# 提纲

# 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

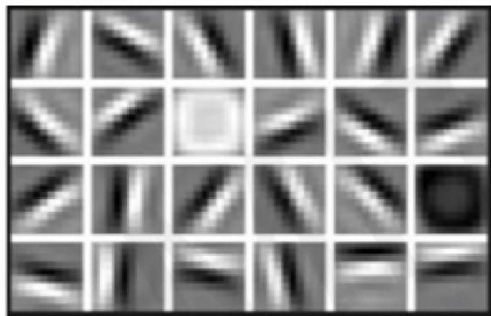Teaching computers how to **learn a task** directly from **raw data**

# 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
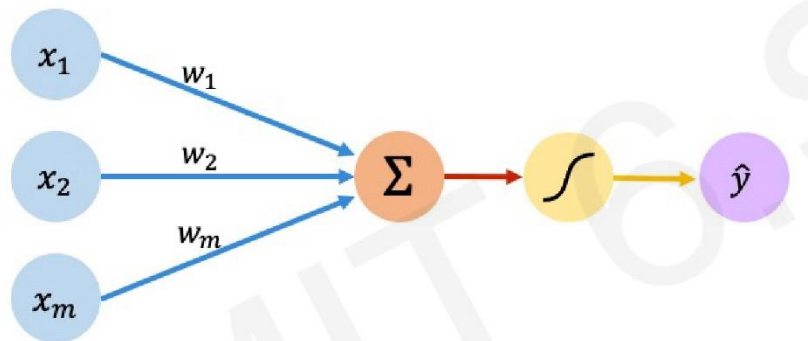
**High Level Features**

Facial Structure

The Perceptron
The structural building block of deep learning

# The Perceptron: Forward Propagation
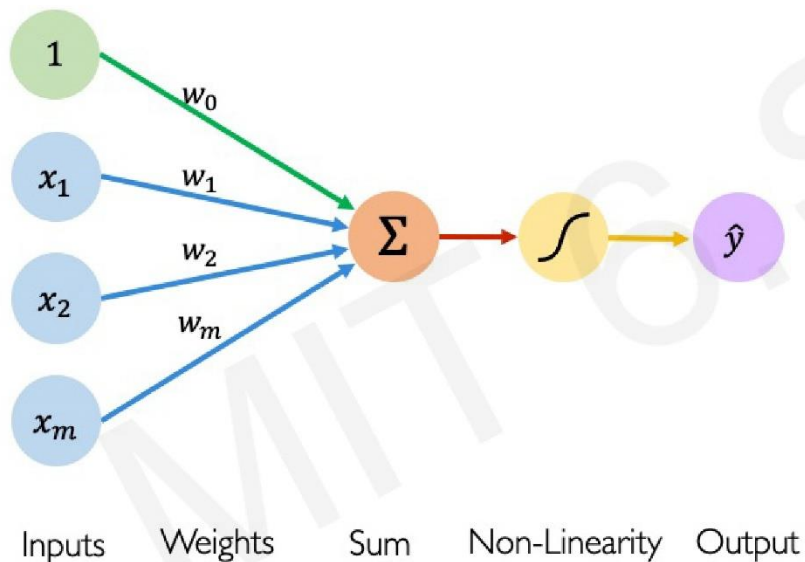


Inputs   Weights   Sum   Non-Linearity   Output

Output $\rightarrow$ 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

# The Perceptron: Forward Propagation



Inputs    Weights    Sum    Non-Linearity    Output

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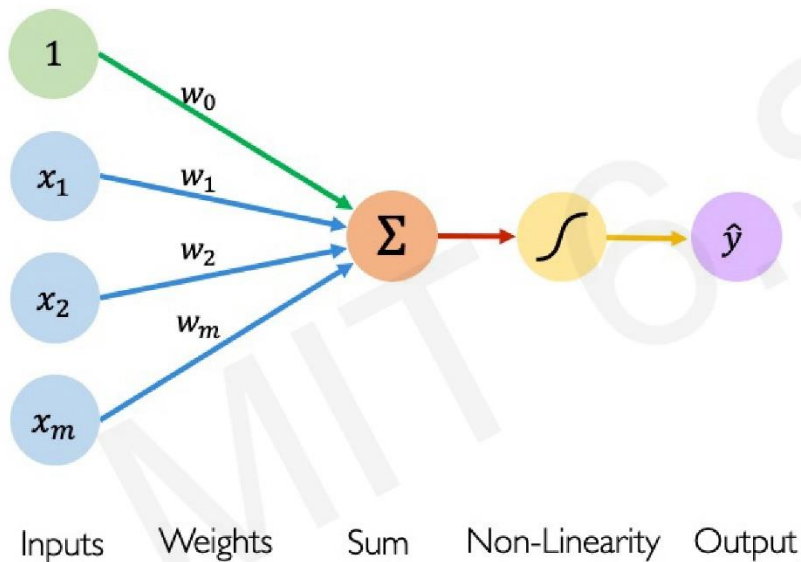
$$\hat{y} = g\left(w_0 + \boldsymbol{X}^T \boldsymbol{W}\right)$$

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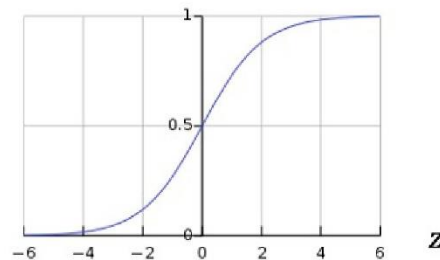
# The Perceptron: Forward Propagation



Inputs   Weights   Sum   Non-Linearity   Output

## Activation Functions

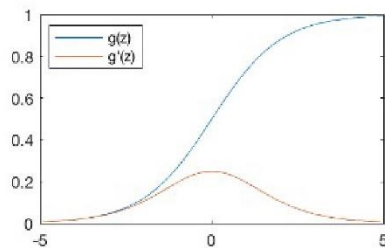$$\hat{y} = g\left( w_0 + X^T W \right)$$

- Example: sigmoid function

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

# Common Activation Functions
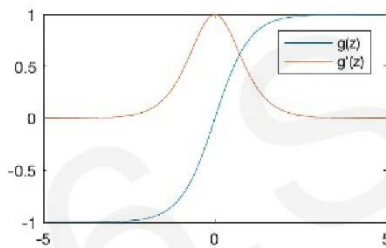
Sigmoid Function



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

$$g'(z) = g(z)(1 - g(z))$$
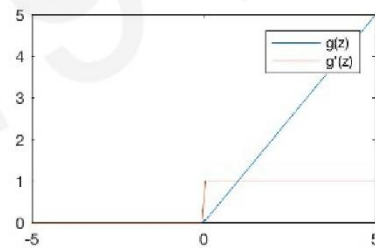
```
tf.math.sigmoid(z)
```

Hyperbolic Tangent



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

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

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

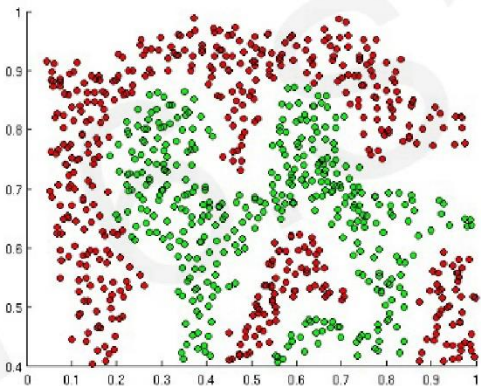Rectified Linear Unit (ReLU)



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

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

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

TensorFlow code blocks

NOTE: All activation functions are non-linear

# Importance of Activation Functions

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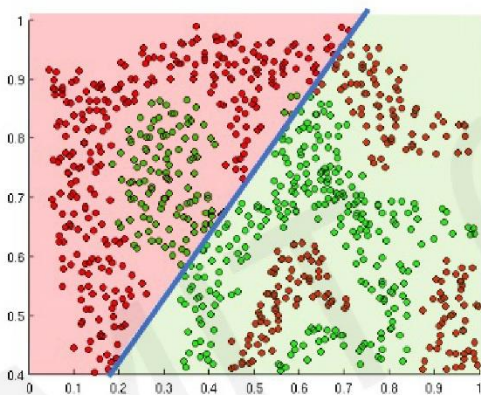


What if we wanted to build a neural network to distinguish green vs red points?

