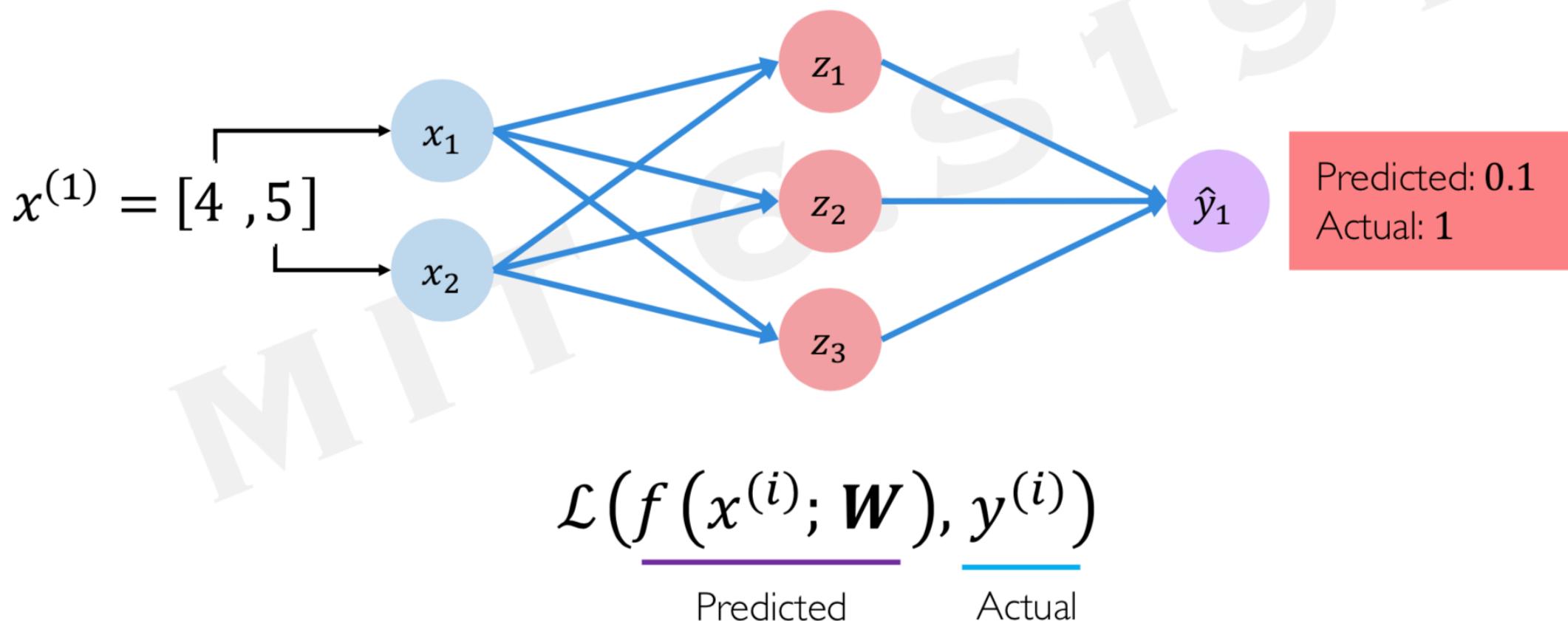
- A NN takes the input space and bends and twists it around until the classes are separated.
- The decision boundary is the visual representation of the threshold where the model will predict 0 or 1

$$\{(x, y) \mid \text{sigmoid}(ax + by + c) = .5\}.$$

- [Visualize transformations](#)

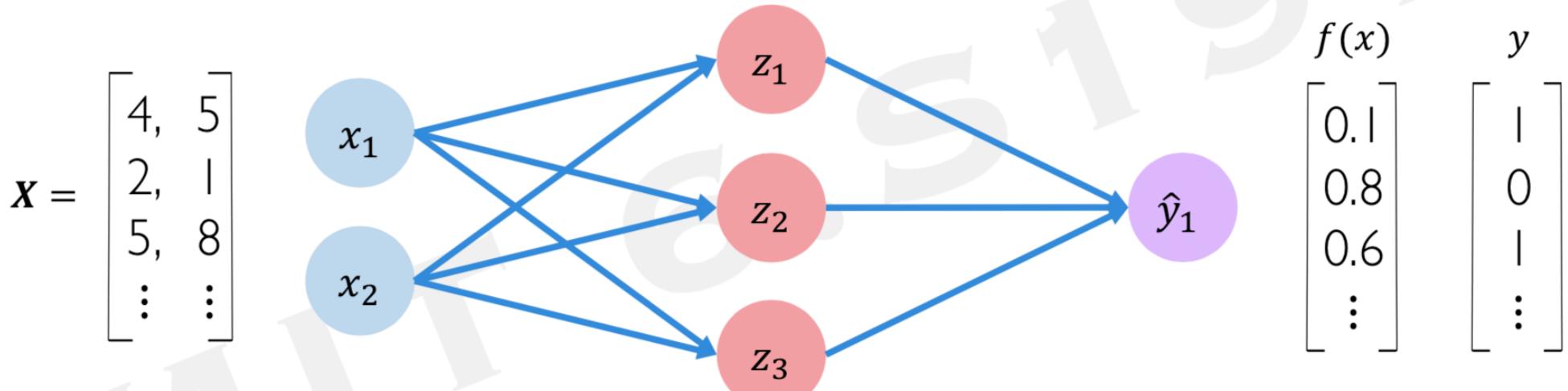
Quantifying Loss

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



Also known as:

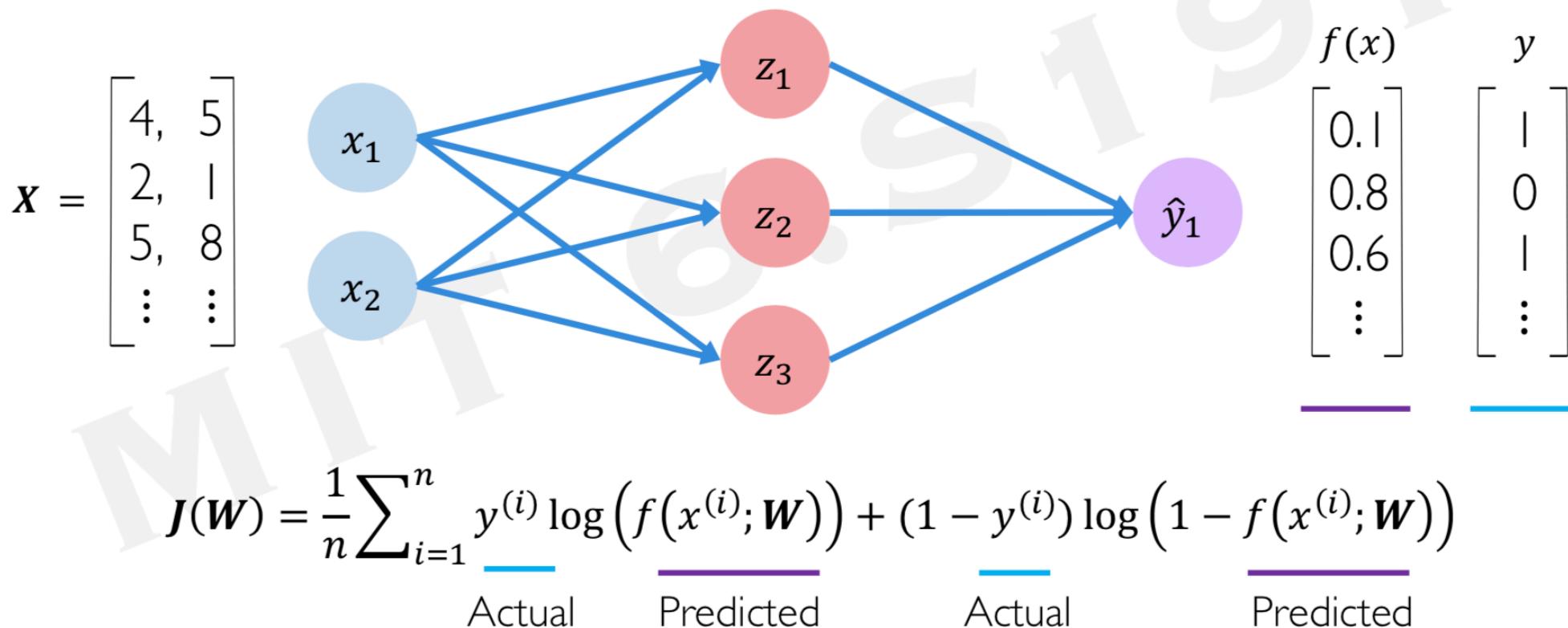
- Objective function
- Cost function
- Empirical Risk

