



Artificial Intelligence and Mobility Lab



GAN and Reinforcement Learning: Theory and Implementation

Prof. Joongheon Kim

Korea University, School of Electrical Engineering

Artificial Intelligence and Mobility Laboratory

<https://joongheon.github.io>

joongheon@korea.ac.kr

- **Introduction**

- Generative Adversarial Network
- Reinforcement Learning
- Deep Reinforcement Learning





Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

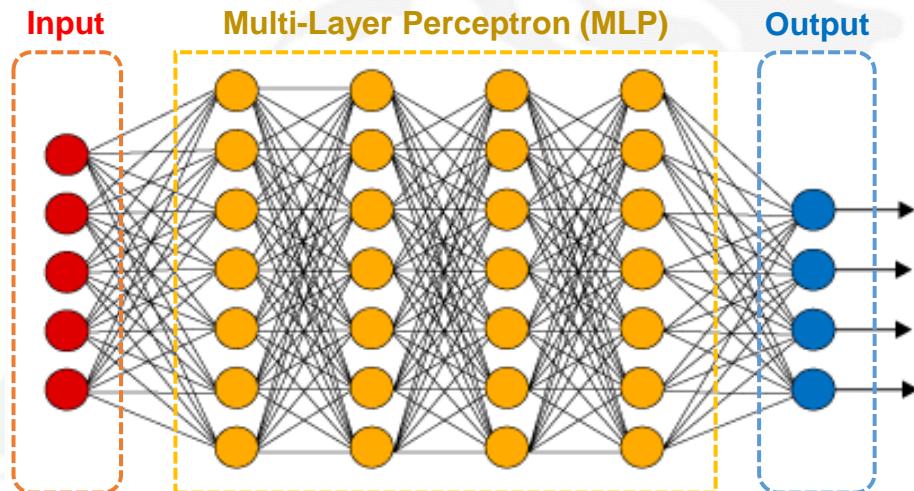
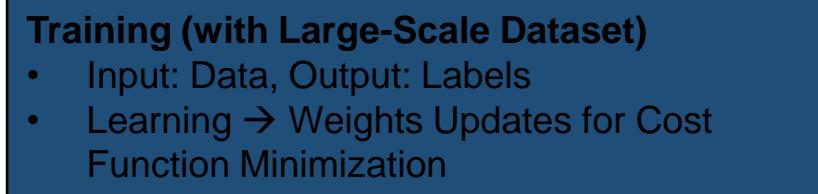
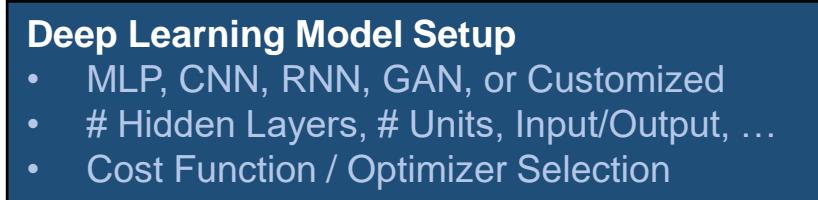
Introduction

Introduction

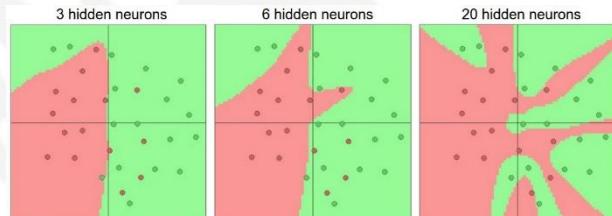
- **Deep Learning Overview**
- Linear Regression
- Binary Classification
- Neural Network Basics

Introduction

- How Deep Learning Works?
 - Deep Learning Computation Procedure



Non-Linear Training (Weights Updates) for Cost Minimization: GD, SGD, Adam, etc.



Introduction

- How Deep Learning Works?
- Deep Learning Computation Procedure

Deep Learning Model Setup

- MLP, CNN, RNN, GAN, or Customized
- # Hidden Layers, # Units, Input/Output, ...
- Cost Function / Optimizer Selection



Training (with Large-Scale Dataset)

- Input: Data, Output: Labels
- Learning → Weights Updates for Cost Function Minimization



Inference / Testing (Real-Word Execution)

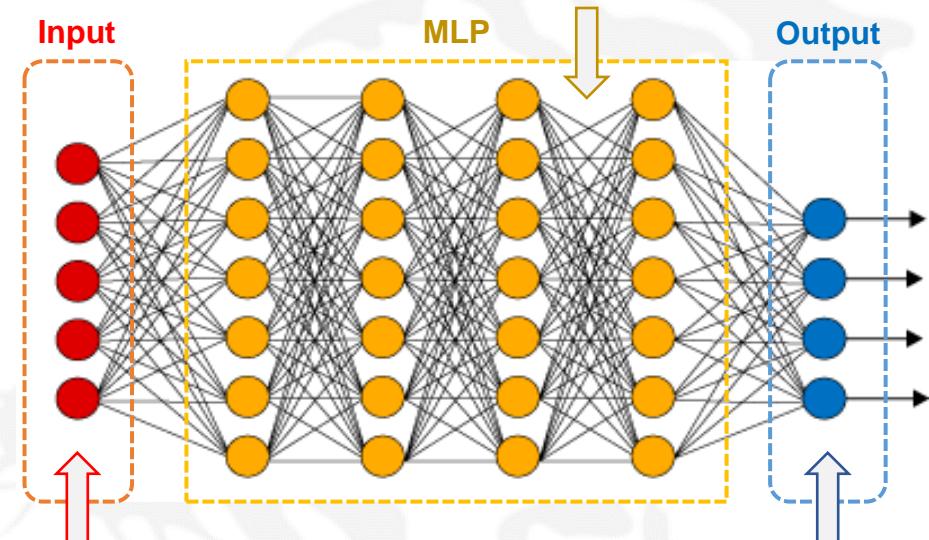
- Input: Real-World Input Data
- Output: Interference Results based on Updated Weights in Deep Neural Networks

All weights in units are trained/set (under cost minimization)

Input

MLP

Output



INPUT: Data

- One-Dimension Vector

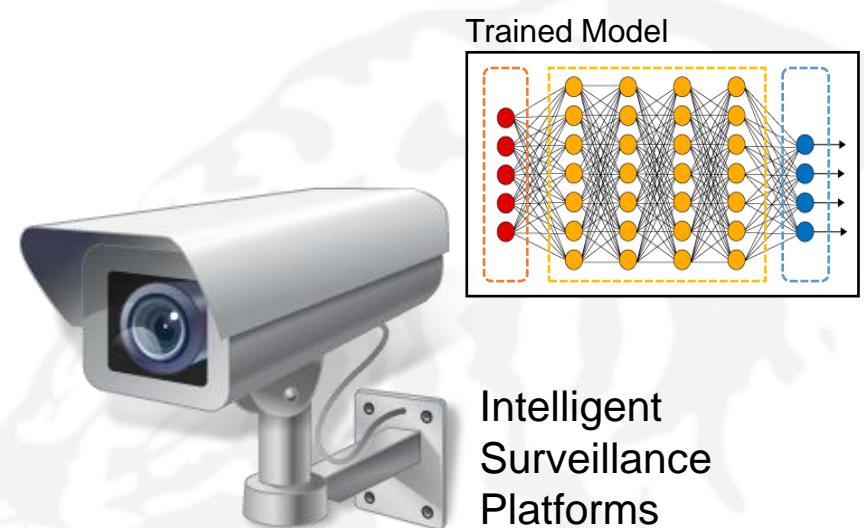
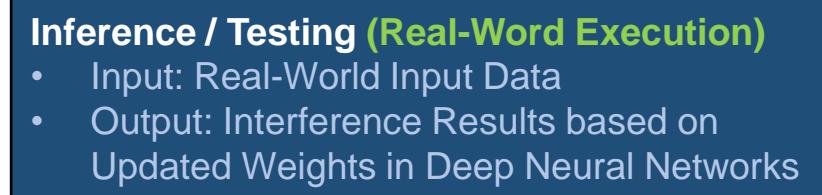
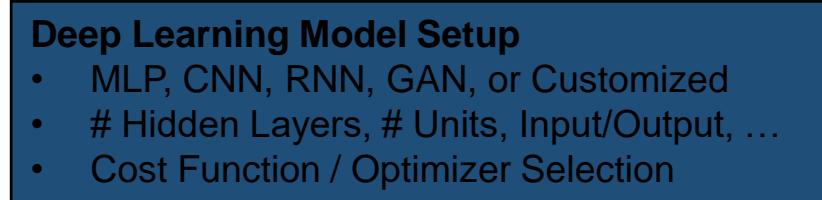
OUTPUT: Labels

- One-Hot Encoding

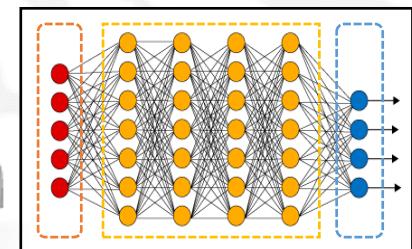
We need a lot of training data for generality
(otherwise, we will suffer from overfitting problem).

Introduction

- How Deep Learning Works?
 - Deep Learning Computation Procedure



Trained Model



Intelligent
Surveillance
Platforms

INPUT: Real-Time Arrivals

OUTPUT: Inference

- Computation Results based on (i) INPUT and (ii) trained weights in units (trained model).

Introduction

- How Deep Learning Works?

- Issue - **Overfitting**

Deep Learning Model Setup

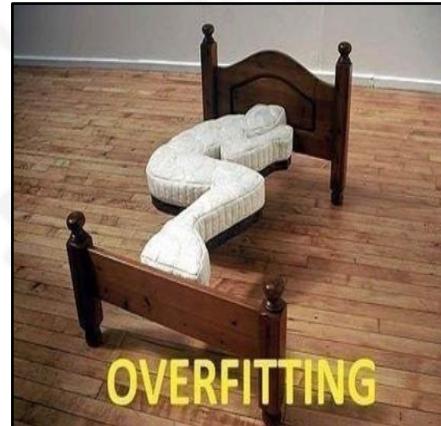
- MLP, CNN, RNN, GAN, or Custom
- # Hidden Layers, # Units, Input/Output, ...
- Cost Function / Optimizer Selection

What if we do not have enough data for training (not enough to derive Gaussian/normal distribution)?

Training (with Large-Scale Dataset)

- Input: Data, Output: Labels
- Learning → Weights Updates for Cost Function Minimization

Situation becomes worse when the model (with insufficient training data) accurately fits on training data.



Inference / Testing (Real-World Execution)

- Input: Real-World Input Data
- Output: Inference Results based on Updated Weights in Deep Neural Networks

To Combat the Overfitting

- More training data
- Autoencoding (or variational auto-encoder (VAE))
- Dropout
- Regularization

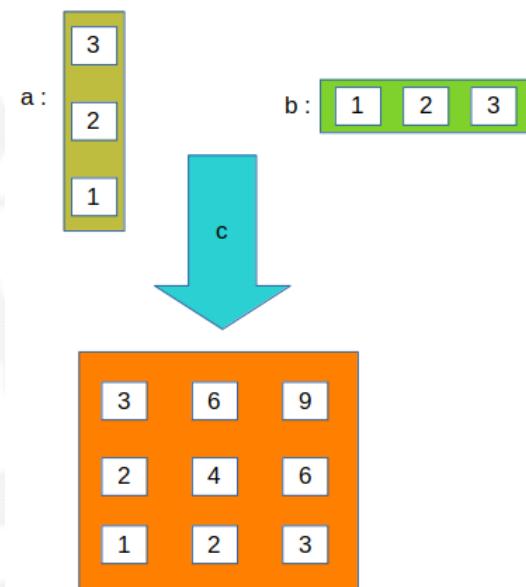
- Quick Start Example

```
ex_add.py
1 import tensorflow as tf
2 # Create nodes in computation graph
3 node1 = tf.constant(3, dtype=tf.int32)
4 node2 = tf.constant(5, dtype=tf.int32)
5 node3 = tf.add(node1, node2)
6
7 # Create session object
8 sess = tf.Session()
9 print("node1 + node2 = ", sess.run(node3))
10 # Close the session
11 sess.close()
```

```
ex_add2.py
1 import tensorflow as tf
2 # Create nodes in computation graph
3 node1 = tf.constant(3, dtype=tf.int32)
4 node2 = tf.constant(5, dtype=tf.int32)
5 node3 = tf.add(node1, node2)
6
7 # Create session object
8 with tf.Session() as sess:
9     print("node1 + node2 = ", sess.run(node3))
```

- Example: Placeholder

```
ex_placeholder.py
1 import tensorflow as tf
2
3 # Create nodes in computation graph
4 a = tf.placeholder(tf.int32, shape=(3,1))
5 b = tf.placeholder(tf.int32, shape=(1,3))
6 c = tf.matmul(a, b)
7
8 # Run computation graph
9 with tf.Session() as sess:
10     print(sess.run(c, feed_dict={a: [[3], [2], [1]], b: [[1, 2, 3]]}))
```





GAN and Reinforcement Learning

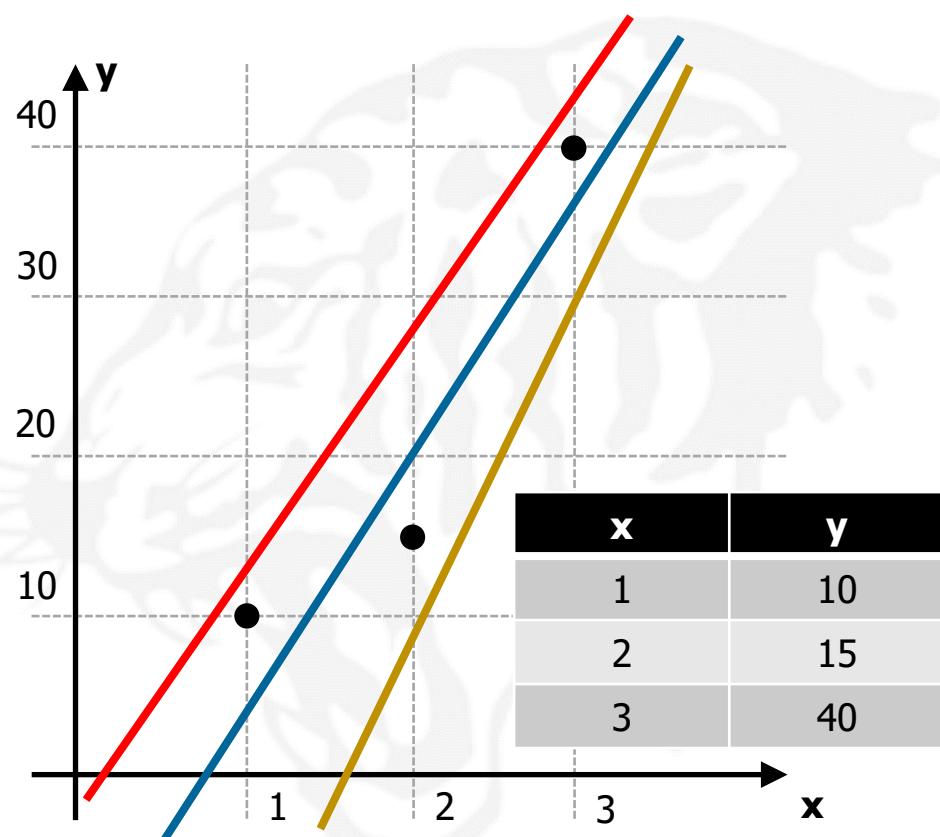
Introduction

Linear Regression

- Deep Learning Overview
- Linear Regression
- Binary Classification
- Neural Network Basics

Linear Regression

- Linear model: $H(x) = Wx + b$
- Which model is the best among the given three?