# Importance of Activation Functions

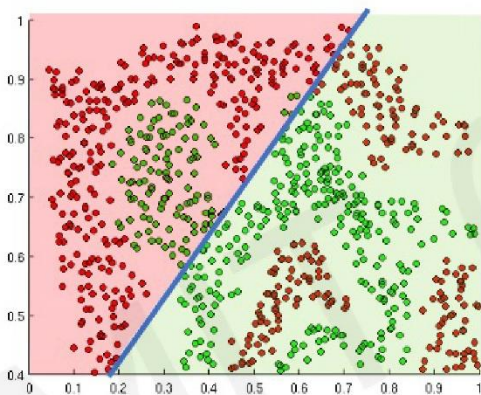*The purpose of activation functions is to **introduce non-linearities** into the network*



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
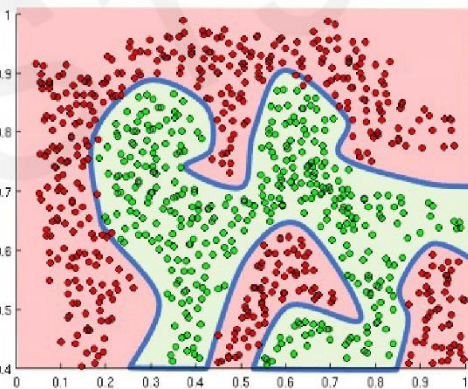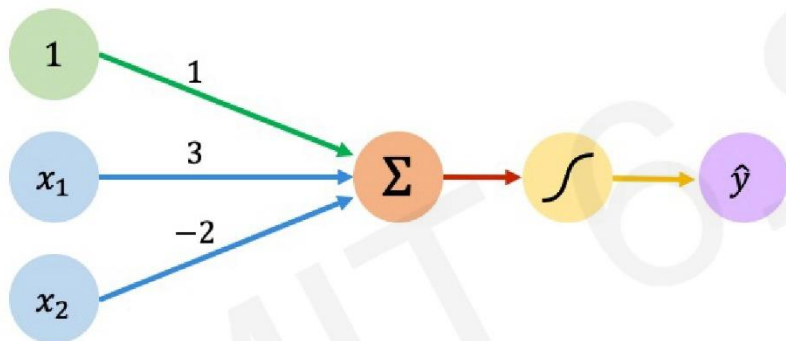
Non-linearities allow us to approximate arbitrarily complex functions
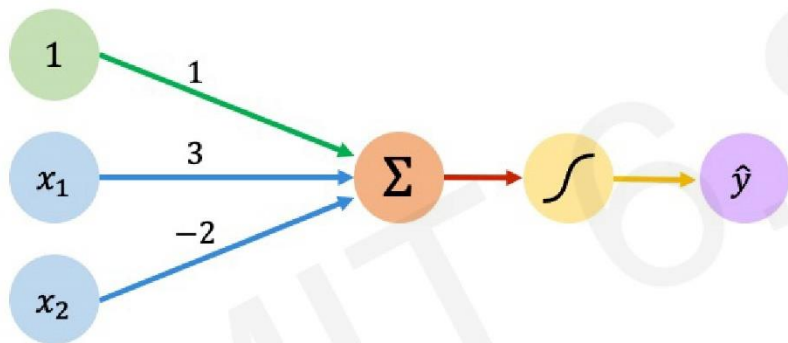
# The Perceptron: Example



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

$$\hat{y} = g(w_0 + \boldsymbol{X}^T \boldsymbol{W})$$
$$= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right)$$
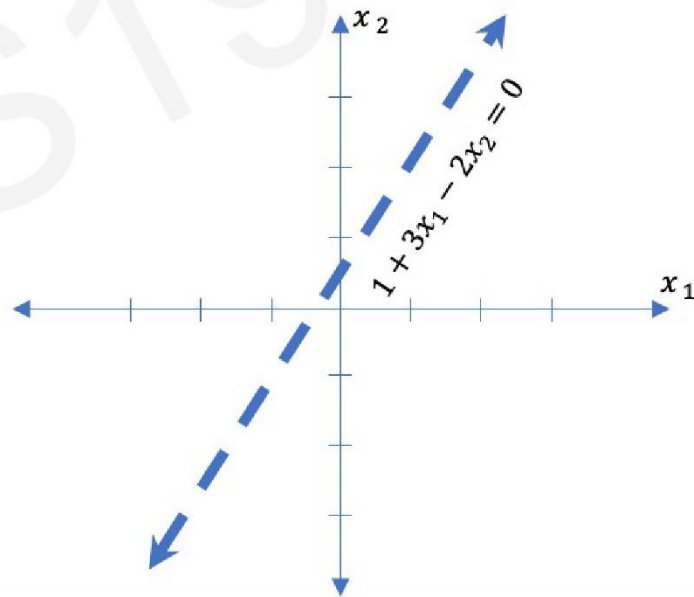$$\hat{y} = g(\underbrace{1 + 3x_1 - 2x_2})$$
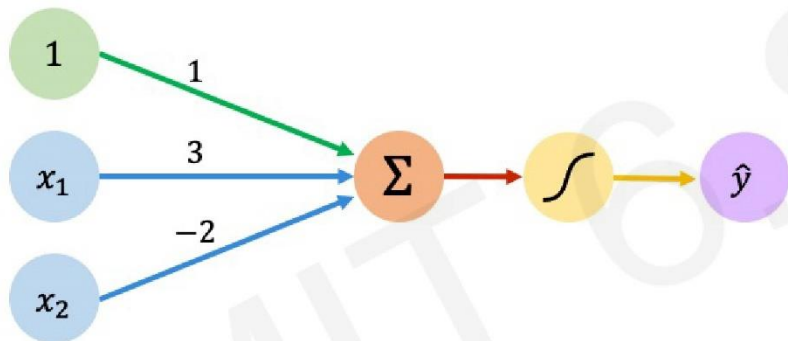
This is just a line in 2D!

# The Perceptron: Example



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

# The Perceptron: Example



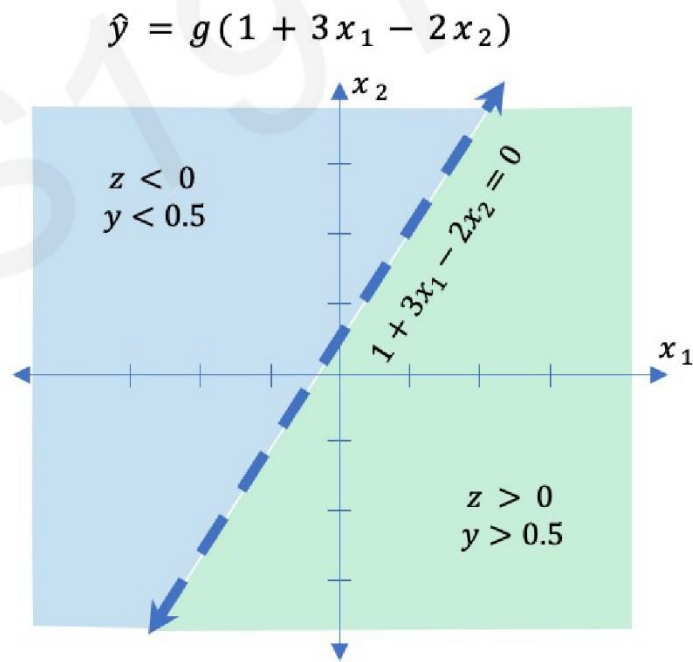$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

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

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

$1 + 3x_1 - 2x_2 = 0$

# The Perceptron: Example



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

# The Perceptron: Simplified
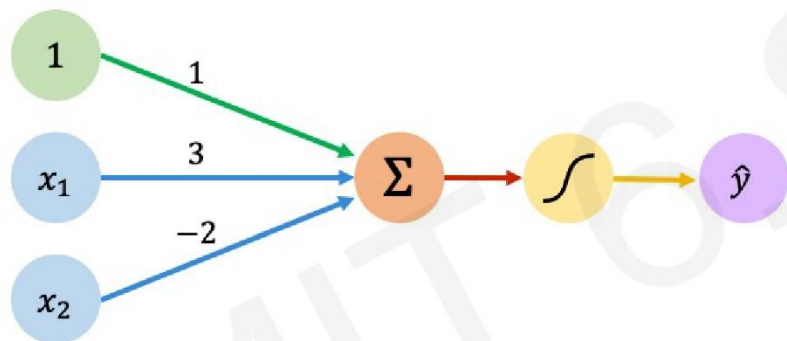
$$\hat{y} = g\left( w_0 + \boldsymbol{X}^T \boldsymbol{W} \right)$$

