# Day 4/5 - Introduction to Neural Nets / Deep Learning for NLP

Advanced Text as Data: Natural Language Processing
Essex Summer School in Social Science Data Analysis

Burt L. Monroe (Instructor) & Sam Bestvater (TA)
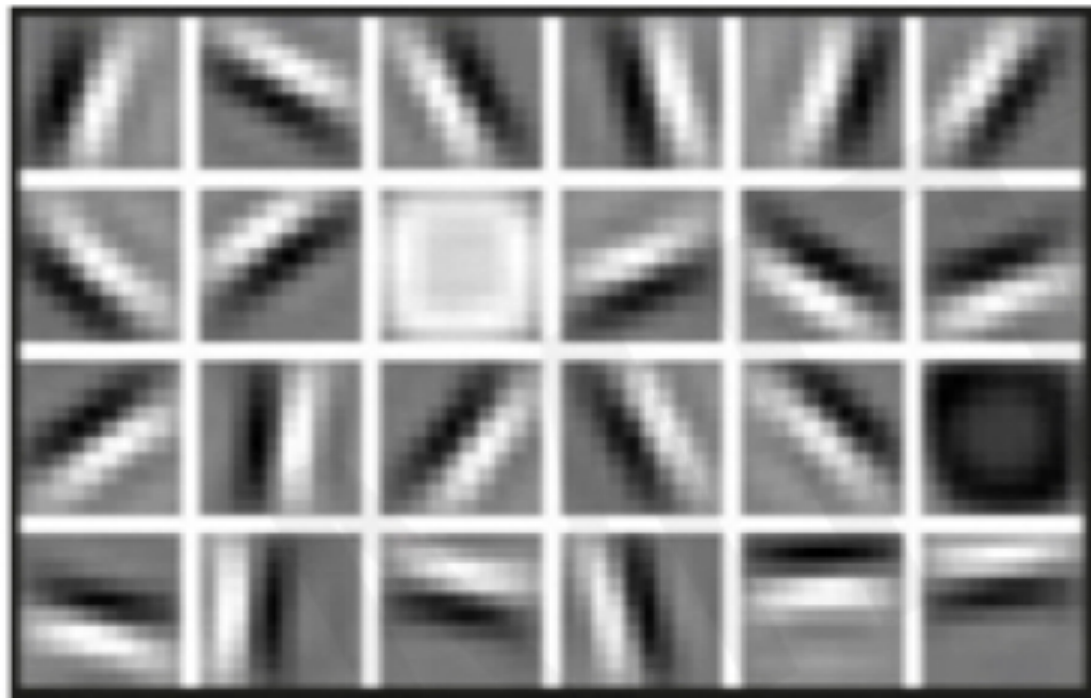Pennsylvania State University

July 29-30, 2021

# Why Deep Learning?

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

Can we learn the **underlying features** directly from data?
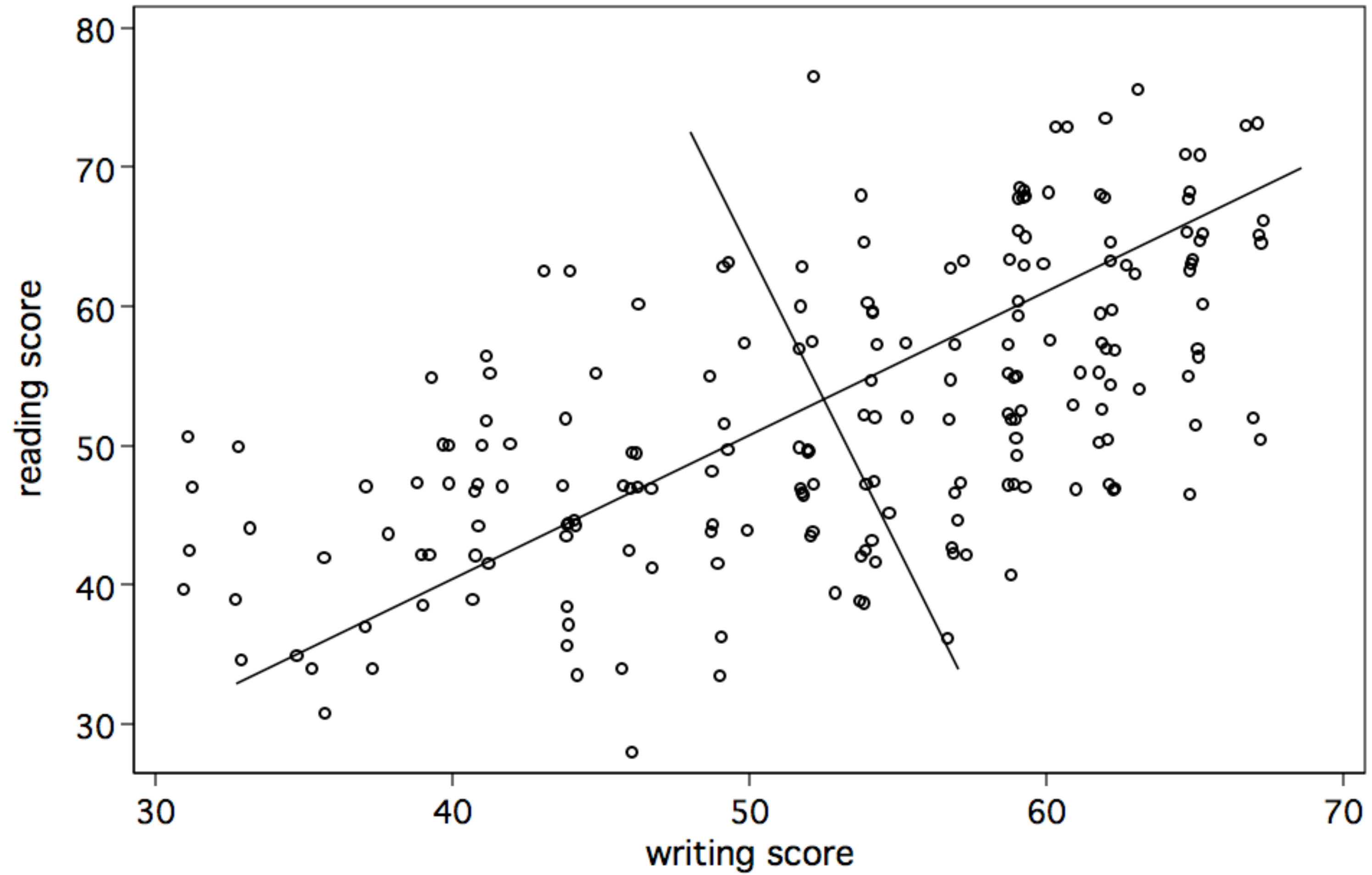
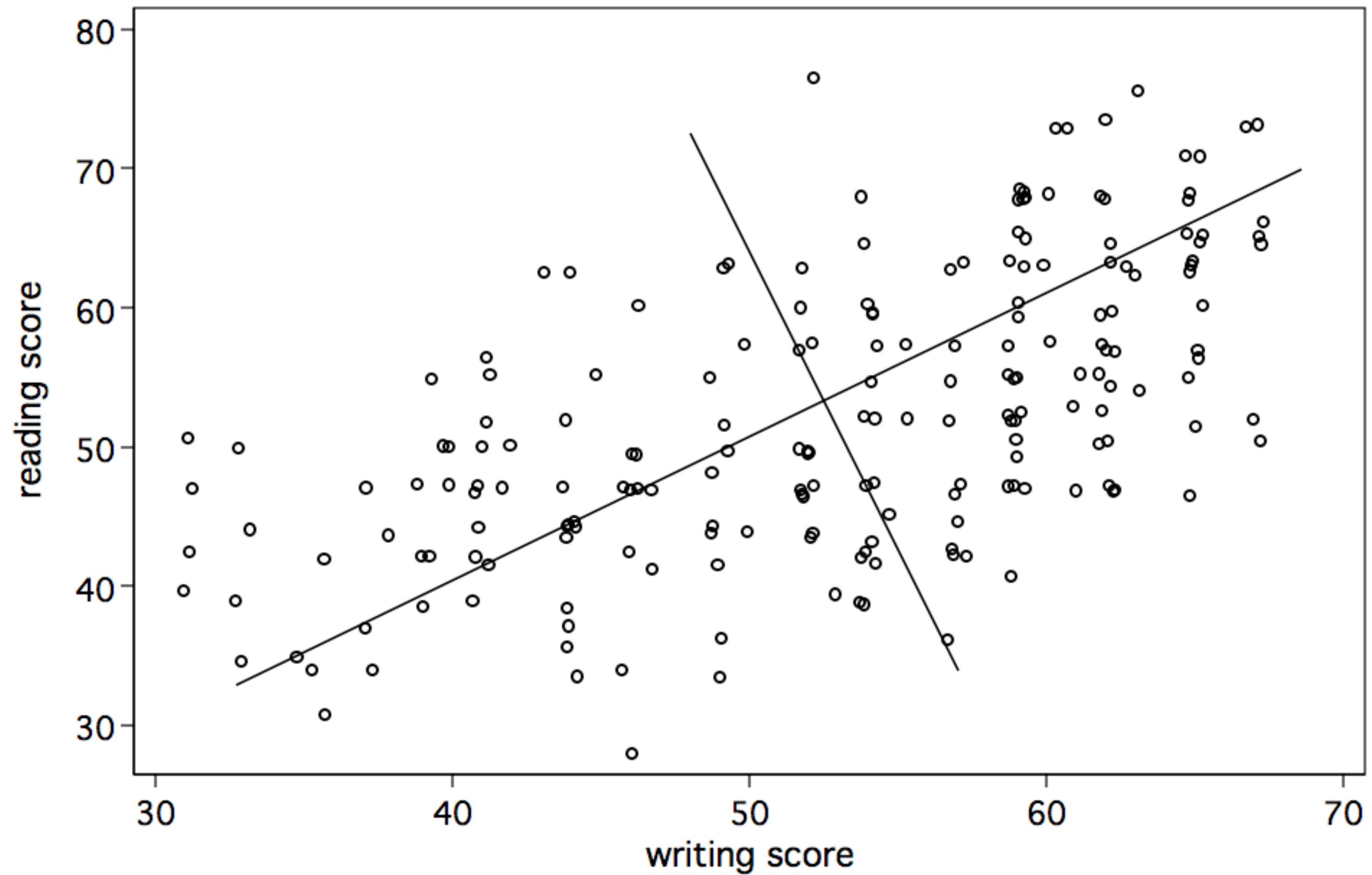| **Low Level Features** | **Mid Level Features** | **High Level Features** |
|:---:|:---:|:---:|
|  |  |  |
| Lines & Edges | Eyes & Nose & Ears | Facial Structure |

You've probably seen a lot of linear approaches to finding new feature representations:

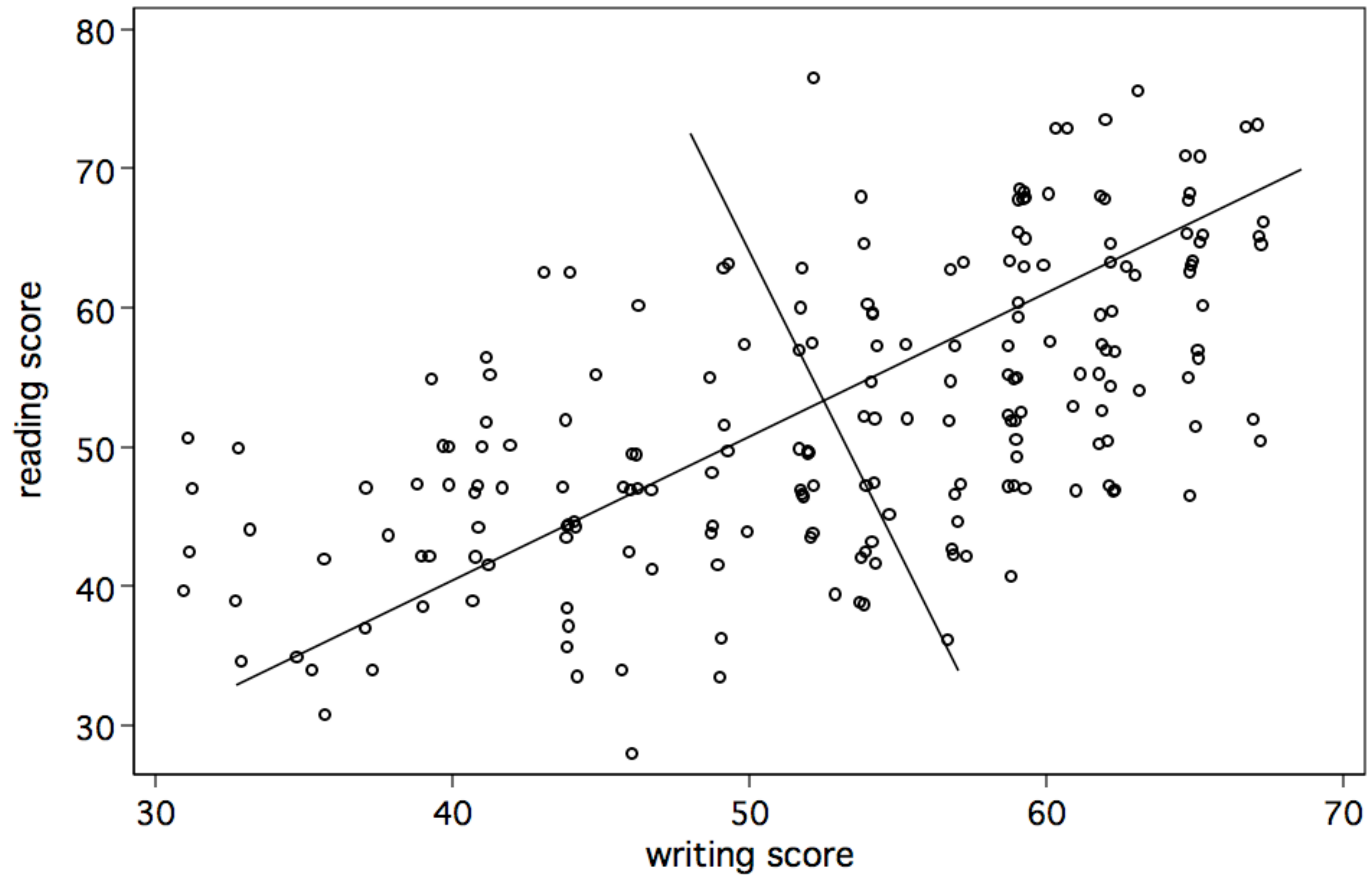You've probably seen a lot of linear approaches to finding new feature representations:

We might do this to find interpretable or intuitive latent concepts.
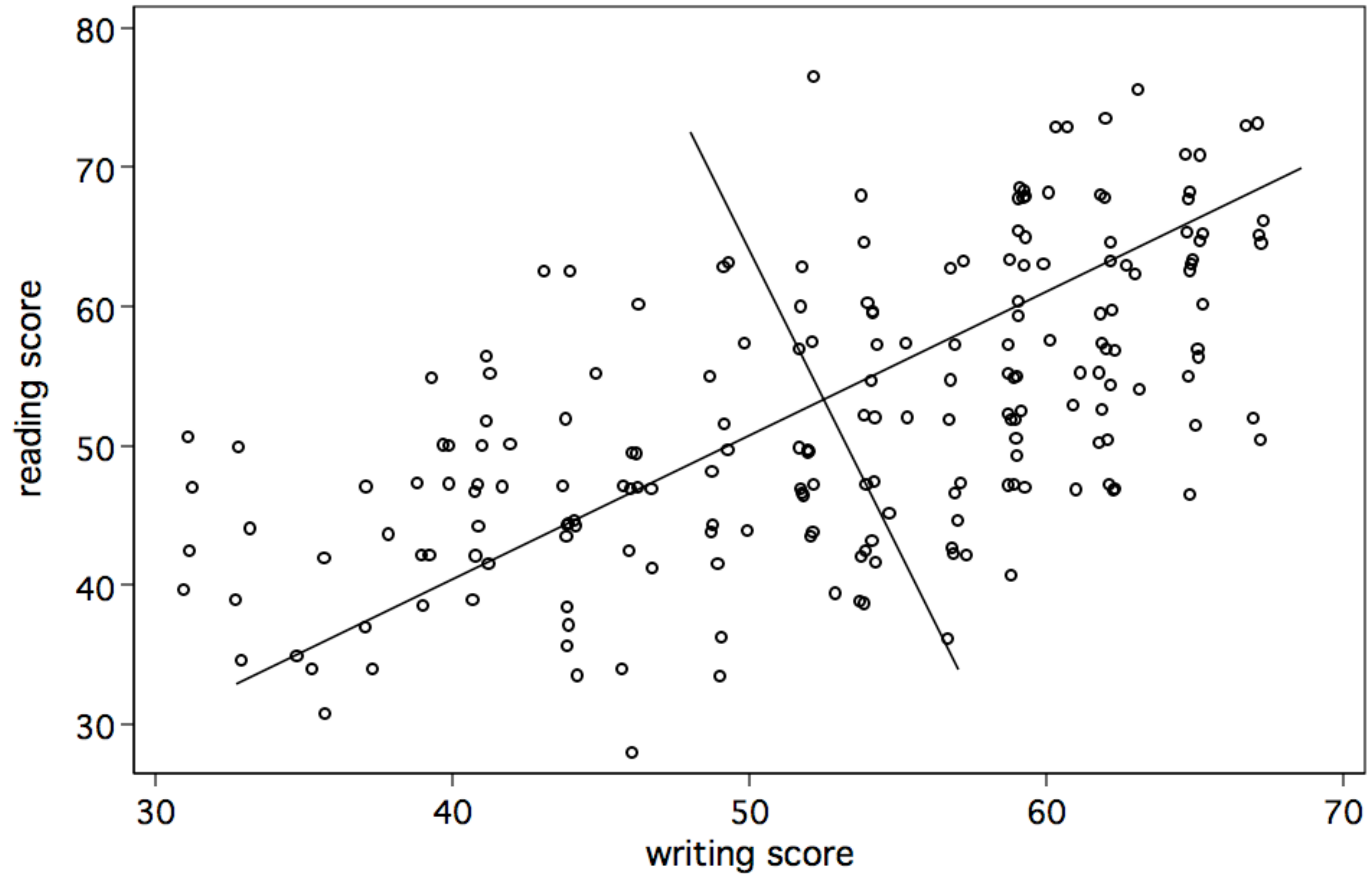
You've probably seen a lot of linear approaches to finding new feature representations:

We might do this to make computing more efficient (e.g., orthogonalization).