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

Predicted Actual

Binary Cross Entropy Loss

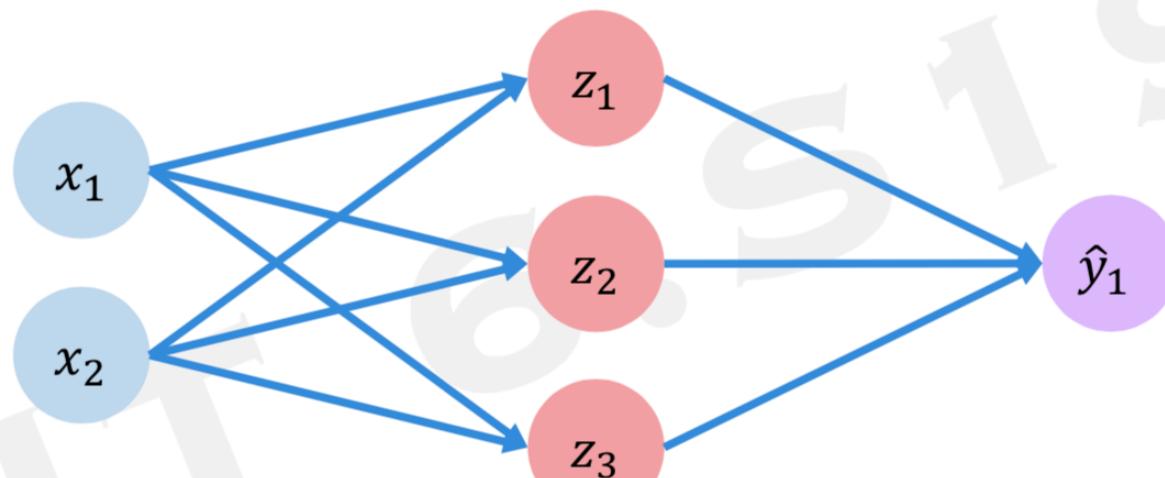
Cross entropy loss can be used with models that output a probability between 0 and 1



Mean Squared Error Loss

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

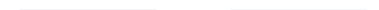
$$\mathbf{x} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \left(\underline{y^{(i)}} - \underline{f(x^{(i)}; \mathbf{W})} \right)^2$$

Actual Predicted

$f(x)$	y
$\begin{bmatrix} 30 \\ 80 \\ 85 \\ \vdots \end{bmatrix}$	$\begin{bmatrix} 90 \\ 20 \\ 95 \\ \vdots \end{bmatrix}$

Final Grades
(percentage)

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:

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

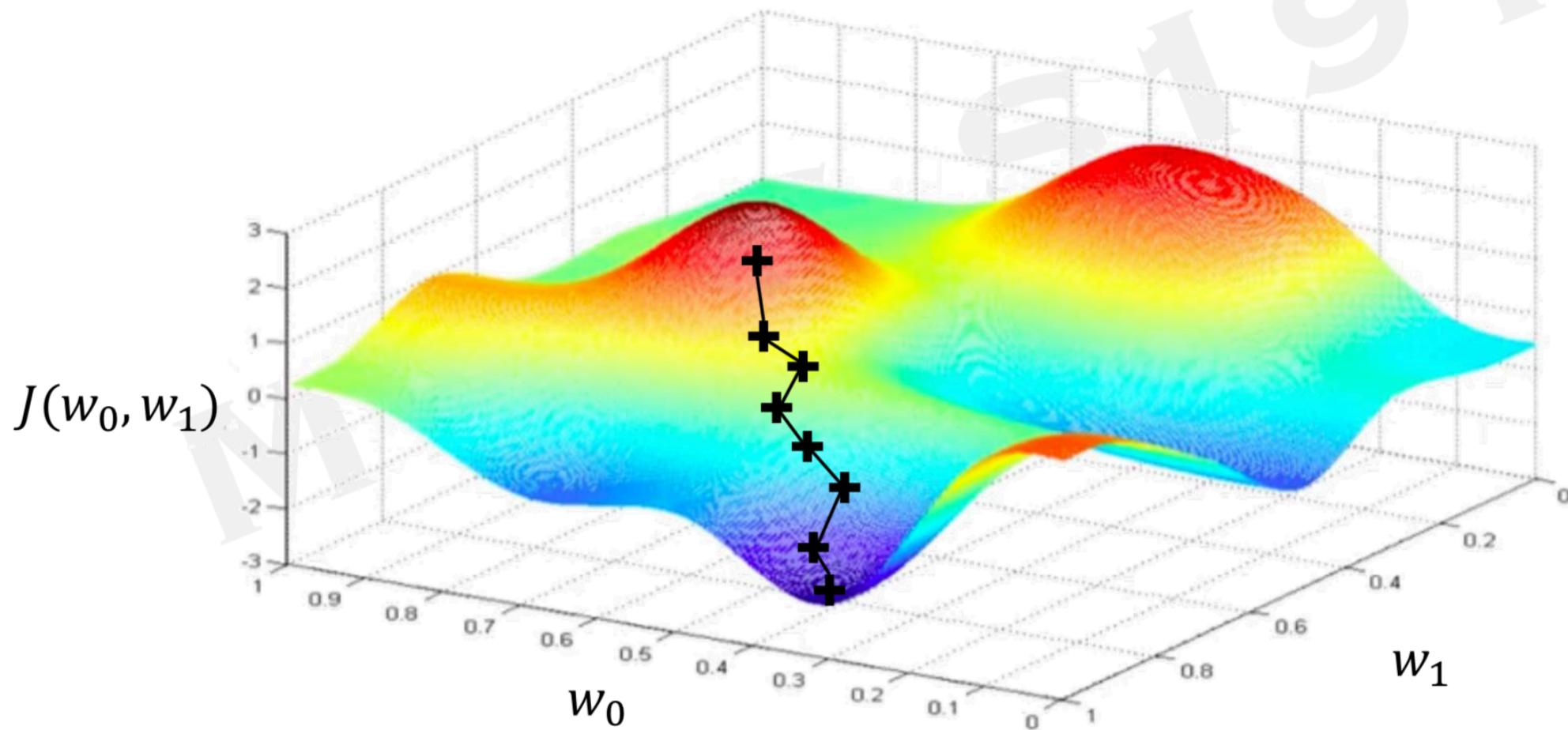
Gradient Descent

Algorithm

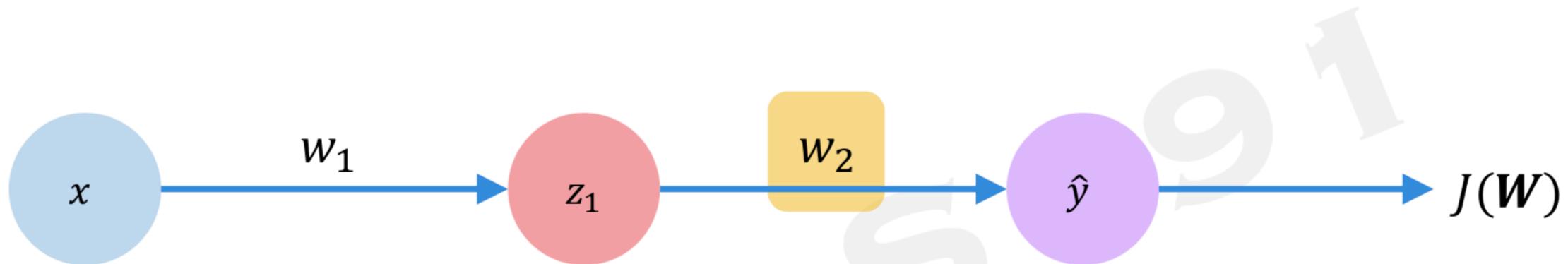
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Gradient Descent