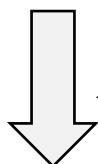


Linear Regression

- Cost Function (or Loss Function)
 - How to fit the line to training data
 - The difference between model values and real measurements:

m : The number of training data

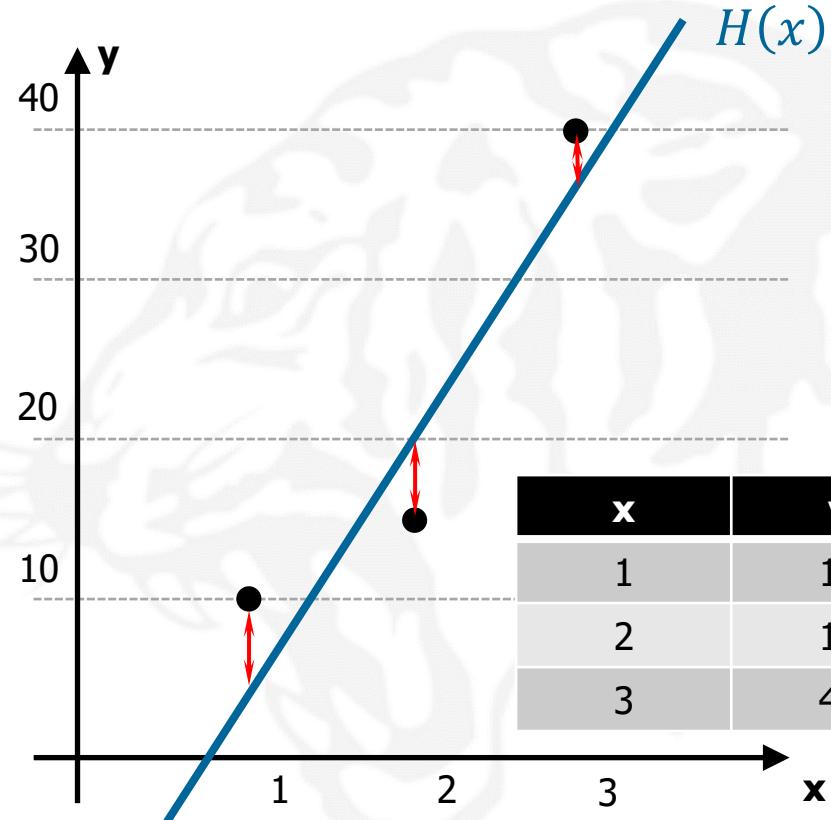
$$\frac{1}{m} \sum_{i=1}^m (H(x^i) - y^i)^2$$



$$H(x) = Wx + b$$

$$\text{Cost}(W, b) =$$

$$\frac{1}{m} \sum_{i=1}^m (H(x^i) - y^i)^2$$



- Cost Function Minimization

- Model: $H(x) = Wx + b$

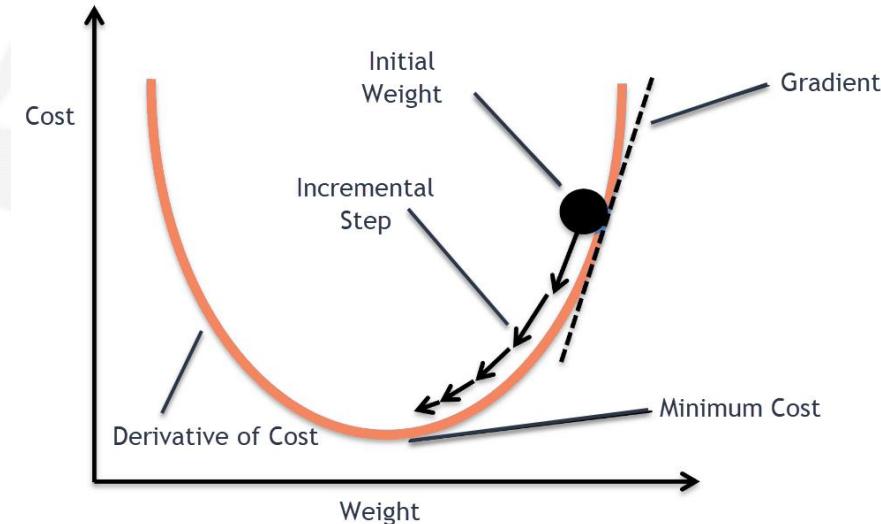
- Cost Function: $Cost(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x^i) - y^i)^2 = \frac{1}{m} \sum_{i=1}^m (Wx^i + b - y^i)^2$

- How to Minimize this Function? → **Gradient Descent Method**

- Angle → Differentiation

$$W \leftarrow W - \alpha \frac{\partial}{\partial W} Cost(W)$$

α : Learning rate



- Multi-Variable Linear Regression

- Model:

$$H(x_1, x_2, \dots, x_n) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

- Cost:

$$Cost(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x_1^i, x_2^i, \dots, x_n^i) - y^i)^2$$

- Multi-Variable Linear Regression

- Model:

$$H(x_1, x_2, \dots, x_n) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$\Rightarrow H(X) = XW + b$$

$$(x_1 \quad x_2 \dots \quad x_n) \cdot \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_n \end{pmatrix} = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

X

W

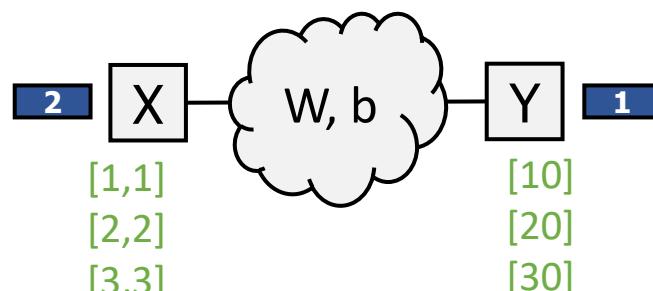


$$H(X) = XW^T + b$$

$$\text{when } W = (w_1 \quad w_2 \dots \quad w_n)$$

Linear Regression Implementation (TensorFlow)

```
1 import tensorflow as tf
2
3 x_data = [[1,1], [2,2], [3,3]]
4 y_data = [[10], [20], [30]]
5 X = tf.placeholder(tf.float32, shape=[None, 2])
6 Y = tf.placeholder(tf.float32, shape=[None, 1])
7
8 W=tf.Variable(tf.random_normal([2,1]))
9 b=tf.Variable(tf.random_normal([1]))
10
11 model = tf.matmul(X,W)+b
12 cost = tf.reduce_mean(tf.square(model - Y))
13 train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
14
15 with tf.Session() as sess:
16     sess.run(tf.global_variables_initializer())
17     # Training
18     for step in range(2001):
19         c, W_, b_, _ = sess.run([cost, W, b, train], feed_dict={X: x_data, Y: y_data})
20         print(step, c, W_, b_)
21     # Testing
22     print(sess.run(model, feed_dict={X: [[4,4]]}))
```



Model, Cost, Train

```
[4.667964 ]] [0.014943]
1991 3.1940912e-05 [[5.325548 ]]
[4.6679726]] [0.01490401]
1992 3.1772186e-05 [[5.3255568]
[4.667981 ]] [0.0148651]
1993 3.1603915e-05 [[5.3255653]
[4.6679897]] [0.01482627]
1994 3.143936e-05 [[5.325574 ]
[4.6679983]] [0.01478756]
1995 3.1277286e-05 [[5.3255825]
[4.668007 ]] [0.01474891]
1996 3.1110336e-05 [[5.325591 ]
[4.6680155]] [0.01471035]
1997 3.095236e-05 [[5.325599 ]
[4.6680236]] [0.01467189]
1998 3.079518e-05 [[5.325608]
[4.668032]] [0.01463356]
1999 3.062387e-05 [[5.325616 ]
[4.6680403]] [0.01459529]
2000 3.0467529e-05 [[5.325624 ]
[4.6680484]] [0.01455714]
[[39.989246]]
```



Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

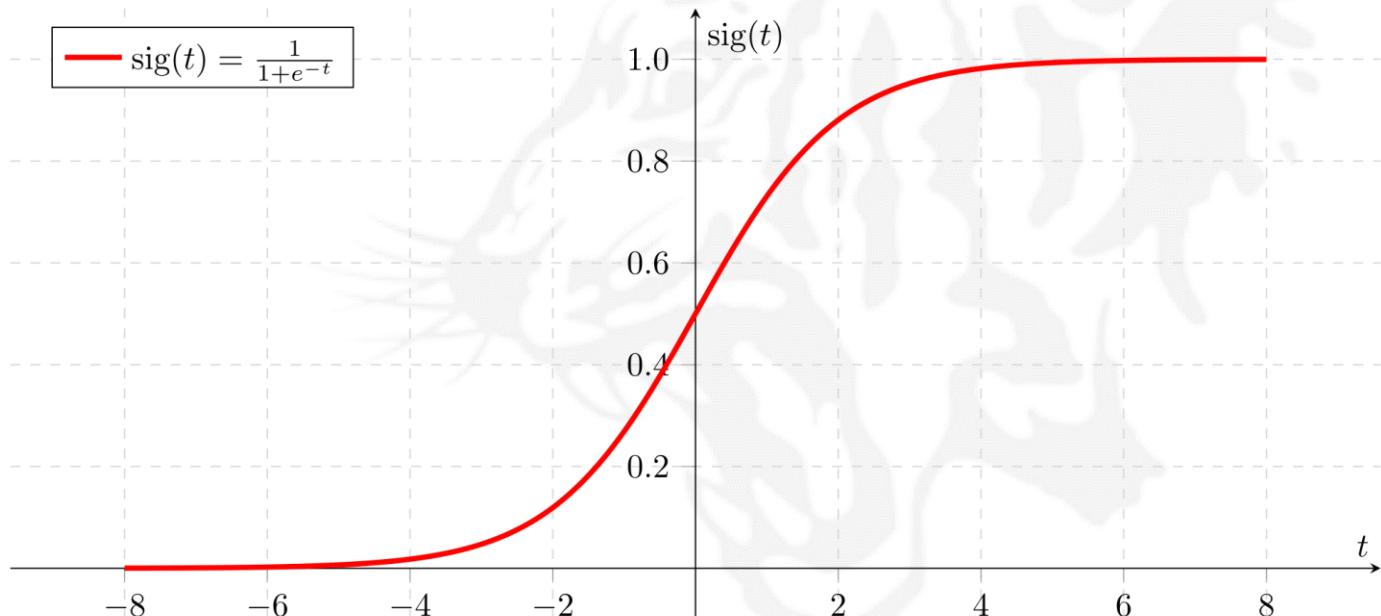
Introduction

Binary Classification

- Deep Learning Overview
- Linear Regression
- **Binary Classification**
- Neural Network Basics

- Binary Classification Examples
 - **Spam Detection:** Spam [1] or Ham [0]
 - **Facebook Feed:** Show [1] or Hide [0]
 - Facebook learns with your like-articles; and shows your favors.
 - **Credit Card Fraudulent Transaction Detection:** Fraud [1] or Legitimate [0]
 - **Tumor Image Detection in Radiology:** Malignant [1] or Benign [0]

- Binary Classification Basic Idea
 - Step 1) Linear regression with $H(x) = Wx + b$
 - Step 2) **Logistic/sigmoid function ($\text{sig}(t)$)** based on the result of Step 1.



Linear Regression Model

$$H(x) = Wx + b \text{ or } H(X) = W^T X$$

Logistic/Sigmoid Function

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

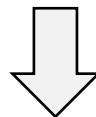
Binary Classification Model

$$g(X) = \frac{1}{1 + e^{-W^T X}}$$

$$\text{Cost}(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x^i) - y^i)^2$$

Linear Regression Model

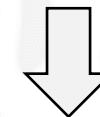
$$H(x) = Wx + b \text{ or } H(X) = W^T X$$



Gradient Descent Method can be used because $\text{Cost}(W, b)$ is convex (local minimum is global minimum).

Binary Classification Model

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



Gradient Descent Method can not be used because $\text{Cost}(W, b)$ is non-convex. **New Cost Function is required.**

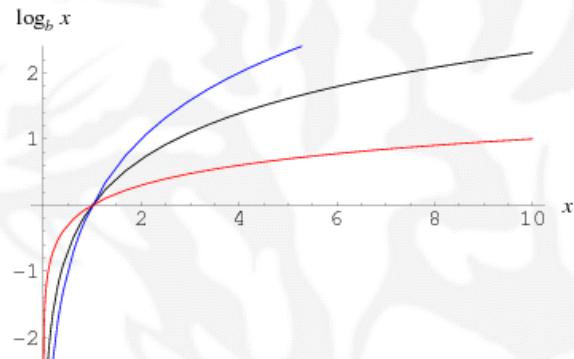
$$\text{Cost}(W) = \frac{1}{m} \sum c(H(x), y)$$

$$c(H(x), y) = \begin{cases} -\log(H(x)), & y = 1 \\ -\log(1 - H(x)), & y = 0 \end{cases}$$

Understanding this Cost Function

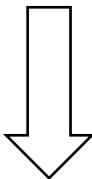
Cost	0	∞	∞	0
$H(x)$	0	0	1	1
y	0	1	0	1

Log Function



$$\text{Cost}(W) = \frac{1}{m} \sum c(H(x), y)$$

$$c(H(x), y) = \begin{cases} -\log(H(x)), & y = 1 \\ -\log(1 - H(x)), & y = 0 \end{cases}$$



$$c(H(x), y) = -y \log(H(x)) - (1 - y) \log(1 - H(x))$$



$$\text{Cost}(W) = -\frac{1}{m} \sum y \log(H(x)) + (1 - y) \log(1 - H(x))$$

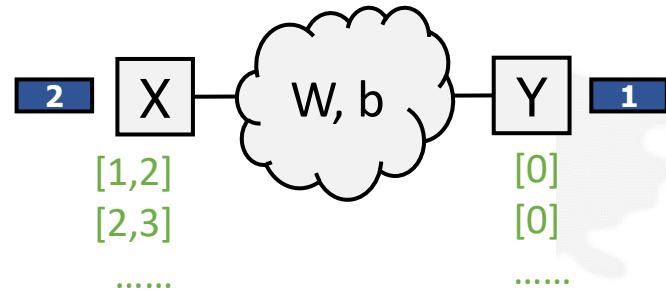


Gradient Descent Method

$$W \leftarrow W - \alpha \frac{\partial}{\partial W} \text{Cost}(W)$$

Binary Classification Implementation (TensorFlow)

```
1 import tensorflow as tf
2
3 x_data = [[1,2], [2,3], [3,1], [4,3], [5,3], [6,2]]
4 y_data = [[0], [0], [0], [1], [1], [1]]
5 X = tf.placeholder(tf.float32, shape=[None, 2])
6 Y = tf.placeholder(tf.float32, shape=[None, 1])
7 W = tf.Variable(tf.random_normal([2,1]))
8 b = tf.Variable(tf.random_normal([1]))
9
10 model = tf.sigmoid(tf.add(tf.matmul(X,W),b))
11 cost = tf.reduce_mean((-1)*Y*tf.log(model) + (-1)*(1-Y)*tf.log(1-model))
12 train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
13
14 prediction = tf.cast(model > 0.5, dtype=tf.float32)
15 accuracy = tf.reduce_mean(tf.cast(tf.equal(prediction, Y), dtype=tf.float32))
16
17 with tf.Session() as sess:
18     sess.run(tf.global_variables_initializer())
19     # Training
20     for step in range(10001):
21         cost_val, train_val = sess.run([cost, train], feed_dict={X: x_data, Y: y_data})
22         print(step, cost_val)
23     # Testing
24     h, c, a = sess.run([model, prediction, accuracy], feed_dict={X: x_data, Y: y_data})
25     print("\nModel: ", h, "\nCorrect: ", c, "\nAccuracy: ", a)
```



Model, Cost, Train



9985 0.1493395
9986 0.14932828
9987 0.14931704
9988 0.14930585
9989 0.14929464
9990 0.1492834
9991 0.14927219
9992 0.14926098
9993 0.14924978
9994 0.14923854
9995 0.14922734
9996 0.14921615
9997 0.14920495
9998 0.14919376
9999 0.1491826
10000 0.1491714

Model: [[0.0306041]
[0.15866211]
[0.30421484]
[0.7816807]
[0.93976283]
[0.9802319]]
Correct: [[0.]
[0.]
[1.]
[1.]
[1.]]
Accuracy: 1.0

```
1 import tensorflow as tf
2
3 x_data = [[1,2], [2,3], [3,1], [4,3], [5,3], [6,2]]
4 y_data = [[0], [0], [0], [1], [1], [1]]
5 X = tf.placeholder(tf.float32, shape=[None, 2])
6 Y = tf.placeholder(tf.float32, shape=[None, 1])
7 W = tf.Variable(tf.random_normal([2,1]))
8 b = tf.Variable(tf.random_normal([1]))
9
10 model = tf.sigmoid(tf.add(tf.matmul(X,W),b))
11 cost = tf.reduce_mean((-1)*Y*tf.log(model) + (-1)*(1-Y)*tf.log(1-model))
12 train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
13
14 prediction = tf.cast(model > 0.5, dtype=tf.float32)
15 accuracy = tf.reduce_mean(tf.cast(tf.equal(prediction, Y), dtype=tf.float32))
16
17 with tf.Session() as sess:
18     sess.run(tf.global_variables_initializer())
19     # Training
20     for step in range(10001):
21         cost_val, train_val = sess.run([cost, train], feed_dict={X: x_data, Y: y_data})
22         print(step, cost_val)
23     # Testing
24     h, c, a = sess.run([model, prediction, accuracy], feed_dict={X: x_data, Y: y_data})
25     print("\nModel: ", h, "\nCorrect: ", c, "\nAccuracy: ", a)
```