You've probably seen a lot of linear approaches to finding new feature representations:

We might do this to reduce dimensionality for generalizability or compression

Domain knowledge may allow us to do successful *feature engineering.*

**Figure 4.3. Feature engineering for reading the time on a clock**

| | | |
|---|---|---|
| Raw data: pixel grid | | |
| Better features: clock hands' coordinates | {x1: 0.7, y1: 0.7} {x2: 0.5, y2: 0.0} | {x1: 0.0, y2: 1.0} {x2: -0.38, 2: 0.32} |
| Even better features: angles of clock hands | theta1: 45 theta2: 0 | theta1: 90 theta2: 140 |

# The Perceptron: Forward Propagation



Inputs   Weights   Sum   Non-Linearity   Output

Output — Linear combination of inputs

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

Non-linear activation function

# The Perceptron: Forward Propagation



Inputs    Weights    Sum    Non-Linearity    Output

Output      Linear combination of inputs

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

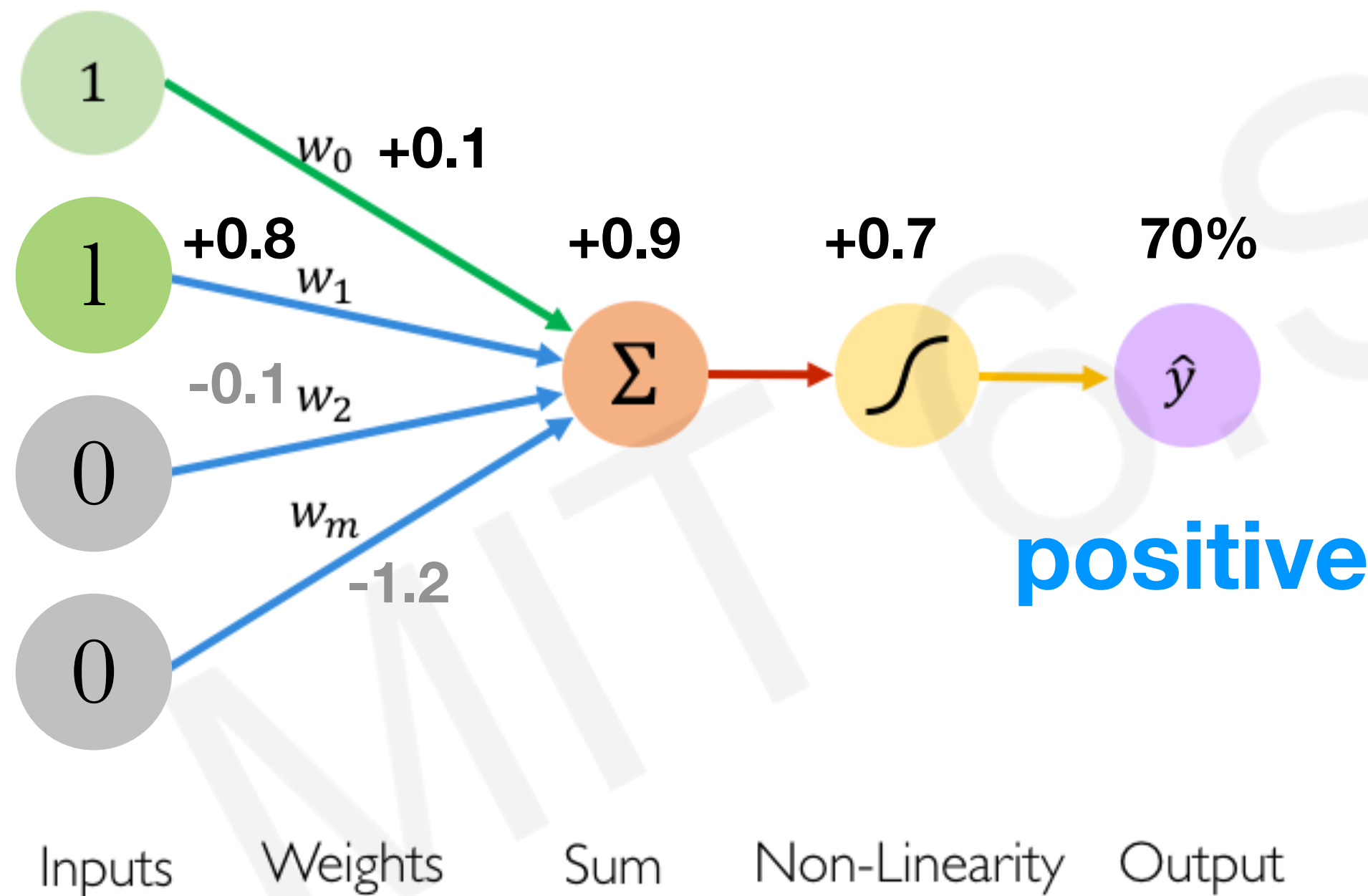Non-linear activation function      Bias

# The Perceptron: Forward Propagation



"It was great."

great

banana

worst

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

Inputs · Weights · Sum · Non-Linearity · Output

Output — Linear combination of inputs — Non-linear activation function — Bias

Massachusetts Institute of Technology

# The Perceptron: Forward Propagation



"The worst."

great    0

banana   0

worst    1

$w_0$ **+0.1**

**+0.8** $w_1$

**-0.1** $w_2$

$w_m$

**-1.2**

**-1.1**    **+0.25**    **25%**

Σ    ∫    $\hat{y}$

**negative**

Inputs    Weights    Sum    Non-Linearity    Output

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

"Acting was great.
Worst script ever."



**great**

**banana**

**worst**

1

great 1 +0.8

banana 0 -0.1

worst 1 -1.2

$w_0$ +0.1

$w_1$

$w_2$

$w_m$

-0.3    +0.43    43%

Σ    ∫    ŷ

**negative**

Inputs   Weights   Sum   Non-Linearity   Output

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

"I want a banana."



$$\hat{y} = g \left( w_0 + \sum_{i=1}^{m} x_i \ w_i \right)$$
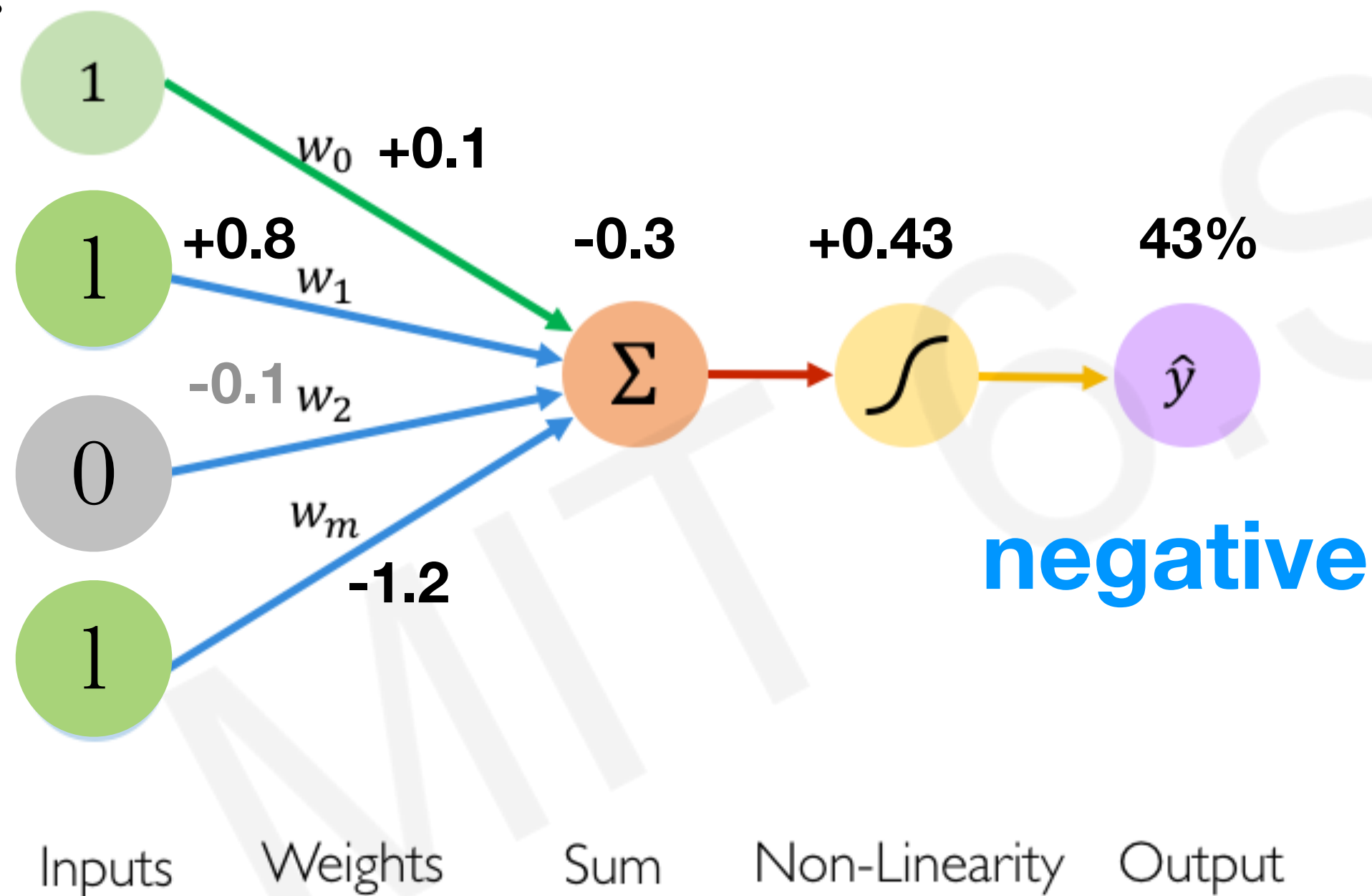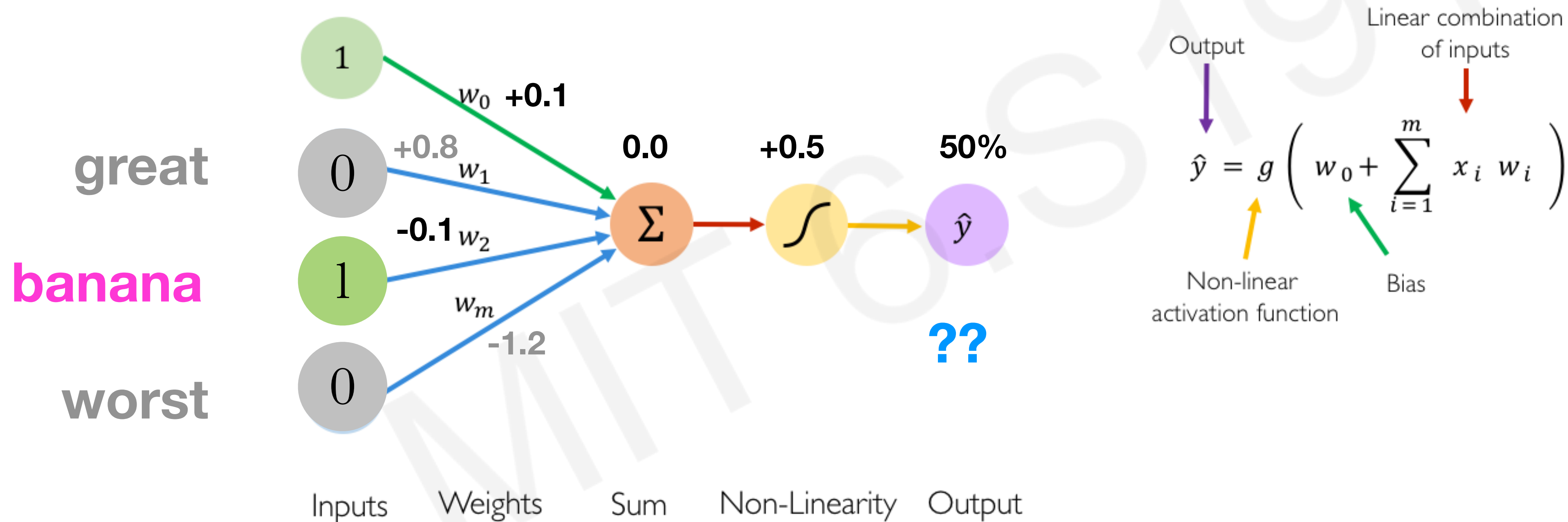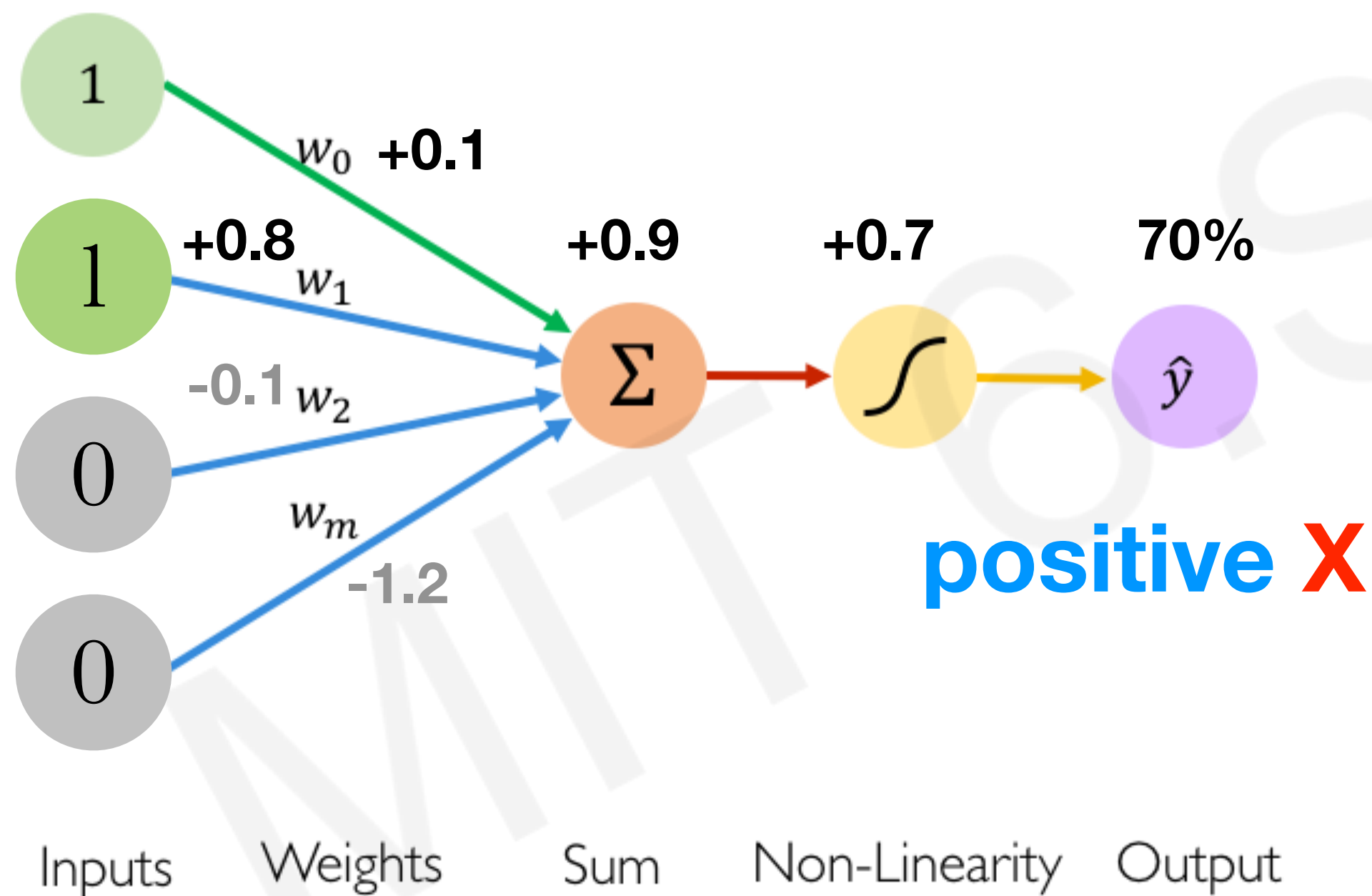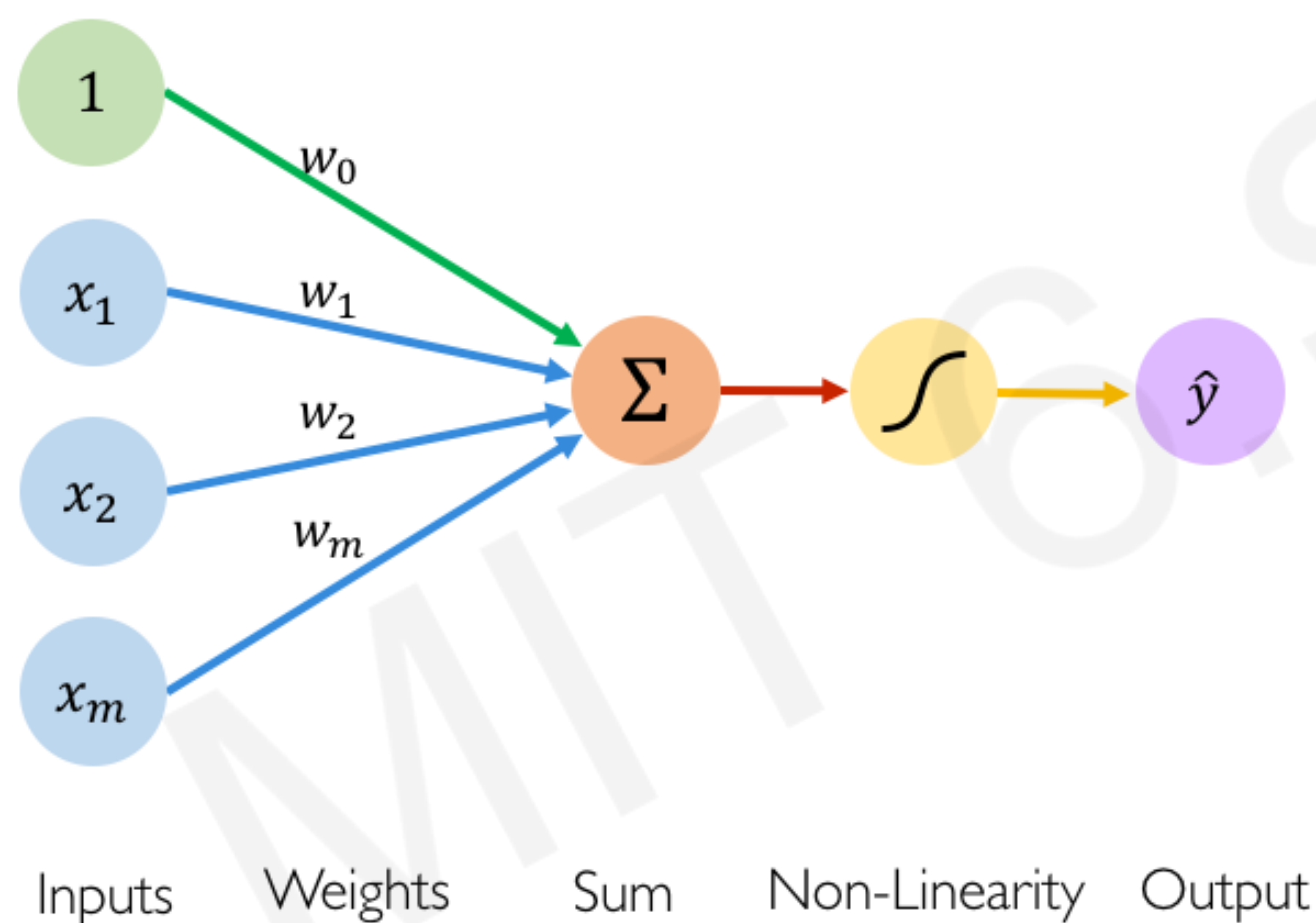
Output — Linear combination of inputs