Repeat until convergence

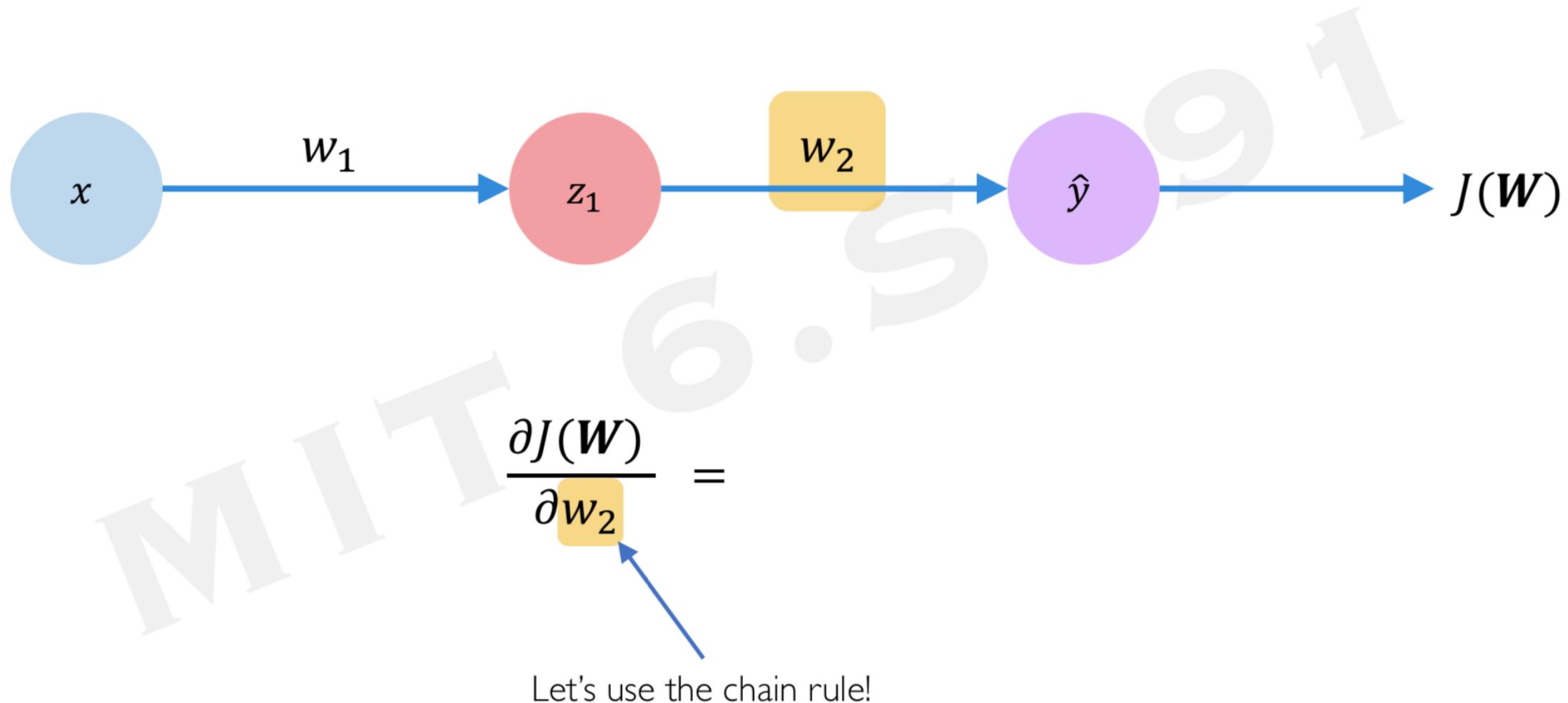


Computing Gradients: Backpropagation

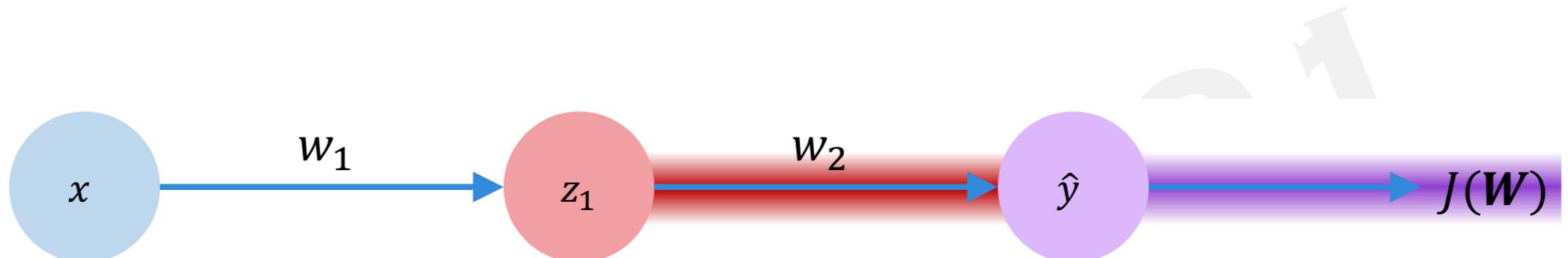


How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

Computing Gradients: Backpropagation

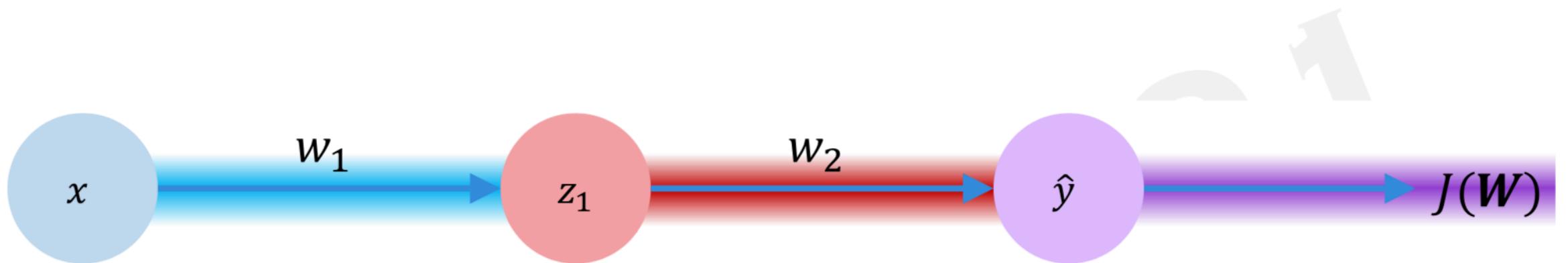


Computing Gradients: Backpropagation



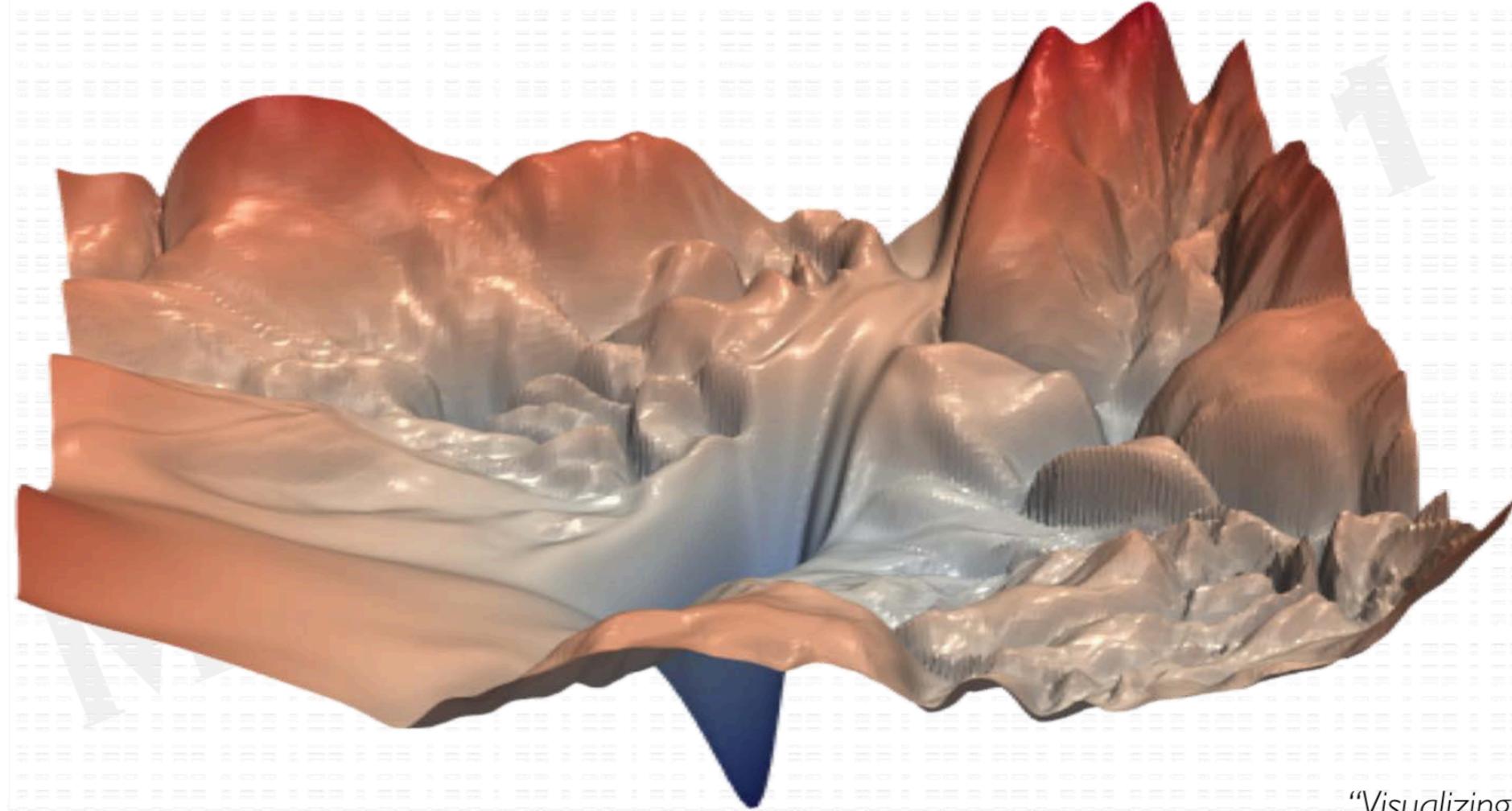
$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} * \underline{\frac{\partial \hat{y}}{\partial w_2}}$$