Artificial Intelligence and Mobility Lab

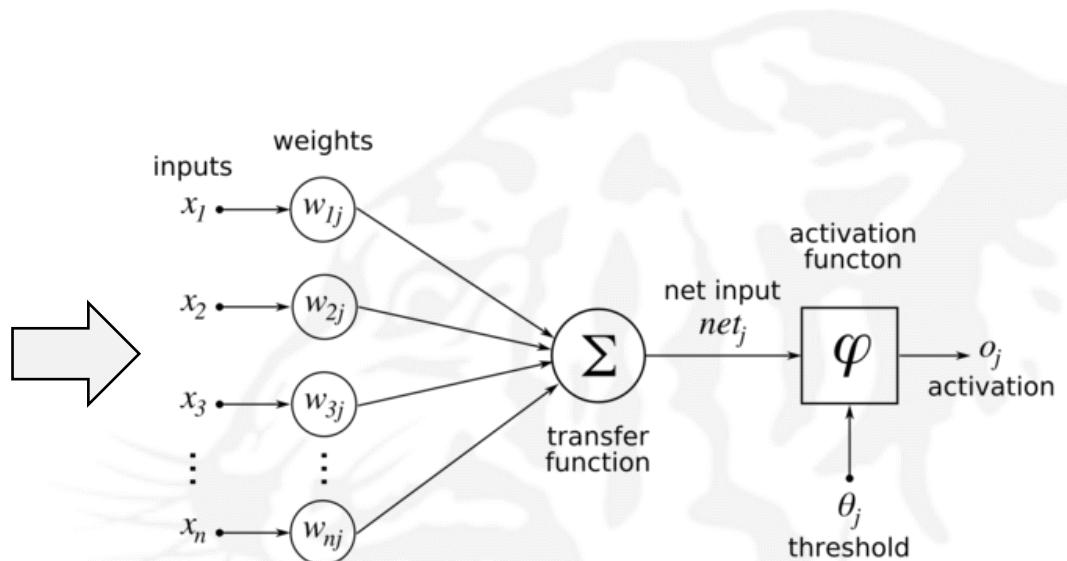
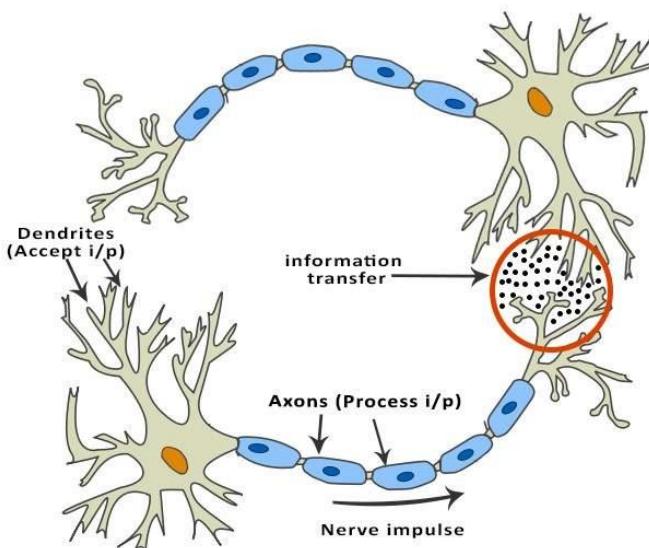
GAN and Reinforcement Learning

Introduction

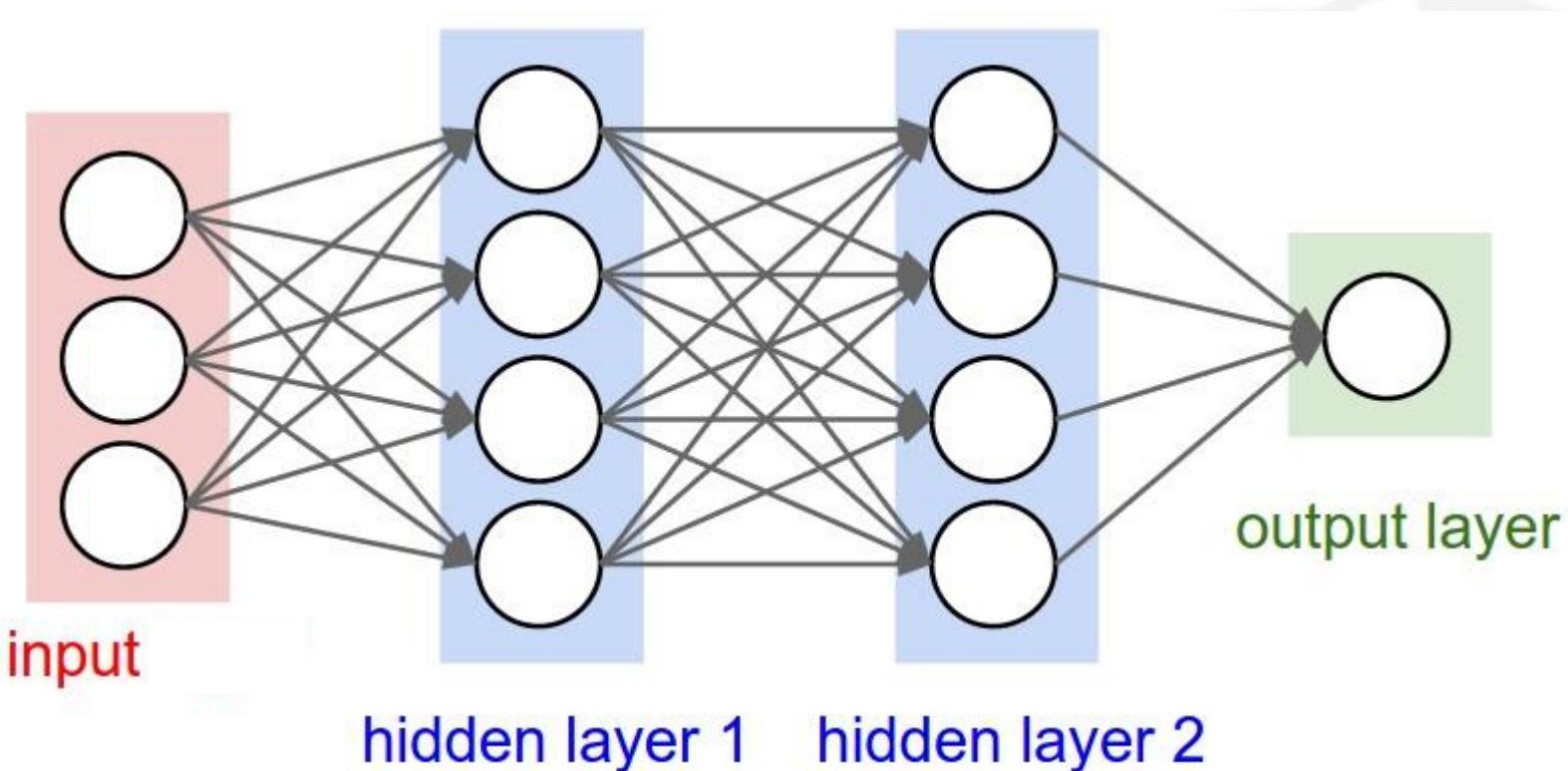
Neural Network Basics

- Deep Learning Overview
- Linear Regression
- **Binary Classification**
- Neural Network Basics

- Human Brain (Neuron)

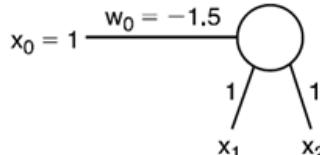


Binary Classification

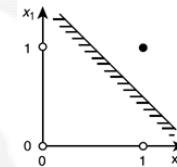


- Application to Logic Gate Design

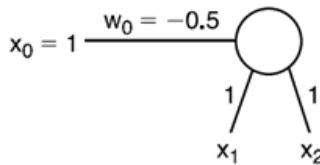
AND
gate



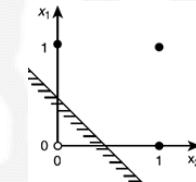
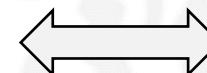
x_1	x_2	$W \cdot X$	y
0	0	-1.5	0
0	1	-0.5	0
1	0	-0.5	0
1	1	0.5	1



OR
gate

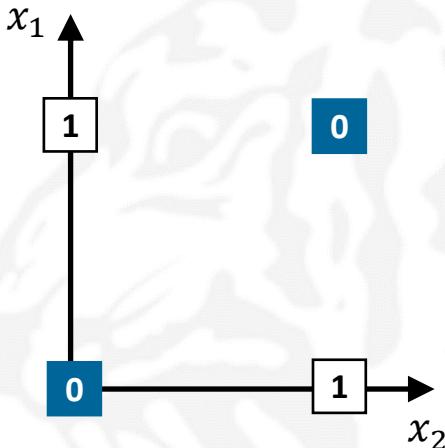


x_1	x_2	$W \cdot X$	y
0	0	-0.5	0
0	1	0.5	1
1	0	0.5	1
1	1	1.5	1



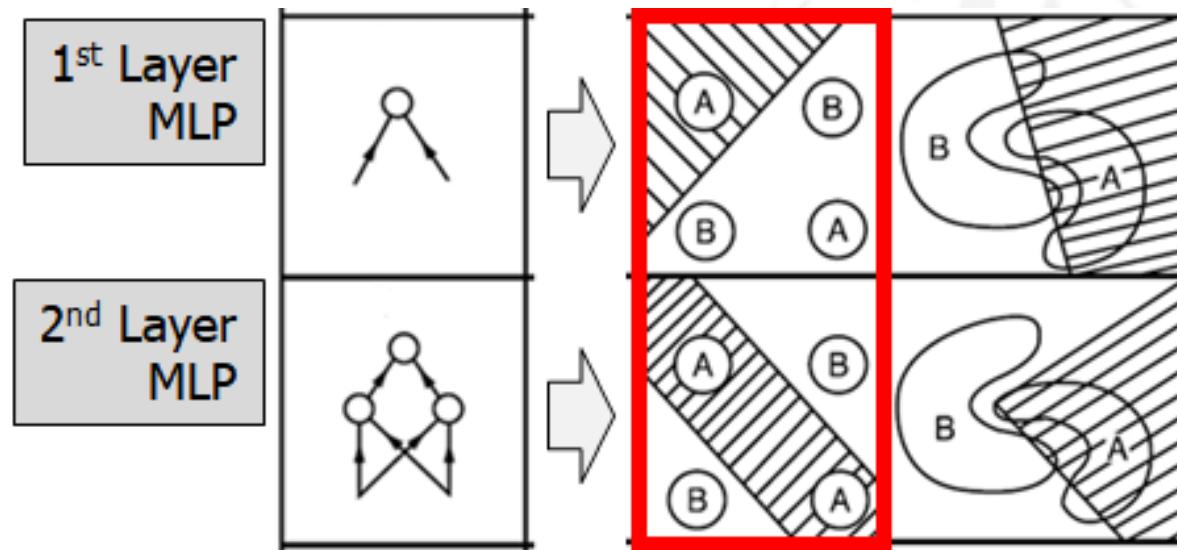
- What about XOR?

x_1	x_2	XOR
0	0	0
0	1	1
1	0	1
1	1	0

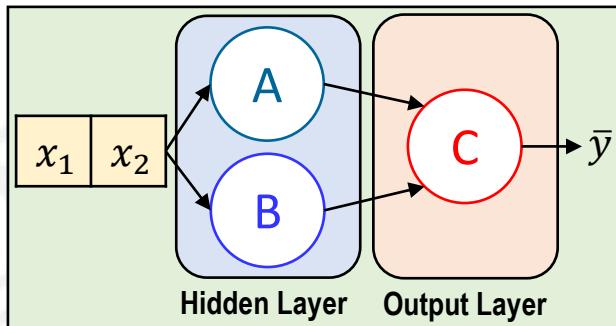
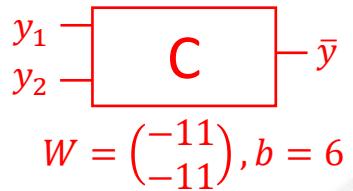
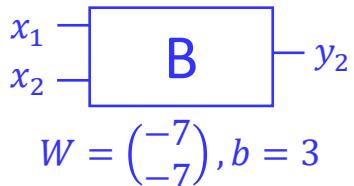
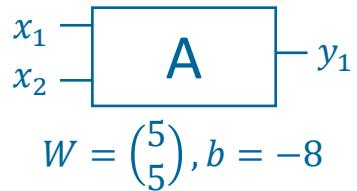


Mathematically proven by
Prof. Marvin Minsky at MIT (1969)

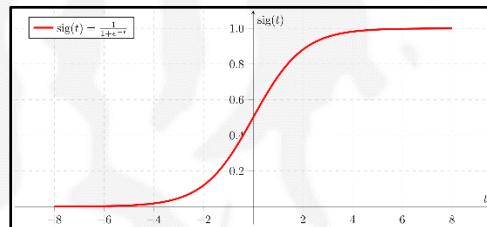
- Multilayer Perceptron (MLP)
 - Proposed by Prof. Marvin Minsky at MIT (1969)
 - Can solve XOR Problem



ANN: Solving XOR with MLP

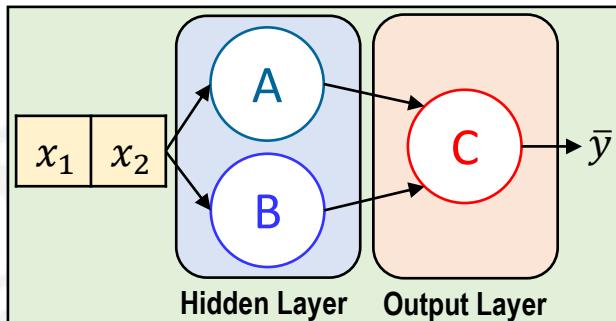
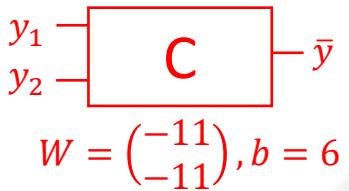
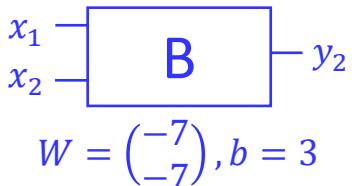
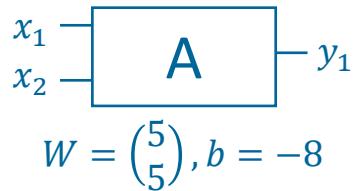


- $(x_1 \ x_2) = (0 \ 0)$
 - $(0 \ 0) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -8$, i.e., $y_1 = \text{Sigmoid}(-8) \cong 0$
 - $(0 \ 0) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = 3$, i.e., $y_2 = \text{Sigmoid}(3) \cong 1$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = -11 + 6 = -5$, i.e., $\bar{y} = \text{Sigmoid}(-5) \cong 0$

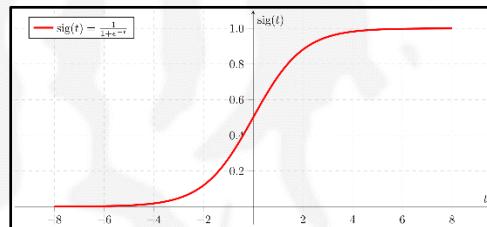


x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1				1
1	0				1
1	1				0

ANN: Solving XOR with MLP

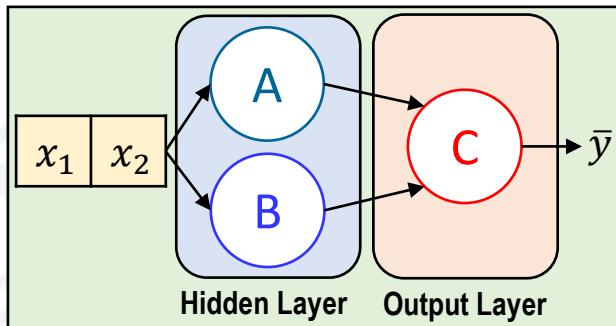
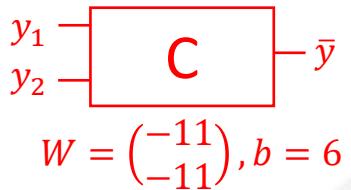
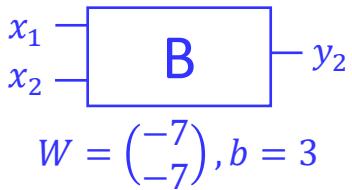
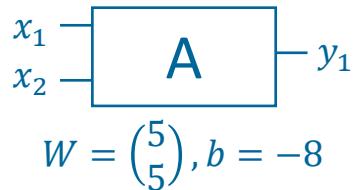


- $(x_1 \ x_2) = (0 \ 1)$
 - $(0 \ 1) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -3$, i.e., $y_1 = Sigmoid(-3) \cong 0$
 - $(0 \ 1) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = -4$, i.e., $y_2 = Sigmoid(-4) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = 6$, i.e., $\bar{y} = Sigmoid(6) \cong 1$

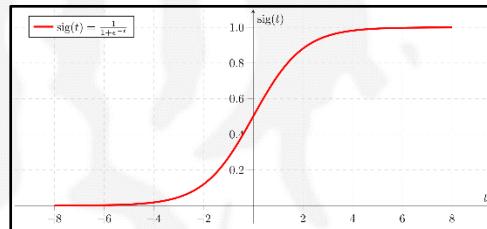


x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0				1
1	1				0

ANN: Solving XOR with MLP

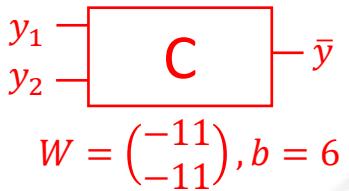
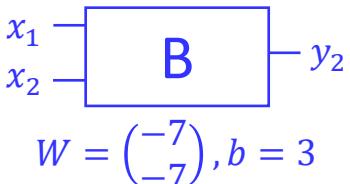
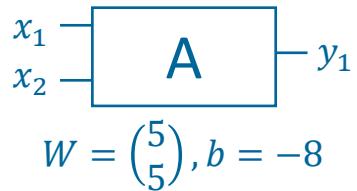


- $(x_1 \ x_2) = (1 \ 0)$
 - $(1 \ 0) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = -3$, i.e., $y_1 = Sigmoid(-3) \cong 0$
 - $(1 \ 0) \begin{pmatrix} -7 \\ 7 \end{pmatrix} + (3) = -4$, i.e., $y_2 = Sigmoid(-4) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ 11 \end{pmatrix} + (6) = 6$, i.e., $\bar{y} = Sigmoid(6) \cong 1$