Non-linear activation function — Bias

Massachusetts Institute of Technology

# The Perceptron: Forward Propagation



"Not great."

great

banana

worst

| 1 | | $w_0$ +0.1 | | | |

+0.8 $w_1$

-0.1 $w_2$

$w_m$

-1.2

+0.9    +0.7    70%

$\Sigma$    $\int$    $\hat{y}$

positive X

Inputs   Weights   Sum   Non-Linearity   Output

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

Output

Linear combination of inputs

Non-linear activation function    Bias

# The Perceptron: Forward Propagation



Inputs    Weights    Sum    Non-Linearity    Output

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

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

where: $\boldsymbol{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\boldsymbol{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$
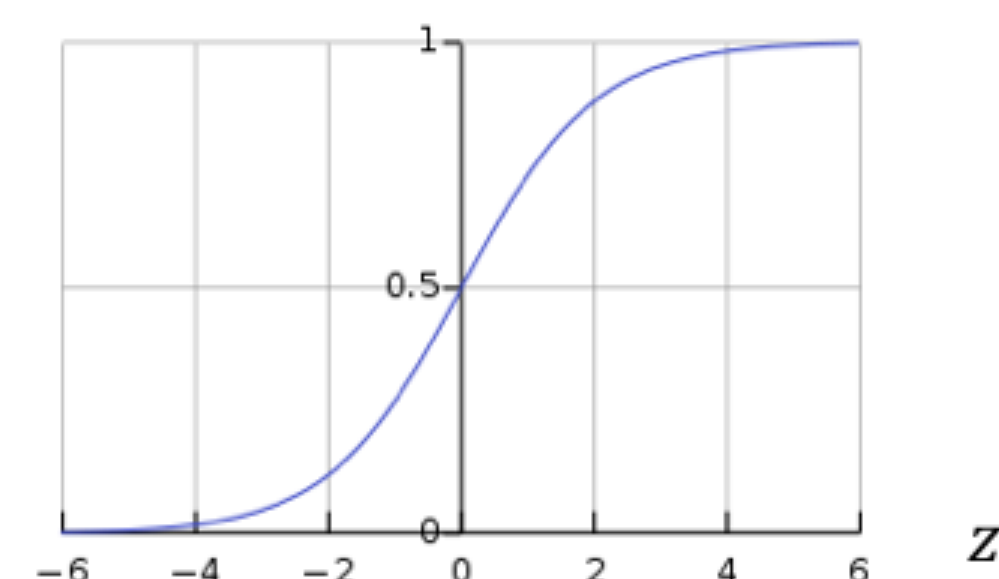
# The Perceptron: Forward Propagation



Inputs    Weights    Sum    Non-Linearity    Output

## Activation Functions

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

- Example: sigmoid function

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

Massachusetts Institute of Technology

This is *exactly the same* as the logit / logistic regression "inverse link function."

# The Perceptron: Forward Propagation



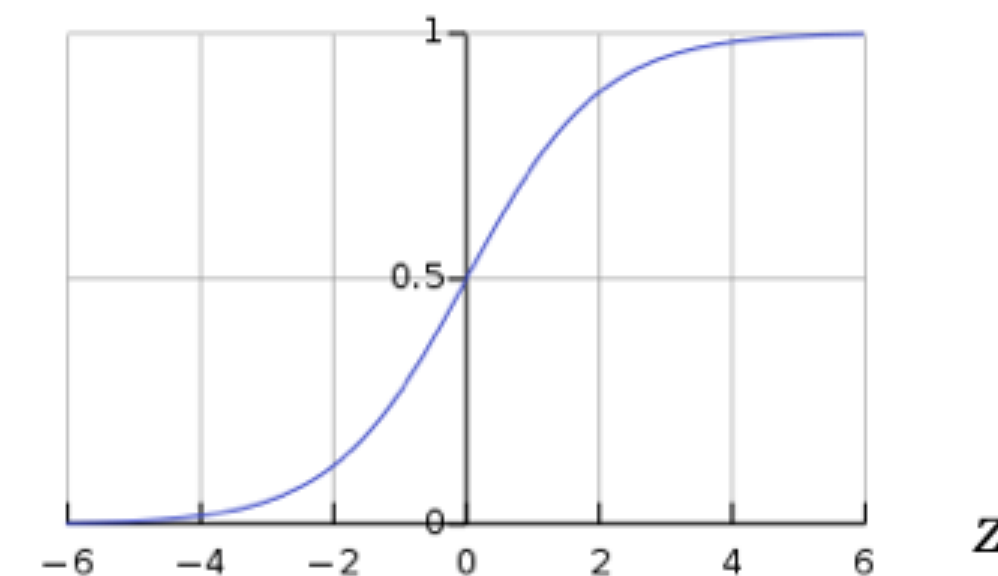Inputs    Weights    Sum    Non-Linearity    Output

## Activation Functions

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

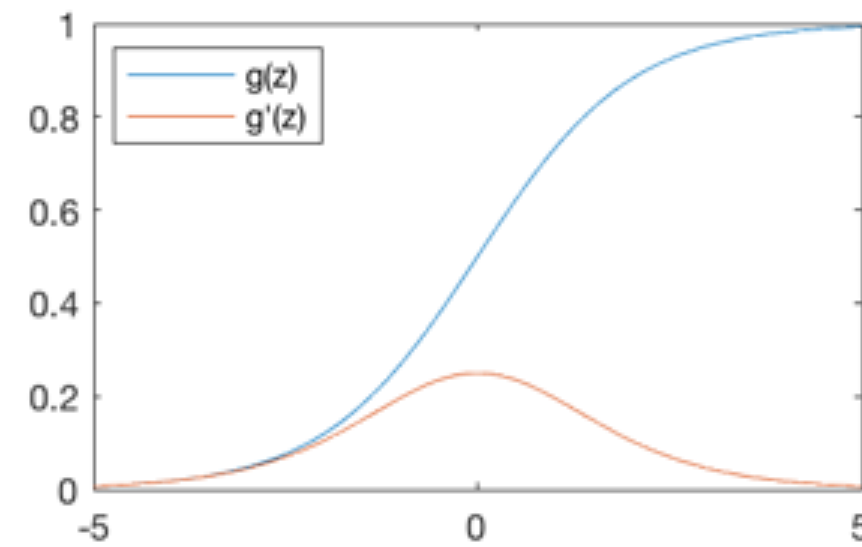- Example: sigmoid function

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

# Machine learning in one slide

| Social science (inference) | Machine learning (prediction) |
|:---:|:---:|
| GLM inverse link function | Activation function |
| $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ | $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ |
| Preferred objective function | |
| Log-likelihood | Cross-entropy |
| $\log \mathcal{L} = \sum_{i=1}^{n} \log P(y_i \mid \mathbf{x}_i, \boldsymbol{\beta})$ | $-\log \mathcal{L} = -\sum_{i=1}^{n} \log P(y_i \mid \mathbf{x}_i, \boldsymbol{\beta})$ |
| Solving algorithm | |
| Newton-Raphson | Gradient descent |
| $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - [\boldsymbol{H} \log \mathcal{L}]^{-1} \nabla \log \mathcal{L}$ | $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - \eta \nabla(-\log \mathcal{L})$ |
| Quantities of interest | |
| $\widehat{\boldsymbol{\beta}}; Var(\widehat{\boldsymbol{\beta}})$ | $\widehat{\boldsymbol{y}}; \sum \mathbf{1}(\widehat{y} = y)/n$ |

# Common Activation Functions

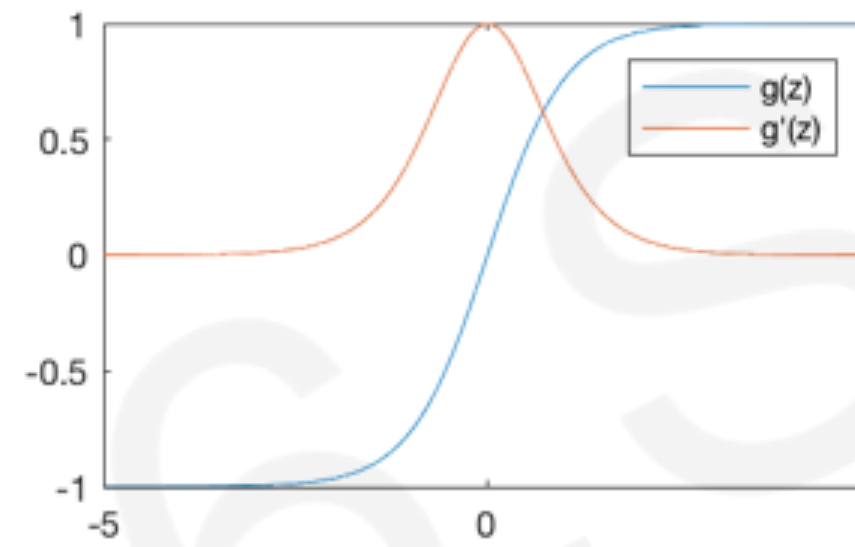| Sigmoid Function | Hyperbolic Tangent | Rectified Linear Unit (ReLU) |
|---|---|---|



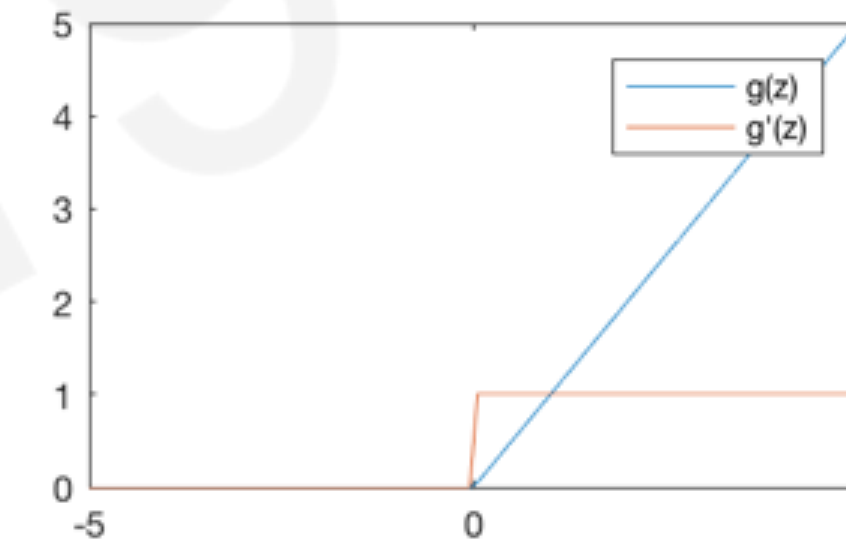$$g(z) = \frac{1}{1 + e^{-z}}$$

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

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

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

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

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

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

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

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

TensorFlow code blocks

NOTE: All activation functions are non-linear

# Common Activation Functions

```python
model = models.Sequential()
model.add(layers.Dense(16, activation = 'relu', input_shape=(5000,)))
model.add(layers.Dense(16, activation = 'relu'))
model.add(layers.Dense(1, activation= 'sigmoid'))

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])

history = model.fit(partial_x_train,
                    partial_y_train,
                    epochs=4,
                    batch_size=512,
                    validation_data=(x_val,y_val))
```

TensorFlow code blocks

NOTE: All activation functions are non-linear
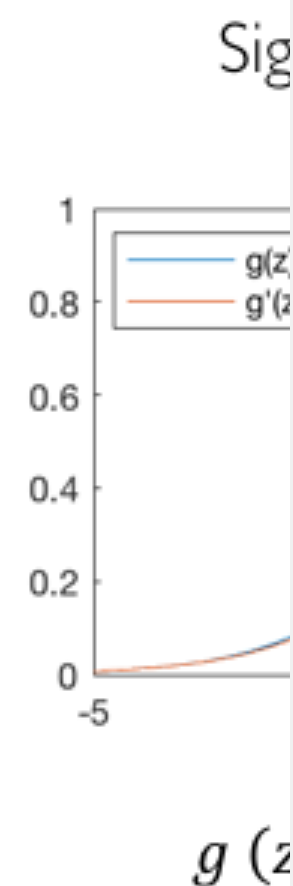
Sig

(ReLU)

```{r}
model <- keras_model_sequential() %>%
  layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
  layer_dense(units = 16, activation = "relu") %>%
  layer_dense(units = 1, activation = "sigmoid")

model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = c("accuracy")
)

model %>% fit(x_train, y_train, epochs = 4, batch_size = 512)
results <- model %>% evaluate(x_test, y_test)
```

