

## Why Deep Learning and Why Now?



Neural Networks date back decades, so why the dominance?

### 1. Big Data

- Larger Datasets
- Easier Collection & Storage

IMAGENET



### 2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



### 3. Software

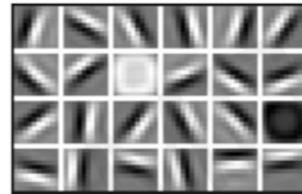
- Improved Techniques
- New Models
- Toolboxes

## 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



MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

1/6/25

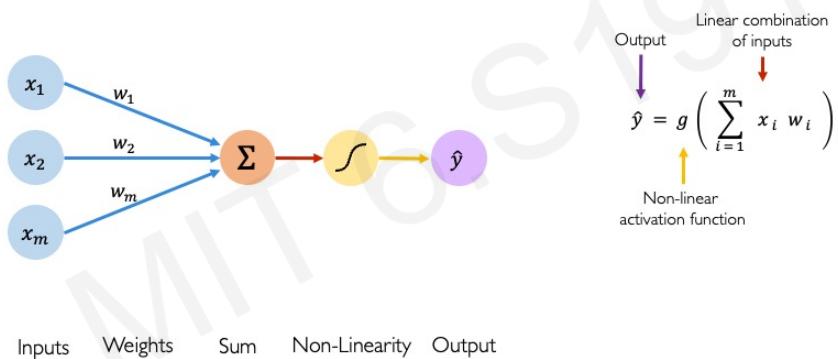
## Why Now?

Neural Networks date back decades, so why the dominance?

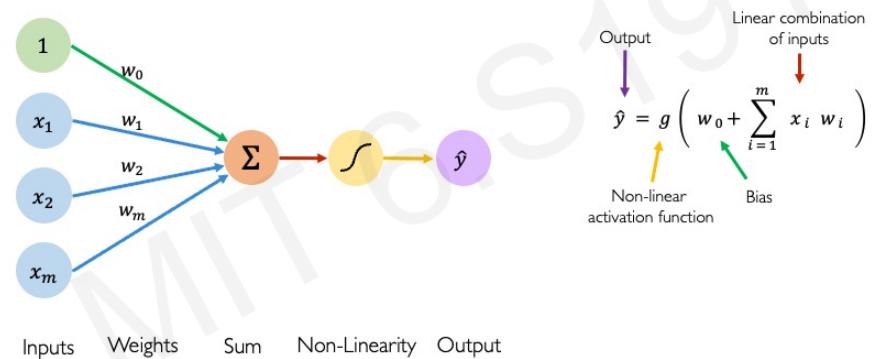
## The Perceptron

The structural building block of deep learning

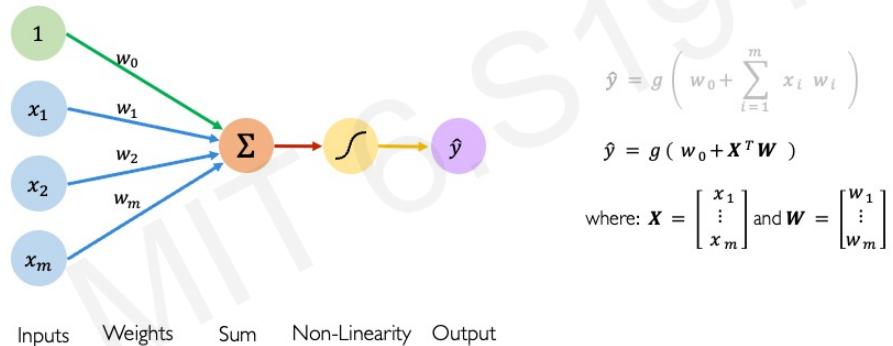
## The Perceptron: Forward Propagation



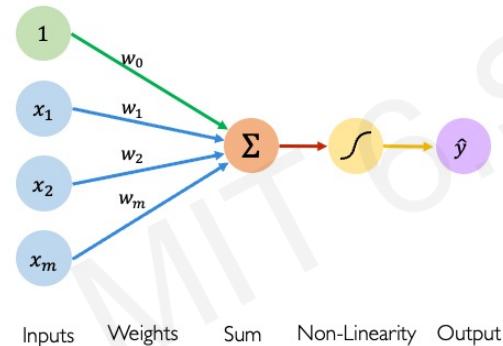
## The Perceptron: Forward Propagation



## The Perceptron: Forward Propagation



## The Perceptron: Forward Propagation

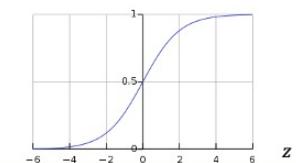


### Activation Functions

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

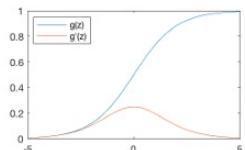
- Example: sigmoid function

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



## Common Activation Functions

Sigmoid Function

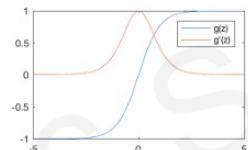


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

`tf.math.sigmoid(z)`  
 `torch.sigmoid(z)`

Hyperbolic Tangent

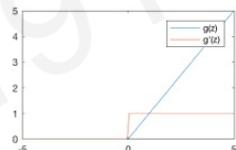


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

`tf.math.tanh(z)`  
 `torch.tanh(z)`

Rectified Linear Unit (ReLU)



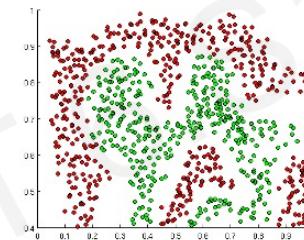
$$g(z) = \max(0, z)$$

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

`tf.nn.relu(z)`  
 `torch.nn.ReLU(z)`

## 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?

TensorFlow code blocks

NOTE: All activation functions are non-linear

MIT Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com) @MITDeepLearning

PyTorch code blocks

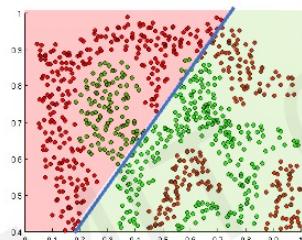
1/6/25 Massachusetts Institute of Technology

MIT Introduction to Deep Learning  
[introtodeelearning.com](http://introtodeeplearning.com) @MITDeepLearning