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



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

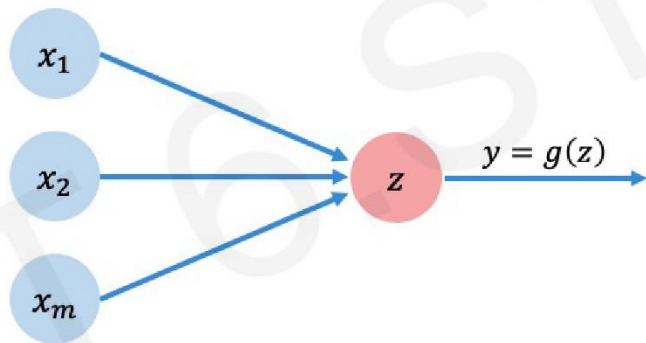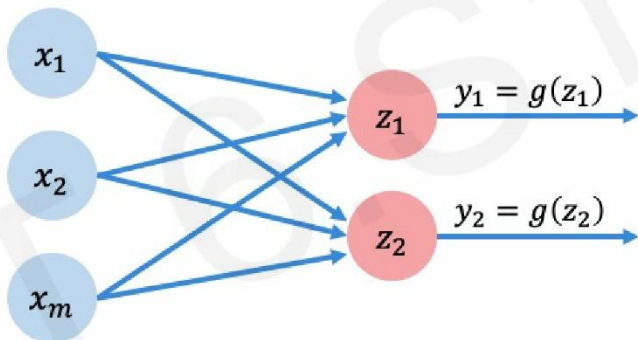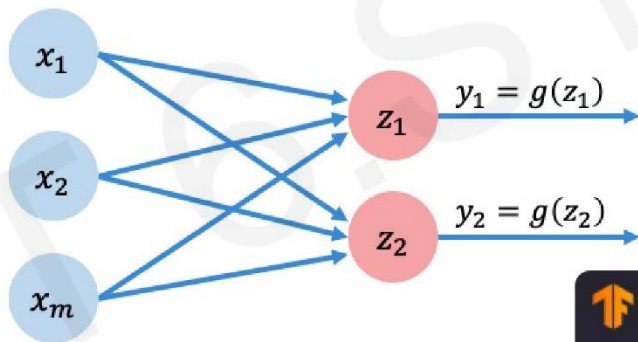# Multi Output Perceptron

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



```
import tensorflow as tf

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

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

# Single Layer Neural Network



$W^{(1)}$

$W^{(2)}$

$g(z_1)$

$g(z_2)$

$g(z_3)$

$g(z_{d_1})$

$x_1$

$x_2$

$x_m$

$z_1$

$z_2$

$z_3$

$z_{d_1}$

$\hat{y}_1$

$\hat{y}_2$

Inputs

Hidden

Final Output

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

# Single Layer Neural Network



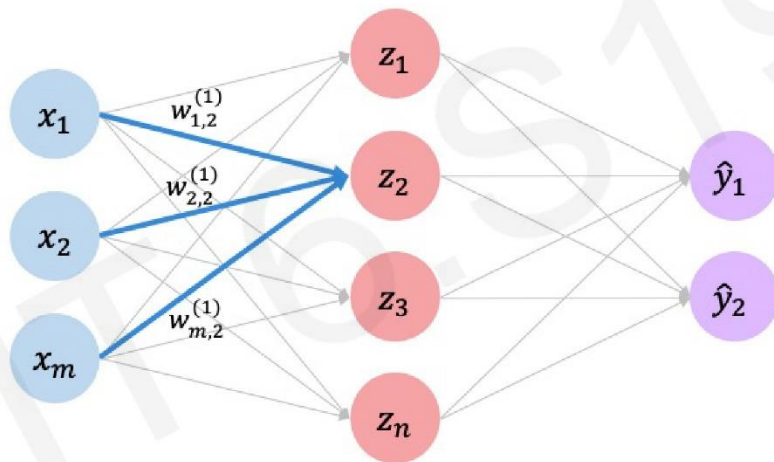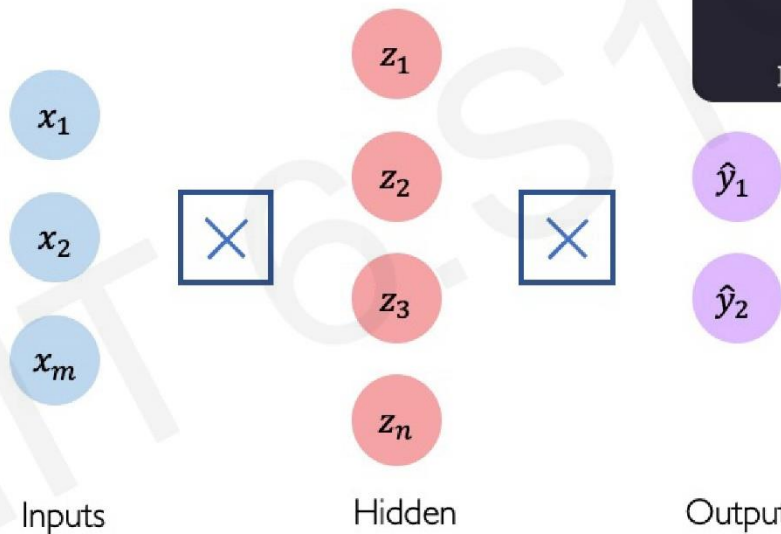$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^{m} x_j\, w_{j,2}^{(1)}$$

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

# Multi Output Perceptron
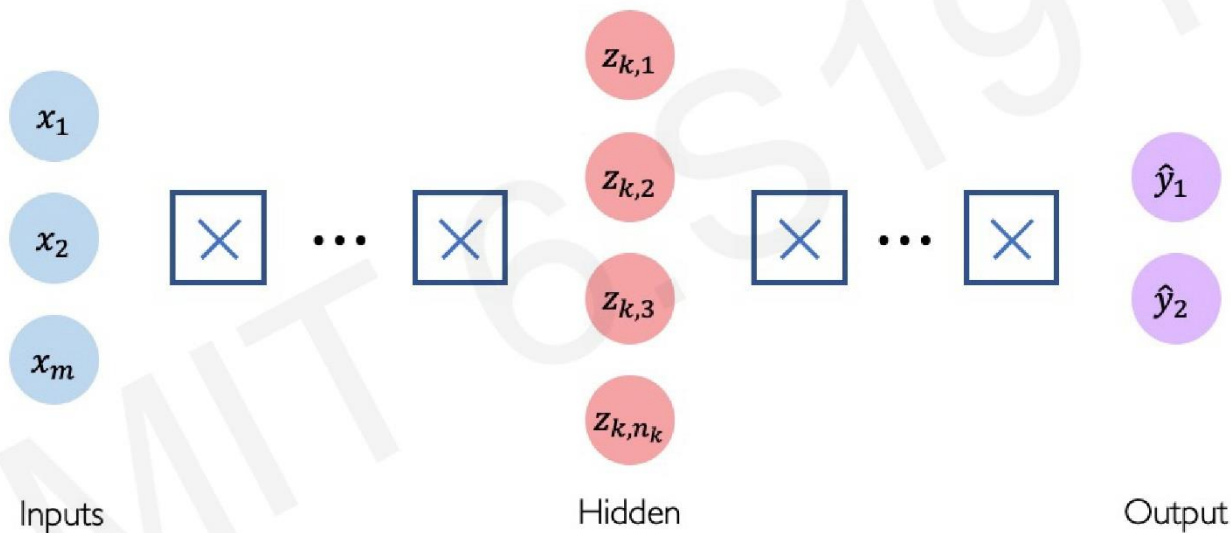
```
import tensorflow as tf