$g(z$

$z)$

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

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

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

tf.math.sigmoid(z)

tf.math.tanh(z)

tf.nn.relu(z)

TensorFlow code blocks

NOTE: All activation functions are non-linear

# Building Neural Networks with Perceptrons

# The Perceptron: Simplified

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



| Inputs | Weights | Sum | Non-Linearity | Output |

# 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



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

# Multi Output Perceptron

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



$$y_1 = g(z_1)$$

$$y_2 = g(z_2)$$

```
import tensorflow as tf

layer = tf.keras.layers.Dense(
    units=2)
```
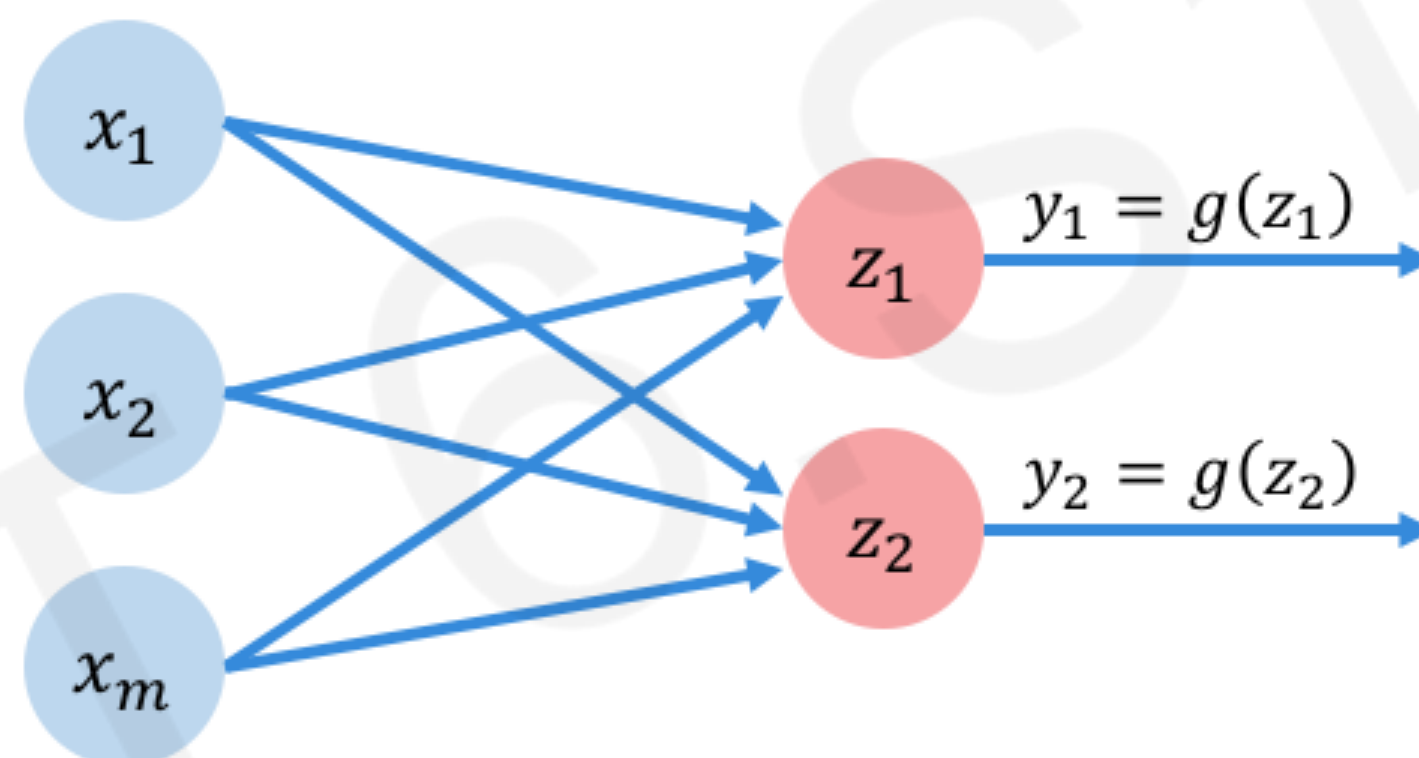
$$z_i = w_{0,i} + \sum_{j=1}^{m} x_j \, w_{j,i}$$

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com     @MITDeepLearning

1/27/20

# Single Layer Neural Network



$W^{(1)}$

$W^{(2)}$

$g(z_1)$

$g(z_2)$

$g(z_3)$

$g(z_{d_1})$

Inputs       Hidden       Final Output

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

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/27/20

There are important contexts — like modeling with Keras — where we would refer to this as a "two-layer" neural network
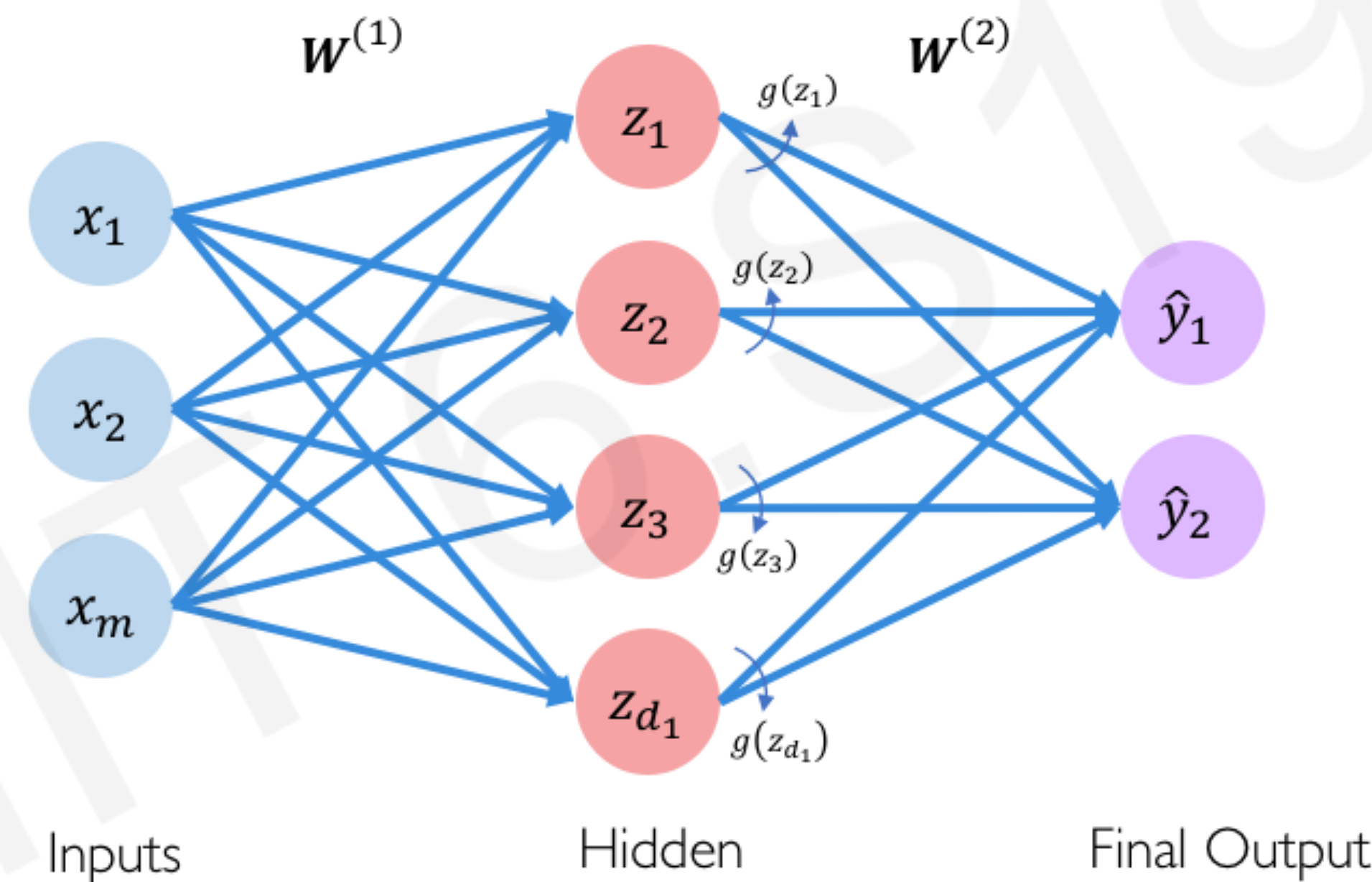
# Single Layer Neural Network



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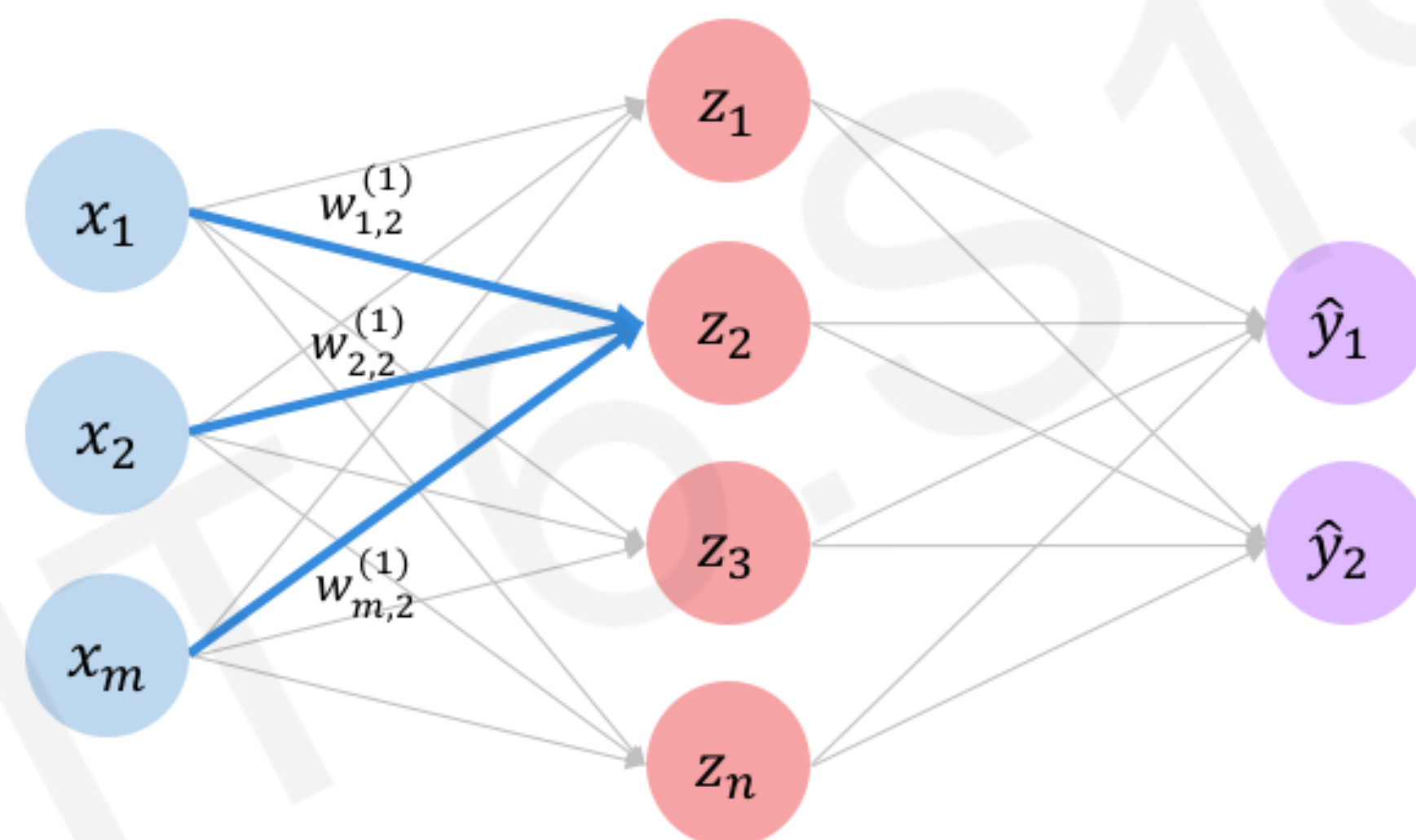
# Single Layer Neural Network



$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^{m} x_j \, w_{j,2}^{(1)}$$

$$= w_{0,2}^{(1)} + x_1 \, w_{1,2}^{(1)} + x_2 \, w_{2,2}^{(1)} + x_m \, w_{m,2}^{(1)}$$

# Multi Output Perceptron



```
import tensorflow as tf

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

$x_1$

$x_2$

$x_m$

$z_1$

$z_2$

$z_3$

$z_n$

$\hat{y}_1$

$\hat{y}_2$

Inputs

Hidden

Output

# Deep Neural Network



Inputs    Hidden    Output

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j})\, w_{j,i}^{(k)}$$

# Deep Neural Network



Inputs            Hidden            Output

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) \, w_{j,i}^{(k)}$$

```python
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(n_1),
    tf.keras.layers.Dense(n_2),
    ⋮
    tf.keras.layers.Dense(2)
])
```

# Deep Neural Network

$x_1$

$x_2$

$x_m$

Input

```python
model = models.Sequential()
model.add(layers.Dense(16, activation = 'relu', input_shape=(5000,)))
model.add(layers.Dense(16, activation = 'relu'))
model.add(layers.Dense(1, activation= 'sigmoid'))

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])

history = model.fit(partial_x_train,
                    partial_y_train,
                    epochs=4,
                    batch_size=512,
                    validation_data=(x_val,y_val))
```

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) \, w_{j,i}^{(k)}$$

# Deep Neural Network

```{r}
model <- keras_model_sequential() %>%
  layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
  layer_dense(units = 16, activation = "relu") %>%
  layer_dense(units = 1, activation = "sigmoid")

model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = c("accuracy")
)

model %>% fit(x_train, y_train, epochs = 4, batch_size = 512)
results <- model %>% evaluate(x_test, y_test)
```

```python
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(n_1),
    tf.keras.layers.Dense(n_2),
    ⋮
    tf.keras.layers.Dense(2)
])
```

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) \, w_{j,i}^{(k)}$$
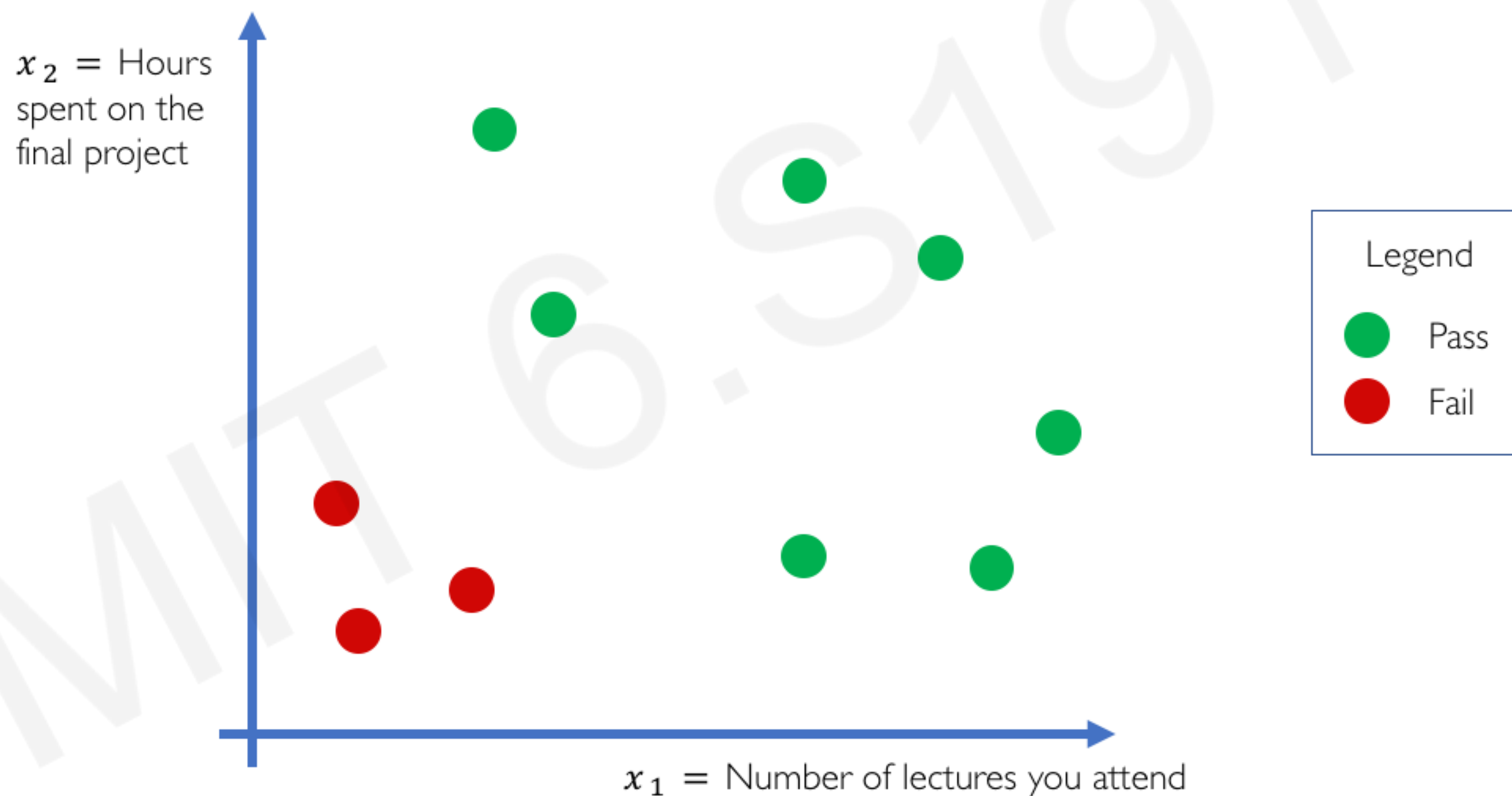
# Example Problem

## 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?



$x_2$ = Hours spent on the final project

$x_1$ = Number of lectures you attend

Legend

● Pass
● Fail

# Example Problem: Will I pass this class?



$x_2 = $ Hours spent on the final project

? $\begin{bmatrix} 4 \\ 5 \end{bmatrix}$

**Legend**

● Pass
● Fail

$x_1 = $ Number of lectures you attend

Massachusetts Institute of Technology

# Example Problem: Will I pass this class?



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

Predicted: 0.1

# Example Problem: Will I pass this class?



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

Predicted: **0.1**
Actual: **1**

Massachusetts
Institute of
Technology

# Quantifying Loss

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



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

$\underline{\text{Predicted}}$      $\underline{\text{Actual}}$

Massachusetts
Institute of
Technology

# Empirical Loss

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



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

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

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

$f(x)$   $y$

$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix} \begin{matrix} \times \\ \times \\ \checkmark \\ \end{matrix} \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

Predicted    Actual

# Binary Cross Entropy Loss