1/6/25

## Importance of Activation Functions

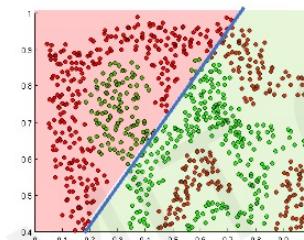
The purpose of activation functions is to **introduce non-linearities** into the network



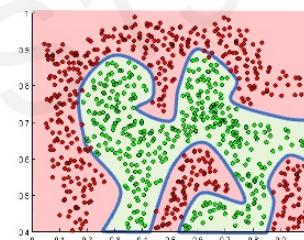
Linear activation functions produce linear decisions no matter the network size

## 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

Massachusetts Institute of Technology

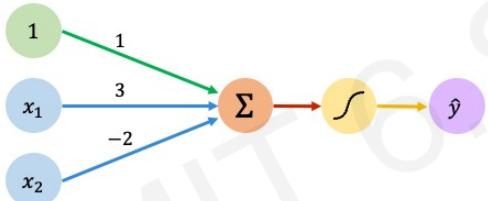
MIT Introduction to Deep Learning  
[introtodeeplearning.com](http://introtodeeplearning.com) @MITDeepLearning

1/6/25 Massachusetts Institute of Technology

MIT Introduction to Deep Learning  
[introtodeelearning.com](http://introtodeeplearning.com) @MITDeepLearning

1/6/25

## The Perceptron: Example

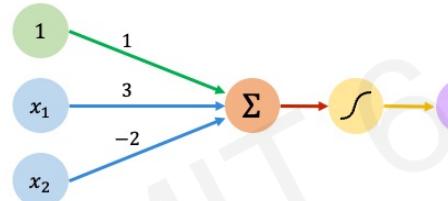


We have:  $w_0 = 1$  and  $\mathbf{w} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

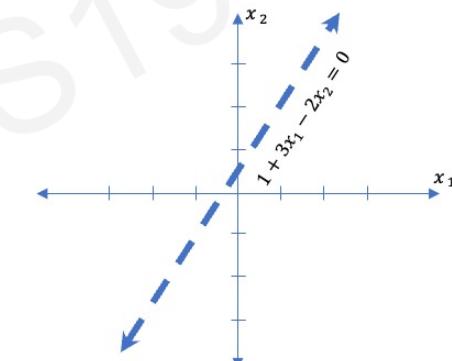
$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{x}^T \mathbf{w}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

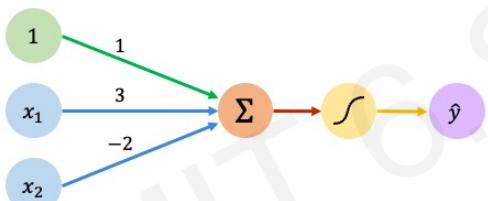
## The Perceptron: Example



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

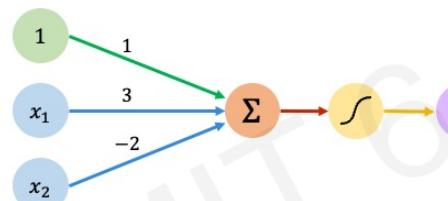
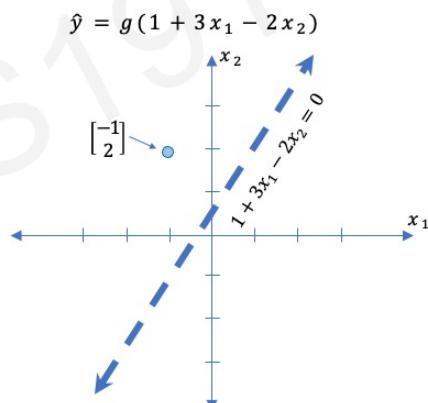


## The Perceptron: Example

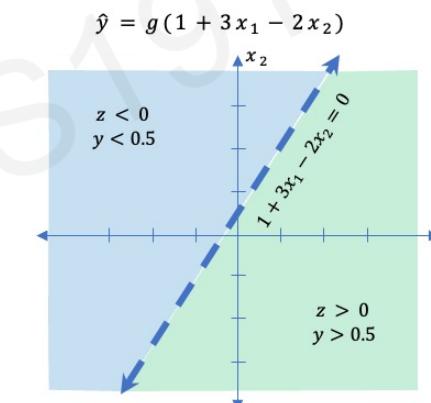


Assume we have input:  $\mathbf{x} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

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



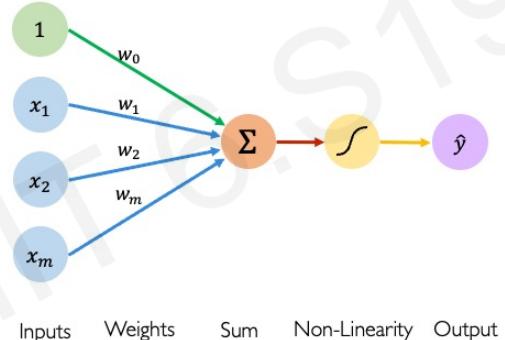
## The Perceptron: Example



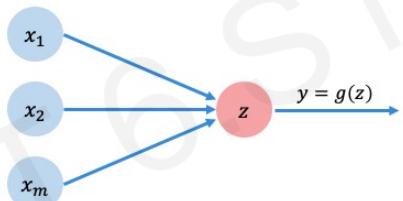
## Building Neural Networks with Perceptrons

### The Perceptron: Simplified

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



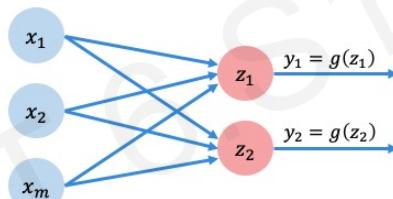
### The Perceptron: Simplified



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

### Multi Output Perceptron

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



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



## Dense layer from scratch

```
class MyDenseLayer(tf.keras.layers.Layer):
    def __init__(self, input_dim, output_dim):
        super(MyDenseLayer, self).__init__()

        # Initialize weights and bias
        self.W = self.add_weight([input_dim, output_dim])
        self.b = self.add_weight([1, output_dim])

    def call(self, inputs):
        # Forward propagate the inputs
        z = tf.matmul(inputs, self.W) + self.b

        # Feed through a non-linear activation
        output = tf.math.sigmoid(z)
        return output
```

```
class MyDenseLayer(nn.Module):
    def __init__(self, input_dim, output_dim):
        super(MyDenseLayer, self).__init__()

        # Initialize weights and bias
        self.W = nn.Parameter(torch.randn(input_dim,
                                         output_dim, requires_grad=True))
        self.b = nn.Parameter(torch.randn(1, output_dim,
                                         requires_grad=True))

    def forward(self, inputs):
        # Forward propagate the inputs
        z = torch.matmul(inputs, self.W) + self.b

        # Feed through a non-linear activation
        output = torch.sigmoid(z)
        return output
```