Computing Gradients: Backpropagation



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

Training Neural Networks is Difficult



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

Loss Functions Can Be Difficult to Optimize

Remember:

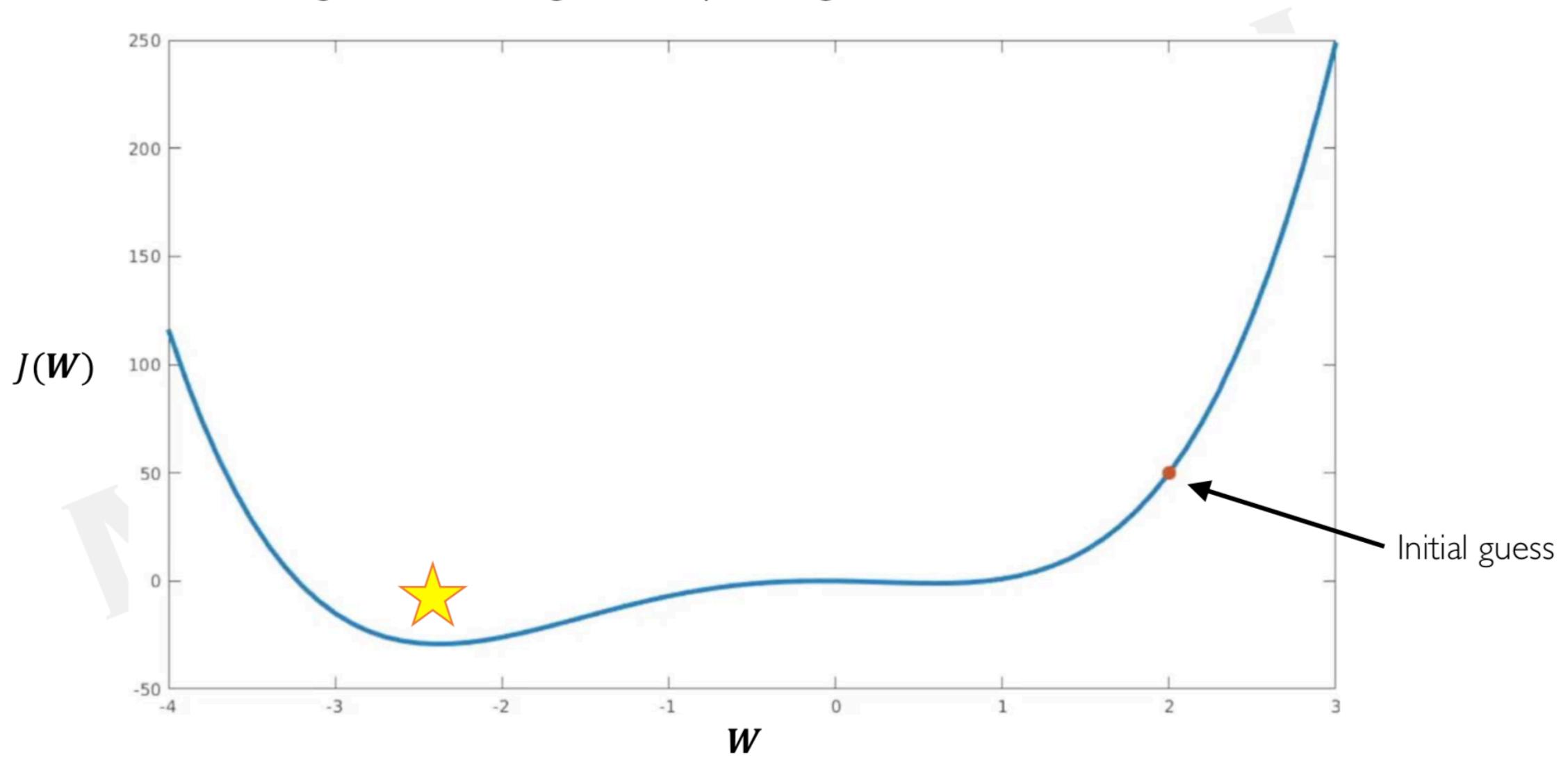
Optimization through gradient descent

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

How can we set the
learning rate?