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



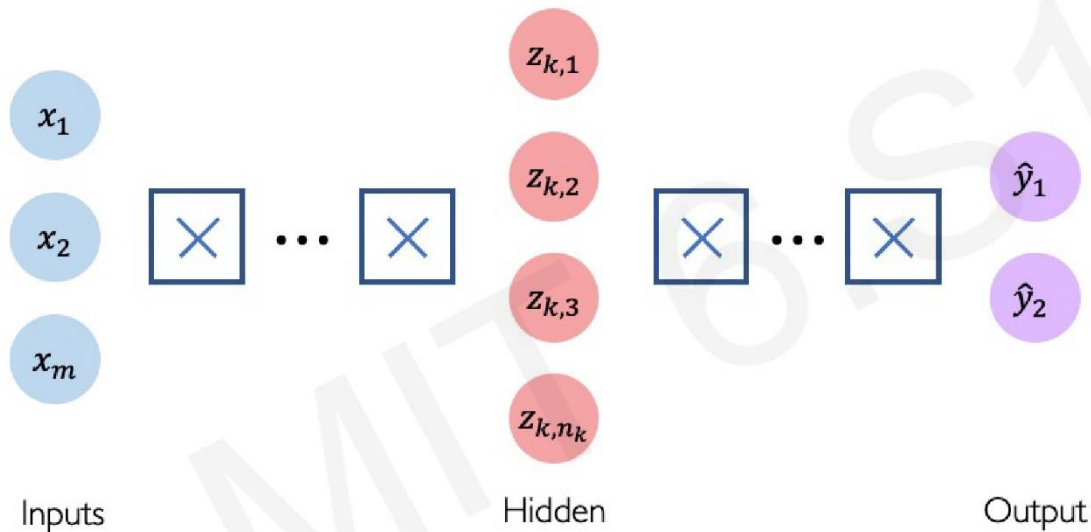Inputs          Hidden          Output

# Deep Neural Network



$x_1$

$x_2$

$x_m$

$z_{k,1}$

$z_{k,2}$

$z_{k,3}$

$z_{k,n_k}$

$\hat{y}_1$

$\hat{y}_2$

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(n₁),
    tf.keras.layers.Dense(n₂),
    ⋮
    tf.keras.layers.Dense(2)
])
```
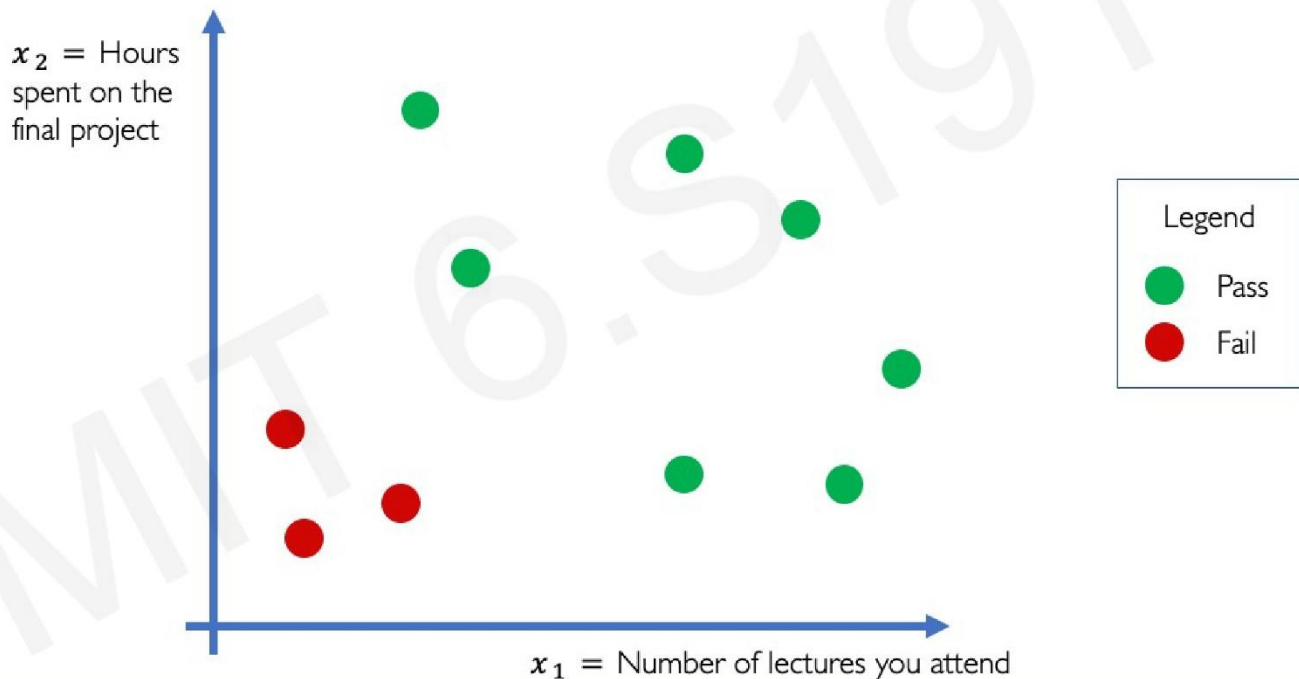
# Example Problem

## Will I pass this class?

Let's start with a simple two feature model

$x_1$ = Number of lectures you attend
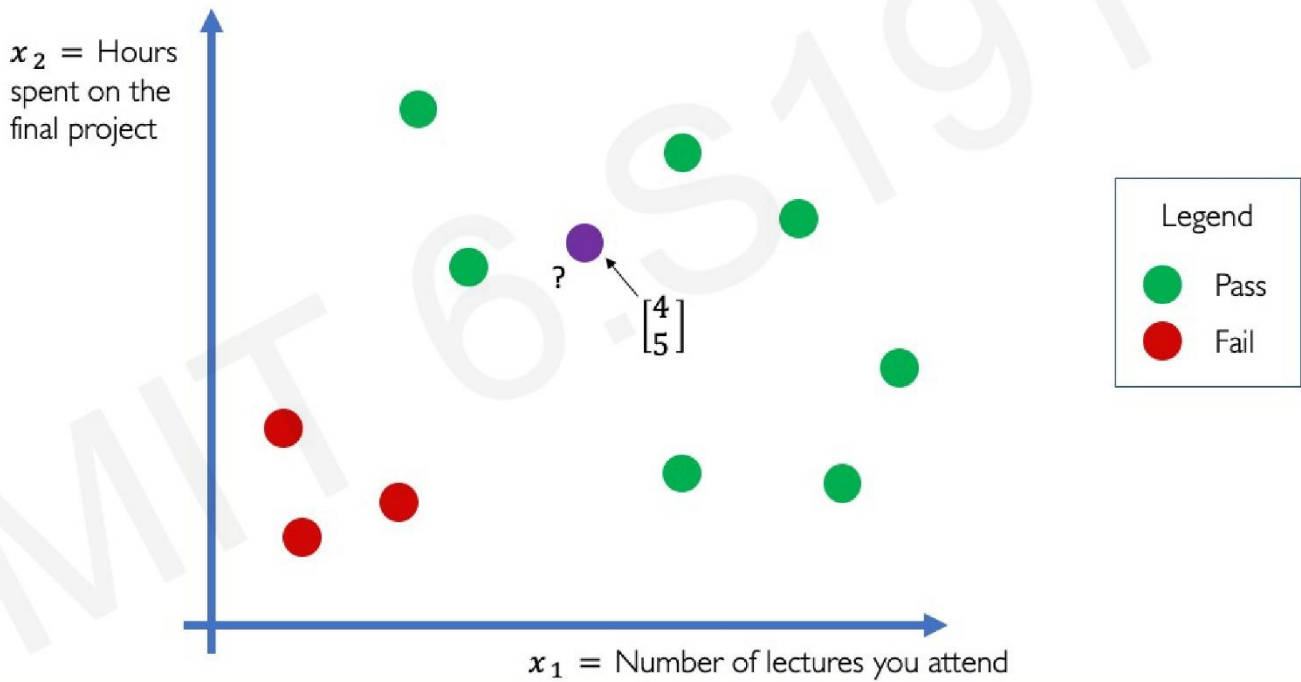
$x_2$ = Hours spent on the final project

Example Problem: Will I pass this class?
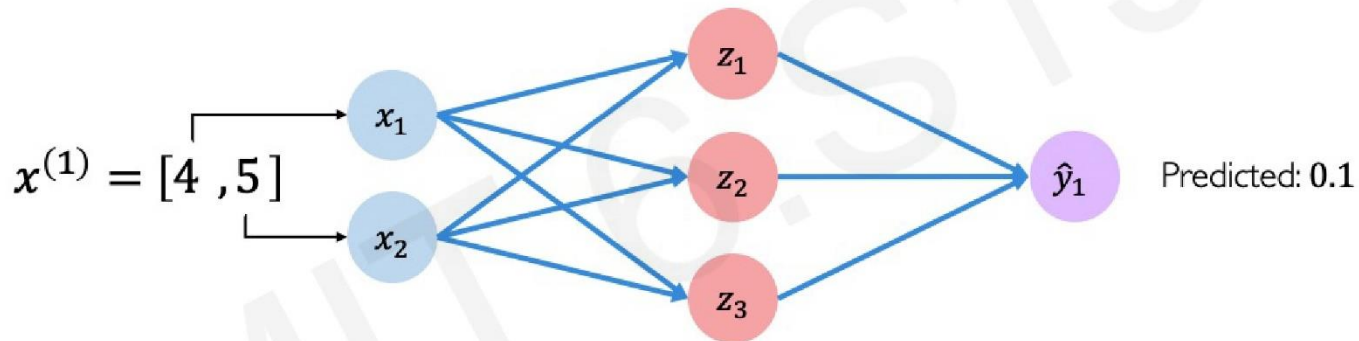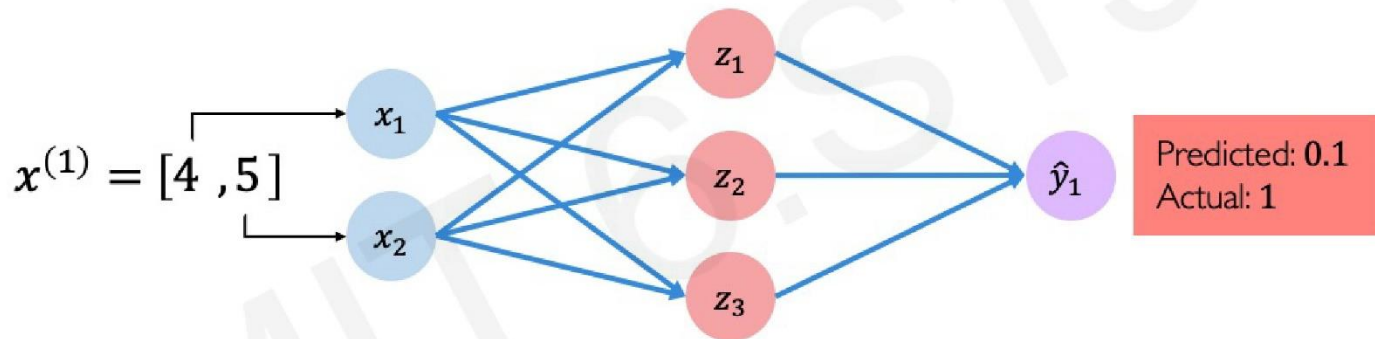
Example Problem: Will I pass this class?

# Example Problem: Will I pass this class?
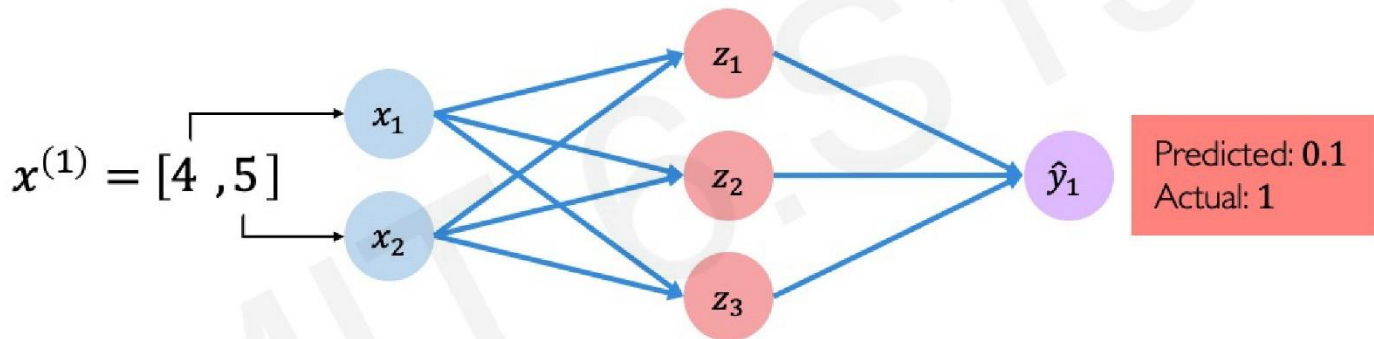


$x^{(1)} = [4, 5]$

Predicted: 0.1

# Example Problem: Will I pass this class?

# Quantifying Loss

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



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

Predicted: 0.1
Actual: 1

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

Predicted    Actual

# Empirical Loss

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



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

Predicted    Actual

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

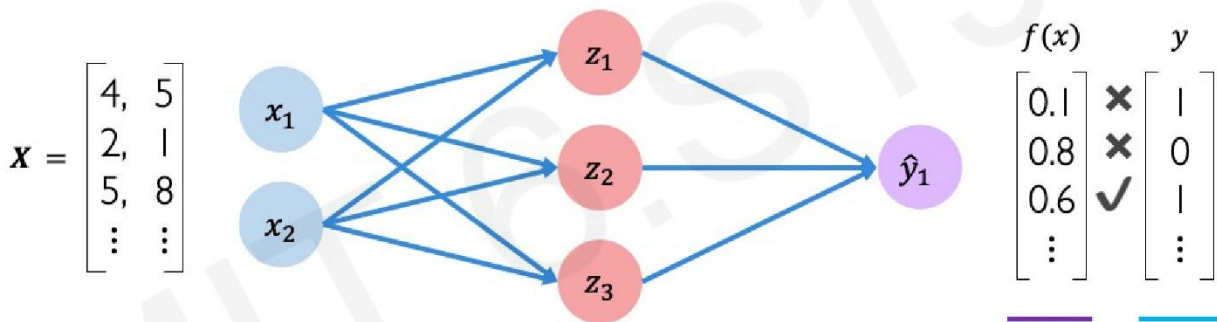# Binary Cross Entropy Loss

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

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

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

$f(x)$  $y$

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

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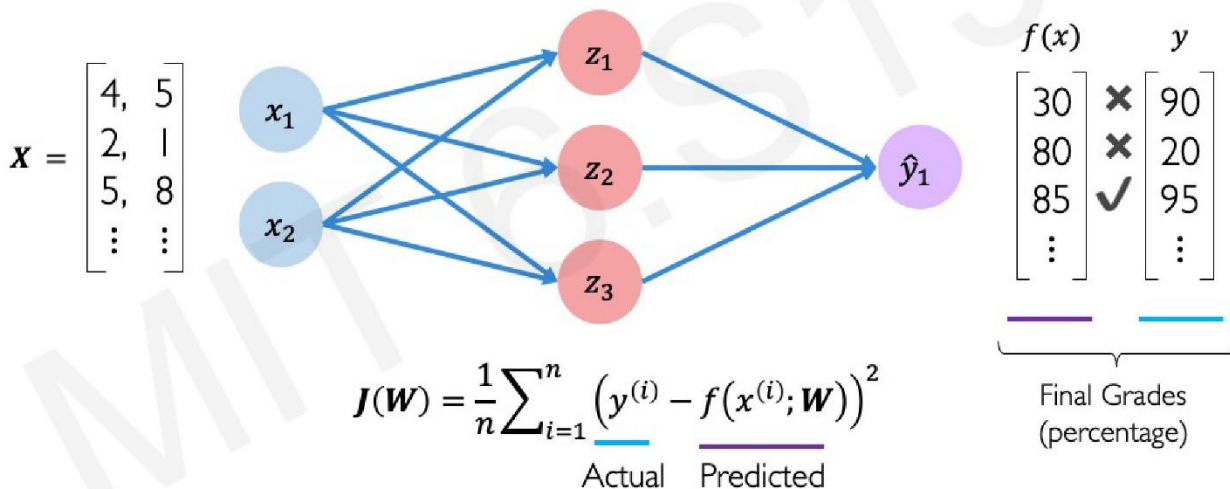
Actual  Predicted  Actual  Predicted

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

# Mean Squared Error Loss

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



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

$x_1$  $x_2$

$z_1$  $z_2$  $z_3$

$\hat{y}_1$

$$f(x) \qquad y$$

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

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)) )
loss = tf.keras.losses.MSE( y, predicted )
```