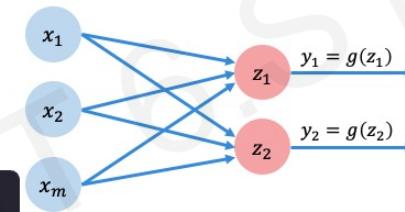
MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning



```
import tensorflow as tf
layer = tf.keras.layers.Dense(
    units=2)
```

MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

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



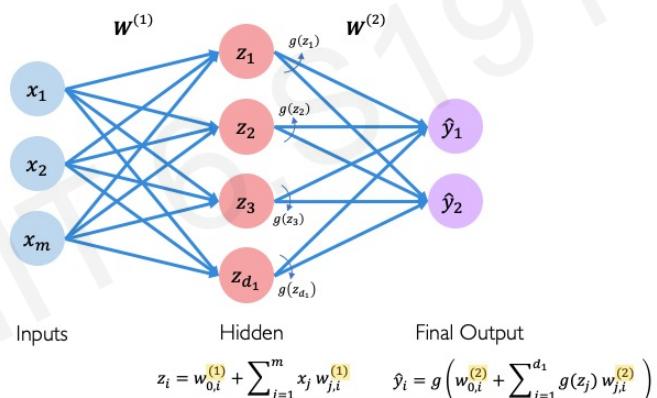
```
import torch.nn as nn
layer = nn.Linear(in_features=m,
                  out_features=2)
```

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

MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

1/6/25

## Single Layer Neural Network

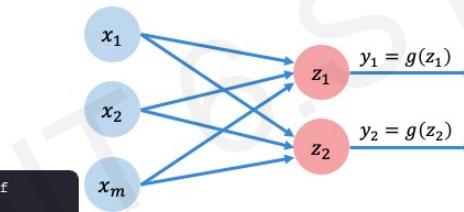


MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

1/6/25

MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

## Multi Output Perceptron



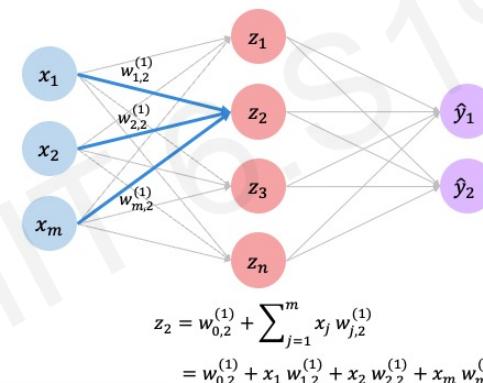
```
import torch.nn as nn
layer = nn.Linear(in_features=m,
                  out_features=2)
```

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

MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

1/6/25

## Single Layer Neural Network



MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

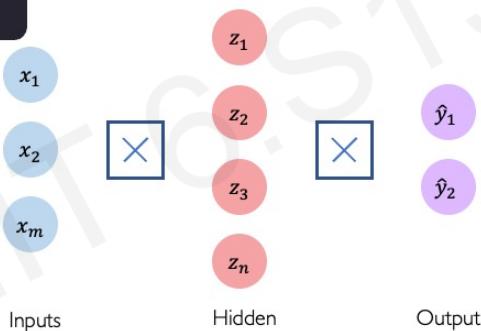
1/6/25

MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

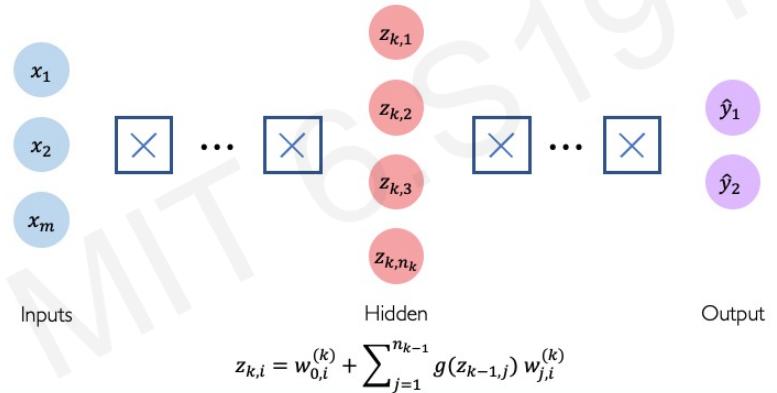
## Multi Output Perceptron

```
import tensorflow as tf
model = tf.keras.Sequential([
    tf.keras.layers.Dense(n),
    tf.keras.layers.Dense(2)
])

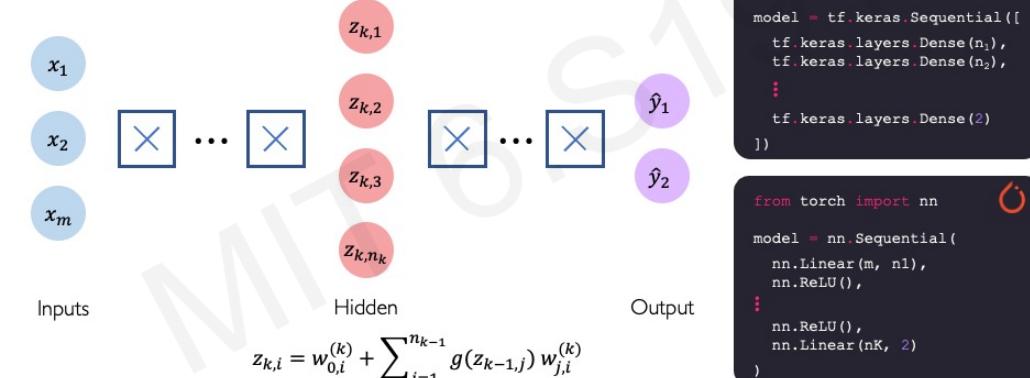
from torch import nn
model = nn.Sequential(
    nn.Linear(m, n),
    nn.ReLU(),
    nn.Linear(n, 2)
)
```