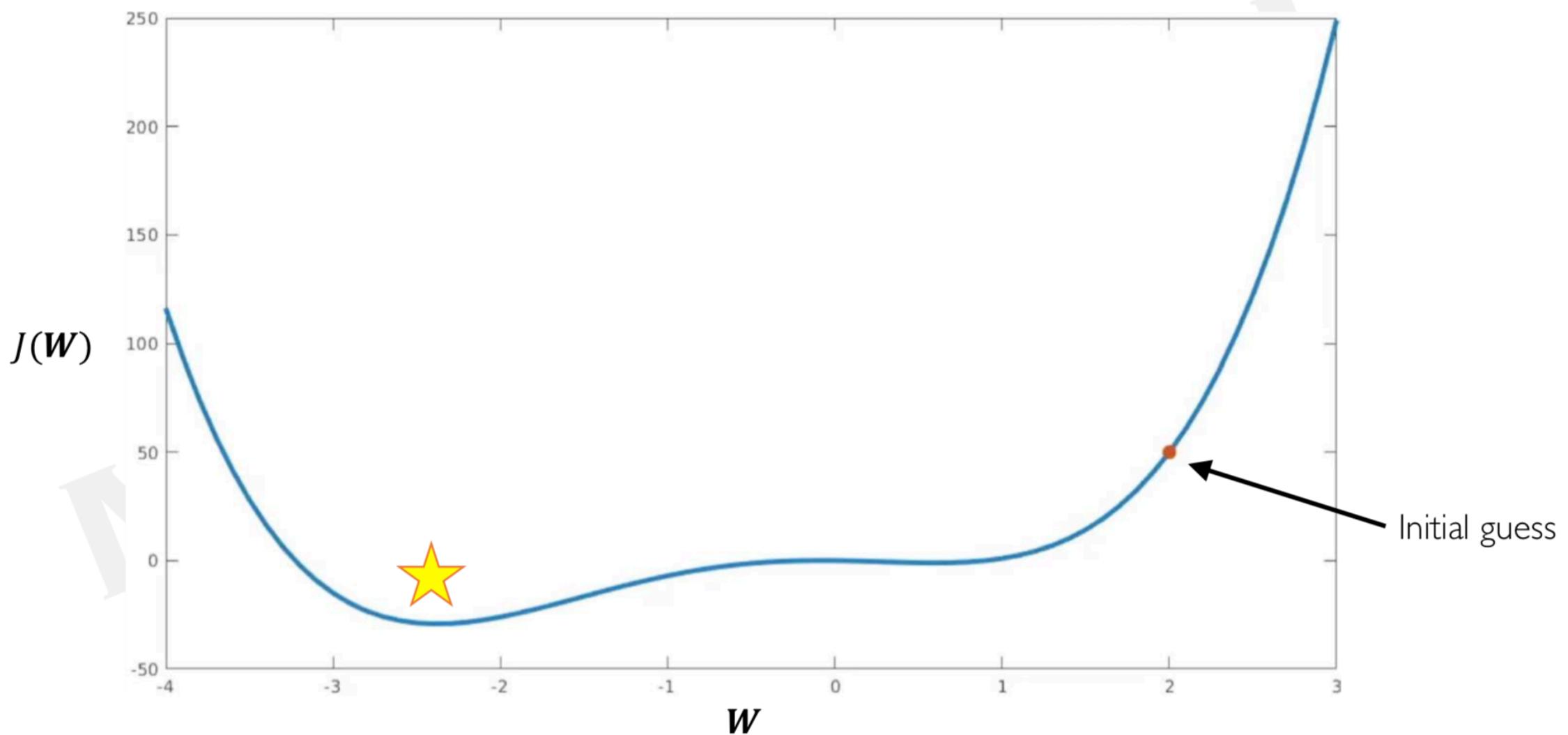
Setting the Learning Rate

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



Setting the Learning Rate

Large learning rates overshoot, become unstable and diverge



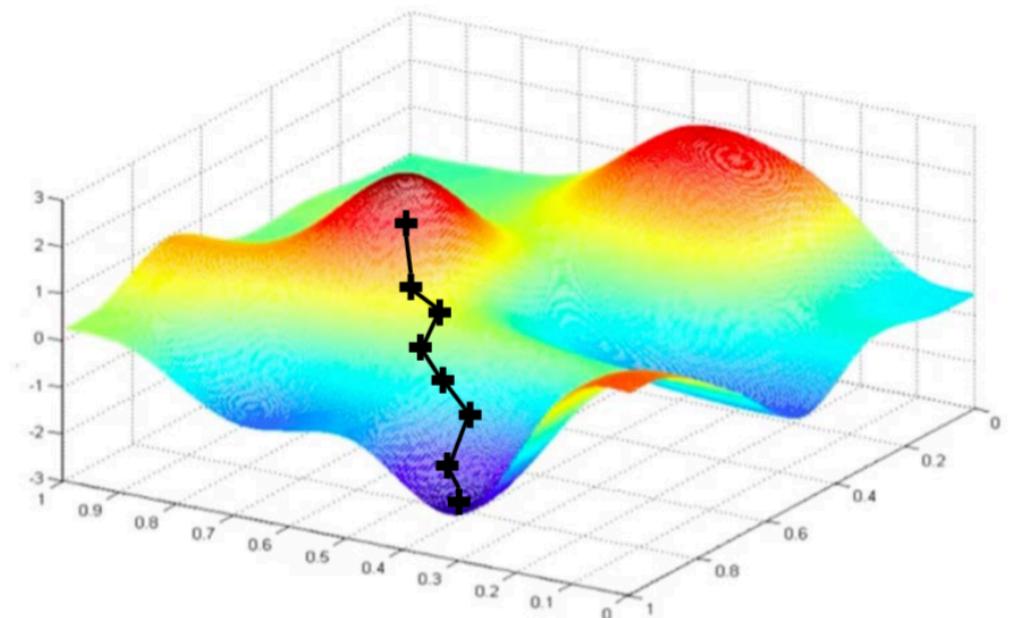
Adaptive Learning Rates

- 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...

Stochastic Gradient Descent

Algorithm

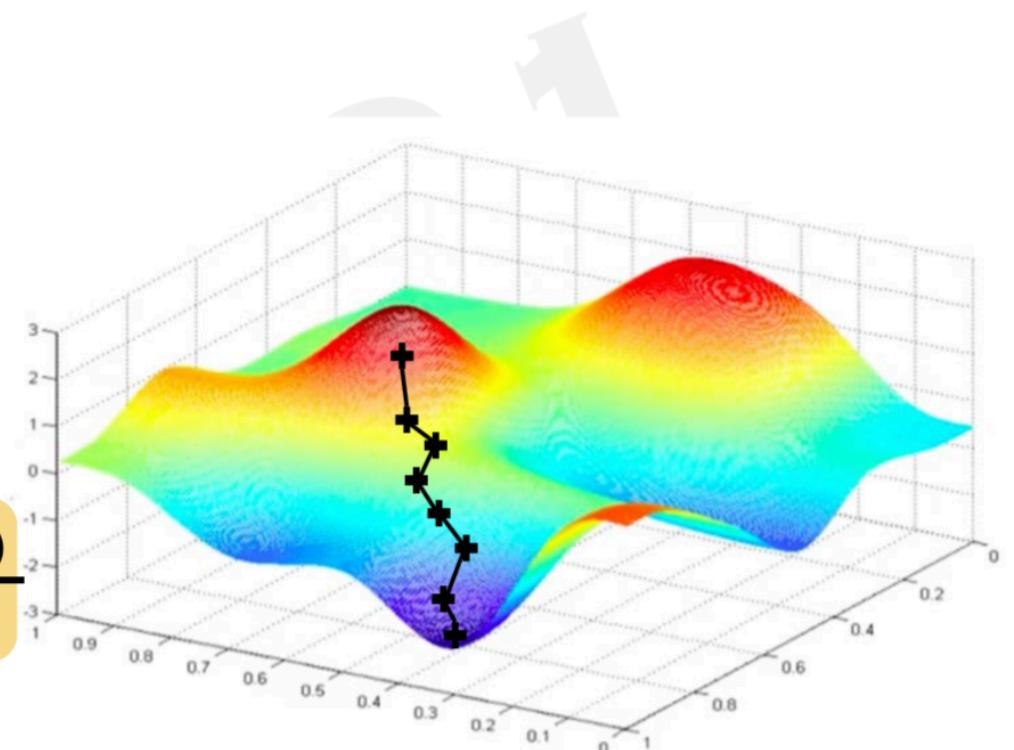
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$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



Stochastic Gradient Descent

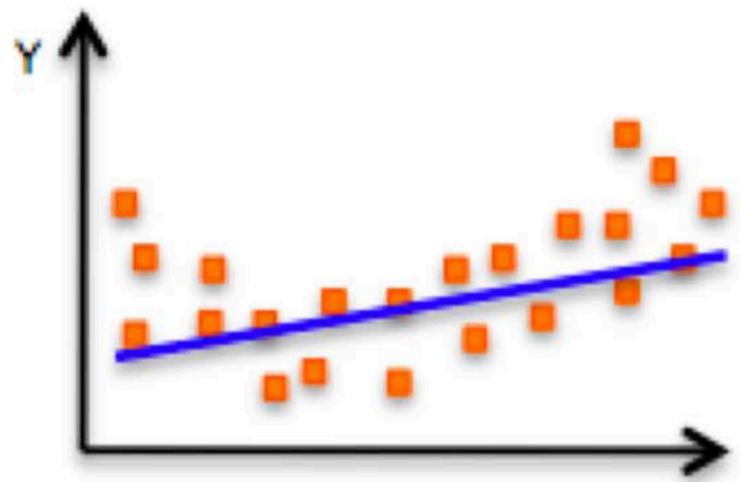
Algorithm

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



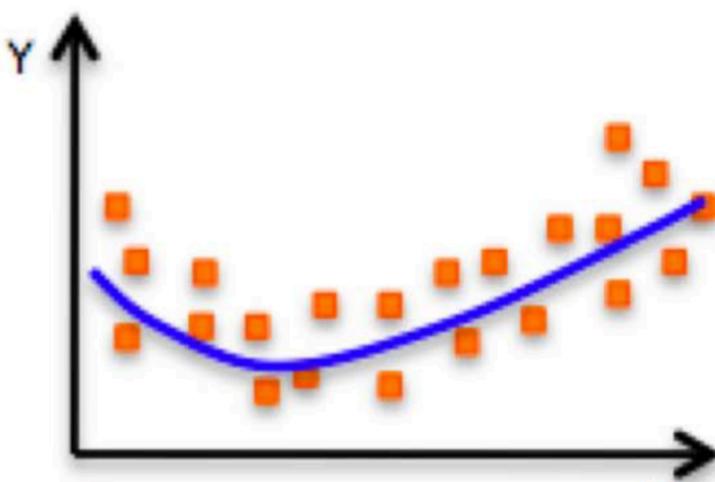
Fast to compute and a much better estimate of the true gradient!