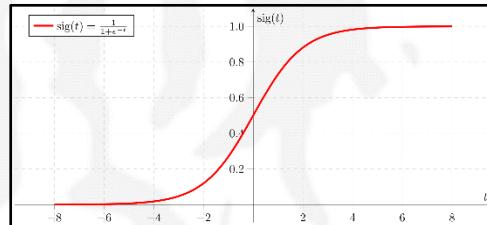
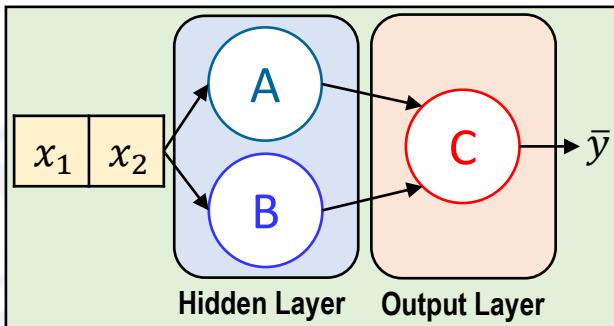


x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0	0	0	1	1
1	1				0

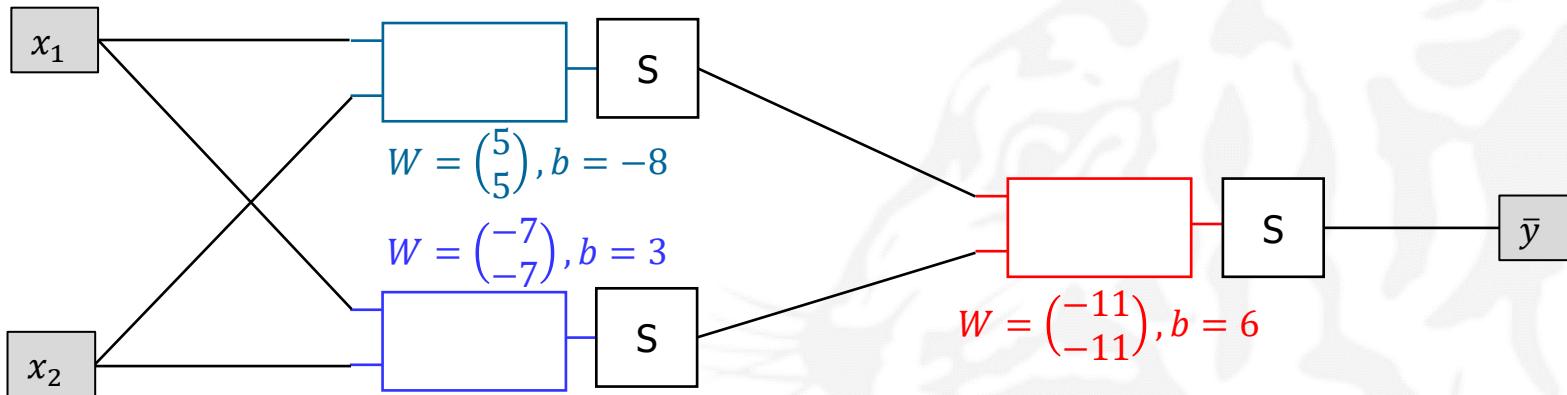
ANN: Solving XOR with MLP



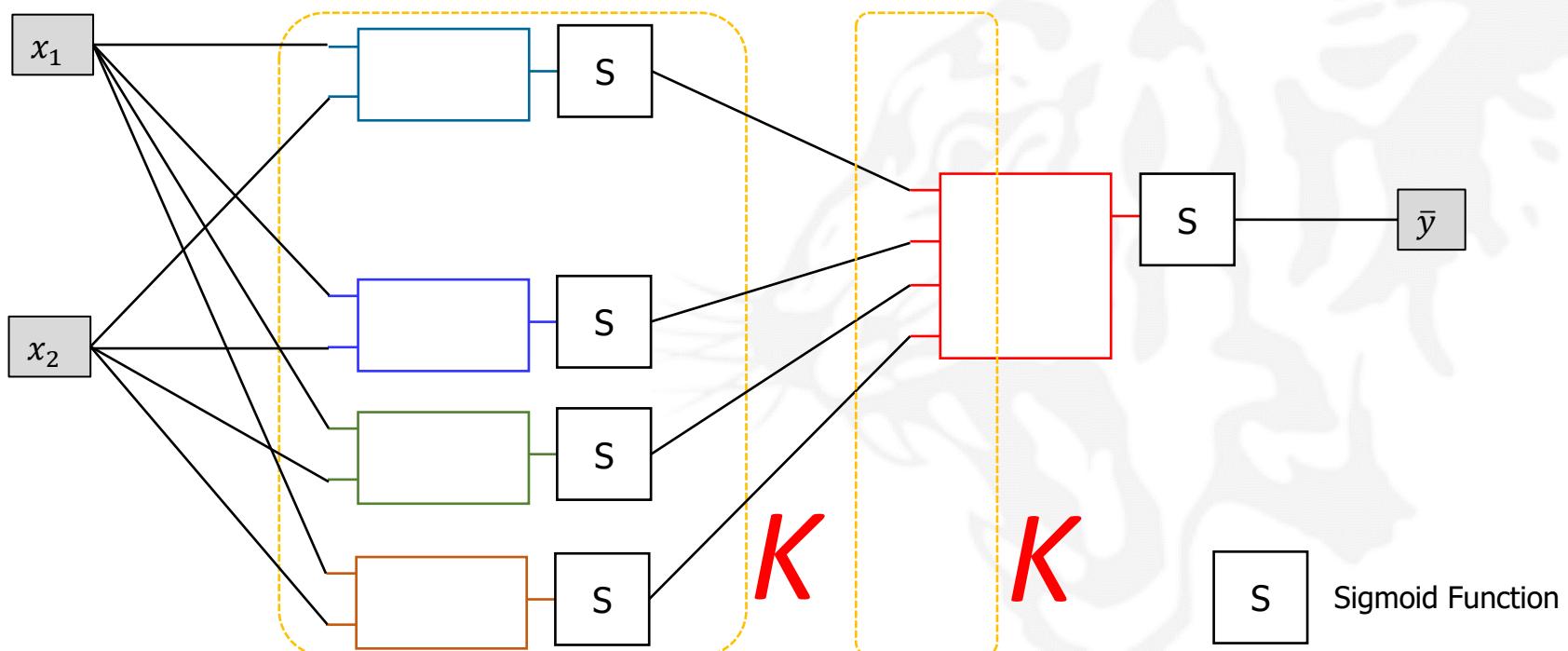
- $(x_1 \ x_2) = (1 \ 1)$
 - $(1 \ 1) \begin{pmatrix} 5 \\ 5 \end{pmatrix} + (-8) = 2$, i.e., $y_1 = Sigmoid(2) \cong 1$
 - $(1 \ 1) \begin{pmatrix} -7 \\ -7 \end{pmatrix} + (3) = -11$, i.e., $y_2 = Sigmoid(-11) \cong 0$
 - $(y_1 \ y_2) \begin{pmatrix} -11 \\ -11 \end{pmatrix} + (6) = -5$, i.e., $\bar{y} = Sigmoid(-5) \cong 0$



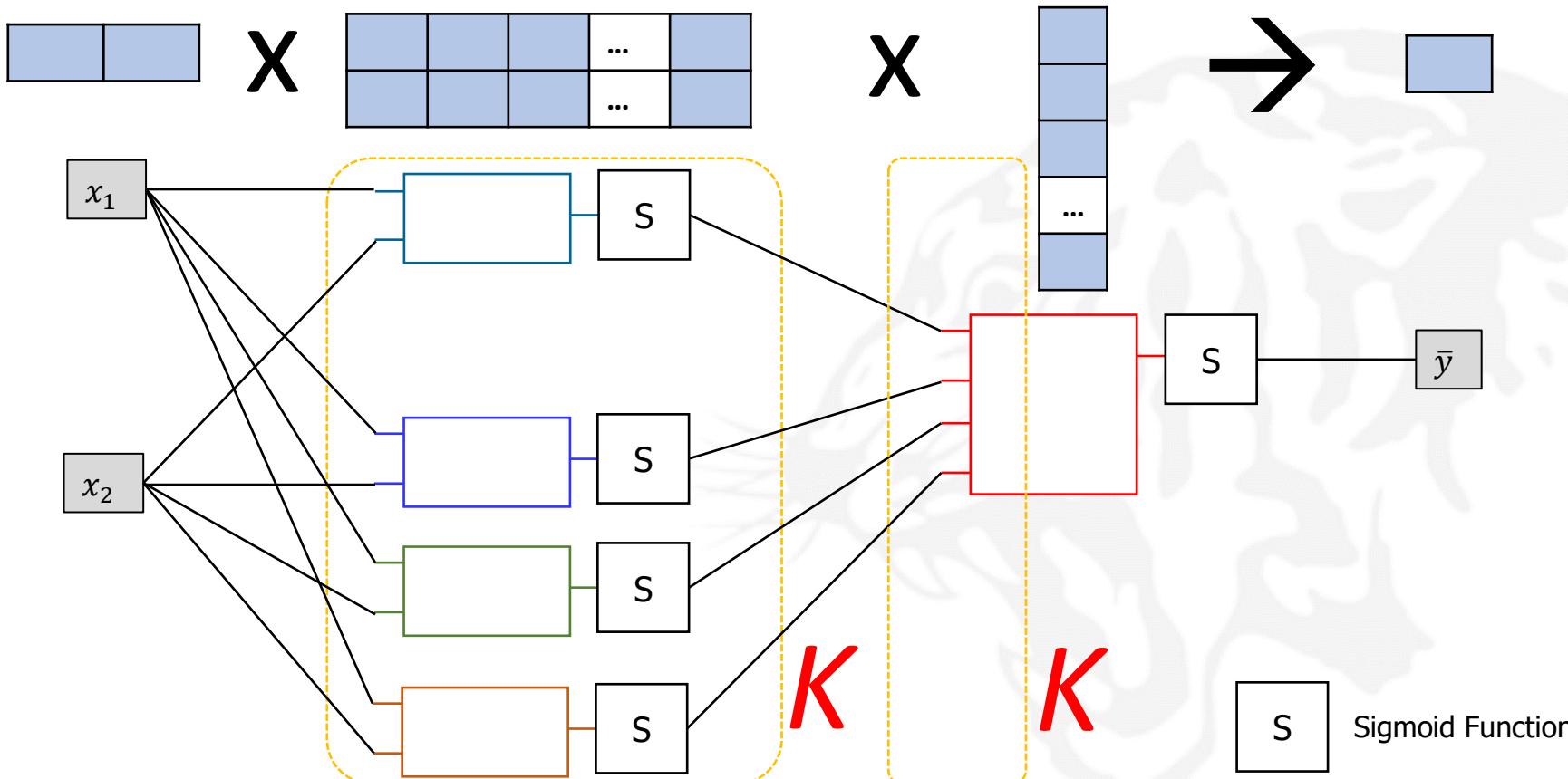
x_1	x_2	y_1	y_2	\bar{y}	XOR
0	0	0	1	0	0
0	1	0	0	1	1
1	0	0	0	1	1
1	1	1	0	0	0



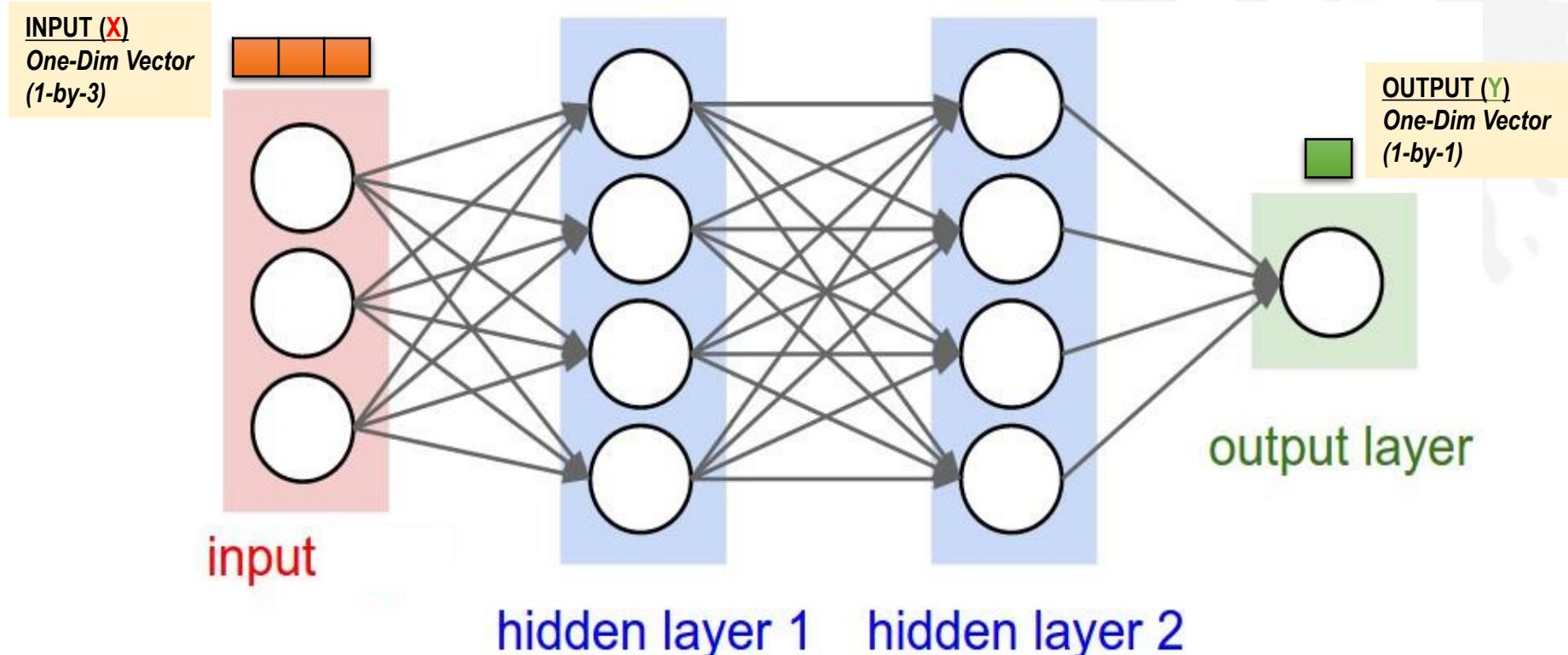
 S Sigmoid Function



ANN: Solving XOR with MLP (Forward Propagation)



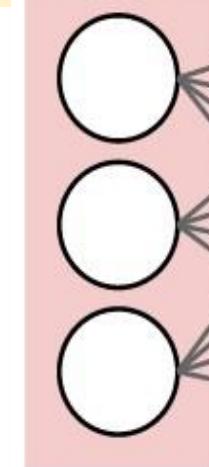
- Toy Model



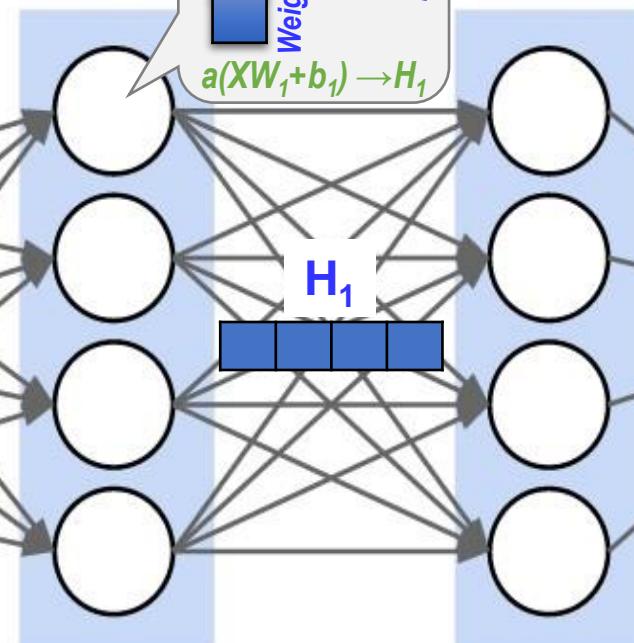
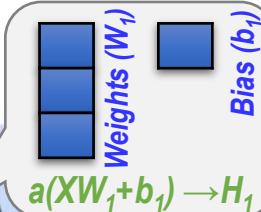
- Toy Model

INPUT (X)

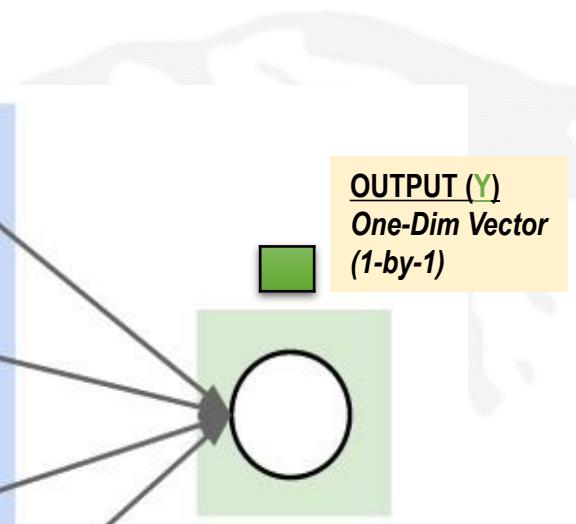
One-Dim Vector
(1-by-3)



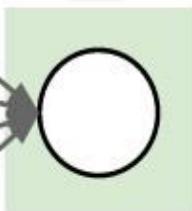
input



H_1



OUTPUT (Y)
One-Dim Vector
(1-by-1)