## Deep Neural Network



## Deep Neural Network



## Applying Neural Networks

## Example Problem

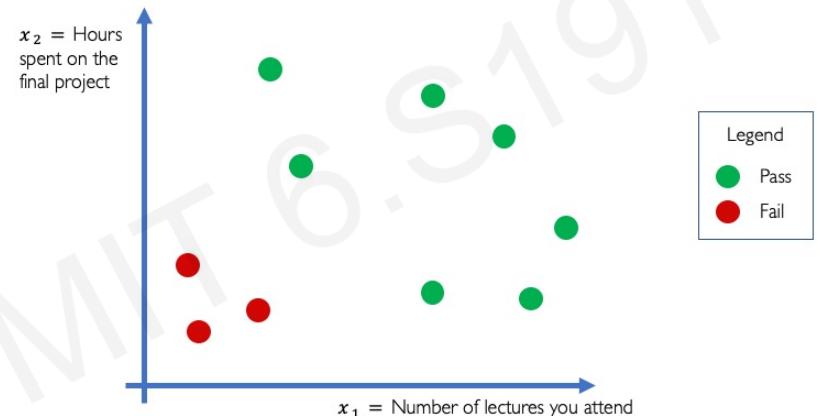
Will I pass this class?

Let's start with a simple two feature model

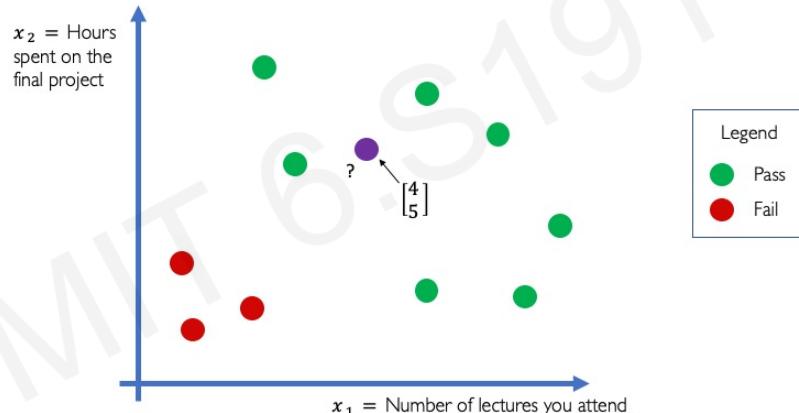
$x_1$  = Number of lectures you attend

$x_2$  = Hours spent on the final project

## Example Problem: Will I pass this class?



## Example Problem: Will I pass this class?



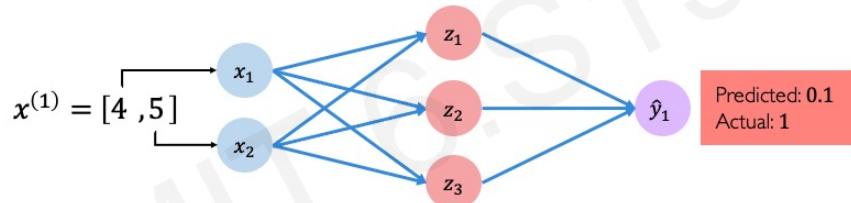
## Example Problem: Will I pass this class?

$$x^{(1)} = [4, 5]$$

$$\begin{array}{ccc} x_1 & \rightarrow & z_1 \\ x_2 & \rightarrow & z_2 \\ & \rightarrow & z_3 \\ & \rightarrow & \hat{y}_1 \end{array}$$

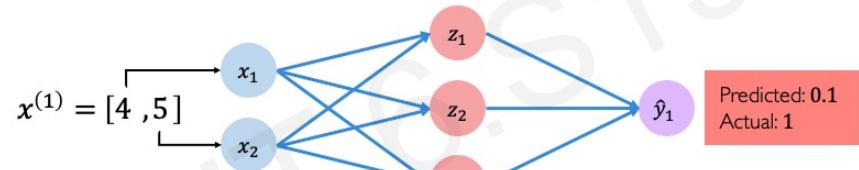
Predicted: 0.1

## Example Problem: Will I pass this class?



## Quantifying Loss

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

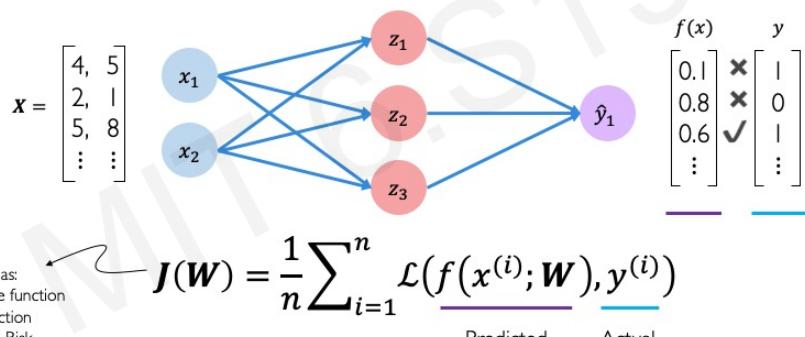


$$\mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted      Actual

## Empirical Loss

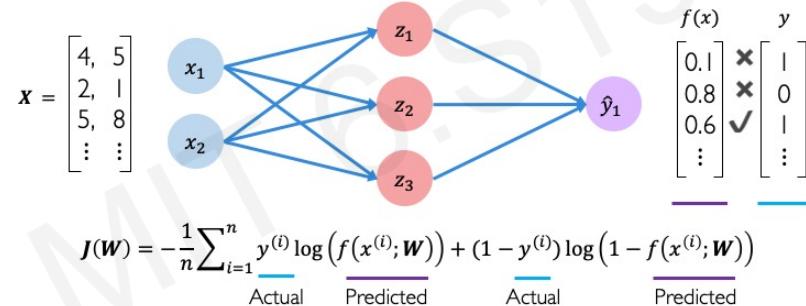
The **empirical loss** measures the total loss over our entire dataset