```python
model = models.Sequential()
model.add(layers.Dense(16, activation = 'relu', input_shape=(5000,)))
model.add(layers.Dense(16, activation = 'relu'))
model.add(layers.Dense(1, activation= 'sigmoid'))

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])


history = model.fit(partial_x_train,
                    partial_y_train,
                    epochs=4,
                    batch_size=512,
                    validation_data=(x_val,y_val))
```
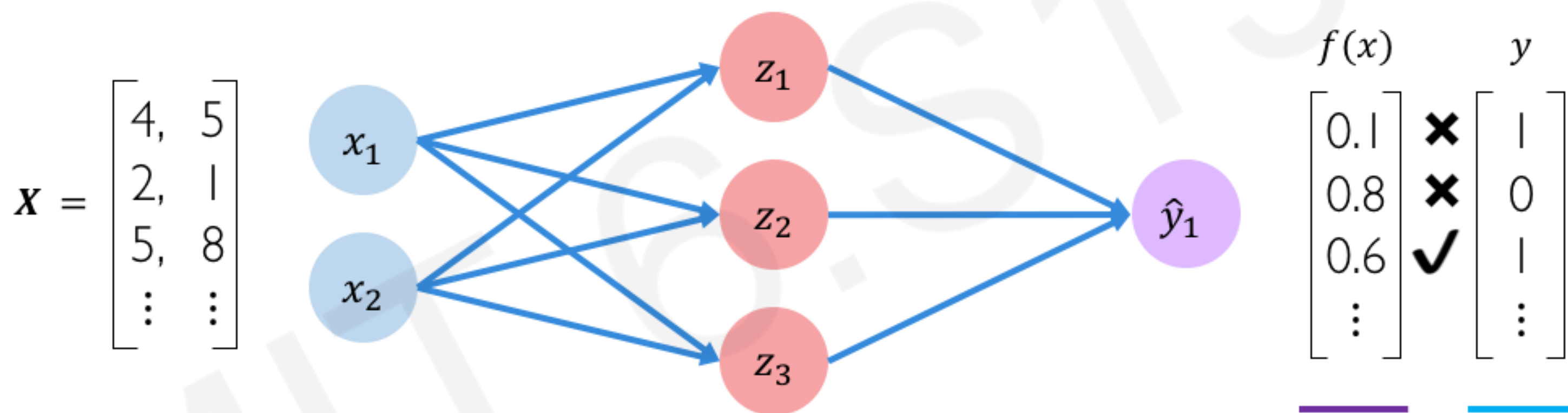
| Actual | Predicted | Actual | Predicted |

```python
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```

# Binary Cross Entropy Loss

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



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

Actual    Predicted    Actual    Predicted

```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```

# Binary Cross Entropy Loss

```{r}
model <- keras_model_sequential() %>%
  layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
  layer_dense(units = 16, activation = "relu") %>%
  layer_dense(units = 1, activation = "sigmoid")

model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = c("accuracy")
)


model %>% fit(x_train, y_train, epochs = 4, batch_size = 512)
results <- model %>% evaluate(x_test, y_test)
```

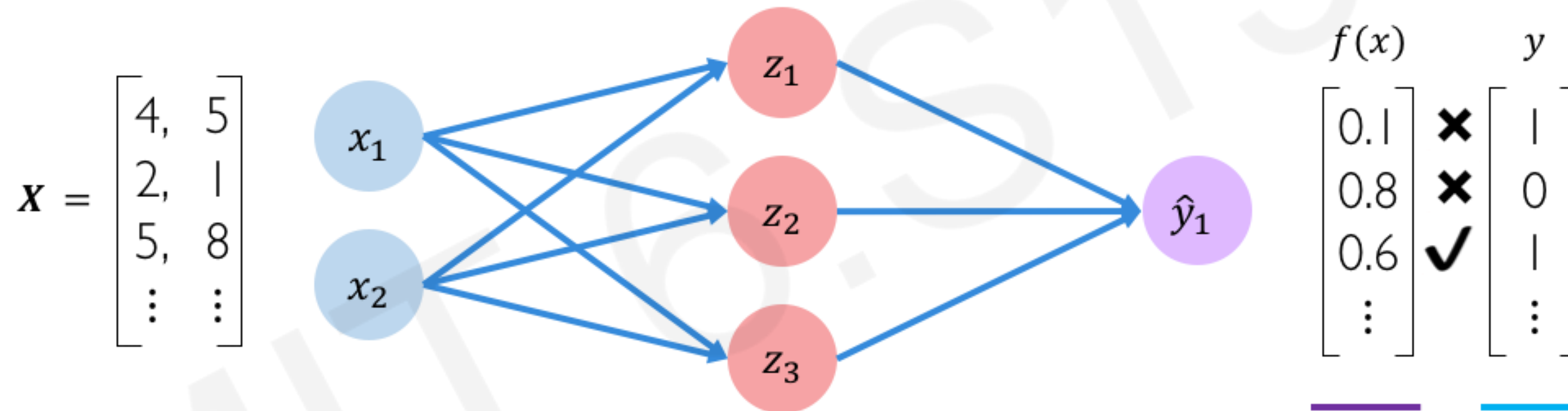Actual     Predicted     Actual     Predicted

```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```

This is exactly equivalent to the negative log-likelihood.
So this is so far identical to logit/logistic regression.

# Binary Cross Entropy Loss

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



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

Actual    Predicted    Actual    Predicted

```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```
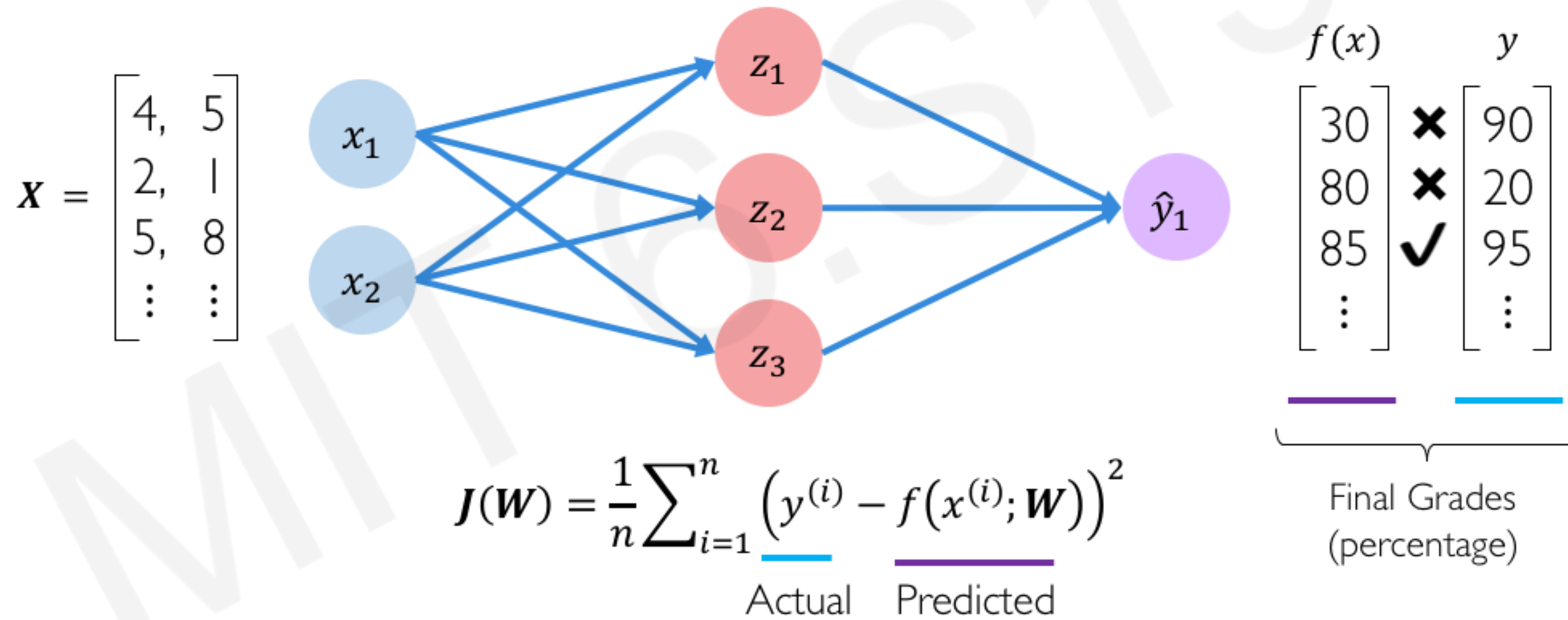
# Machine learning in one slide

| Social science (inference) | Machine learning (prediction) |
|:---:|:---:|
| GLM inverse link function | Activation function |
| $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ | $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ |
| Preferred objective function | |
| Log-likelihood | Cross-entropy |
| $\log \mathcal{L} = \sum_{i=1}^{n} \log P(y_i|\mathbf{x}_i, \boldsymbol{\beta})$ | $-\log \mathcal{L} = -\sum_{i=1}^{n} \log P(y_i|\mathbf{x}_i, \boldsymbol{\beta})$ |
| Solving algorithm | |
| Newton-Raphson | Gradient descent |
| $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - [\boldsymbol{H} \log \mathcal{L}]^{-1} \nabla \log \mathcal{L}$ | $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - \eta \nabla(-\log \mathcal{L})$ |
| Quantities of interest | |
| $\widehat{\boldsymbol{\beta}}; Var(\widehat{\boldsymbol{\beta}})$ | $\widehat{y}; \sum \mathbf{1}(\widehat{y} = y)/n$ |

# Mean Squared Error Loss

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

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

$x_1$  $x_2$  $z_1$  $z_2$  $z_3$  $\hat{y}_1$

$f(x)$   $y$

$$\begin{bmatrix} 30 \\ 80 \\ 85 \\ \vdots \end{bmatrix} \begin{matrix} \times \\ \times \\ \checkmark \end{matrix} \begin{bmatrix} 90 \\ 20 \\ 95 \\ \vdots \end{bmatrix}$$

Final Grades
(percentage)

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

Actual    Predicted