Training Neural Networks

# Loss Optimization

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

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

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

# Loss Optimization

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

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

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

Remember:

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

Loss Optimization

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

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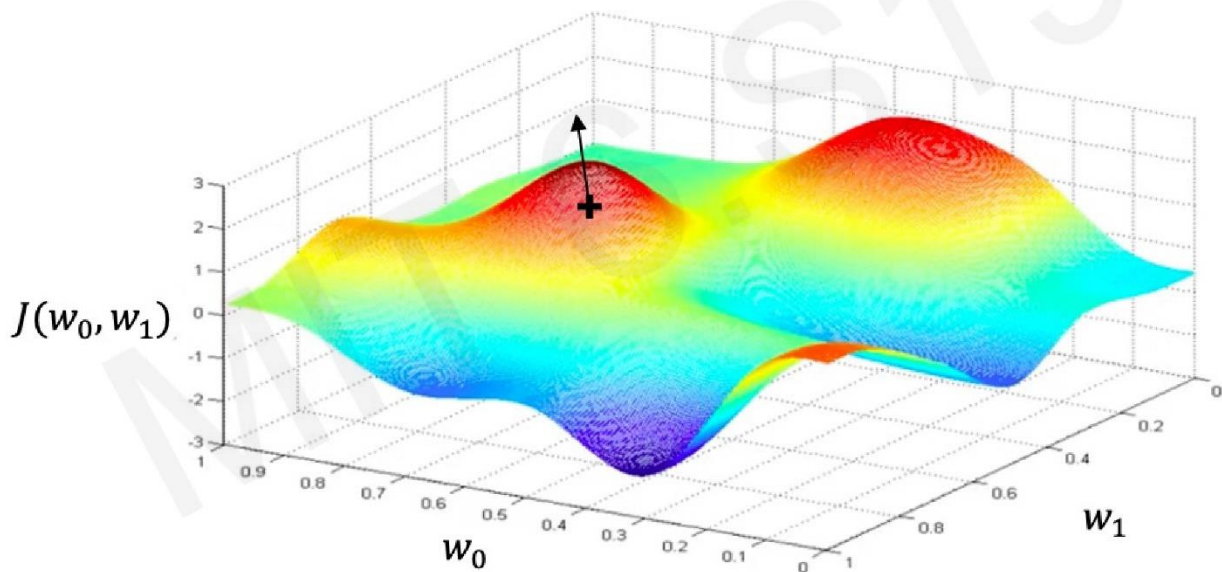*Our loss is a function of the network weights!*