- Also known as:  
 • Objective function  
 • Cost function  
 • Empirical Risk

## Binary Cross Entropy Loss

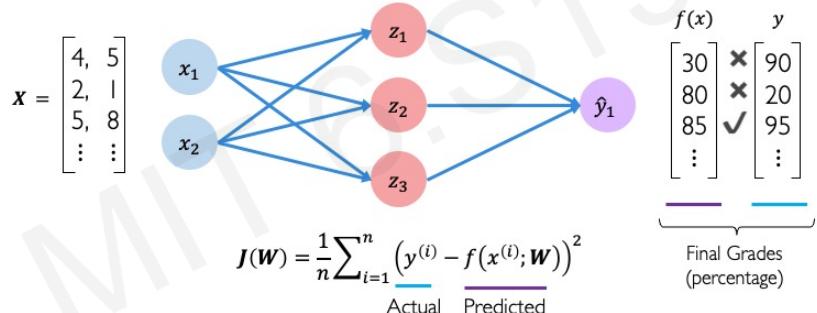
**Cross entropy loss** can be used with models that output a probability between 0 and 1



```
TensorFlow: loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(y, predicted))
PyTorch: loss = torch.nn.functional.cross_entropy(predicted, y)
```

# Mean Squared Error Loss

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



```
loss = tf.reduce_mean(tf.square(tf.subtract(y, predicted)))  
loss = tf.keras.losses.MSE(y, predicted)  
  
loss = torch.nn.functional.mse_loss(predicted, y)
```

# Training Neural Networks

## Loss Optimization

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

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

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

## Loss Optimization

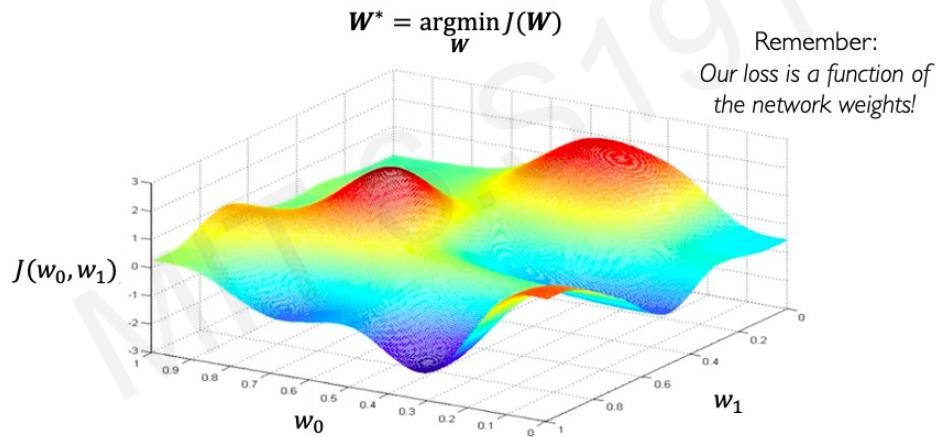
We want to find the network weights that achieve the lowest loss