The Problem of Overfitting

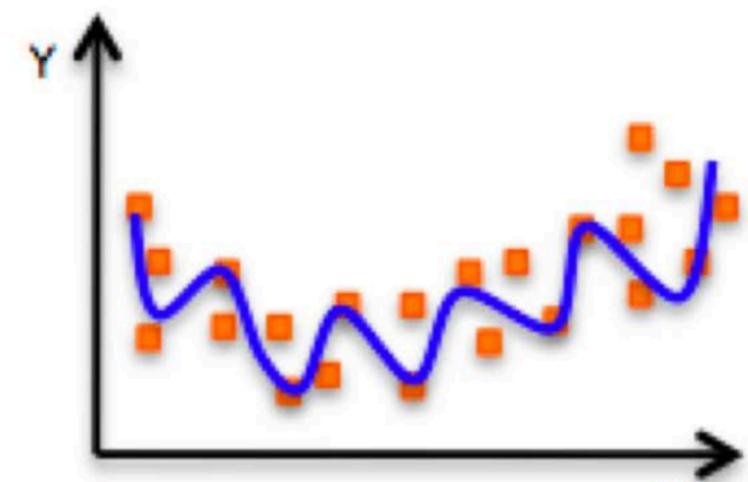


Underfitting

Model does not have capacity
to fully learn the data



Ideal fit



Overfitting

Too complex, extra parameters,
does not generalize well

Bias and Variance trade off

Bias and Variance trade off

Bias

- Bias is error introduced in your model due to over simplification of machine learning algorithm
- It can lead to under fitting. When you train your model at that time model makes simplified assumptions to make the target function easier to understand

Bias and Variance trade off

Bias

- Bias is error introduced in your model due to over simplification of machine learning algorithm
- It can lead to under fitting. When you train your model at that time model makes simplified assumptions to make the target function easier to understand

Variance

- Variance is error introduced in your model due to complex machine learning algorithm
- your model learns noise also from the training data set and performs bad on test data set.

Bias and Variance trade off

- The goal of any supervised machine learning algorithm is to have low bias and low variance to achieve good prediction performance.
- The **k-nearest neighbours** algorithm has low bias and high variance, but the trade-off can be changed by increasing the value of k which increases the number of neighbours that contribute to the prediction and in turn increases the bias of the model.
- The **support vector machine** algorithm has low bias and high variance, but the trade-off can be changed by increasing the C parameter that influences the number of violations of the margin allowed in the training data which increases the bias but decreases the variance.

Bias and Variance trade off

- There is no escaping the relationship between bias and variance in machine learning. Increasing the bias will decrease the variance. Increasing the variance will decrease the bias.

Regularization

What is it?

Technique that constrains our optimization problem to discourage complex models

Regularization

What is it?

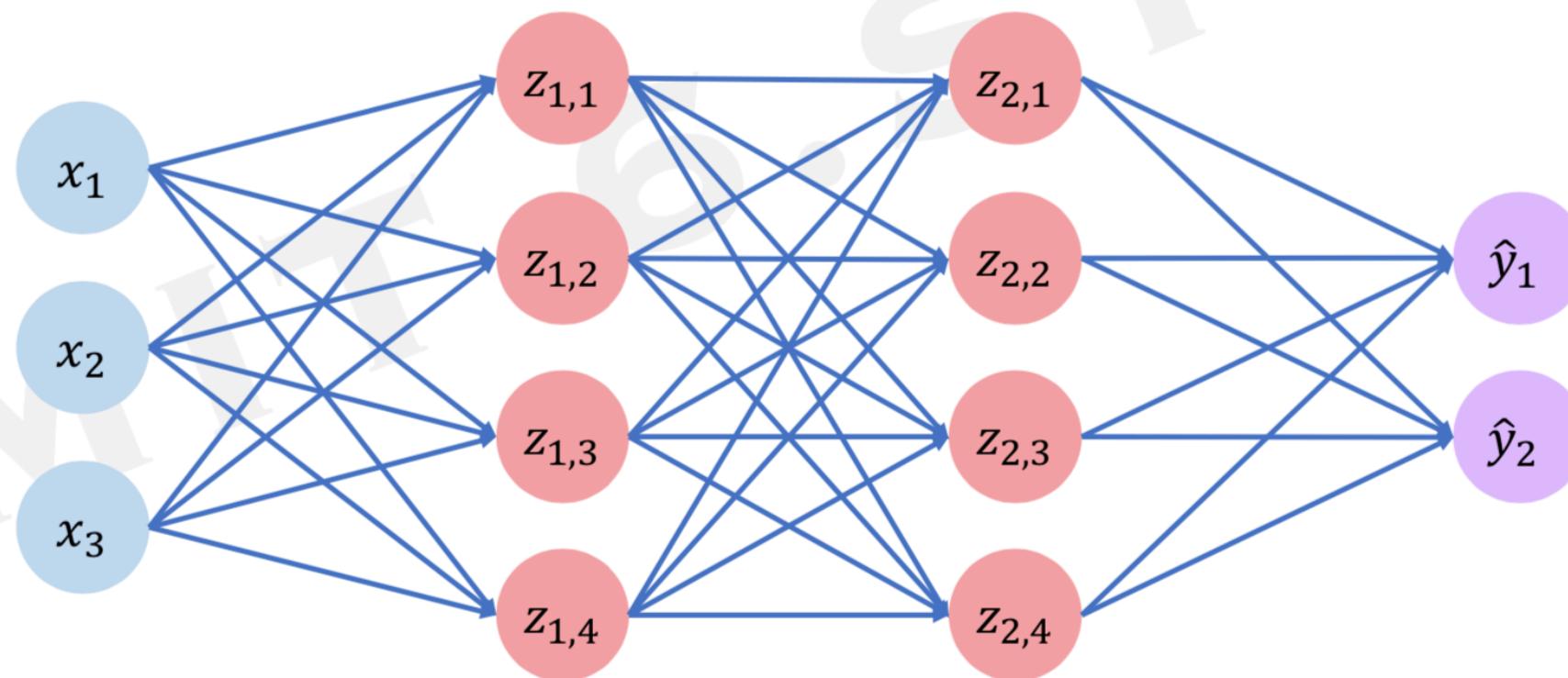
Technique that constrains our optimization problem to discourage complex models

Why do we need it?

Improve generalization of our model on unseen data

Regularization I: Dropout

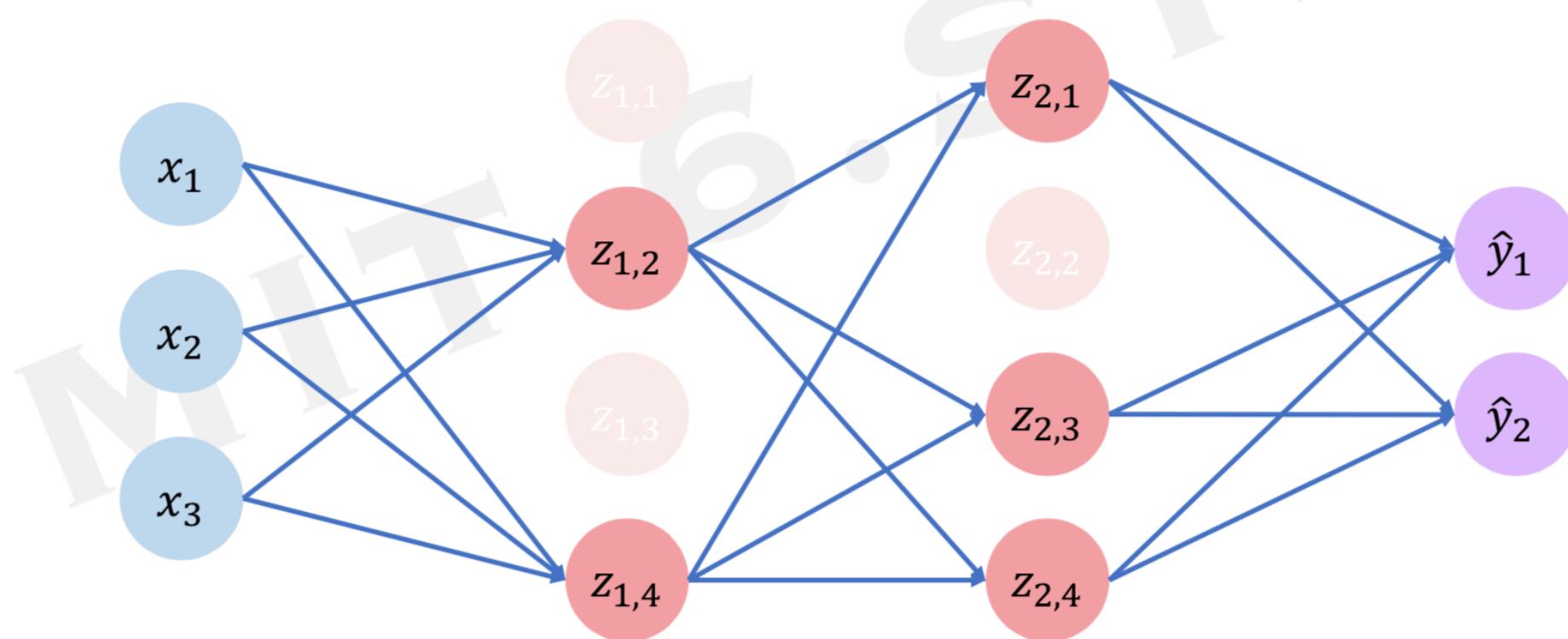
- During training, randomly set some activations to 0