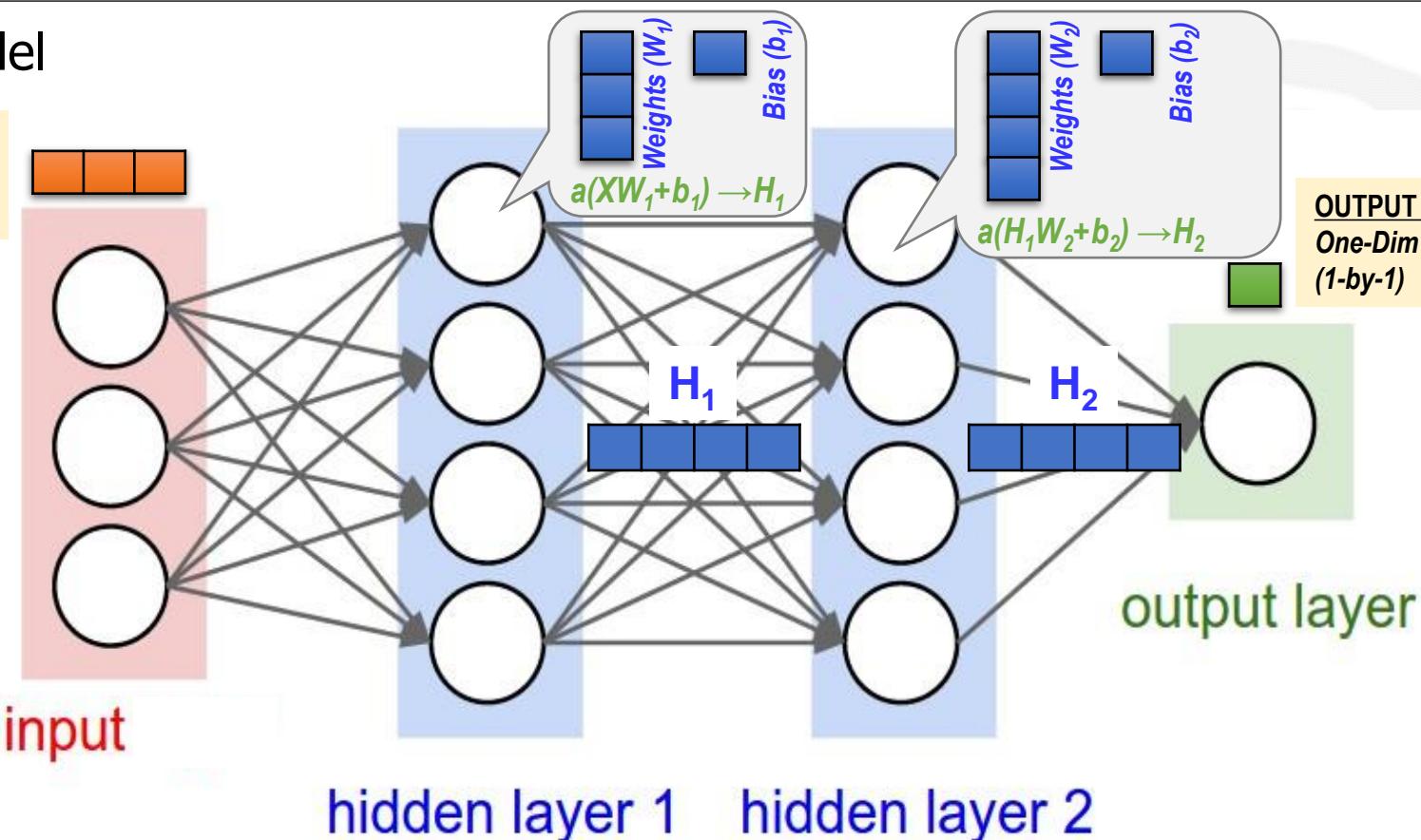


output layer

Conventional Deep Neural Network Training and Inference

- Toy Model

INPUT (X)
One-Dim Vector
(1-by-3)



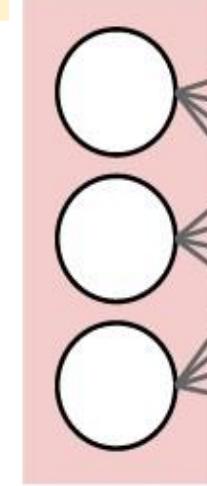
Conventional Deep Neural Network Training and Inference



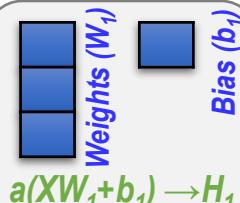
- Toy Model

INPUT (X)

One-Dim Vector
(1-by-3)

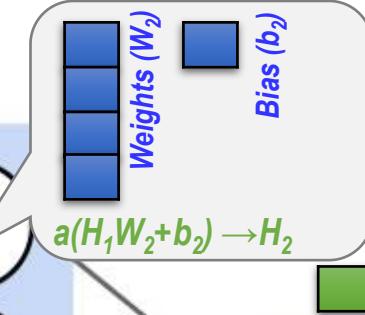


input



$$a(XW_1 + b_1) \rightarrow H_1$$

H_1



$$a(H_1 W_2 + b_2) \rightarrow H_2$$

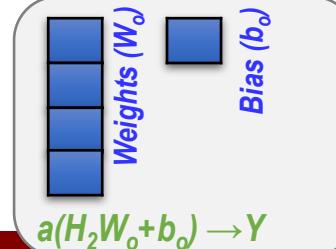
H_2

OUTPUT (Y)

One-Dim Vector
(1-by-1)



output layer



$$a(H_2 W_0 + b_0) \rightarrow Y$$

hidden layer 1 hidden layer 2

TensorFlow for ANN (XOR with Binary Classification)



A	0 1.0826586
M	2000 0.6931472
	4000 0.6931472
	6000 0.6931472
	8000 0.6931472
	10000 0.6931472
	12000 0.6931472
	14000 0.6931472
	16000 0.6931472
	18000 0.6931472
	20000 0.6931472
	[0.5]
	[0.5]
	[0.5]
	[0.5]
	[0.]
	[0.]
	[0.]
	0.5

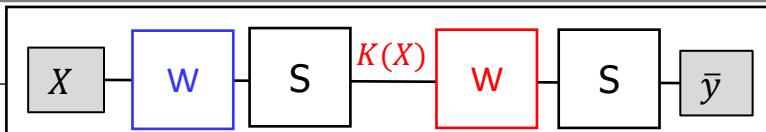
```
1 import tensorflow as tf
2
3 x_data = [[0,0], [0,1], [1,0], [1,1]]
4 y_data = [[0], [1], [1], [0]]
5 X = tf.placeholder(tf.float32, shape=[None, 2])
6 Y = tf.placeholder(tf.float32, shape=[None, 1])
7 W = tf.Variable(tf.random_normal([2,1]))
8 b = tf.Variable(tf.random_normal([1]))
9 model = tf.sigmoid(tf.matmul(X,W)+b)
10 cost = tf.reduce_mean((-1)*Y*tf.log(model) + (-1)*(1-Y)*tf.log(1-model))
11 train = tf.train.GradientDescentOptimizer(0.1).minimize(cost)
12
13 prediction = tf.cast(model > 0.5, dtype=tf.float32)
14 accuracy = tf.reduce_mean(tf.cast(tf.equal(prediction, Y), dtype=tf.float32))
15
16 with tf.Session() as sess:
17     sess.run(tf.global_variables_initializer())
18     # Training
19     for step in range(20001):
20         c, _ = sess.run([cost, train], feed_dict={X: x_data, Y: y_data})
21         if step % 2000 == 0:
22             print(step, c)
23     # Testing
24     m, p, a = sess.run([model, prediction, accuracy], feed_dict={X: x_data, Y: y_data})
25     print(m,p,a)
```

TensorFlow for ANN (XOR with ANN)

```

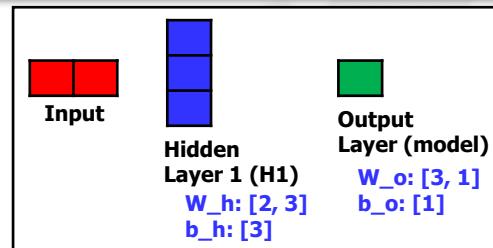
1 import tensorflow as tf
2
3 x_data = [[0,0], [0,1], [1,0], [1,1]]
4 y_data = [[0], [1], [1], [0]]
5 X = tf.placeholder(tf.float32, shape=[None, 2])
6 Y = tf.placeholder(tf.float32, shape=[None, 1])
7 W_h = tf.Variable(tf.random_normal([2,3]))
8 b_h = tf.Variable(tf.random_normal([3]))
9 H1 = tf.sigmoid(tf.matmul(X,W_h)+b_h)
10 W_o = tf.Variable(tf.random_normal([3,1]))
11 b_o = tf.Variable(tf.random_normal([1]))
12 model = tf.sigmoid(tf.matmul(H1,W_o)+b_o)
13 cost = tf.reduce_mean((-1)*Y*tf.log(model) + (-1)*(1-Y)*tf.log(1-model))
14 train = tf.train.GradientDescentOptimizer(0.1).minimize(cost)
15
16 prediction = tf.cast(model > 0.5, dtype=tf.float32)
17 accuracy = tf.reduce_mean(tf.cast(tf.equal(prediction, Y), dtype=tf.float32))
18
19 with tf.Session() as sess:
20     sess.run(tf.global_variables_initializer())
21     # Training
22     for step in range(20001):
23         c, _ = sess.run([cost, train], feed_dict={X: x_data, Y: y_data})
24         if step % 2000 == 0:
25             print(step, c)
26     # Testing
27     m, p, a = sess.run([model, prediction, accuracy], feed_dict={X: x_data, Y: y_data})
28     print(m,p,a)

```



$$K(X) = \text{Sigmoid}(XW_1 + B_1)$$

$$\bar{Y} = H(X) = \text{Sigmoid}(K(X)W + b)$$



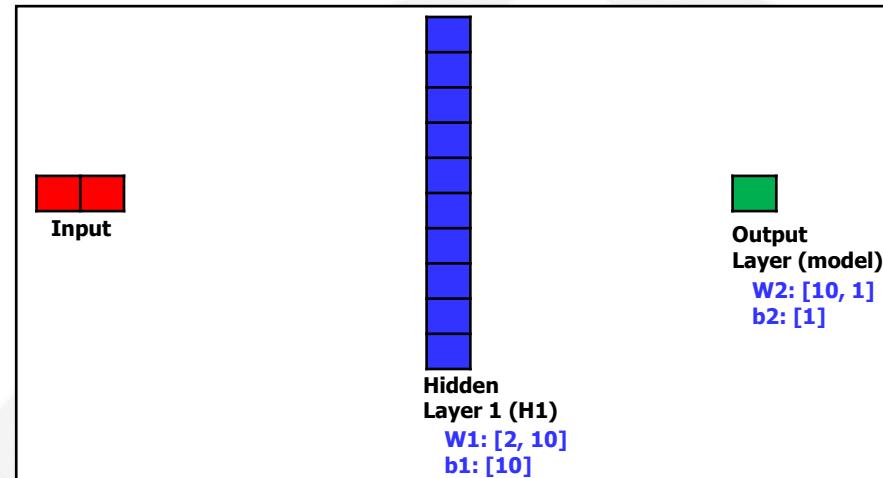
```

0 0.86801255
2000 0.27334553
4000 0.046823796
6000 0.02246793
8000 0.0143708335
10000 0.010443565
12000 0.008153362
14000 0.0066633224
16000 0.0056207716
18000 0.0048525333
20000 0.004264111
[[0.00480547]
 [0.99502313]
 [0.9966924 ]
 [0.00392833]] [[0.]
 [1.]
 [1.]
 [0.]] 1.0

```

- **Wide ANN for XOR**

```
W1 = tf.Variable(tf.random_normal([2, 10]))  
b1 = tf.Variable(tf.random_normal([10]))  
H1 = tf.sigmoid(tf.matmul(X, W1) + b1)  
  
W2 = tf.Variable(tf.random_normal([10, 1]))  
b2 = tf.Variable(tf.random_normal([1]))  
model = tf.sigmoid(tf.matmul(H1, W2) + b2)
```



- Deep ANN for XOR

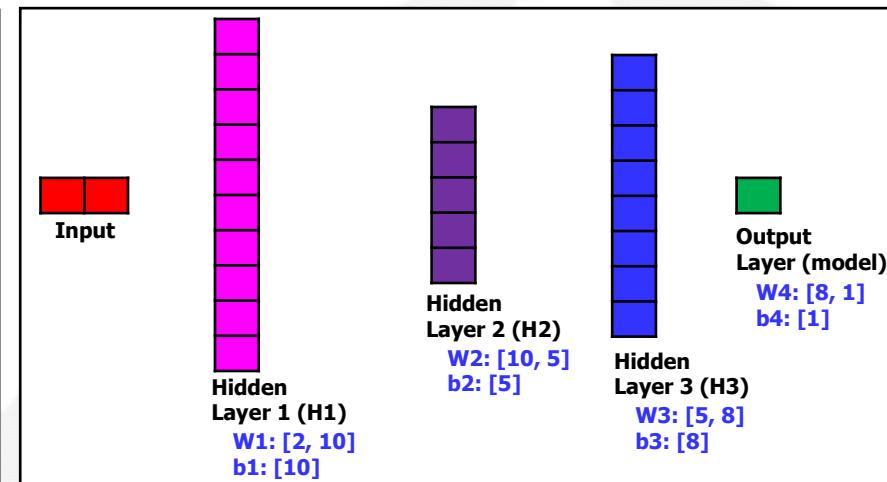
```

W1 = tf.Variable(tf.random_normal([2, 10]))
b1 = tf.Variable(tf.random_normal([10]))
H1 = tf.sigmoid(tf.matmul(X, W1) + b1)

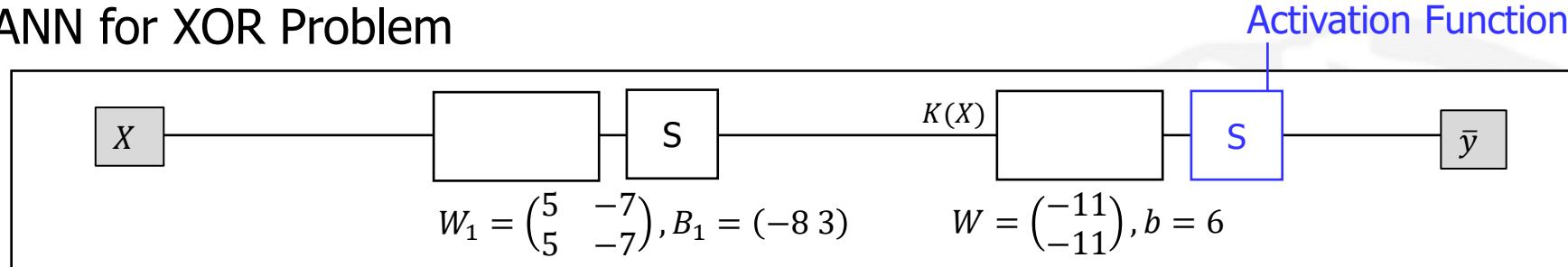
W2 = tf.Variable(tf.random_normal([10, 5]))
b2 = tf.Variable(tf.random_normal([5]))
H2 = tf.sigmoid(tf.matmul(H1, W2) + b2)

W3 = tf.Variable(tf.random_normal([5, 8]))
b3 = tf.Variable(tf.random_normal([8]))
H3 = tf.sigmoid(tf.matmul(H2, W3) + b3)

W4 = tf.Variable(tf.random_normal([8, 1]))
b4 = tf.Variable(tf.random_normal([1]))
model = tf.sigmoid(tf.matmul(H3, W4) + b4)
    
```