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

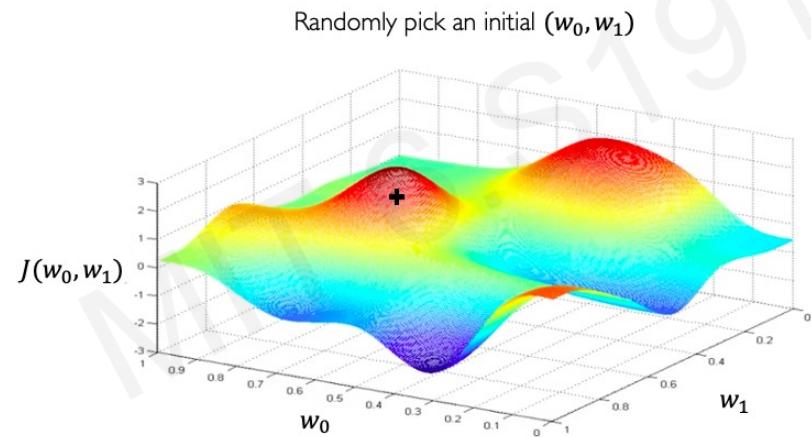
$$W^* = \operatorname{argmin}_W J(W)$$

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

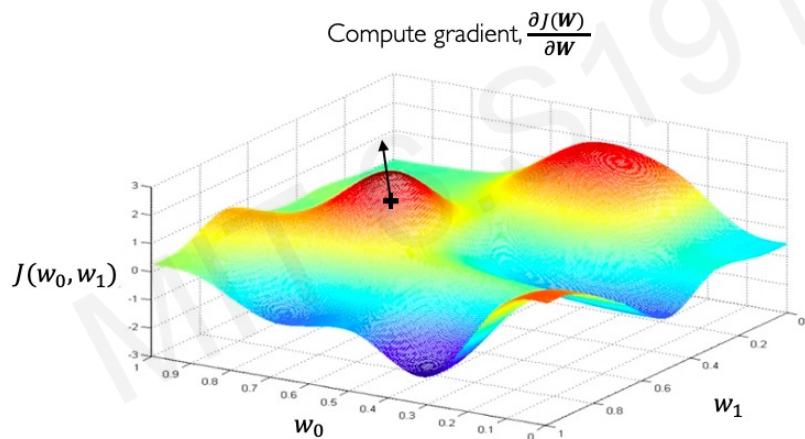
## Loss Optimization



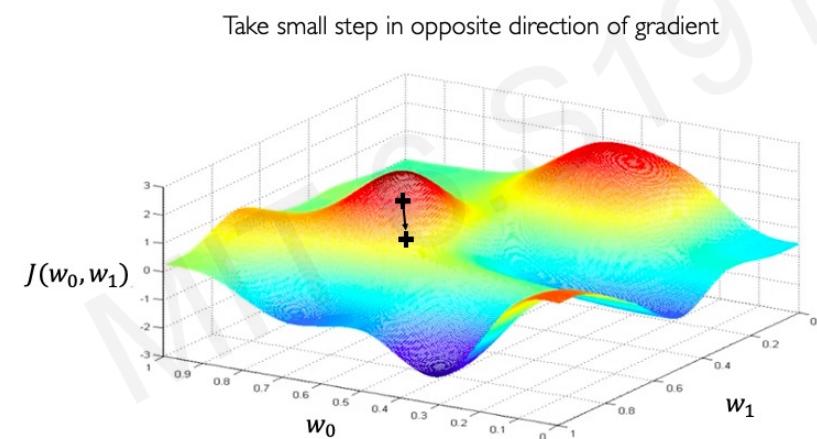
## Loss Optimization



## Loss Optimization

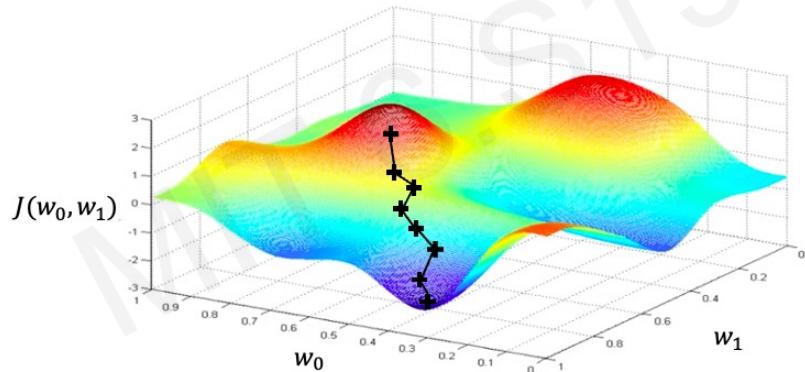


## Loss Optimization



# 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

# 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

```
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True:    # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient
```

# 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

```
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True:    # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

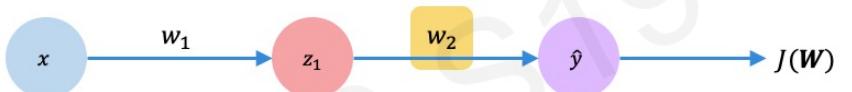
    weights = weights - lr * gradient
```

## Computing Gradients: Backpropagation



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

## Computing Gradients: Backpropagation



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

Let's use the chain rule!

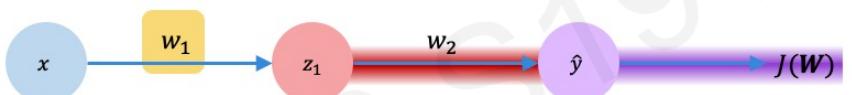
## Computing Gradients: Backpropagation



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

—      —

## Computing Gradients: Backpropagation



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

—      —

Apply chain rule!

Apply chain rule!

## Computing Gradients: Backpropagation



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

\_\_\_\_\_    \_\_\_\_\_    \_\_\_\_\_

## Computing Gradients: Backpropagation



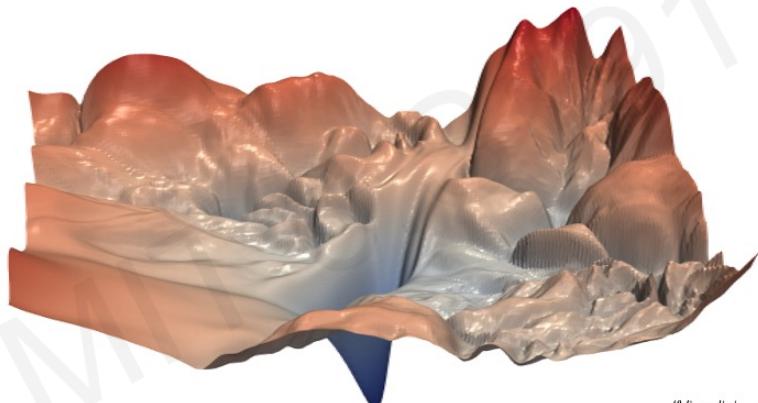
$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

\_\_\_\_\_    \_\_\_\_\_    \_\_\_\_\_

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

## Neural Networks in Practice: Optimization

## 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

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

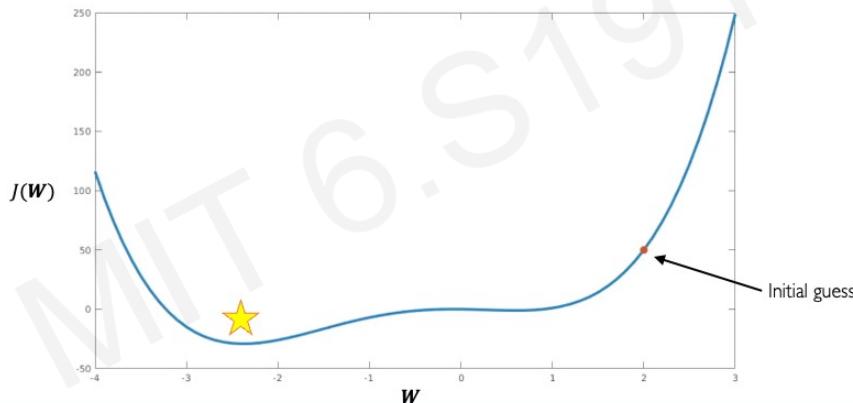
**Remember:**  
Optimization through gradient descent

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

How can we set the learning rate?