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

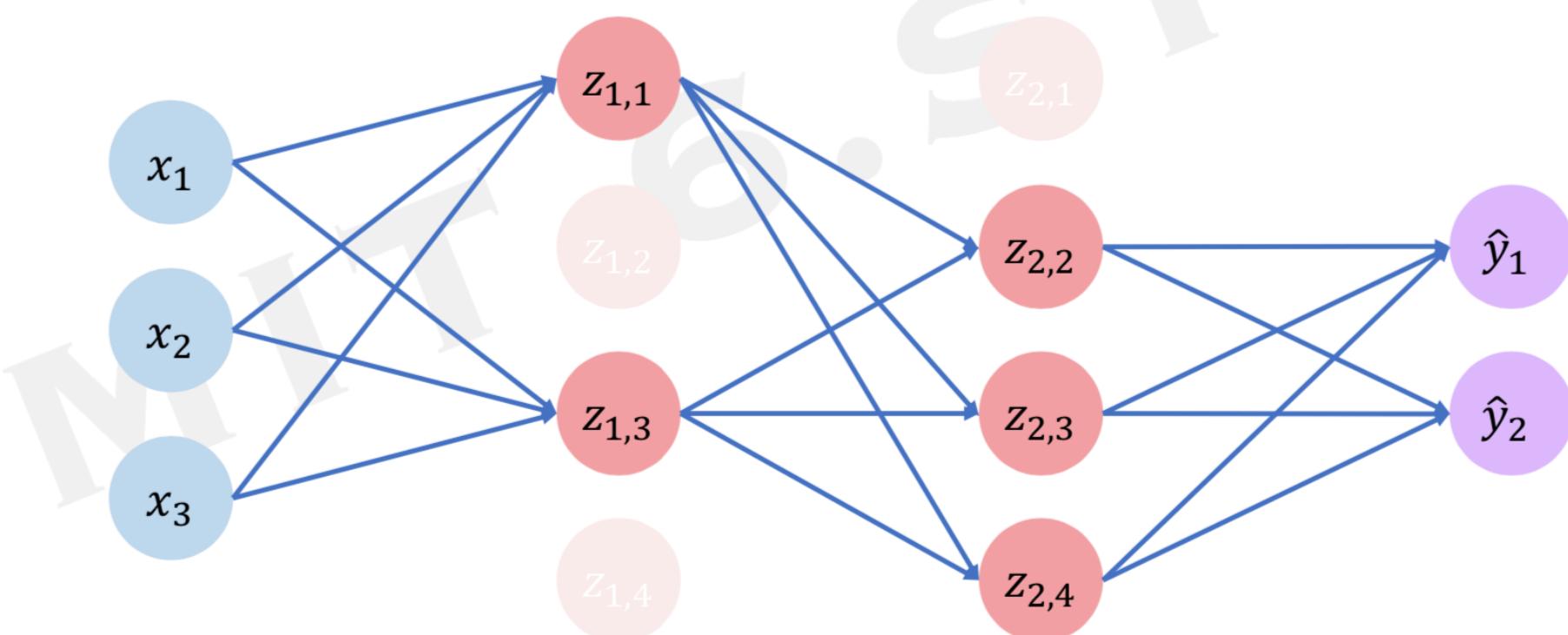
 `tf.keras.layers.Dropout(p=0.5)`



Regularization I: Dropout

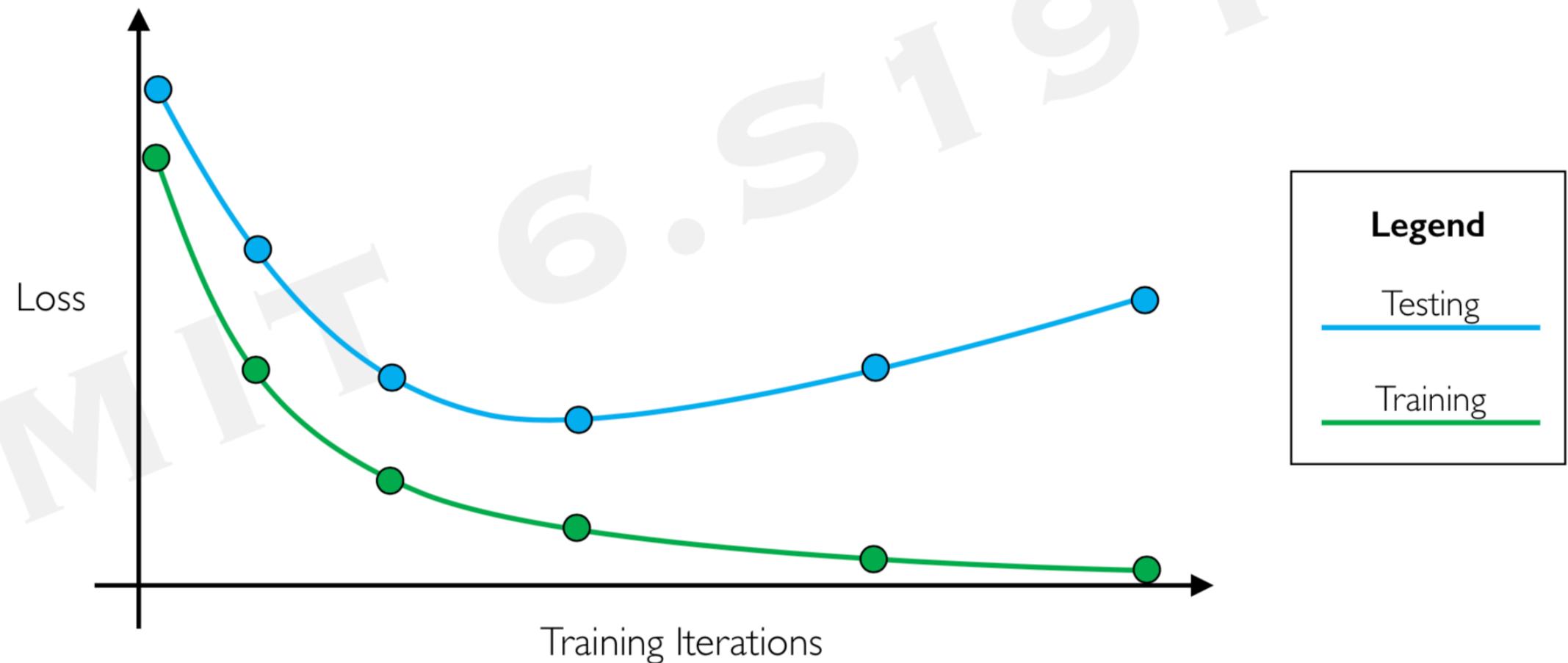
- During training, randomly set some activations to 0
 - Typically ‘drop’ 50% of activations in layer
 - Forces network to not rely on any 1 node