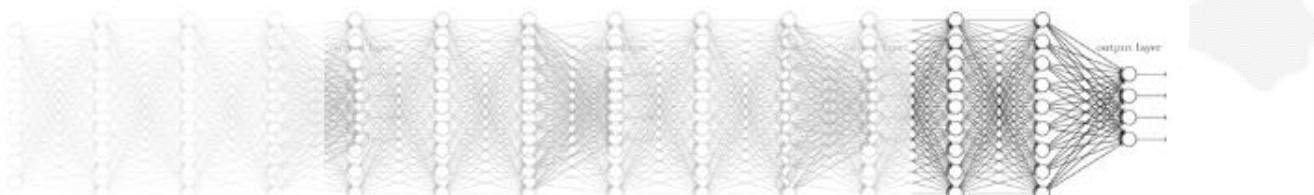
• ANN for XOR Problem



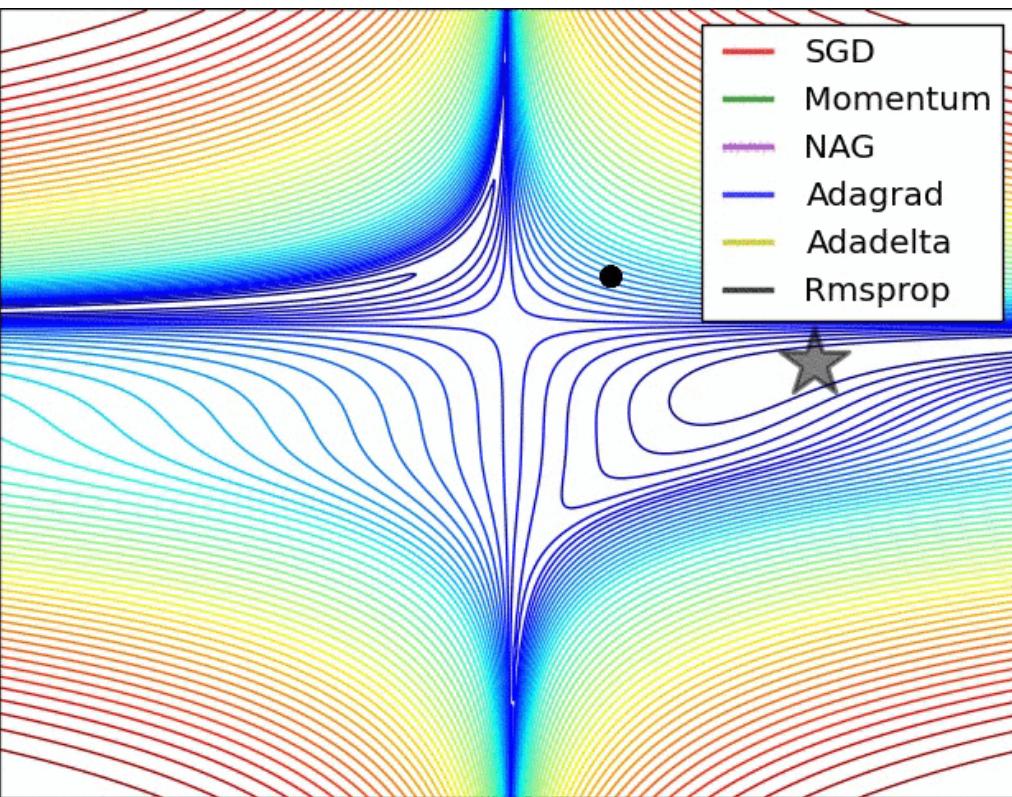
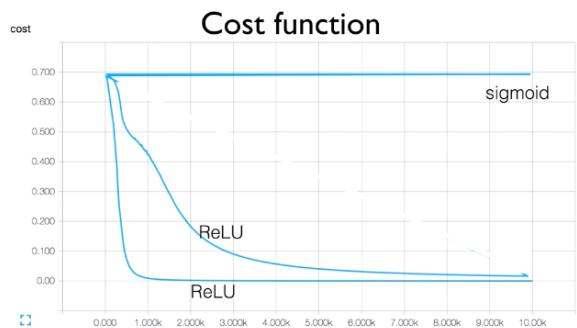
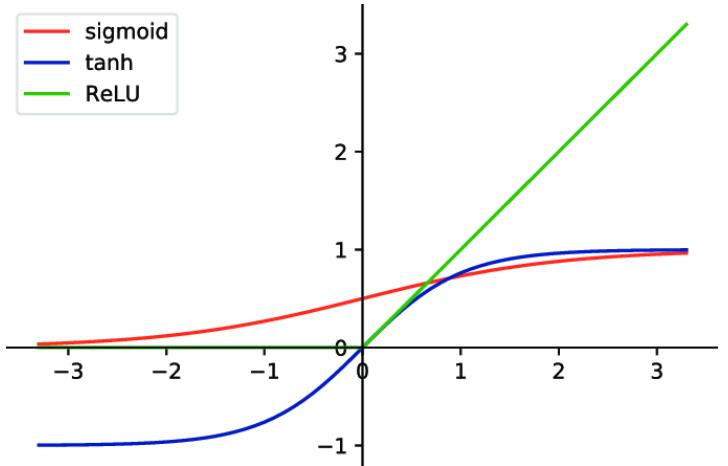
• Observation)

- There exists cases when the accuracy is low even if the # layers is high. Why?
- Answer)

- The result of one ANN is the result of sigmoid function (**between 0 and 1**).
- The numerous multiplication of this result converges to near zero.
→ **Gradient Vanishing Problem**



ANN: ReLU (Rectified Linear Unit)



Outline

- Introduction
 - Theory
 - TensorFlow
- **Generative Adversarial Network**
- Reinforcement Learning
- Deep Reinforcement Learning





Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

Generative Adversarial Networks (GAN)

GAN Theory

- **GAN Theory**
- GAN Implementation

- GAN: Generative Adversarial Network
- Training both of **generator** and **discriminator**; and then generates samples which are similar to the original samples

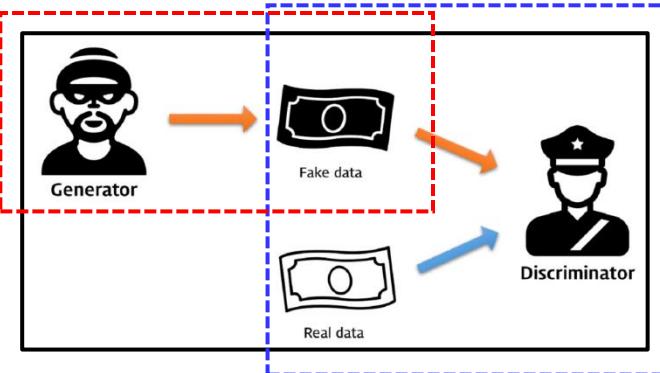


Generators

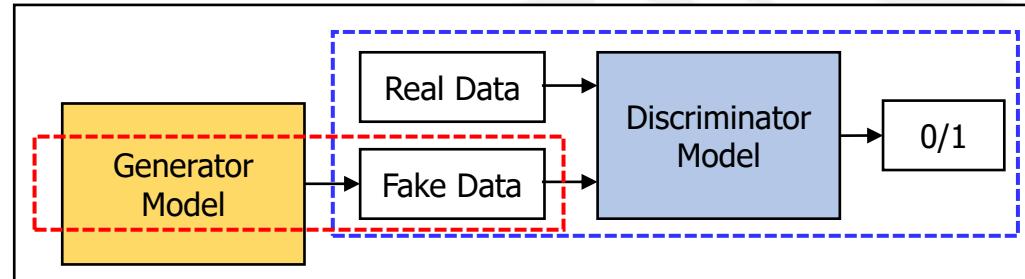
Performance
Improvements via
Competition



Discriminator



GAN architecture



Discriminator Model

- The discriminative model learns **how to classify** input to its class (fake → fake class, real → real class).
- Binary classifier.

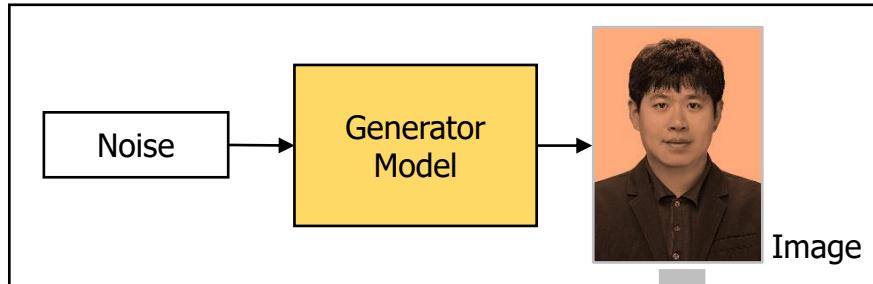
Supervised Learning

Generator Model

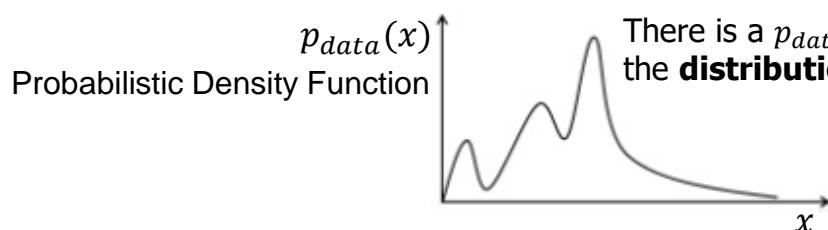
- The generative model learns **the distribution of training data**.

Unsupervised Learning

- Generative Model

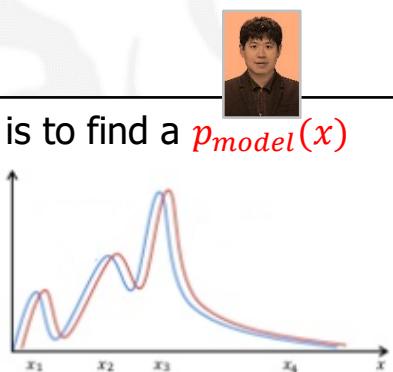


The generator model learns **the distribution of training data**.

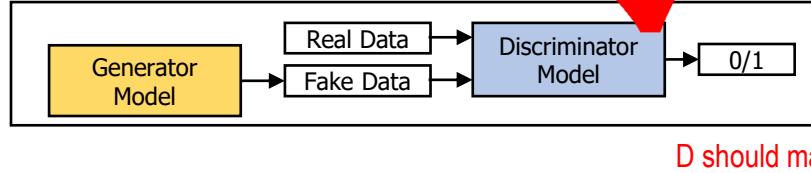


There is a $p_{data}(x)$ that represents
the distribution of actual images (training data).

The goal of the generative model is to find a $p_{model}(x)$
that approximates $p_{data}(x)$ well.



GAN architecture



Objective of D

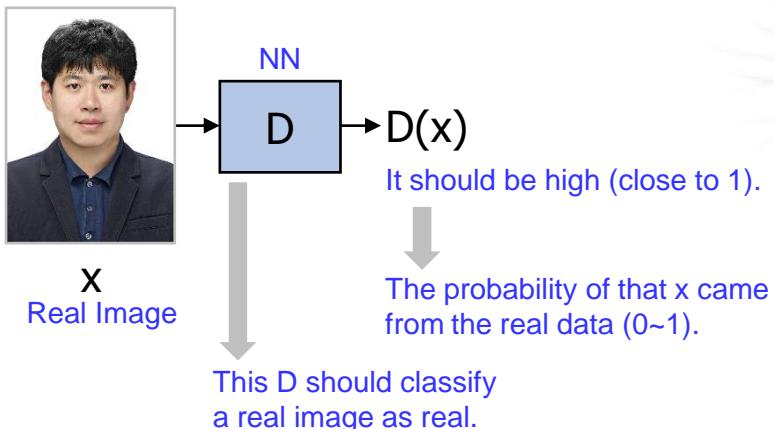
$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

D should maximize $V(D, G)$

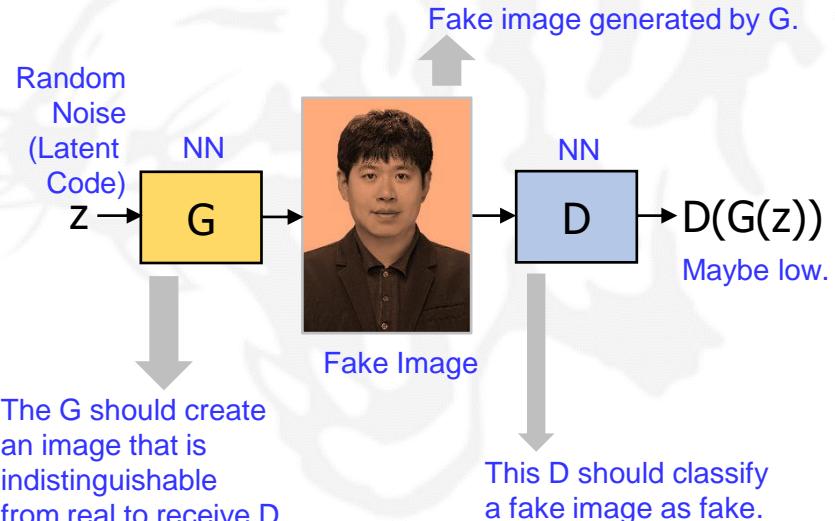
Maximize when $D(x)=1$

Maximize when $D(G(z))=0$

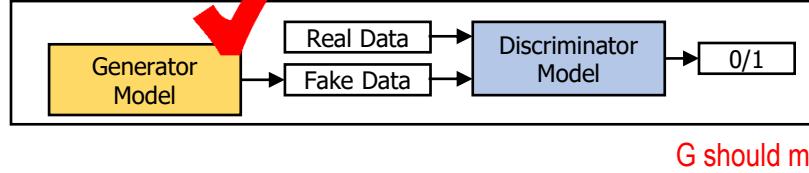
Training with REAL



Training with FAKE



GAN architecture

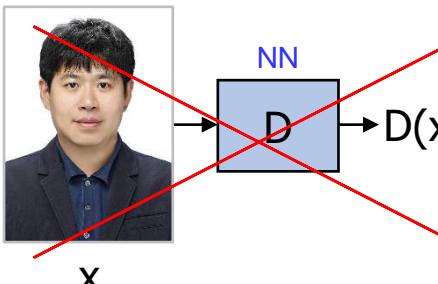


Objective of G

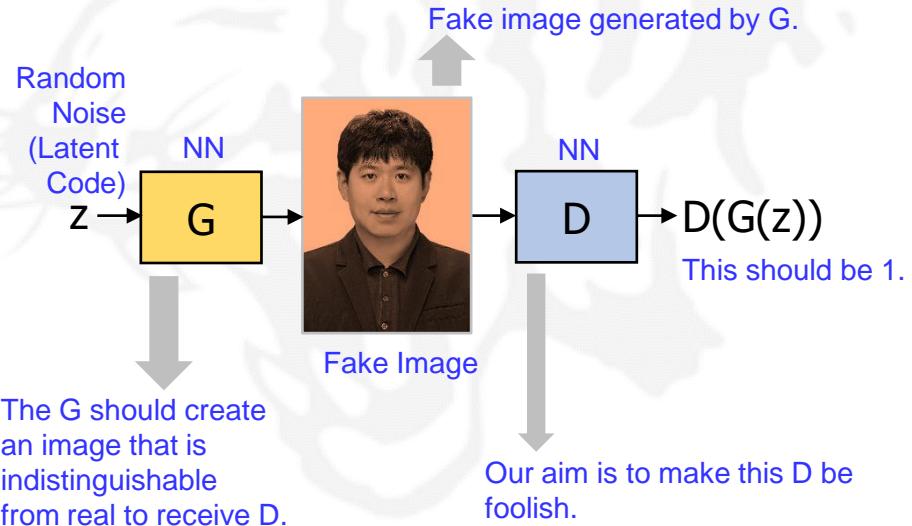
$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

G should minimize $V(D, G)$
 G is independent to this part
 Sample latent code x from Gaussian distribution
 Minimum when $D(G(z))=1$

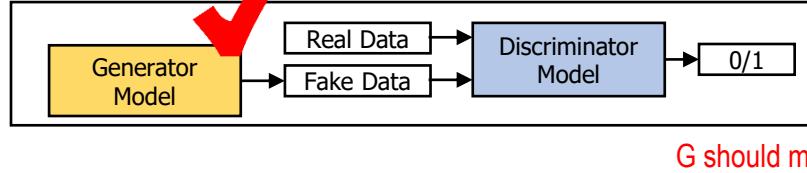
Training with REAL



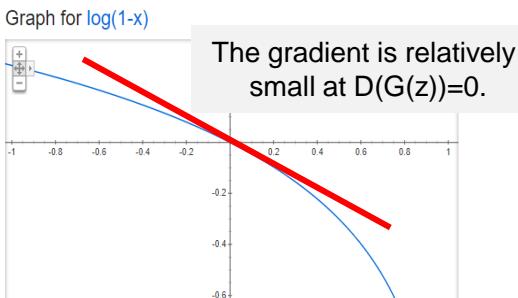
Training with FAKE



GAN architecture



- At the beginning of training, the D can clearly classify the generated image as fake because the quality of the image is very low.
- This means $D(G(z))$ is almost zero at early stage of training.