```
loss = tf.reduce_mean( tf.square(tf.subtract(y, predicted)) )
```

# Training Neural Networks

# Loss Optimization

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

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

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

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

Massachusetts
Institute of
Technology

1/27/20

# Loss Optimization

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

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

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

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

# Loss Optimization

$$W^* = \underset{W}{\arg\min}\, J(W)$$

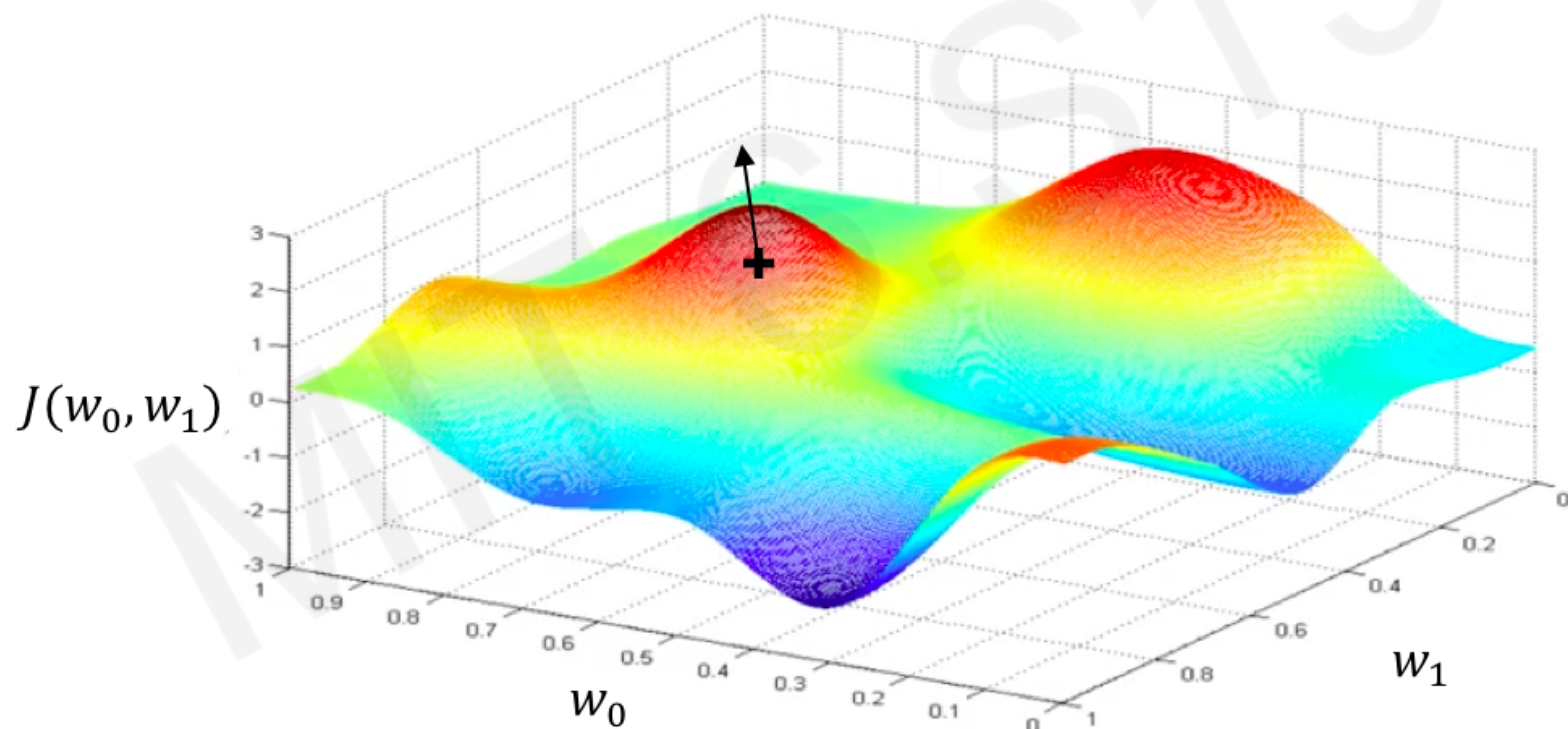Remember:
*Our loss is a function of the network weights!*



$J(w_0, w_1)$

$w_0$

$w_1$

# Loss Optimization

Randomly pick an initial $(w_0, w_1)$

# Loss Optimization

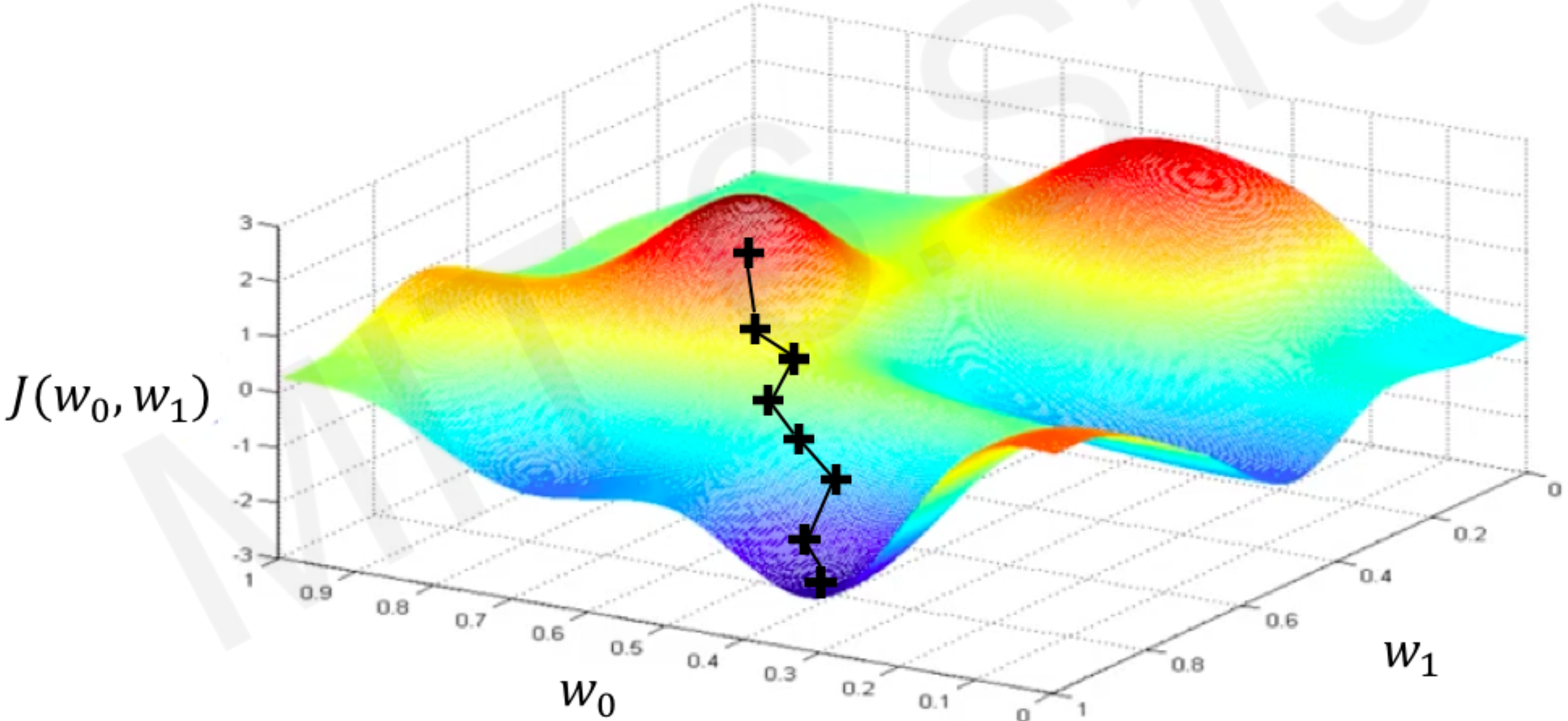Compute gradient, $\dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$



$J(w_0, w_1)$

$w_0$

$w_1$

# Loss Optimization

Take small step in opposite direction of gradient



$J(w_0, w_1)$

$w_0$

$w_1$

# Gradient Descent

Repeat until convergence



$J(w_0, w_1)$

$w_0$

$w_1$

Contrast with Newton-Raphson, which also uses *second derivative* (Hessian)

# Gradient Descent
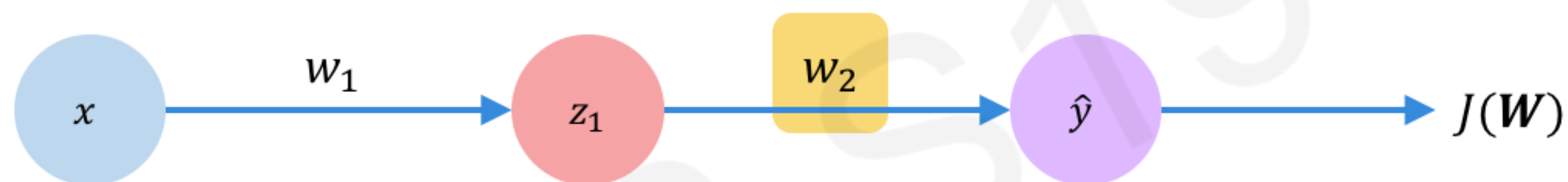
Repeat until convergence

# Machine learning in one slide

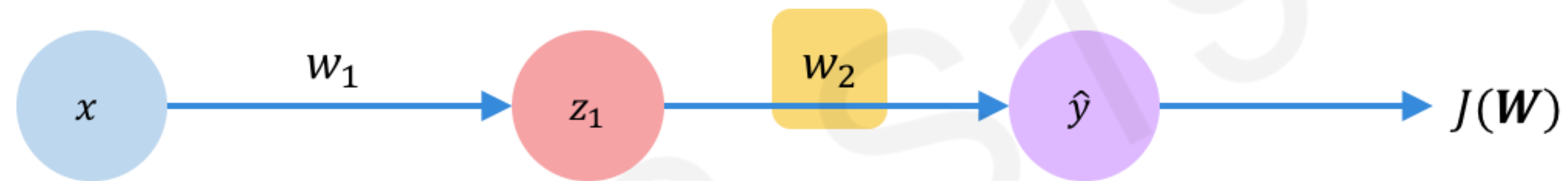| Social science (inference) | Machine learning (prediction) |
| --- | --- |
| GLM inverse link function | Activation function |
| $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ | $\mathbb{E}(y) = f(\mathbf{x}'\boldsymbol{\beta})$ |
| Preferred objective function | |
| Log-likelihood | Cross-entropy |
| $\log \mathcal{L} = \sum_{i=1}^{n} \log P(y_i\|\mathbf{x}_i, \boldsymbol{\beta})$ | $-\log \mathcal{L} = -\sum_{i=1}^{n} \log P(y_i\|\mathbf{x}_i, \boldsymbol{\beta})$ |
| Solving algorithm | |
| Newton-Raphson | Gradient descent |
| $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - [\boldsymbol{H}\log\mathcal{L}]^{-1}\nabla\log\mathcal{L}$ | $\boldsymbol{\beta}_t := \boldsymbol{\beta}_{t-1} - \eta\nabla(-\log\mathcal{L})$ |
| Quantities of interest | |
| $\widehat{\boldsymbol{\beta}}; Var(\widehat{\boldsymbol{\beta}})$ | $\widehat{y}; \sum \mathbf{1}(\widehat{y} = y)/n$ |

# Gradient Descent

**Algorithm**

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

2. Loop until convergence:

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

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

5. Return weights

# Computing Gradients: Backpropagation



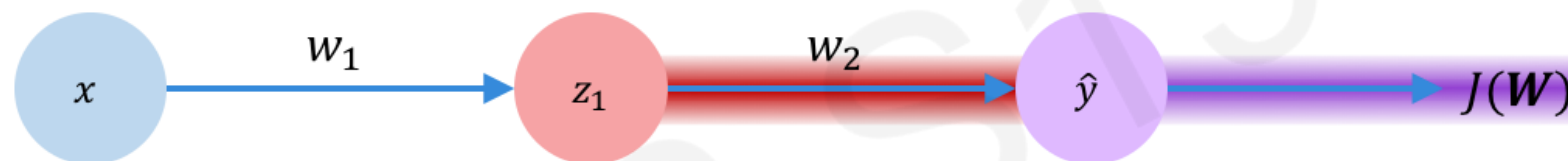How does a small change in one weight (ex. $w_2$) affect the final loss $J(W)$?

# Computing Gradients: Backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} =$$
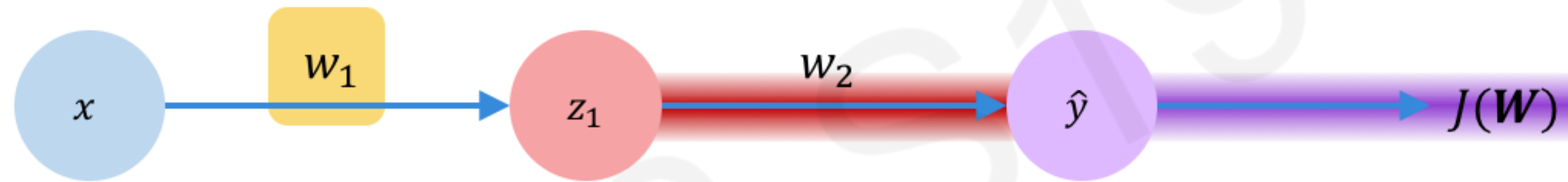
Let's use the chain rule!

Massachusetts
Institute of
Technology

# Computing Gradients: Backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$
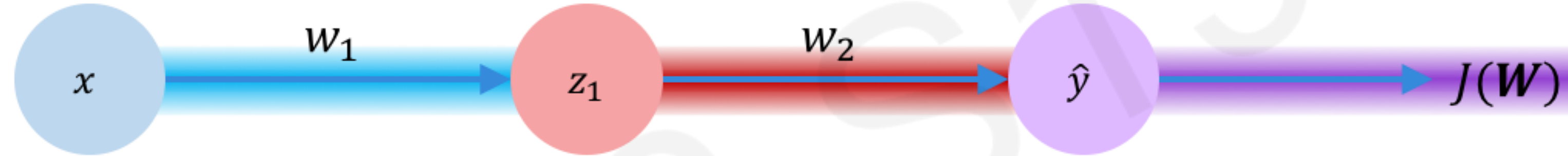
# Computing Gradients: Backpropagation



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

Apply chain rule!          Apply chain rule!

Massachusetts
Institute of
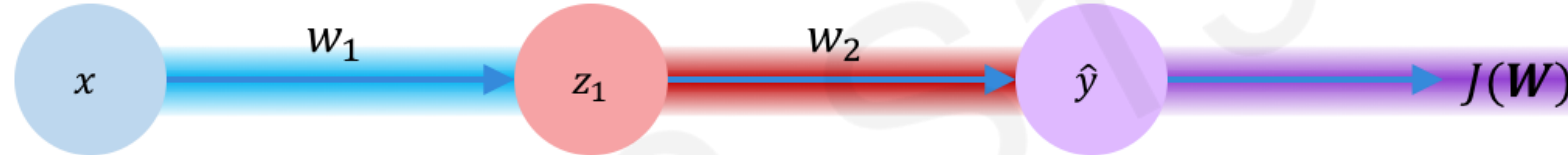Technology

# Computing Gradients: Backpropagation



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

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

1/27/20

Massachusetts
Institute of
Technology

# 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

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

# Loss Functions Can Be Difficult to Optimize

## Remember:

Optimization through gradient descent
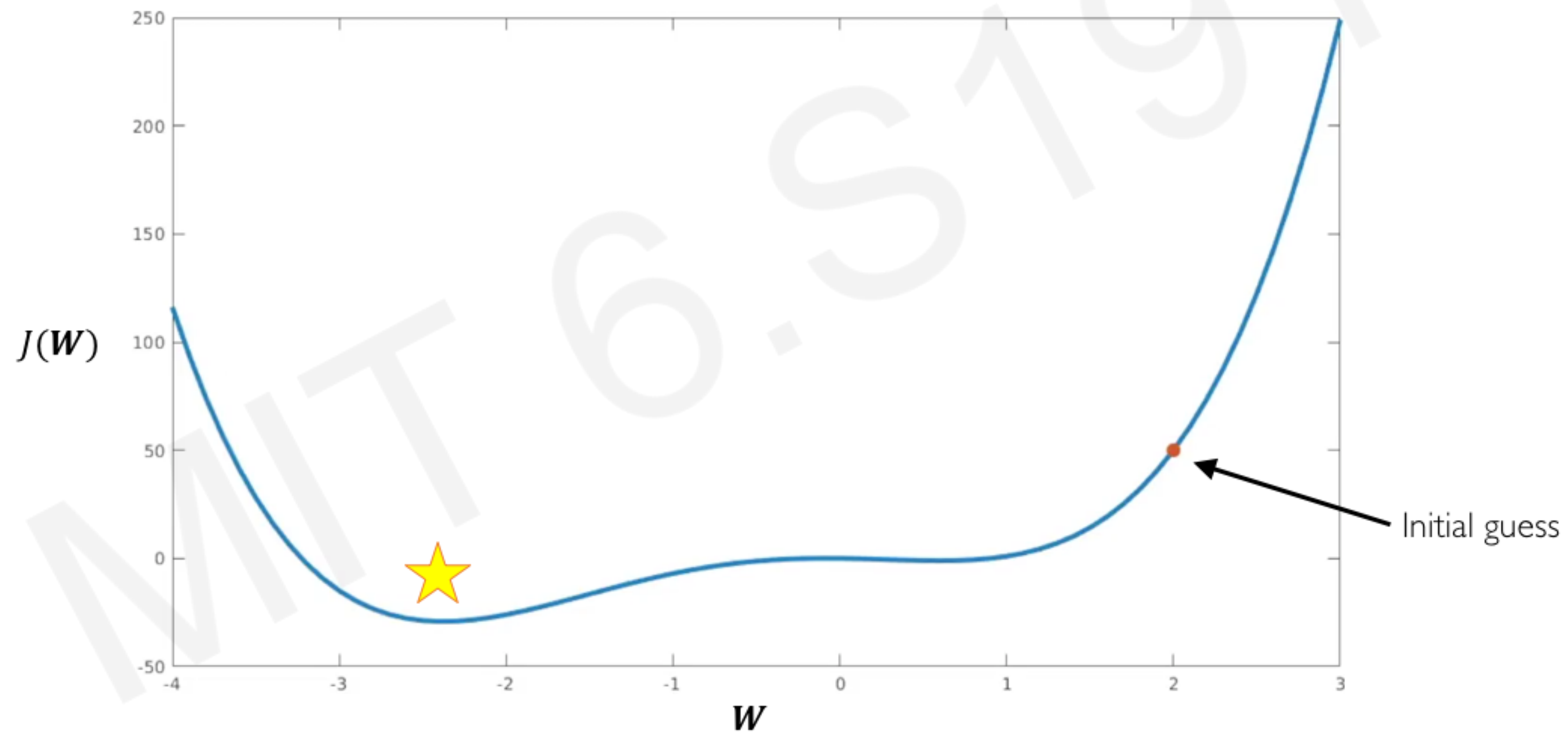
$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

How can we set the
learning rate?

Massachusetts
Institute of
Technology

6.S191 Introduction to Deep Learning
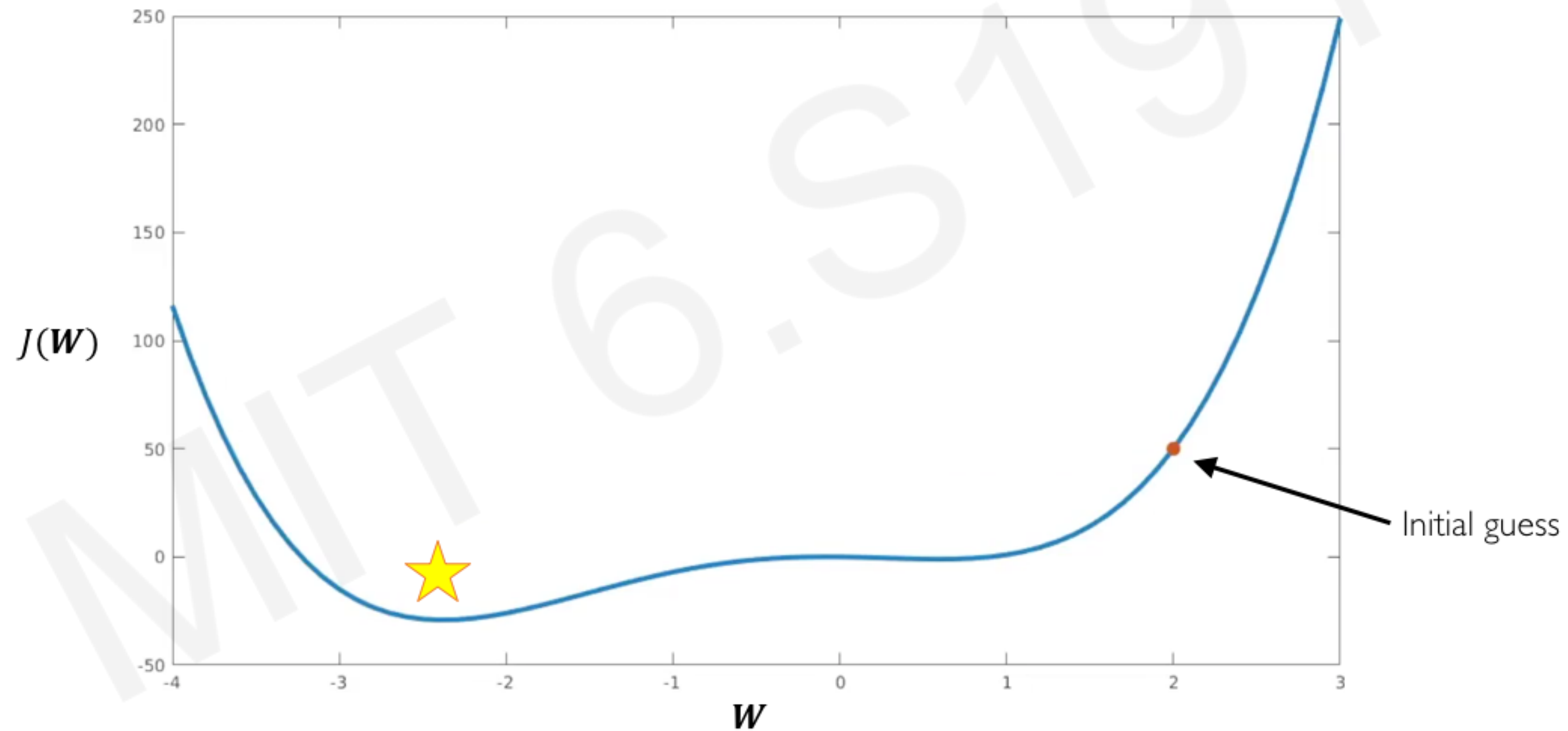introtodeeplearning.com  @MITDeepLearning

1/27/20

# Setting the Learning Rate

*Small learning rate* converges slowly and gets stuck in false local minima
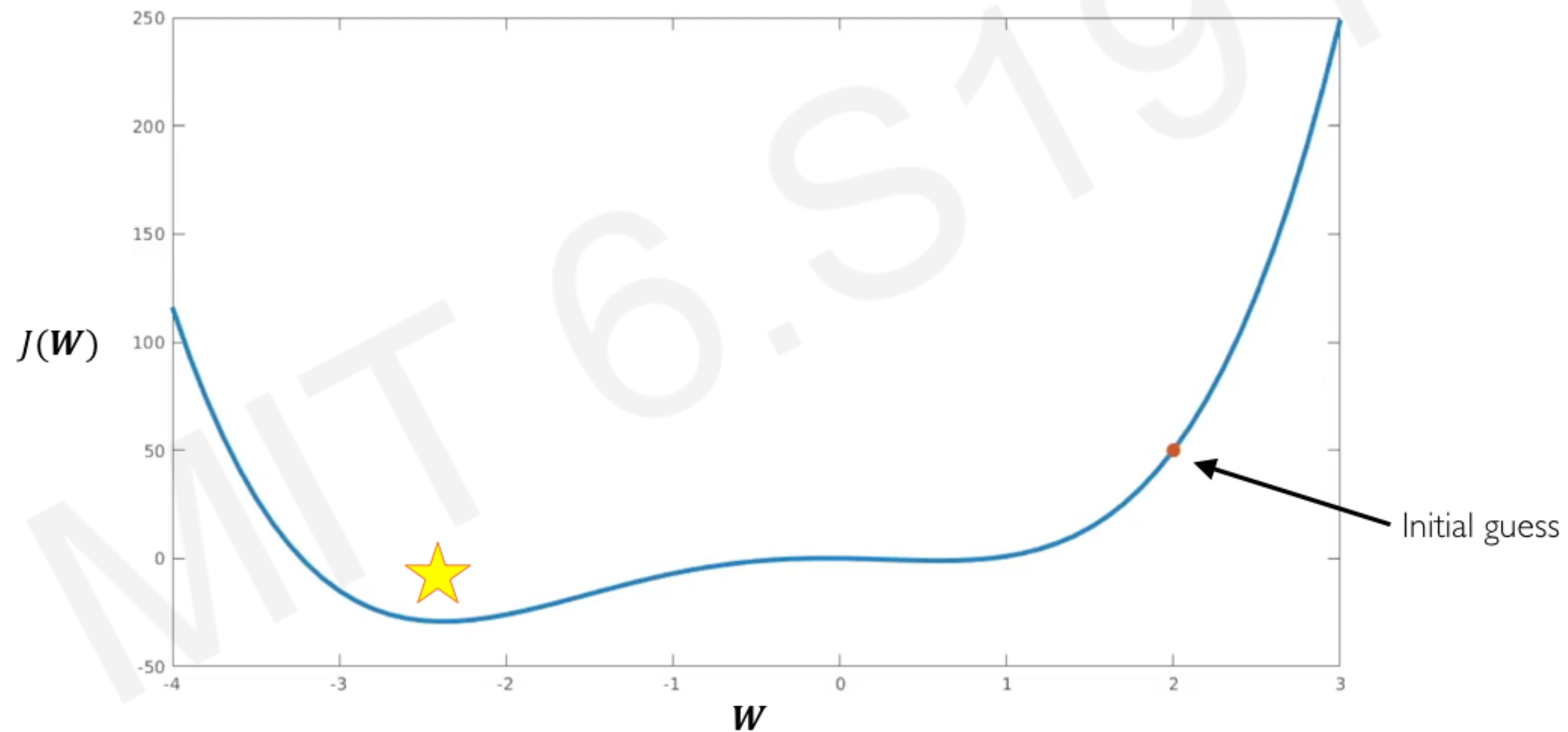
# Setting the Learning Rate

*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 | TF Implementation | Reference |
|-----------|-------------------|-----------|
| • SGD | `tf.keras.optimizers.SGD` | Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952. |
| • Adam | `tf.keras.optimizers.Adam` | Kingma et al. "Adam: A Method for Stochastic Optimization." 2014. |
| • Adadelta | `tf.keras.optimizers.Adadelta` | Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012. |
| • Adagrad | `tf.keras.optimizers.Adagrad` | Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011. |
| • RMSProp | `tf.keras.optimizers.RMSProp` | |

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

# Gradient Descent Algorithms