$J(w_0, w_1)$
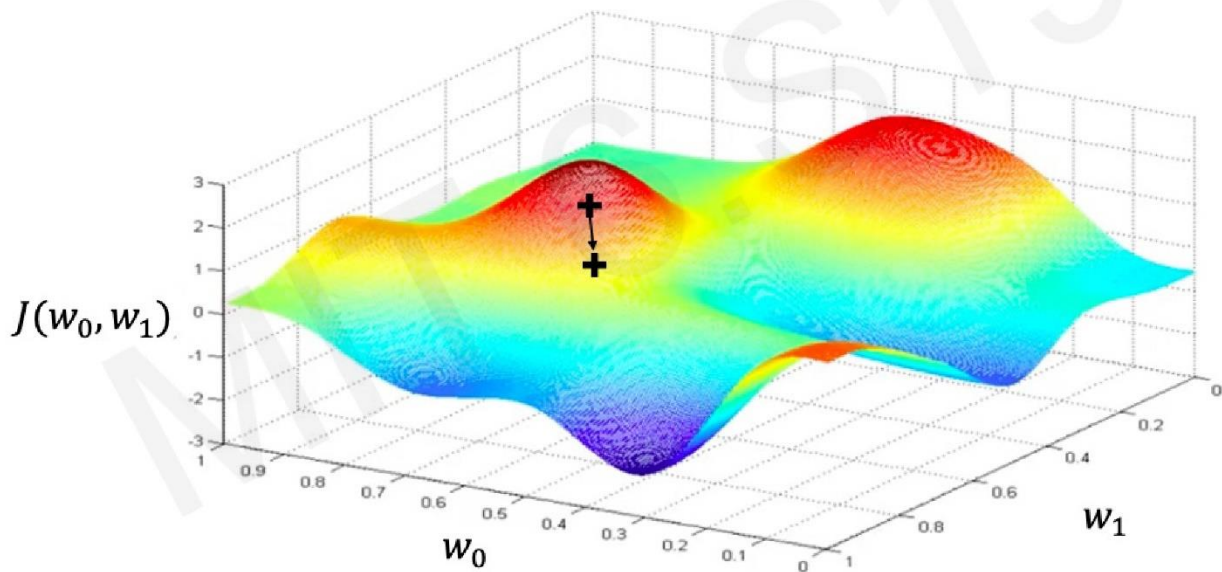
Loss Optimization

Randomly pick an initial $(w_0, w_1)$

Loss Optimization

Loss Optimization

Take small step in opposite direction of gradient

# 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

# 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

```python
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, $\dfrac{\partial J(W)}{\partial W}$

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

5. Return weights

```python
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 $J(W)$?

# Computing Gradients: Backpropagation



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

Let's use the chain rule!

# Computing Gradients: Backpropagation



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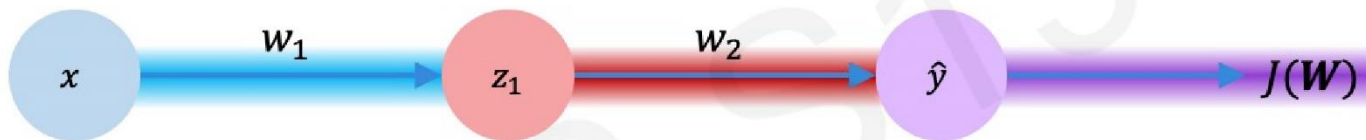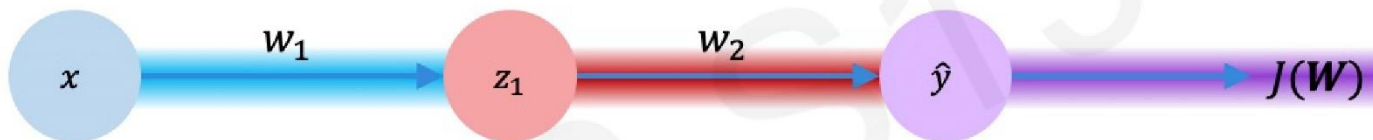
# Computing Gradients: Backpropagation



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

Apply chain rule!          Apply chain rule!

# 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

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

# 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?

# 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/

Neural Networks in Practice:
Mini-batches
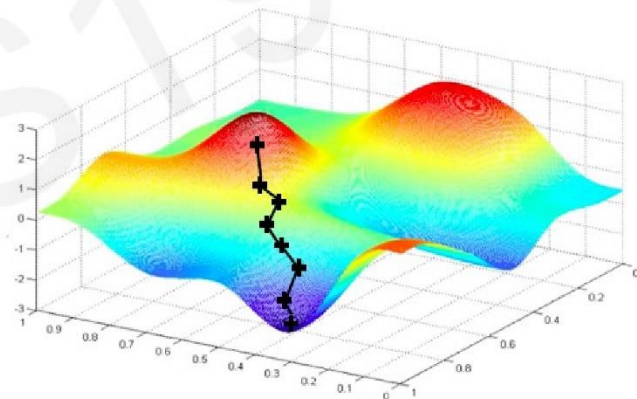
# Gradient Descent

## Algorithm

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

2. Loop until convergence:

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

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

5. Return weights

# Gradient Descent

**Algorithm**

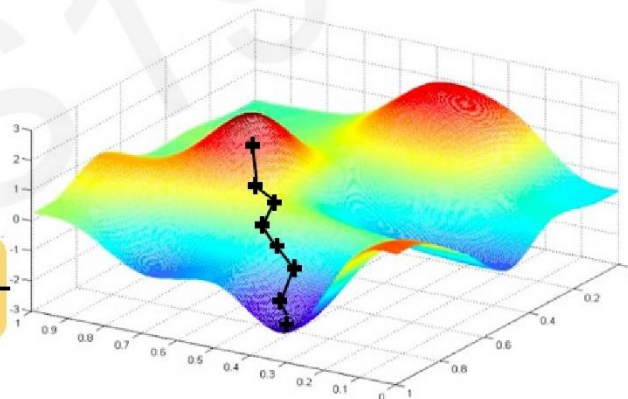1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

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

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

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

6. Return weights

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

**Underfitting**
Model does not have capacity to fully learn the data

Ideal fit

**Overfitting**
Too complex, extra parameters, does not generalize well

# 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 1: Dropout

- During training, randomly set some activations to 0

# Regularization 1: 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 1: 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



| Legend |
|--------|
| Testing |
| Training |

# Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



Loss

Training Iterations

**Legend**

Testing

Training

# 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



**Legend**

Testing

Training

# 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

# 提纲

Sequences in the Wild

Sequence Modeling Applications

# The Perceptron Revisited

Feed-Forward Networks Revisited

# Feed-Forward Networks Revisited



$$\boldsymbol{x_t} \in \mathbb{R}^m \qquad\qquad \boldsymbol{\hat{y}_t} \in \mathbb{R}^n$$

Handling Individual Time Steps

# Neurons with Recurrence



$$\hat{y}_t = f(x_t, h_{t-1})$$

output · input · past memory

# Neurons with Recurrence

output vector

input vector

$\hat{y}_t$

recurrent cell

$h_t$

$x_t$

$\hat{y}_0$ $\hat{y}_1$ $\hat{y}_2$

$h_0$ $h_1$

$x_0$ $x_1$ $x_2$

$$\hat{y}_t = f(x_t, h_{t-1})$$

output    input    past memory

# Recurrent Neural Networks (RNNs)

output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(x_t, h_{t-1})$$

cell state | function with weights W | input | old state

Note: the same function and set of parameters are used at every time step

RNNs have a **state**, $h_t$, that is updated **at each time step** as a sequence is processed

# RNN State Update and Output



output vector $\hat{y}_t$

RNN $h_t$

input vector $x_t$

**Update Hidden State**

$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

**Input Vector**

$$x_t$$

# RNNs: Computational Graph Across Time



$= $ Represent as computational graph unrolled across time

# RNNs: Computational Graph Across Time

Forward pass

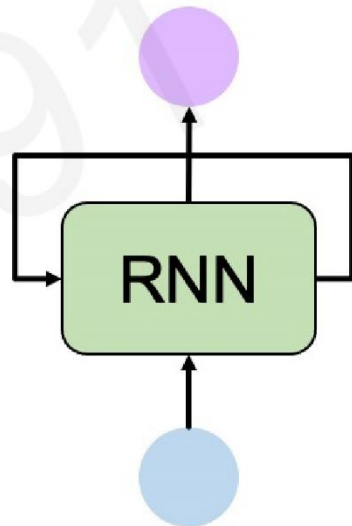Re-use the **same weight matrices** at every time step

# Sequence Modeling: Design Criteria

To model sequences, we need to:

1.  Handle **variable-length** sequences

2.  Track **long-term** dependencies

3.  Maintain information about **order**

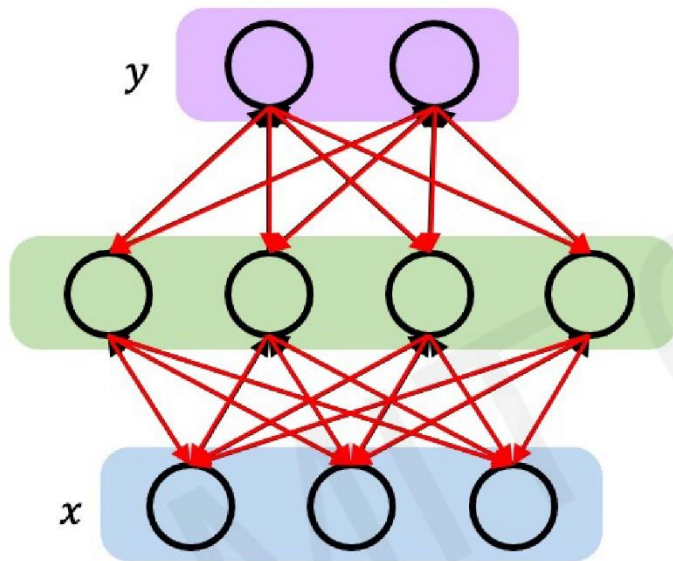4.  **Share parameters** across the sequence

**Recurrent Neural Networks (RNNs)** meet
these sequence modeling design criteria

RNN

A Sequence Modeling Problem:
Predict the Next Word

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

given these words

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

given these words          predict the
                           next word

# A Sequence Modeling Problem: Predict the Next Word

"This morning I took my cat for a walk."

given these words     predict the next word

## Representing Language to a Neural Network



"deep" ✗→ ✗→ "learning"

*Neural networks cannot interpret words*

$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \end{bmatrix}$ → → $\begin{bmatrix} 0.9 \\ 0.2 \\ 0.4 \end{bmatrix}$

*Neural networks require numerical inputs*

# Encoding Language for a Neural Network

Neural networks cannot interpret words

Neural networks require numerical inputs

Embedding: transform indexes into a vector of fixed size.