Objective of G

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

G should minimize $V(D, G)$

G is independent to this part

Sample latent code x from Gaussian distribution

Minimum when $D(G(z))=1$

$$\min_G E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

$$\max_G E_{z \sim p_z(z)}[\log(D(G(z)))]$$



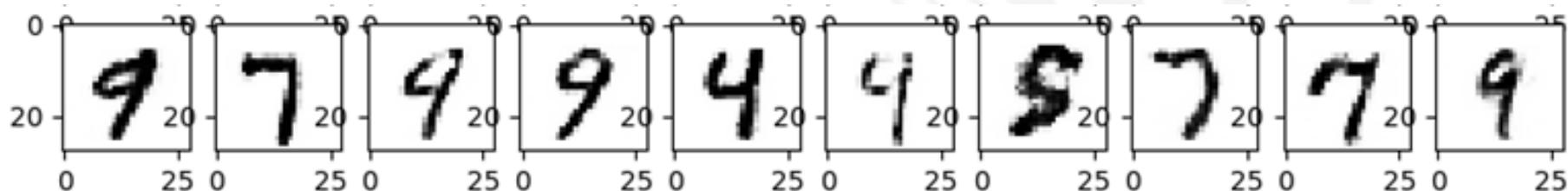
Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

Generative Adversarial Networks (GAN)

GAN Implementation

- GAN Theory
- GAN Implementation



```
1 # MNIST data
2 from tensorflow.examples.tutorials.mnist import input_data
3 mnist = input_data.read_data_sets("data_MNIST", one_hot=True)
4
5 import matplotlib.pyplot as plt
6 import numpy as np
7 import tensorflow as tf
8
9 # Training Params
10 num_steps = 100000
11 batch_size = 128
12
13 # Network Params
14 dim_image = 784 # 28*28 pixels
15 nHL_G = 256
16 nHL_D = 256
17 dim_noise = 100 # Noise data points
18
19 # A custom initialization (Xavier Glorot init)
20 def glorot_init(shape):
21     return tf.random_normal(shape=shape, stddev=1. / tf.sqrt(shape[0] / 2.))
```

```
23 W = {
24     'HL_G' : tf.Variable(glorot_init([dim_noise, nHL_G])),
25     'OL_G' : tf.Variable(glorot_init([nHL_G, dim_image])),
26     'HL_D' : tf.Variable(glorot_init([dim_image, nHL_D])),
27     'OL_D' : tf.Variable(glorot_init([nHL_D, 1])),
28 }
29 b = {
30     'HL_G' : tf.Variable(tf.zeros([nHL_G])),
31     'OL_G' : tf.Variable(tf.zeros([dim_image])),
32     'HL_D' : tf.Variable(tf.zeros([nHL_D])),
33     'OL_D' : tf.Variable(tf.zeros([1])),
34 }
35
36 # Neural Network: Generator
37 def nn_G(x):
38     HL = tf.nn.relu(tf.add(tf.matmul(x, W['HL_G']), b['HL_G']))
39     OL = tf.nn.sigmoid(tf.add(tf.matmul(HL, W['OL_G']), b['OL_G']))
40     return OL
41
42 # Neural Network: Discriminator
43 def nn_D(x):
44     HL = tf.nn.relu(tf.add(tf.matmul(x, W['HL_D']), b['HL_D']))
45     OL = tf.nn.sigmoid(tf.add(tf.matmul(HL, W['OL_D']), b['OL_D']))
46     return OL
47
48 # Network Inputs
49 IN_G = tf.placeholder(tf.float32, shape=[None, dim_noise])
50 IN_D = tf.placeholder(tf.float32, shape=[None, dim_image])
```

```

52 # Build Generator Neural Network
53 sample_G = nn_G(IN_G) --> G(z)
54                                         D(x)           D(G(z))
55 # Build Discriminator Neural Network (one from noise input, one from generated samples)
56 D_real = nn_D(IN_D)
57 D_fake = nn_D(sample_G)
58 vars_G = [W['HL_G'], W['OL_G'], b['HL_G'], b['OL_G']]
59 vars_D = [W['HL_D'], W['OL_D'], b['HL_D'], b['OL_D']]
60
61 # Cost, Train
62 cost_G = -tf.reduce_mean(tf.log(D_fake))
63 cost_D = -tf.reduce_mean(tf.log(D_real) + tf.log(1. - D_fake))
64 train_G = tf.train.AdamOptimizer(0.0002).minimize(cost_G, var_list=vars_G)
65 train_D = tf.train.AdamOptimizer(0.0002).minimize(cost_D, var_list=vars_D)

```

Objective of G $\max_G E_{z \sim p_z(z)}[\log(D(G(z)))]$

Objective of D

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

D should maximize $V(D, G)$

Sample x from real data distribution

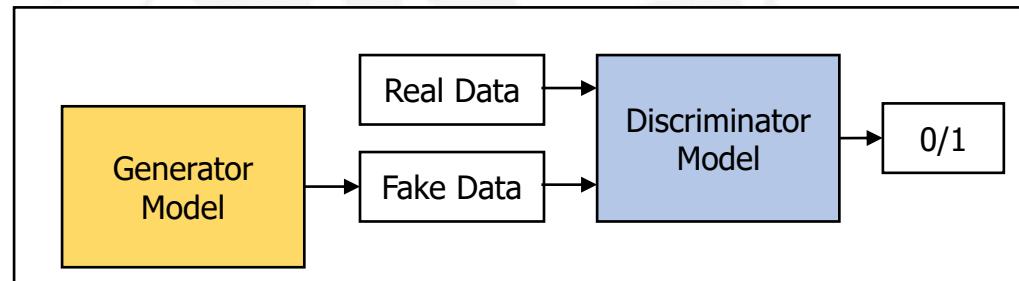
Sample latent code z from Gaussian distribution

Maximize when $D(x)=1$

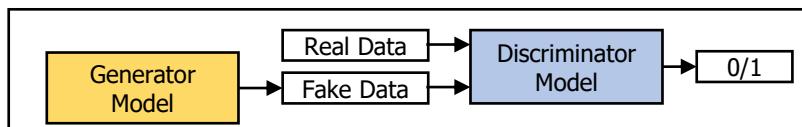
Maximize when $D(G(z))=0$

```
52 # Build Generator Neural Network
53 sample_G = nn_G(IN_G)
54
55 # Build Discriminator Neural Network (one from noise input, one from generated samples)
56 D_real = nn_D(IN_D)
57 D_fake = nn_D(sample_G)
58 vars_G = [W['HL_G'], W['OL_G'], b['HL_G'], b['OL_G']]
59 vars_D = [W['HL_D'], W['OL_D'], b['HL_D'], b['OL_D']]
60
61 # Cost, Train
62 cost_G = -tf.reduce_mean(tf.log(D_fake))
63 cost_D = -tf.reduce_mean(tf.log(D_real) + tf.log(1. - D_fake))
64 train_G = tf.train.AdamOptimizer(0.0002).minimize(cost_G, var_list=vars_G)
65 train_D = tf.train.AdamOptimizer(0.0002).minimize(cost_D, var_list=vars_D)
```

GAN architecture



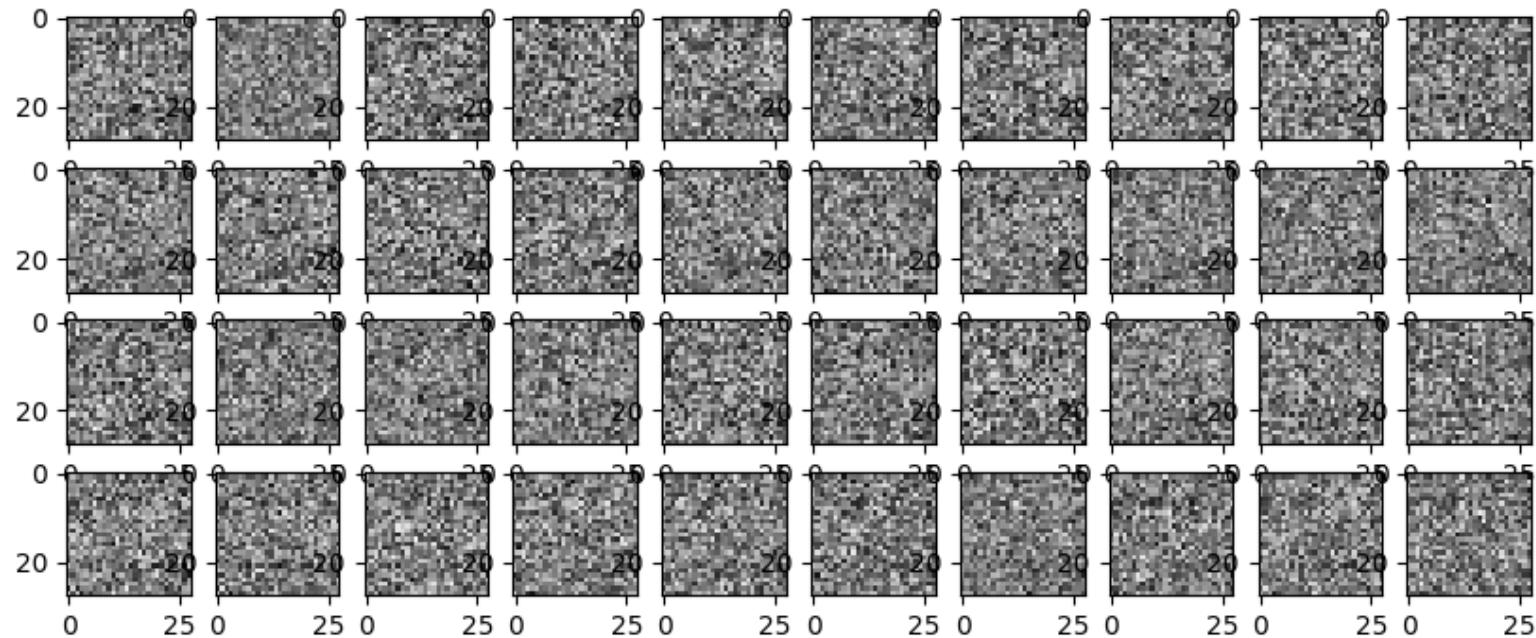
GAN architecture



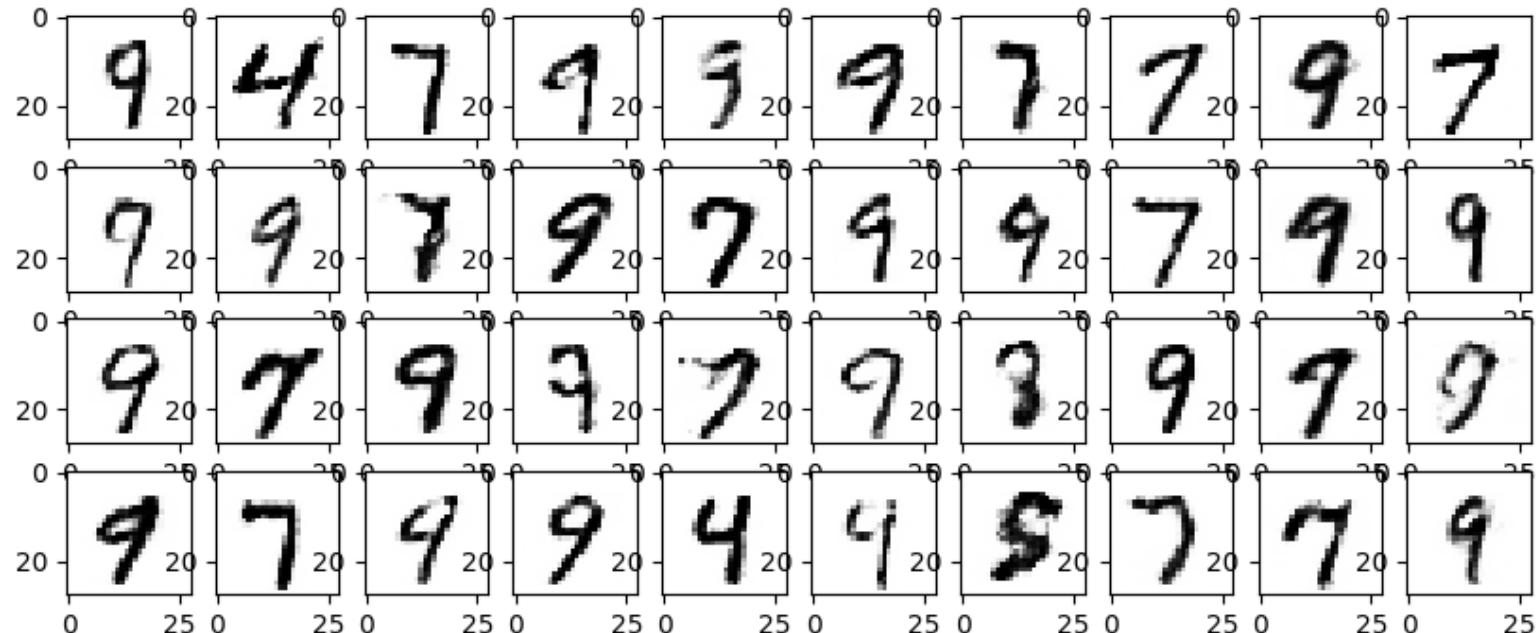
```

67 # Session
68 with tf.Session() as sess:
69     sess.run(tf.global_variables_initializer())
70     for i in range(1, num_steps+1):
71         # Get the next batch of MNIST data
72         batch_images, _ = mnist.train.next_batch(batch_size)
73         # Generate noise to feed to the generator G
74         z = np.random.uniform(-1., 1., size=[batch_size, dim_noise])
75         # Train
76         sess.run([train_G, train_D], feed_dict = {IN_D: batch_images, IN_G: z})
77         f, a = plt.subplots(4, 10, figsize=(10, 4))
78         for i in range(10):
79             z = np.random.uniform(-1., 1., size=[4, dim_noise])
80             g = sess.run([sample_G], feed_dict={IN_G: z})
81             g = np.reshape(g, newshape=(4, 28, 28, 1))
82             # Reverse colors for better display
83             g = -1 * (g - 1)
84             for j in range(4):
85                 # Generate image from noise. Extend to 3 channels for matplotlib figure.
86                 img = np.reshape(np.repeat(g[j][:, :, np.newaxis], 3, axis=2), newshape=(28, 28, 3))
87                 a[j][i].imshow(img)
88         f.show()
89         plt.draw()
90         plt.waitforbuttonpress()
  
```

num_steps: 1



num_steps: 100000



Outline

- Introduction
 - Theory
 - TensorFlow
- Generative Adversarial Network
- **Reinforcement Learning**
- Deep Reinforcement Learning





Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

Reinforcement Learning (RL)

RL Theory

- **RL Theory**
- RL Implementation

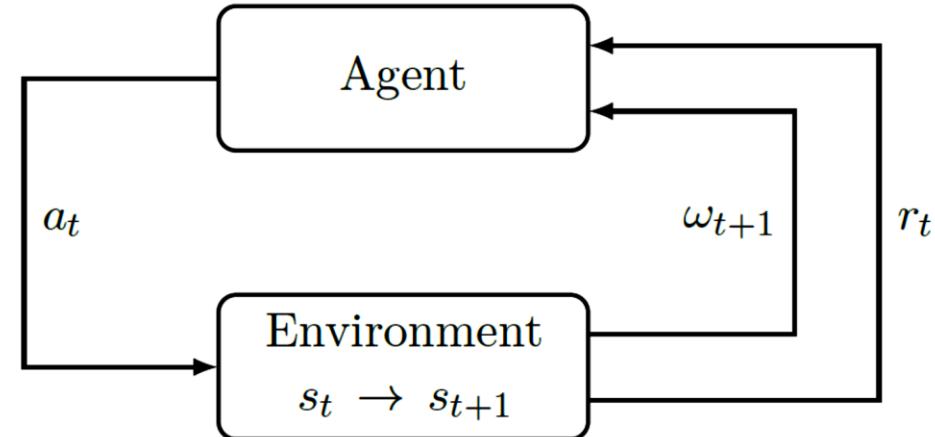
- Action Sequence (also called **Policy**, later in this presentation)!

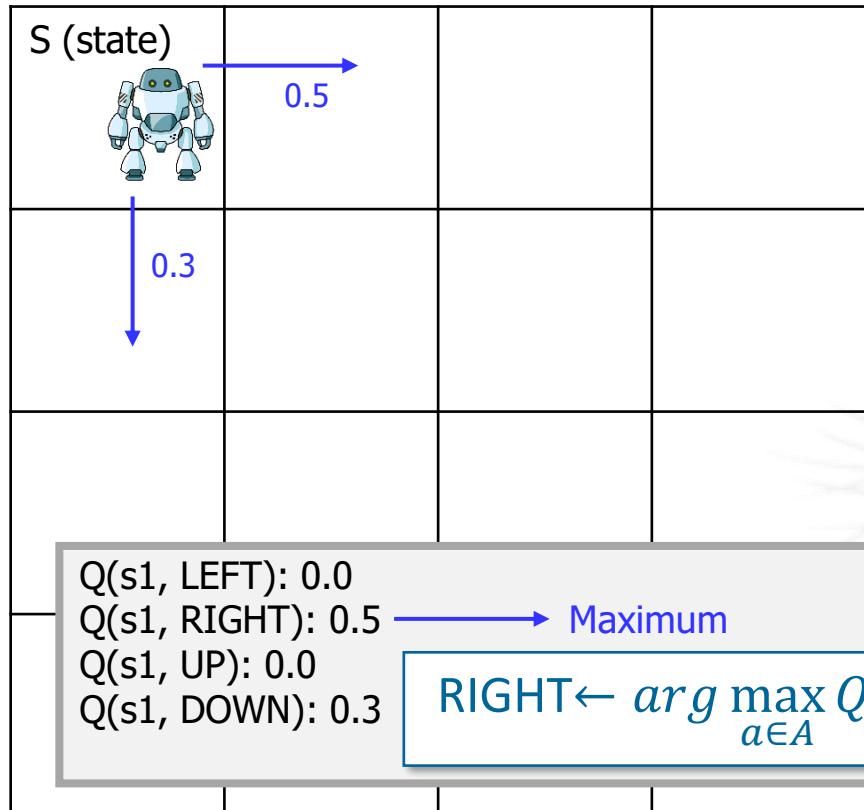


• RL Setting

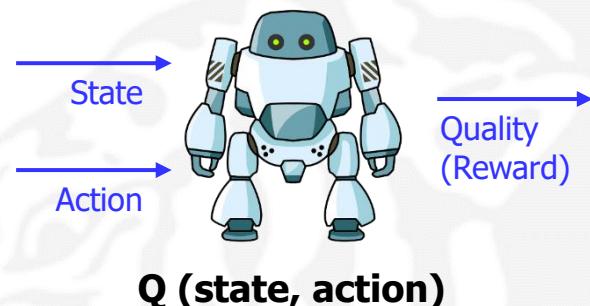
- The general RL problem is formalized as a **discrete time stochastic control process** where **an agent interacts with its environment** as follows:

1. The agent starts in a given state within its environment $s_0 \in S$ by gathering an initial observation $\omega_0 \in \Omega$.
2. At each time step t ,
The agent has to take an action $a_t \in A$.
It follows three consequences:
 - 1) Obtains a reward $r_t \in R$
 - 2) State transitions to $s_{t+1} \in S$
 - 3) Obtains an observation $\omega_{t+1} \in \Omega$





- Q-Function (State-action value)