```python
model = models.Sequential()
model.add(layers.Dense(16, activation = 'relu', input_shape=(5000,)))
model.add(layers.Dense(16, activation = 'relu'))
model.add(layers.Dense(1, activation= 'sigmoid'))

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])

history = model.fit(partial_x_train,
                    partial_y_train,
                    epochs=4,
                    batch_size=512,
                    validation_data=(x_val,y_val))
```

Rate

Online

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

# Gradient Descent Algorithms

Alg

```{r}
model <- keras_model_sequential() %>%
  layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
  layer_dense(units = 16, activation = "relu") %>%
  layer_dense(units = 1, activation = "sigmoid")

model %>% compile(
  optimizer = "adam",
  loss = "binary_crossentropy",
  metrics = c("accuracy")
)

model %>% fit(x_train, y_train, epochs = 4, batch_size = 512)
results <- model %>% evaluate(x_test, y_test)
```
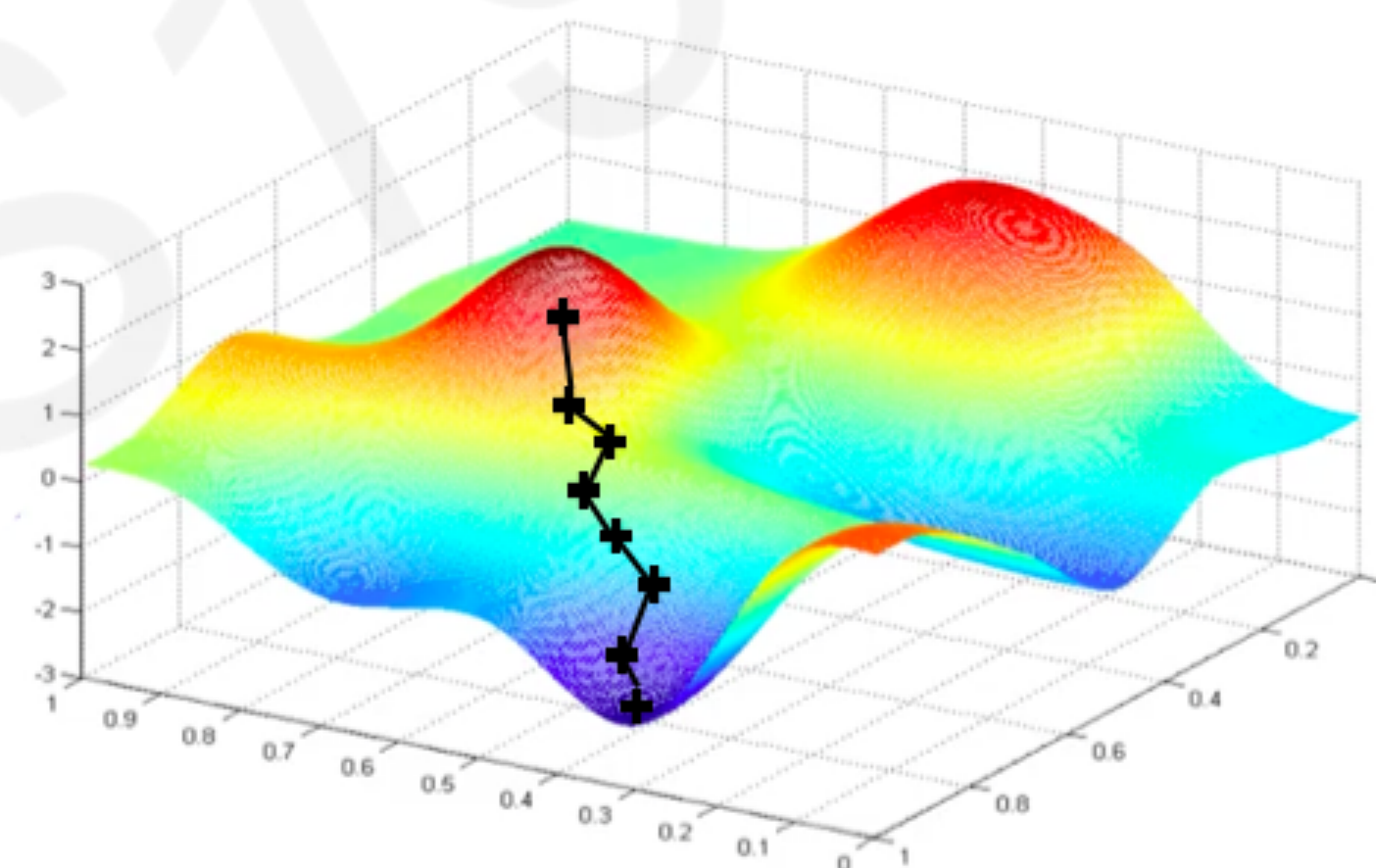
Additional details: http://ruder.io/optimizing-gradient-descent/

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, $\dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

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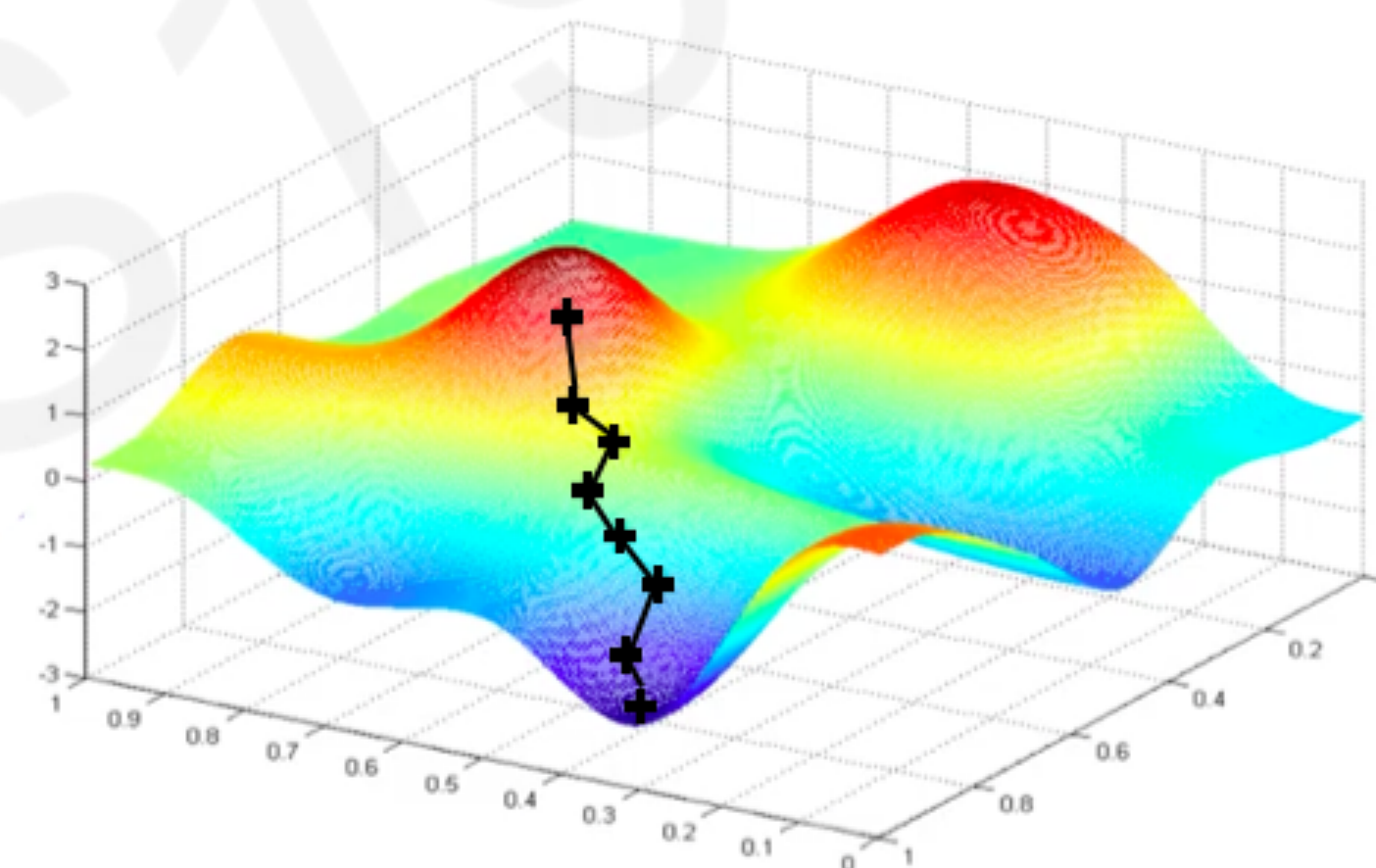
5. Return weights

# Gradient Descent

**Algorithm**

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

2. Loop until convergence:

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

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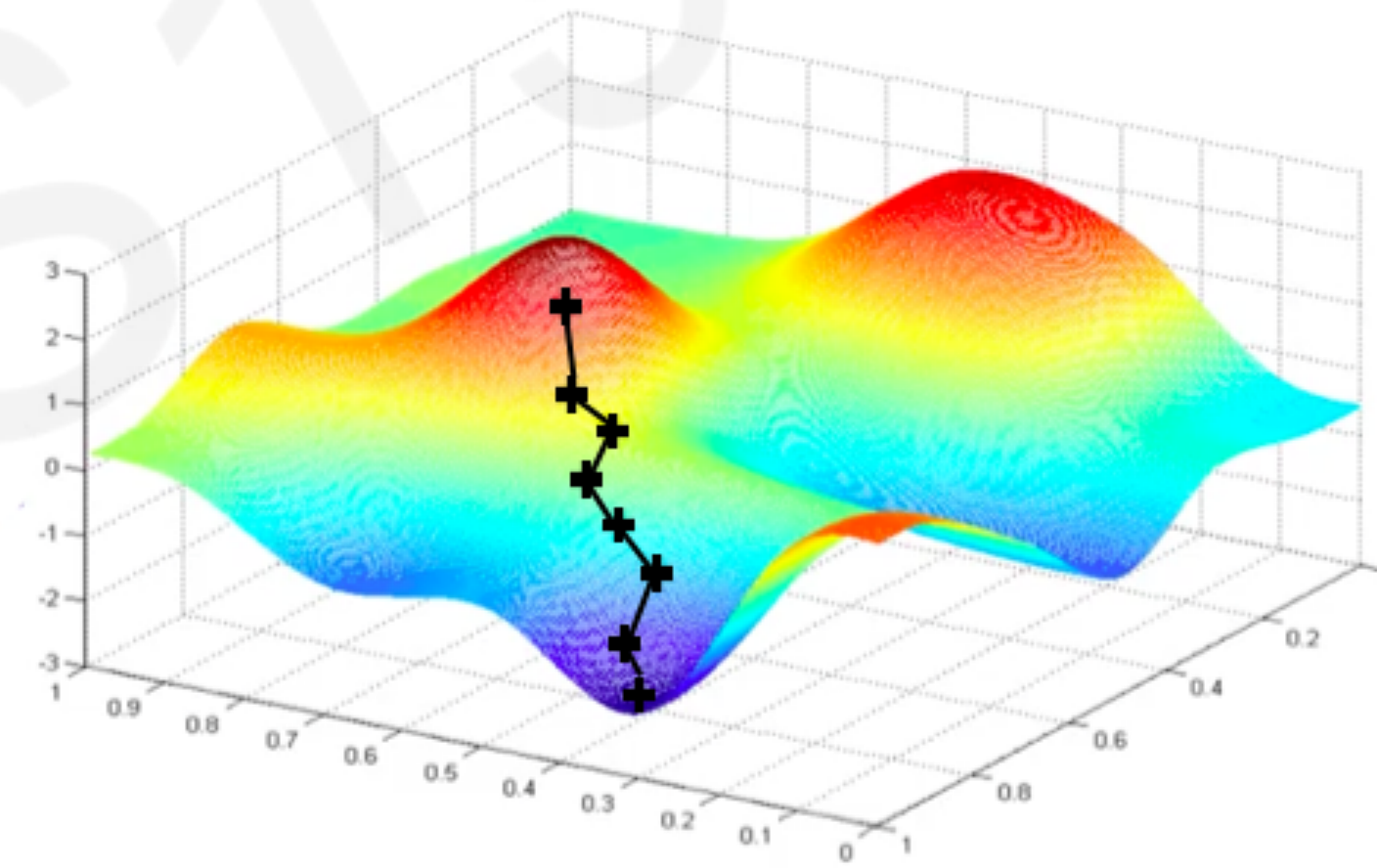
5. Return weights

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, $\dfrac{\partial J_i(W)}{\partial W}$

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

6. Return weights

# Stochastic Gradient Descent
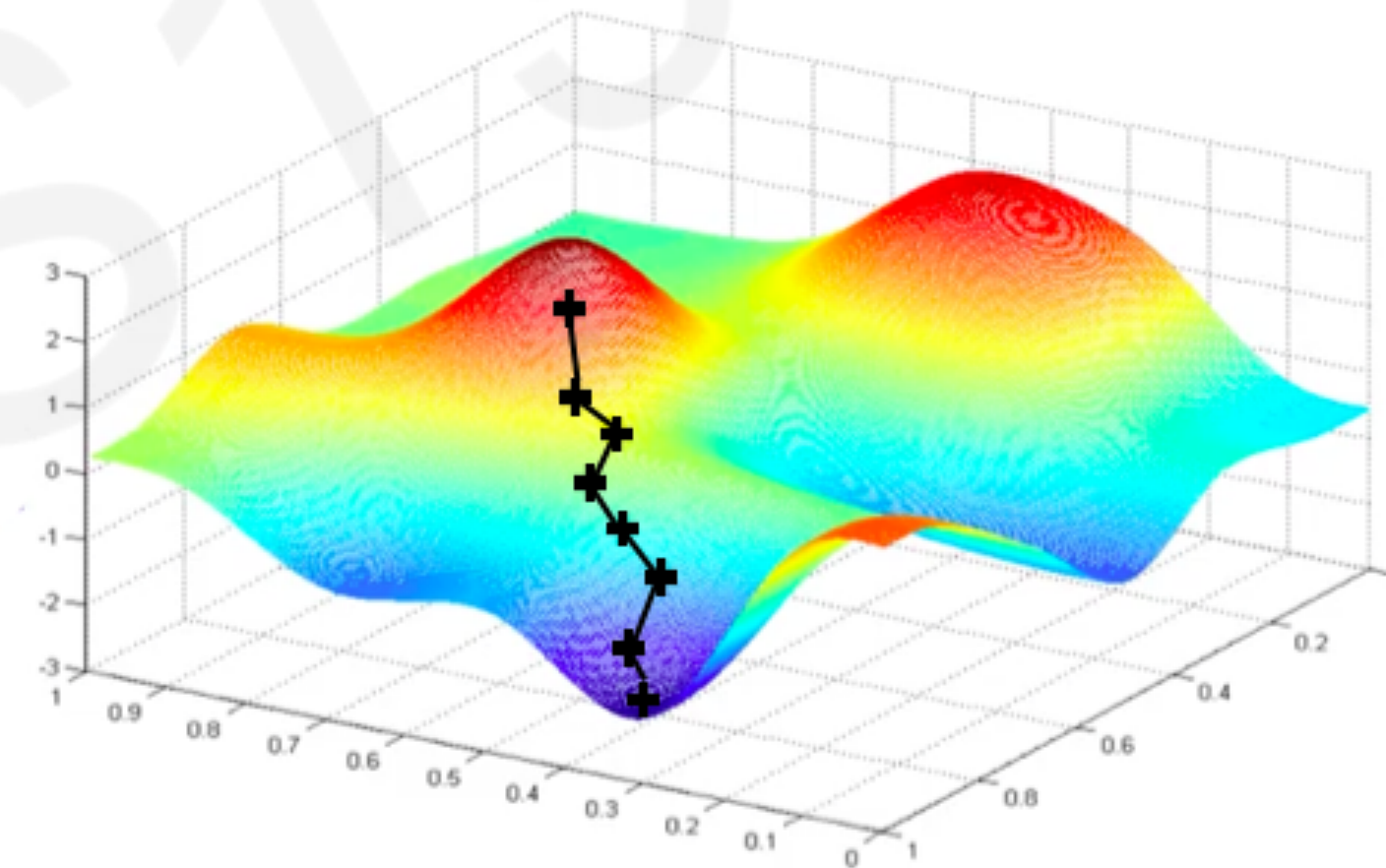
## Algorithm

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

2. Loop until convergence:

3.      Pick single data point $i$

4.      Compute gradient, $\dfrac{\partial J_i(W)}{\partial W}$

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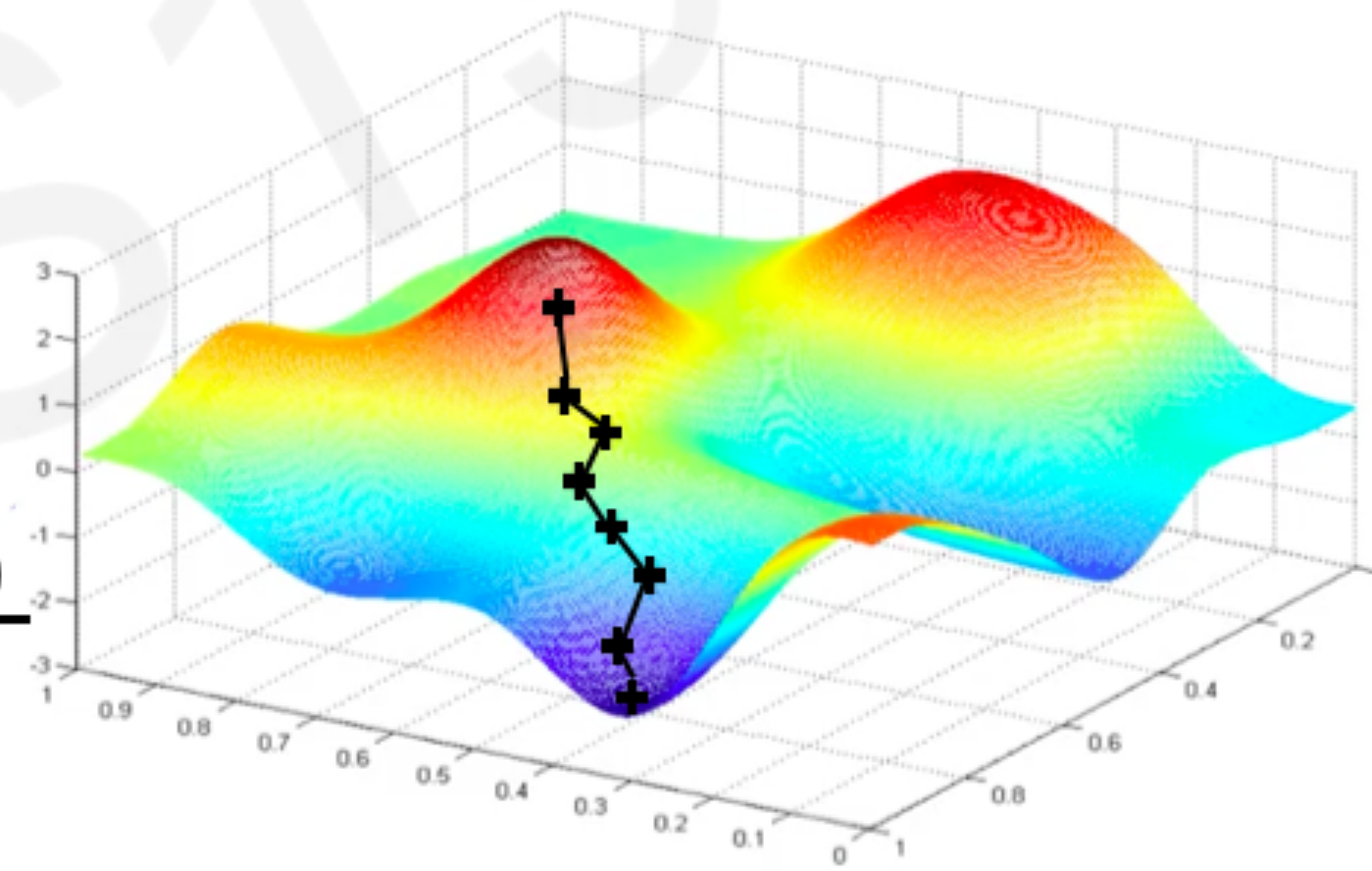
6. Return weights

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, $\frac{\partial J(W)}{\partial W} = \frac{1}{B} \sum_{k=1}^{B} \frac{\partial J_k(W)}{\partial W}$

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

6. Return weights

# Stochastic Gradient Descent
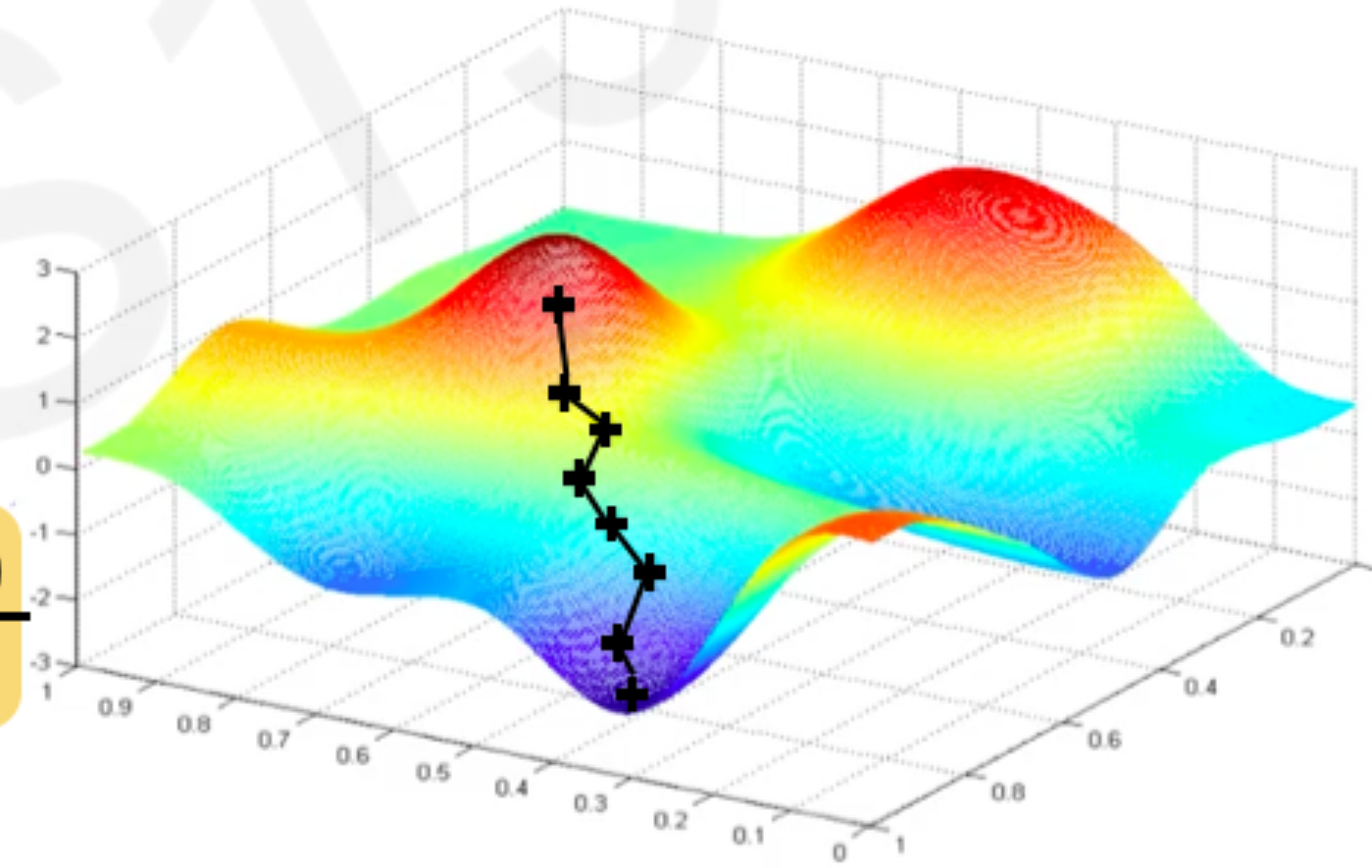