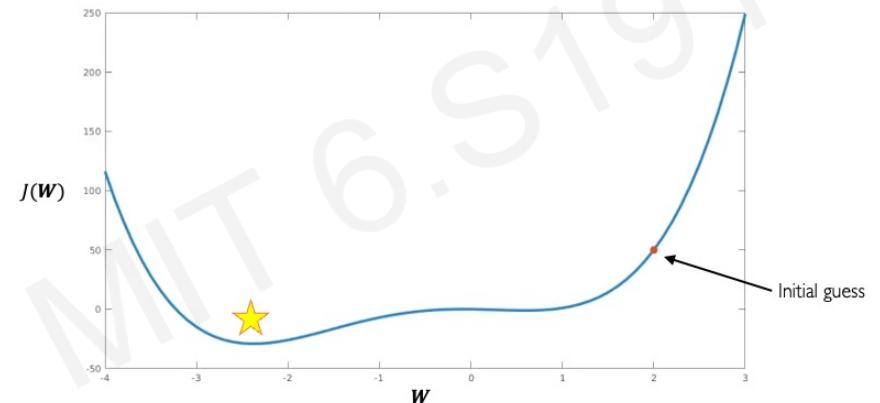
## Setting the Learning Rate

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



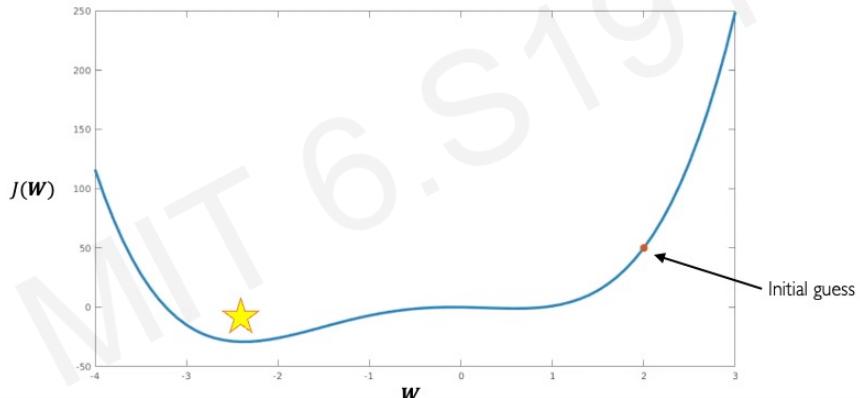
## Setting the Learning Rate

*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 1:

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

## TF Implementation

<code>tf.keras.optimizers.SGD</code>
<code>tf.keras.optimizers.Adam</code>
<code>tf.keras.optimizers.Adadelta</code>
<code>tf.keras.optimizers.Adagrad</code>
<code>tf.keras.optimizers.RMSProp</code>

## Torch Implementation

<code>torch.optim.SGD</code>
<code>torch.optim.Adam</code>
<code>torch.optim.Adadelta</code>
<code>torch.optim.Adagrad</code>
<code>torch.optim.RMSProp</code>

## Reference

- Kiefer & Wolfowitz, 1952.  
Kingma et al., 2014.  
Zeiler et al., 2012.  
Duchi et al., 2011.

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

## Putting it all together

```
import tensorflow as tf
model = tf.keras.Sequential([...])
# pick your favorite optimizer
optimizer = tf.keras.optimizer.SGD()
while True: # loop forever
    # forward pass through the network
    prediction = model(x)
    with tf.GradientTape() as tape:
        # compute the loss
        loss = compute_loss(y, prediction)
    # update the weights using the gradient
    grads = tape.gradient(loss, model.trainable_variables)
    optimizer.apply_gradients(zip(grads, model.trainable_variables))
```

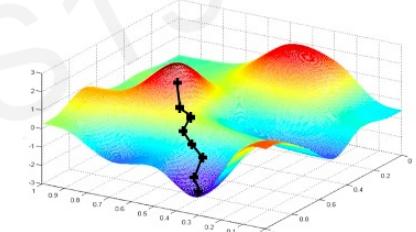
Can replace with  
any TensorFlow  
optimizer!

# 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

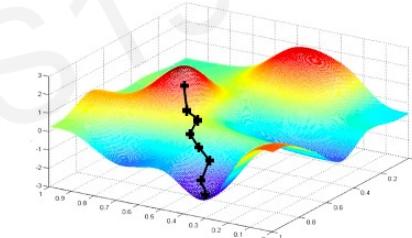


# Gradient Descent

## 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

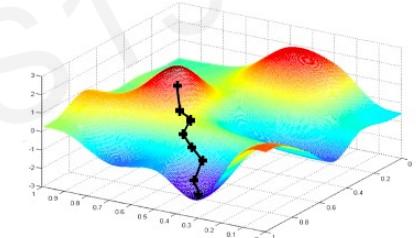
Can be very  
computationally  
intensive to compute!



# 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

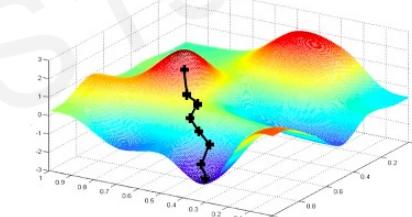


# Stochastic Gradient Descent

## 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

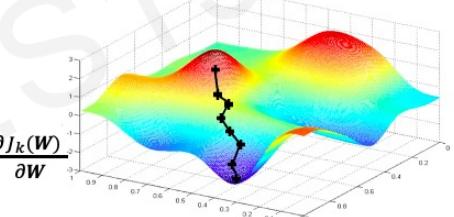
Easy to compute but  
very noisy (stochastic)!



# 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



# Stochastic Gradient Descent

## Algorithm

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

2. Loop until convergence:

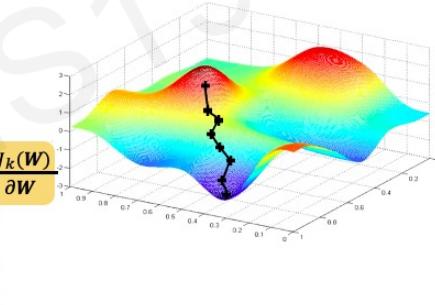
3. Pick batch of  $B$  data points

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!



# Mini-batches while training

## More accurate estimation of gradient

Smoother convergence

Allows for larger learning rates

# Mini-batches while training

## More accurate estimation of gradient

Smoother convergence

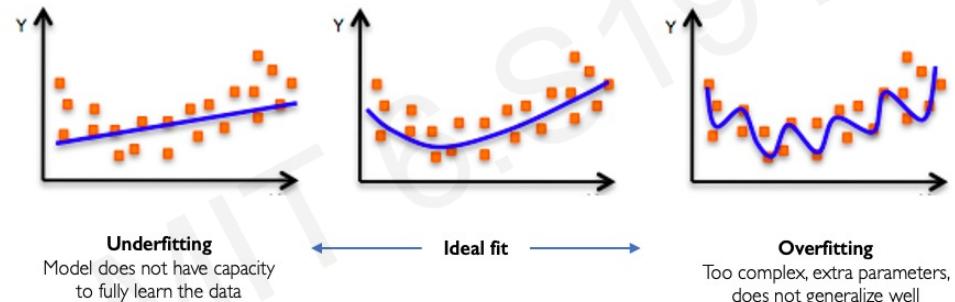
Allows for larger learning rates

## Mini-batches lead to fast training!

Can parallelize computation + achieve significant speed increases on GPU's

# Neural Networks in Practice: Overfitting

# The Problem of Overfitting



# Regularization

## What is it?

Technique that constrains our optimization problem to discourage complex models

# Regularization

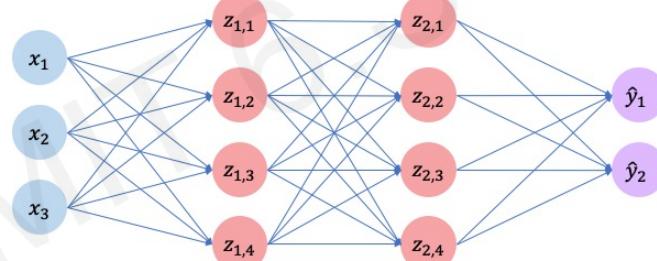
What is it?  
Technique that constrains our optimization problem to discourage complex models