 `tf.keras.layers.Dropout(p=0.5)`



Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



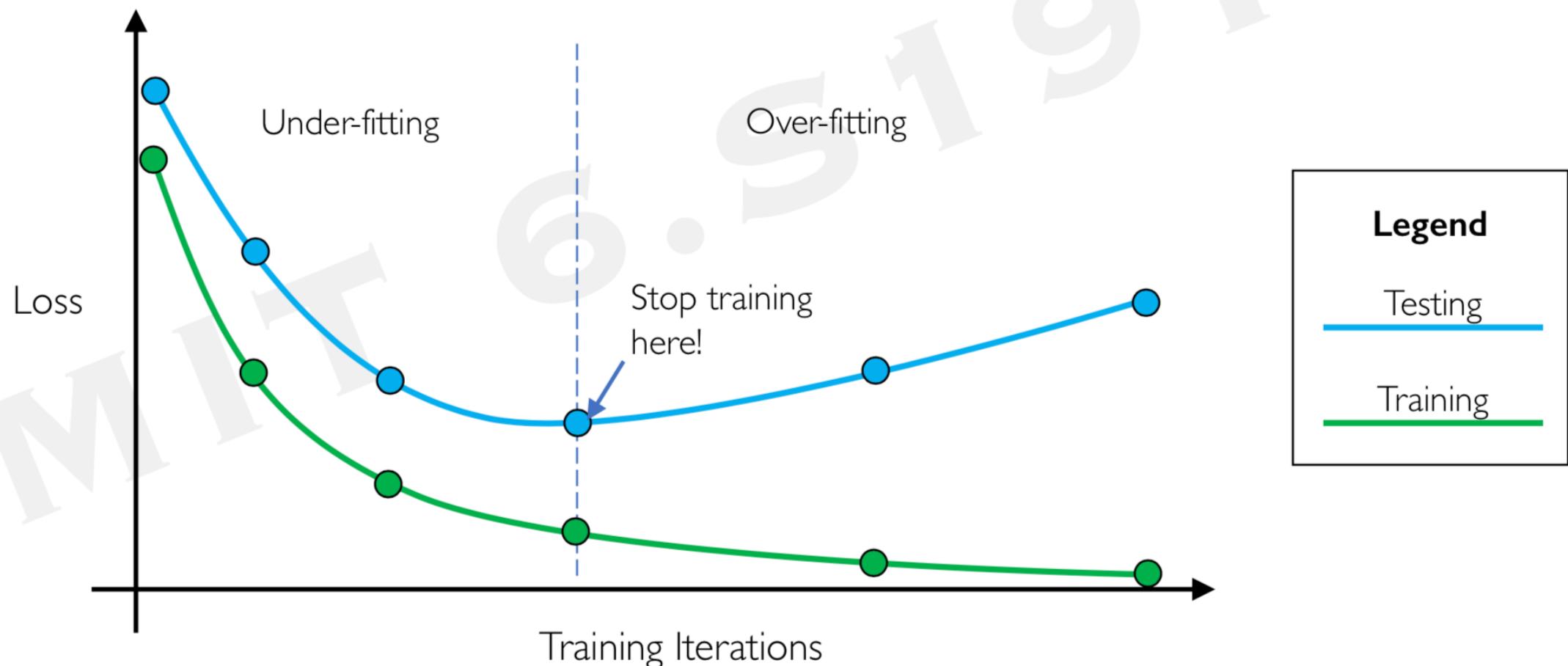
Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



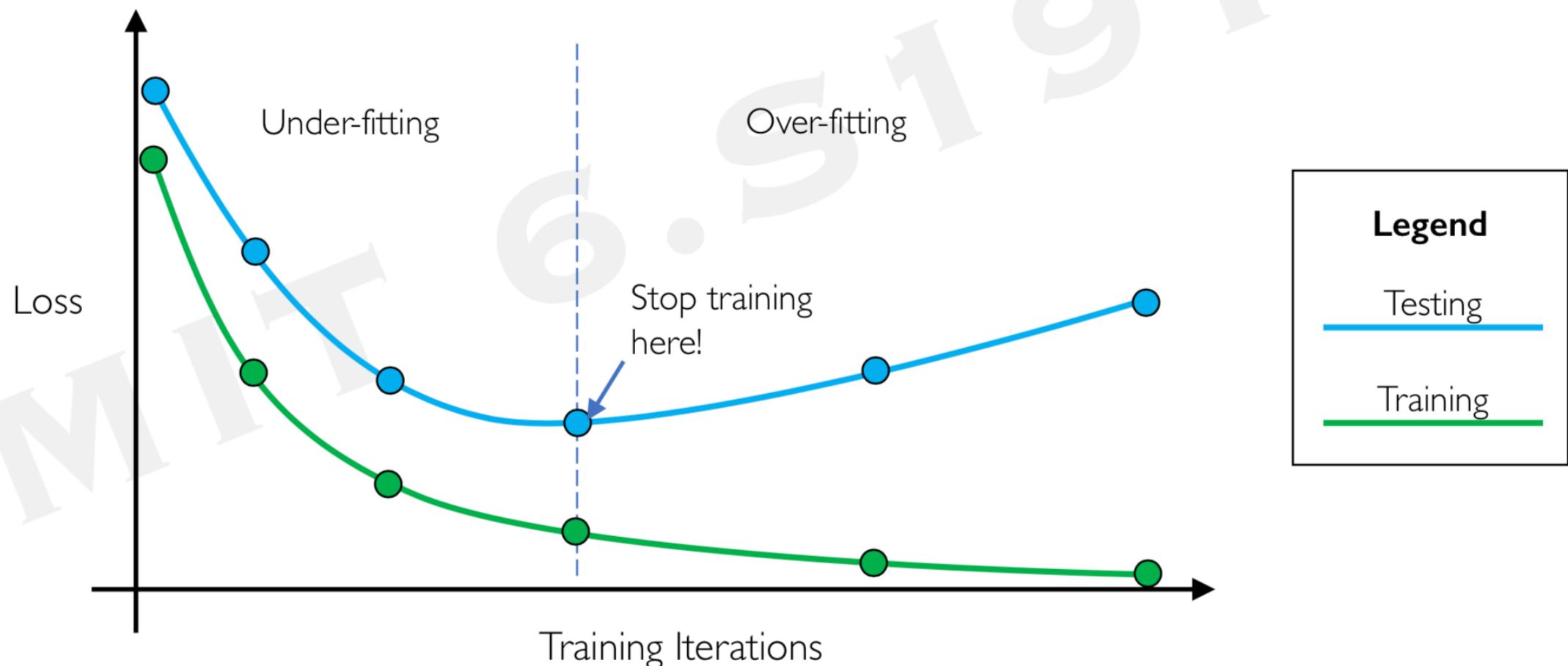
Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



Regularization 2: Early Stopping

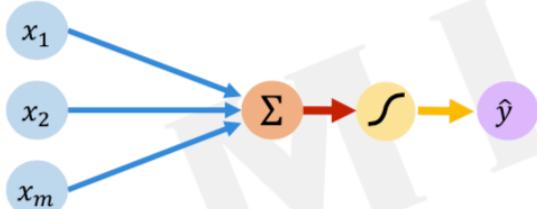
- Stop training before we have a chance to overfit



Core Foundation Review

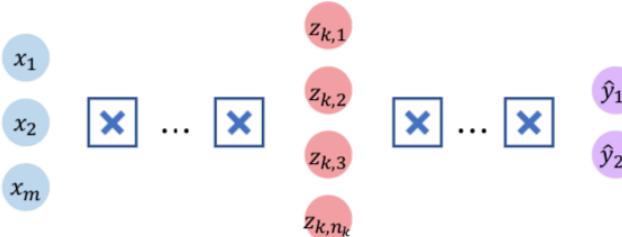
The Perceptron

- Structural building blocks
- Nonlinear activation functions



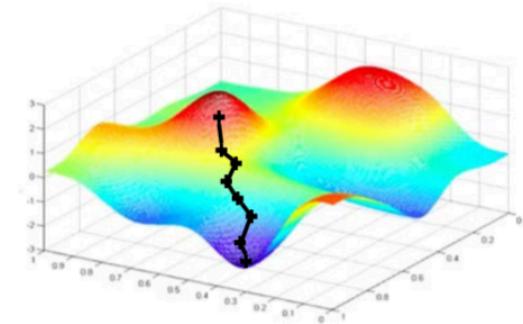
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization



References

- MIT 6.S191 Introduction to Deep Learning
- [Visualizing the Learning of a Neural Network Geometrically](#), Scott Rome
- [Top-30-data-science-interview-questions](#), Nitin Panwar