**Algorithm**

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

2. Loop until convergence:

3.      Pick batch of $B$ data points

4.      Compute gradient, $\dfrac{\partial J(W)}{\partial W} = \dfrac{1}{B} \sum_{k=1}^{B} \dfrac{\partial J_k(W)}{\partial W}$

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

6. Return weights

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

# Stochastic Gradient Descent

```python
model = models.Sequential()
model.add(layers.Dense(16, activation = 'relu', input_shape=(5000,)))
model.add(layers.Dense(16, activation = 'relu'))
model.add(layers.Dense(1, activation= 'sigmoid'))

model.compile(optimizer='adam',
              loss='binary_crossentropy',
              metrics=['accuracy'])

history = model.fit(partial_x_train,
                    partial_y_train,
                    epochs=4,
                    batch_size=512,
                    validation_data=(x_val,y_val))
```

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

# Stochastic Gradient Descent

**Algo**

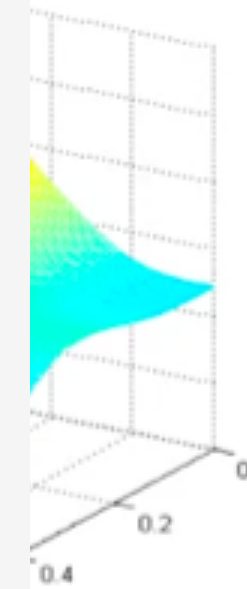1. In

2. L

3.

4.

5.

6. R

```{r}
model <- keras_model_sequential() %>%
    layer_dense(units = 16, activation = "relu", input_shape = c(5000)) %>%
    layer_dense(units = 16, activation = "relu") %>%
    layer_dense(units = 1, activation = "sigmoid")

model %>% compile(
    optimizer = "adam",
    loss = "binary_crossentropy",
    metrics = c("accuracy")
)

model %>% fit(x_train, y_train, epochs = 4, batch_size = 512)
results <- model %>% evaluate(x_test, y_test)
```

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

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

Massachusetts
Institute of
Technology

1/27/20

# 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

Massachusetts
Institute of
Technology

So, SGD is different from Newton-Raphson in derivative information used,
*and* in its optimization over small subsets of the data at a time and in parallel.

# Mini-batches while training

More accurate estimation of gradient
Smoother convergence
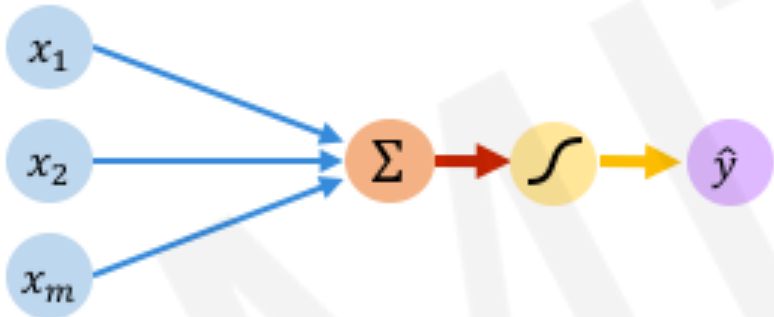Allows for larger learning rates

## Mini-batches lead to fast training!
Can parallelize computation + achieve significant speed increases on GPU's

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

Massachusetts
Institute of
Technology

1/27/20

https://playground.tensorflow.org

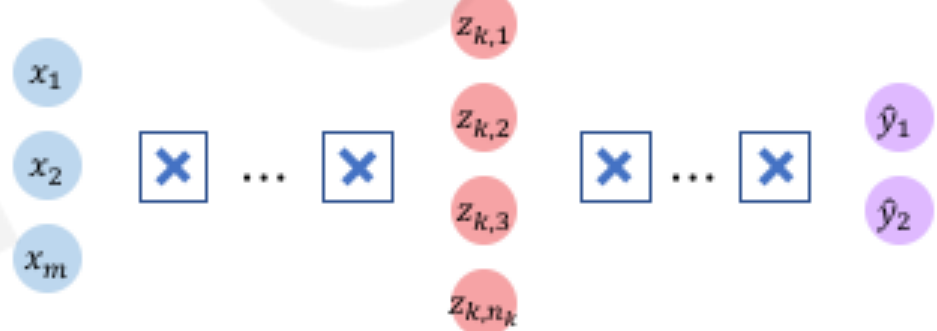# Core Foundation Review

## 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