**1. Vocabulary:** Corpus of words

**2. Indexing:** Word to index

**3. Embedding:** Index to fixed-sized vector

# Handle Variable Sequence Lengths

The food was great

vs.

We visited a restaurant for lunch

vs.

We were hungry but cleaned the house before eating

# Capture Differences in Sequence Order

The food was good, not bad at all.

vs.

The food was bad, not good at all.

# Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences

2. Track **long-term** dependencies

3. Maintain information about **order**

4. **Share parameters** across the sequence

**Recurrent Neural Networks (RNNs)** meet
these sequence modeling design criteria

Backpropagation Through Time (BPTT)

# Recall: Backpropagation in Feed Forward Models



**Backpropagation algorithm:**

1. Take the derivative (gradient) of the loss with respect to each parameter

2. Shift parameters in order to minimize loss

# RNNs: Backpropagation Through Time

RNNs: Backpropagation Through Time

Standard RNN Gradient Flow

Standard RNN Gradient Flow

Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + **repeated gradient computation!**

# Standard RNN Gradient Flow: Exploding Gradients



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + **repeated gradient computation**!

Many values > 1:
**exploding gradients**

**Gradient clipping** to scale big gradients

# Standard RNN Gradient Flow:Vanishing Gradients



Computing the gradient wrt $h_0$ involves **many factors of $W_{hh}$** + repeated gradient computation!

Many values > 1:
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1:
**vanishing gradients**

1. Activation function
2. Weight initialization
3. Network architecture

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

↓

Errors due to further back time steps
have smaller and smaller gradients

↓

Bias parameters to capture short-term
dependencies

# The Problem of Long-Term Dependencies

"The clouds are in the ____"

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

↓

Errors due to further back time steps
have smaller and smaller gradients

↓

Bias parameters to capture short-term
dependencies

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

Errors due to further back time steps have smaller and smaller gradients

Bias parameters to capture short-term dependencies

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

↓

Errors due to further back time steps have smaller and smaller gradients

↓

Bias parameters to capture short-term dependencies

"The clouds are in the ___"

$\hat{y}_0$ $\hat{y}_1$ $\hat{y}_2$ $\hat{y}_3$ $\hat{y}_4$

$x_0$ $x_1$ $x_2$ $x_3$ $x_4$

"I grew up in France, … and I speak fluent___"

# The Problem of Long-Term Dependencies

**Why are vanishing gradients a problem?**

Multiply many **small numbers** together

↓

Errors due to further back time steps have smaller and smaller gradients

↓

Bias parameters to capture short-term dependencies



"The clouds are in the ___"

"I grew up in France, … and I speak fluent___"

Trick #1: Activation Functions

# Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

Initialize **biases** to zero

$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.

# Trick #3: Gated Cells

Idea: use **gates** to selectively **add** or **remove** information within **each recurrent unit with**

Pointwise multiplication

Sigmoid neural net layer

$\sigma$

gated cell

LSTM, GRU, etc.

Gates optionally let information through the cell

**Long Short Term Memory (LSTMs)** networks rely on a gated cell to track information throughout many time steps.

# Long Short Term Memory (LSTMs)

Gated LSTM cells **control information flow:**
1) Forget    2) Store    3) Update    4) Output



LSTM cells are able to track information throughout many timesteps

`tf.keras.layers.LSTM(num_units)`

# LSTMs: Key Concepts

1. Maintain a **cell state**

2. Use **gates** to control the **flow of information**

   - **Forget** gate gets rid of irrelevant information

   - **Store** relevant information from current input

   - Selectively **update** cell state

   - **Output** gate returns a filtered version of the cell state

3. Backpropagation through time with partially **uninterrupted gradient flow**

# RNN Applications & Limitations

# Example Task: Music Generation

F#        G        C        A

**Input:** sheet music

**Output:** next character in sheet music

Listening to
3rd movement

E        F#        G        C

Software Lab!

# Example Task: Sentiment Classification



sentiment
<positive>

**Input:** sequence of words

**Output:** probability of having positive sentiment

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

I   love   this   class!

# Example Task: Sentiment Classification

# Limitations of Recurrent Models

# Goal of Sequence Modeling

RNNs: recurrence to model sequence dependencies

# Goal of Sequence Modeling

RNNs: recurrence to model sequence dependencies

## Limitations of RNNs

Encoding bottleneck

Slow, no parallelization

Not long memory

# Goal of Sequence Modeling

Can we eliminate the need for recurrence entirely?

**Desired Capabilities**

Continuous stream

Parallelization

Long memory

# Goal of Sequence Modeling

Can we eliminate the need for recurrence entirely?

**Desired Capabilities**

Continuous stream

Parallelization

Long memory

$\hat{y}_0 \quad \hat{y}_1 \quad \hat{y}_2 \qquad \hat{y}_{t-2} \quad \hat{y}_{t-1} \quad \hat{y}_t$   output

· · ·   feature vector

$x_0 \quad x_1 \quad x_2 \qquad x_{t-2} \quad x_{t-1} \quad x_t$   input

$t$

# Goal of Sequence Modeling

Idea 1: Feed everything into dense network

✔ **No recurrence**

✘ **Not scalable**

✘ **No order**

✘ **No long memory**

Idea: Identify and attend to what's important

Can we eliminate the need for recurrence entirely?

| $\hat{y}_0$ | $\hat{y}_1$ | $\hat{y}_2$ | | $\hat{y}_{t-2}$ | $\hat{y}_{t-1}$ | $\hat{y}_t$ | output |

feature vector

| $x_0$ | $x_1$ | $x_2$ | | $x_{t-2}$ | $x_{t-1}$ | $x_t$ | input |

Intuition Behind Self-Attention

Attending to the most important parts of an input.

1. Identify which parts to attend to

2. Extract the features with high attention

Similar to a search problem!

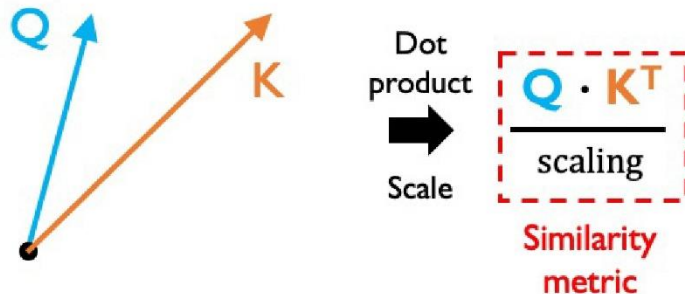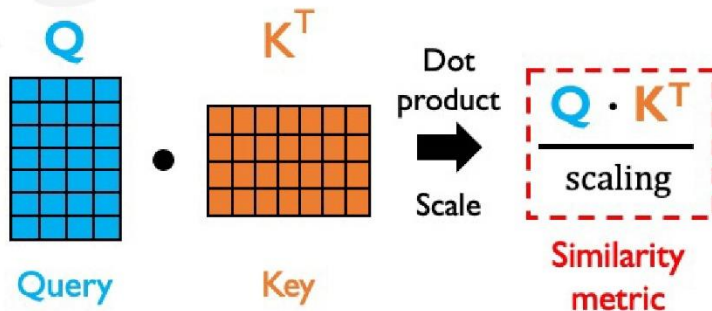A Simple Example: Search

Understanding Attention with Search

# Learning Self-Attention with Neural Networks

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract **features with high attention**

$x$

He    tossed    the    tennis    ball    to    serve

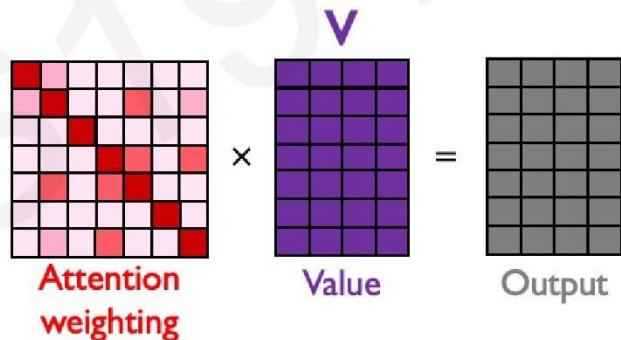Data is fed in all at once! Need to encode position information to understand order.

# Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

I. Encode **position** information

2. Extract **query, key, value** for search

3. Compute **attention weighting**

4. Extract **features with high attention**



Position-aware encoding

Data is fed in all at once! Need to encode position information to understand order.
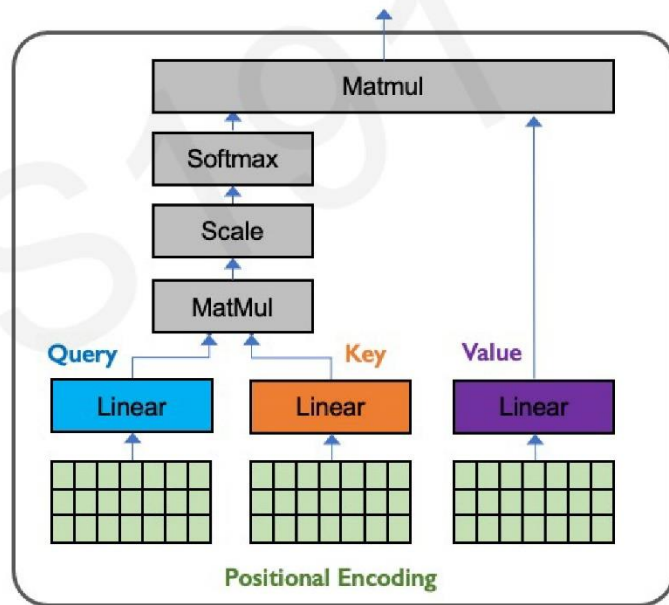
# Learning Self-Attention with Neural Networks

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute attention weighting
4. Extract features with high attention

# Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

**Attention score:** compute pairwise similarity between each query and key

How to compute similarity between two sets of features?

Q

K

Dot product

Scale

$$\frac{Q \cdot K^T}{\text{scaling}}$$

**Similarity metric**

Also known as the "cosine similarity"

# Learning Self-Attention with Neural Networks

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

**Attention score:** compute pairwise similarity between each query and key

How to compute similarity between two sets of features?



Also known as the "cosine similarity"

# Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information

2. Extract **query, key, value** for search

3. Compute **attention weighting**

4. Extract features with high attention

Attention weighting: where to attend to! How similar is the key to the query?

$$softmax\left(\frac{Q \cdot K^T}{scaling}\right)$$

**Attention weighting**

# Learning Self-Attention with Neural Networks

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract **features with high attention**

Last step: self-attend to extract features



$$softmax\left(\frac{Q \cdot K^T}{scaling}\right) \cdot V = A(Q, K, V)$$

# Learning Self-Attention with Neural Networks

**Goal:** identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract **features with high attention**

These operations form a self-attention head that can plug into a larger network. Each head attends to a different part of input.

Matmul

Softmax

Scale

MatMul

**Query** Linear

**Key** Linear

**Value** Linear

**Positional Encoding**

$$softmax\left(\frac{Q \cdot K^T}{scaling}\right) \cdot V$$

Applying Multiple Self-Attention Heads

# Self-Attention Applied

**Language Processing**

An armchair in the shape of an avocado

BERT, GPT-3

Devlin et al., *NAACL* 2019
Brown et al., *NeurIPS* 2020

**Biological Sequences**

AlphaFold2

Jumper et al., *Nature* 2021

**Computer Vision**

Vision Transformers

Dosovitskiy et al., *ICLR* 2020

# Deep Learning for Sequence Modeling: Summary

1. RNNs are well suited for **sequence modeling** tasks

2. Model sequences via a **recurrence relation**

3. Training RNNs with **backpropagation through time**

4. Models for **music generation**, classification, machine translation, and more

5. Self-attention to model **sequences without recurrence**

# 提纲

"To know what is
where by looking."

To discover from images what is present in the world, where things are, what actions are taking place, to predict and anticipate events in the world

Impact: Facial Detection & Recognition

Impact: Self-Driving Cars

What Computers "See"

Images are Numbers

What the computer sees

An image is just a matrix of numbers [0,255]!
i.e., 1080x1080x3 for an RGB image

# Tasks in Computer Vision



Input Image → Pixel Representation → classification →

| | |
|---|---|
| Lincoln | 0.8 |
| Washington | 0.1 |
| Jefferson | 0.05 |
| Obama | 0.05 |

- **Regression**: output variable takes continuous value
- **Classification**: output variable takes class label. Can produce probability of belonging to a particular class

# High Level Feature Detection

Let's identify key features in each image category



Nose,
Eyes,
Mouth



Wheels,
License Plate,
Headlights



Door,
Windows,
Steps

# Learning Feature Representations

Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



Facial structure

Learning Visual Features

# Fully Connected Neural Network

**Input:**
- 2D image
- Vector of pixel values

$x_1$

$x_2$

$\vdots$

$x_p$

**Fully Connected:**
- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?

# Using Spatial Structure

**Input:** 2D image.
Array of pixel values

**Idea:** connect patches of input to neurons in hidden layer.

Neuron connected to region of input. Only "sees" these values.

# Using Spatial Structure



Connect patch in input layer to a single neuron in subsequent layer.
Use a sliding window to define connections.
*How can we* **weight** *the patch to detect particular features?*

# Applying Filters to Extract Features

1) Apply a set of weights – a filter – to extract **local features**

2)  Use **multiple filters** to extract different features

3) Spatially **share** parameters of each filter
(features that matter in one part of the input should matter elsewhere)

# Feature Extraction with Convolution



- Filter of size 4x4 : 16 different weights
- Apply this same filter to 4x4 patches in input
- Shift by 2 pixels for next patch

This "patchy" operation is **convolution**

1) Apply a set of weights – a filter – to extract **local features**

2) Use **multiple filters** to extract different features

3) **Spatially share** parameters of each filter

The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:

# Producing Feature Maps



Original

Sharpen

Edge Detect

"Strong" Edge Detect

Convolutional Neural Networks (CNNs)

# CNNs for Classification

Input image     Convolution (feature maps)     Maxpooling     Fully-connected layer

1. **Convolution**: Apply filters to generate feature maps.

2. **Non-linearity**: Often ReLU.

3. **Pooling**: Downsampling operation on each feature map.

`tf.keras.layers.Conv2D`

`tf.keras.activations.*`

`tf.keras.layers.MaxPool2D`

**Train model with image data.**
**Learn weights of filters in convolutional layers.**

# Convolutional Layers: Local Connectivity

`tf.keras.layers.Conv2D`

**For a neuron in hidden layer:**
- Take inputs from patch
- Compute weighted sum
- Apply bias

4x4 filter: matrix of weights $w_{ij}$

$$\sum_{i=1}^{4}\sum_{j=1}^{4} w_{ij}\, x_{i+p,j+q} + b$$

for neuron (p,q) in hidden layer

1) applying a window of weights
2) computing linear combinations
3) activating with non-linear function

# Introducing Non-Linearity

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**

**Rectified Linear Unit (ReLU)**

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

`tf.keras.layers.ReLU`

# Pooling



max pool with 2x2 filters
and stride 2

```
tf.keras.layers.MaxPool2D(
    pool_size=(2,2),
    strides=2
)
```

1) Reduced dimensionality
2) Spatial invariance

How else can we downsample and preserve spatial invariance?

# Representation Learning in Deep CNNs



Low level features

Edges, dark spots

Conv Layer 1

Mid level features

Eyes, ears, nose

Conv Layer 2

High level features

Facial structure

Conv Layer 3

CNNs for Classification: Feature Learning

1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**

# CNNs for Classification: Class Probabilities



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

An Architecture for Many Applications

Classification
Object detection
Segmentation
Probabilistic control

# Classification: Breast Cancer Screening



International evaluation of an AI system for breast cancer screening    nature

CNN-based system outperformed expert radiologists at detecting breast cancer from mammograms

Breast cancer case missed by radiologist but detected by AI

Object Detection

# Semantic Segmentation: Biomedical Image Analysis



Brain Tumors
Dong+ *MIUA* 2017.

Malaria Infection
Soleimany+ *arXiv* 2019.

# Continuous Control: Navigation from Vision



**Raw Perception**
*I*
(ex. camera)

**Coarse Maps**
*M*
(ex. GPS)

**Possible Control Commands**

End-to-End Framework for Autonomous Navigation

Entire model is trained end-to-end **without any human labelling or annotations**

前馈神经网络
MLP

循环神经网络
RNN

卷积神经网络
CNN