## Why do we need it?

Improve generalization of our model on unseen data

# Regularization I: Dropout

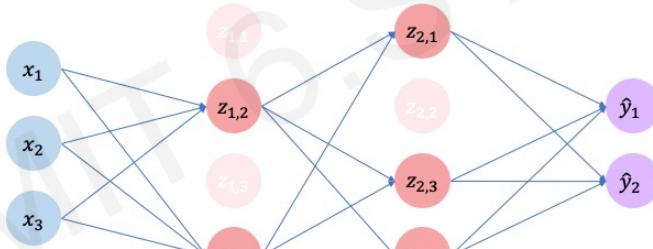
- During training, randomly set some activations to 0



## 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

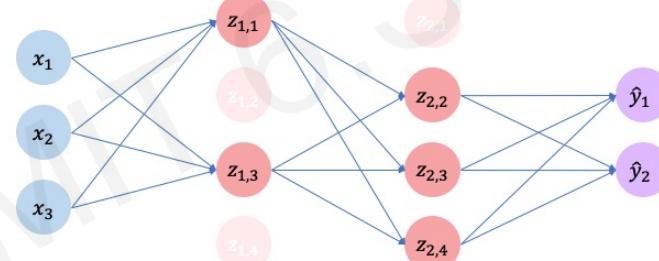
 `tf.keras.layers.Dropout(p=0.5)`  
 `torch.nn.Dropout(p=0.5)`



## 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)`  
 `torch.nn.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



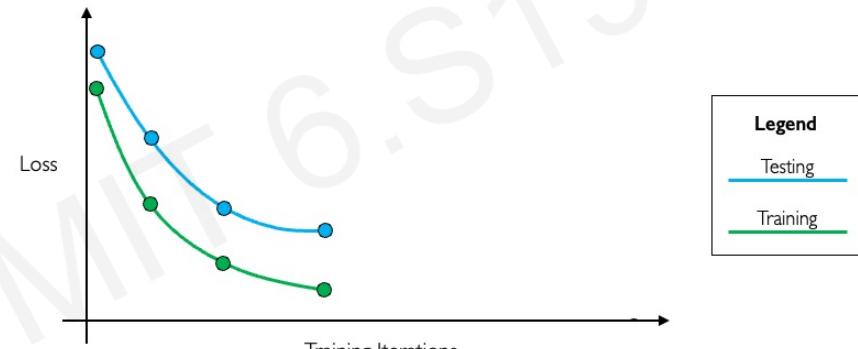
## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



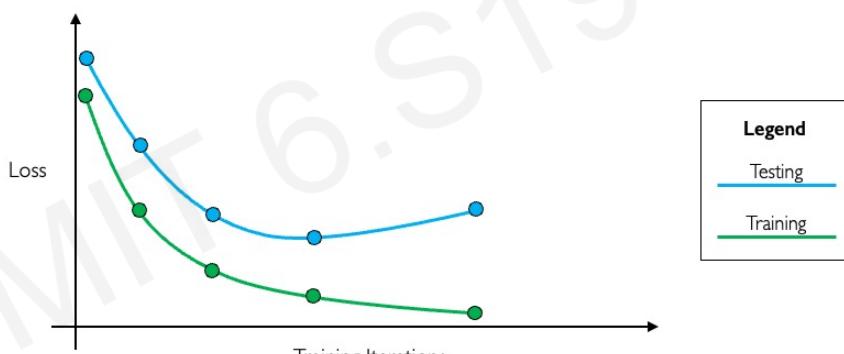
## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



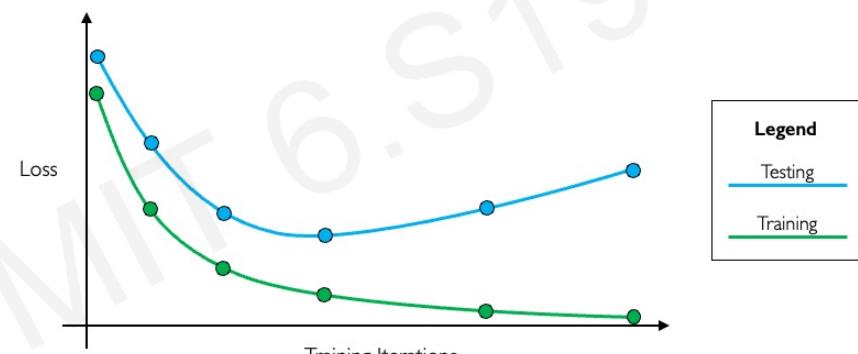
## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



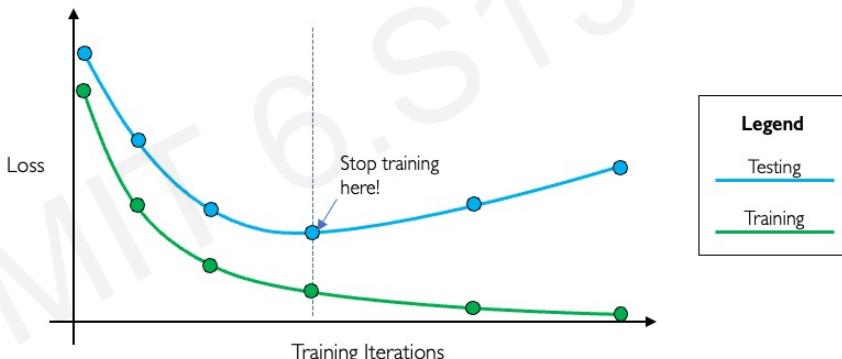
## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

## Regularization 2: Early Stopping

- Stop training before we have a chance to overfit

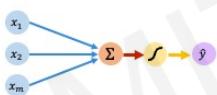


MIT Introduction to Deep Learning  
introtodeeplearning.com @MITDeepLearning

## 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