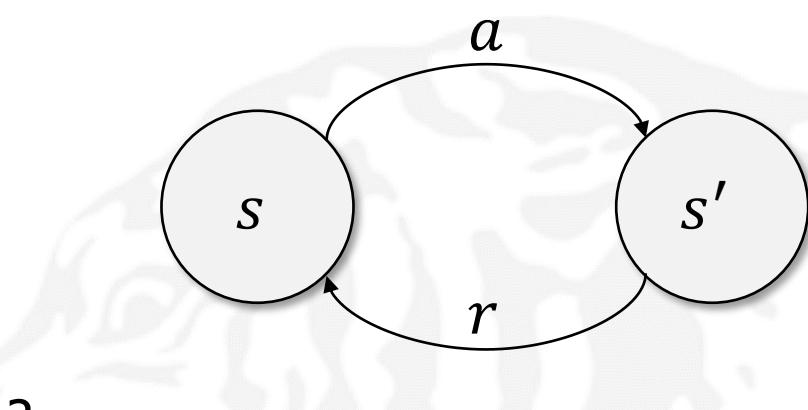
Optimal Policy π and Max Q

- $\text{Max } Q = \max_{a'} Q(s, a')$
- $\pi^*(s) = \arg \max_a Q(s, a)$

- My condition
 - I am now in state s
 - When I do action a , I will go to s' .
 - When I do action a , I will get reward r
 - Q in s' , it means $Q(s', a')$ exists.
- How can we express $Q(s, a)$ using $Q(s', a')$?

$$Q(s, a) = r + \max_{a'} Q(s', a')$$

- **Solution with Dynamic Programming**
 - **Pseudo-polynomial complexity**



Recurrence (e.g., factorial)

$F(x)\{$

```
if (x != 1){ x * F(x-1) }
if (x == 1){ F(x) = 1 }
}
```

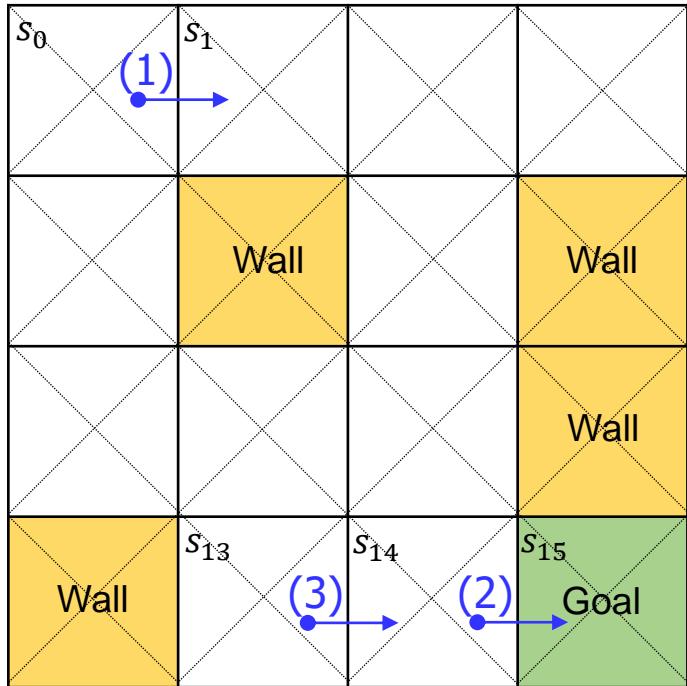
}

$$3! = F(3) = 3 * F(2)$$

$$= 3 * 2 * F(1)$$

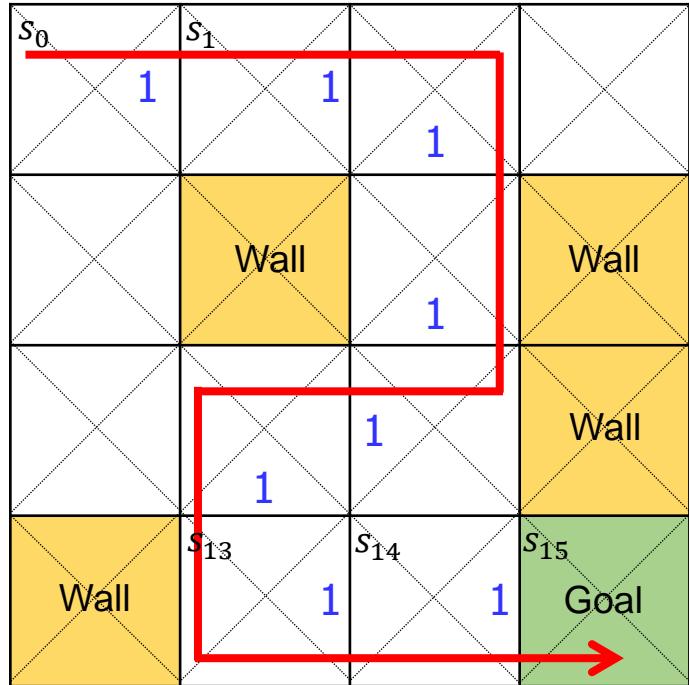
$$= 3 * 2 * 1 = 6$$

- 16 states and 4 actions (U, D, L, R)



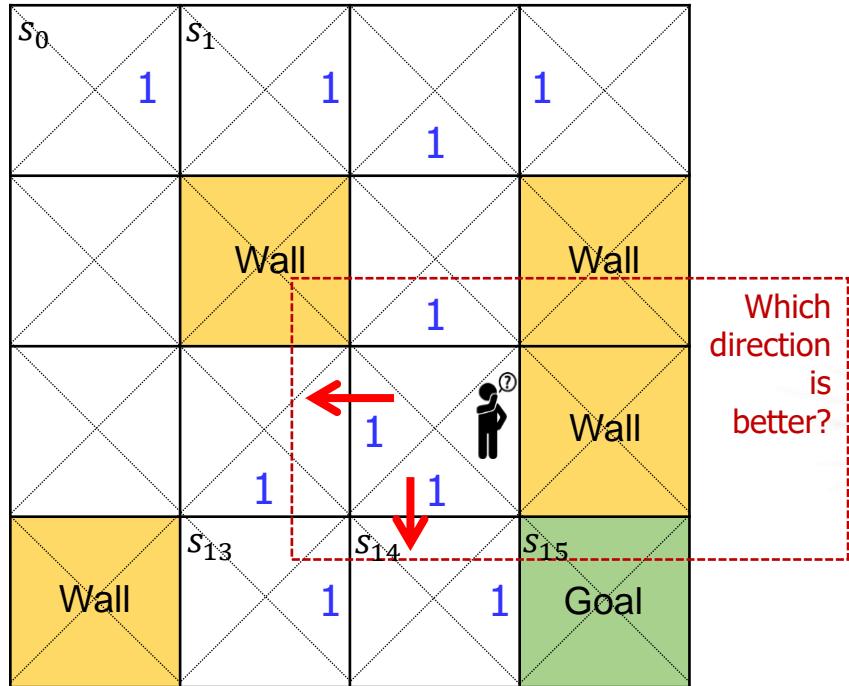
- Initial Status
 - All 64 Q values are 0,
 - Reward are all zero except $r_{s_{15},L} = 1$
- For (1), from s_0 to s_1
 - $Q(s_0, a_R) = r + \max_a Q(s_1, a) = 0 + \max\{0,0,0,0\} = 0$
- For (2), from s_{14} to s_{15} (goal)
 - $Q(s_{14}, a_R) = r + \max_a Q(s_{15}, a) = 1 + \max\{0,0,0,0\} = 1$
- For (3), from s_{13} to s_{14}
 - $Q(s_{13}, a_R) = r + \max_a Q(s_{14}, a) = 0 + \max\{0,0,1,0\} = 1$

- 16 states and 4 actions (U, D, L, R)



Q-Learning

- 16 states and 4 actions (U, D, L, R)



Learning $Q(s, a)$ with Discounted Reward

$$Q(s, a) = r + \gamma \cdot \arg \max_a Q(s', a')$$

$$0 < \gamma \leq 1$$

- For each s, a , initialize table entry $Q(s, a) \leftarrow 0$
- Observe current state s
- Do forever
 - Select an action a and execute it
 - Receive immediate reward r
 - Observe the new state s'
 - Update the table entry for $Q(s, a)$ as follows:

$$Q(s, a) \leftarrow r + \max_{a'} Q(s', a')$$

- $s \leftarrow s'$



Finding the Best Restaurant

- Try the best one during weekdays.
- Try new ones during weekends.



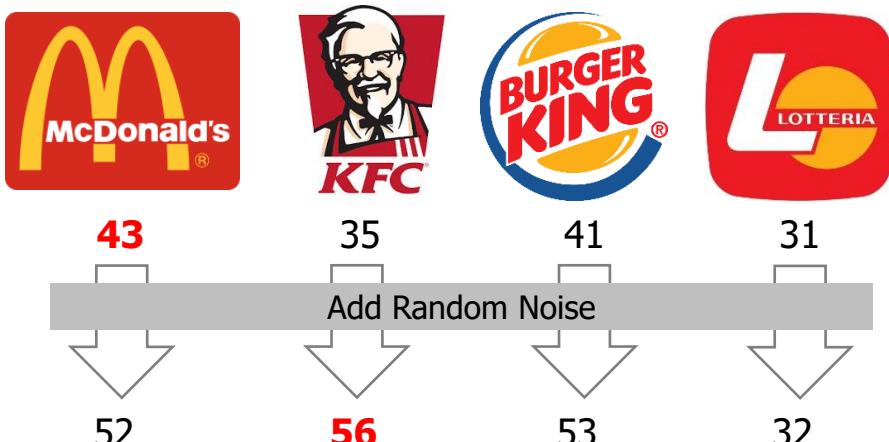
ϵ -Greedy

```
e=0.1
IF (random < e)
    a = random;
ELSE
    a = argmax(Q(s,a));
```

Decaying ϵ -Greedy

```
for i in range (1000); e=0.1 / (i+1);
IF (random < e)
    a = random;
ELSE
    a = argmax(Q(s,a));
```

Finding the Best Restaurant



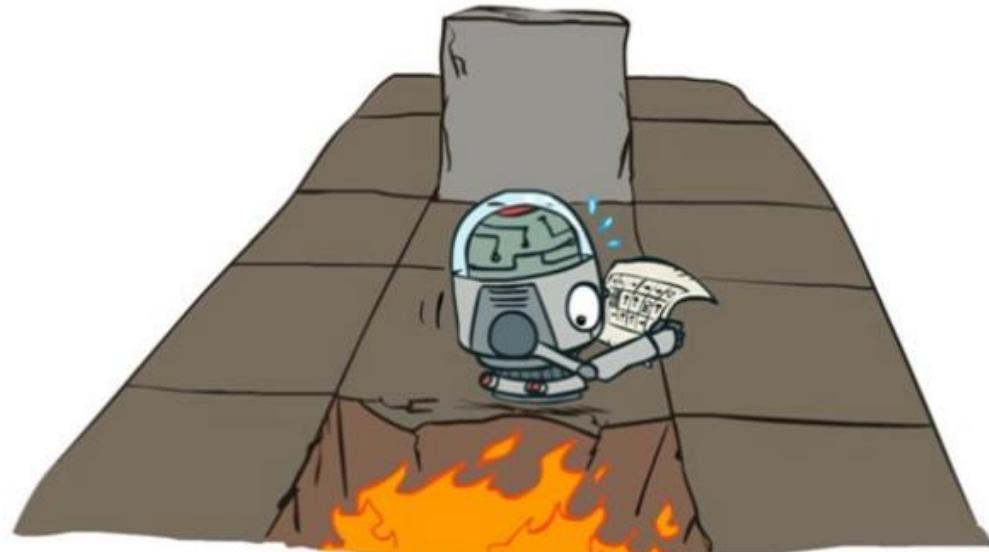
Add Random Noise

```
a = argmax(Q(s,a) + random_values);
```

Add Decaying Random Noise

```
for i in range (1000);
    a = argmax(Q(s,a) + random/(i+1));
```

- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$
 - S : Set of states
 - A : Set of actions
 - R : Reward function
 - T : Transition function
 - γ : Discount factor

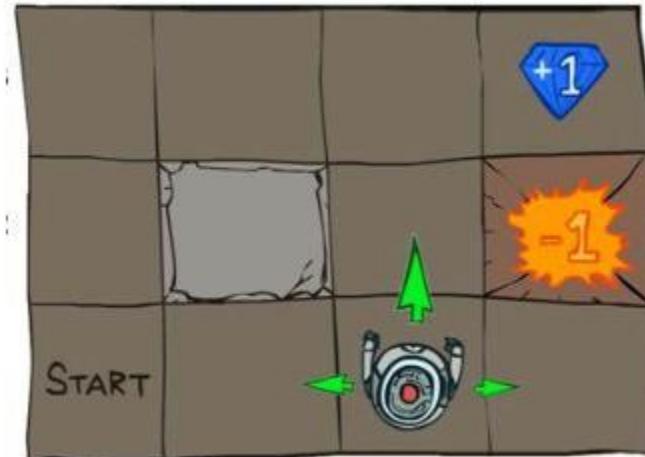


How can we use MDP to model agent in a maze?

- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$

- **S : Set of states**

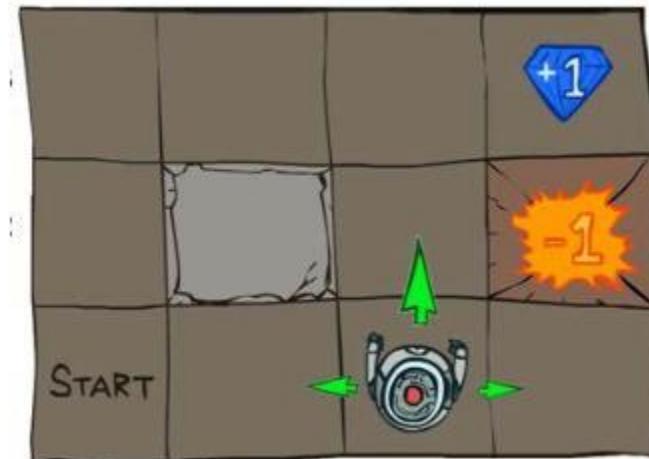
- A : Set of actions
- R : Reward function
- T : Transition function
- γ : Discount factor



S : location (x, y) if the maze is a 2D grid

- s_0 : starting state
- s : current state
- s' : next state
- s_t : state at time t

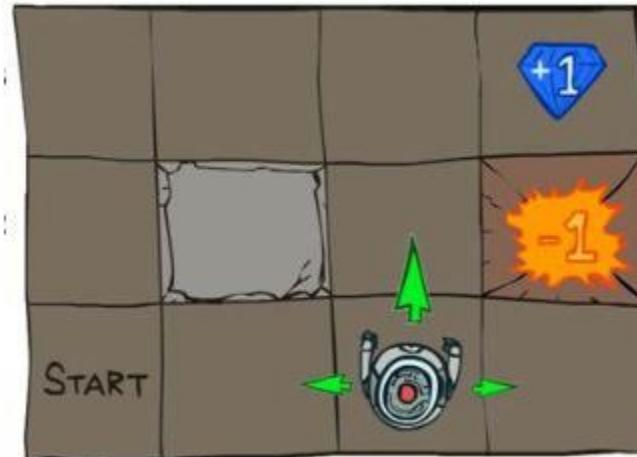
- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$
 - S : Set of states
 - **A: Set of actions**
 - R : Reward function
 - T : Transition function
 - γ : Discount factor



S : location (x, y) if the maze is a 2D grid
 A : move up, down, left, or right

- $s \rightarrow s'$

- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$
 - S : Set of states
 - A : Set of actions
 - **R : Reward function**
 - T : Transition function
 - γ : Discount factor

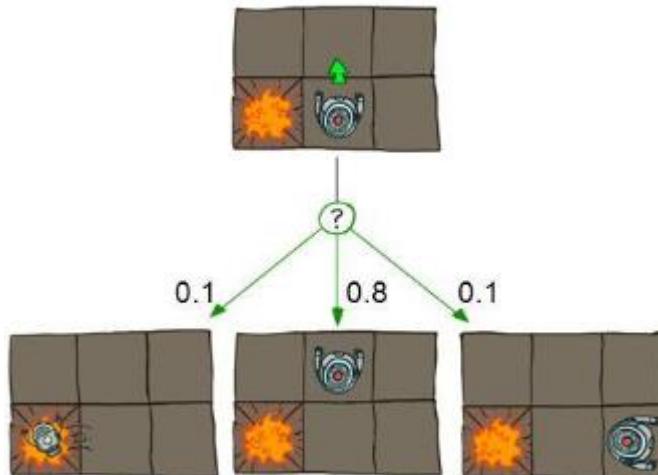


S : location (x, y) if the maze is a 2D grid
 A : move up, down, left, or right
 R : how good was the chosen action?

- $r = R(s, a, s')$
- -1 for moving (battery used)
- +1 for jewel? +100 for exit?

- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$

- S : Set of states
- A : Set of actions
- R : Reward function
- **T : Transition function**
- γ : Discount factor



Stochastic Transition

S : location (x, y) if the maze is a 2D grid
 A : move up, down, left, or right
 R : how good was the chosen action?
 T : where is the robot's new location?
• $T = P(s'|s, a)$

- Markov Decision Process (MDP) Components: $\langle S, A, R, T, \gamma \rangle$
 - S : Set of states
 - A : Set of actions
 - R : Reward function
 - T : Transition function
 - γ : **Discount factor**



1

Worth Now



γ

Worth Next Step



γ^2

Worth In Two Steps

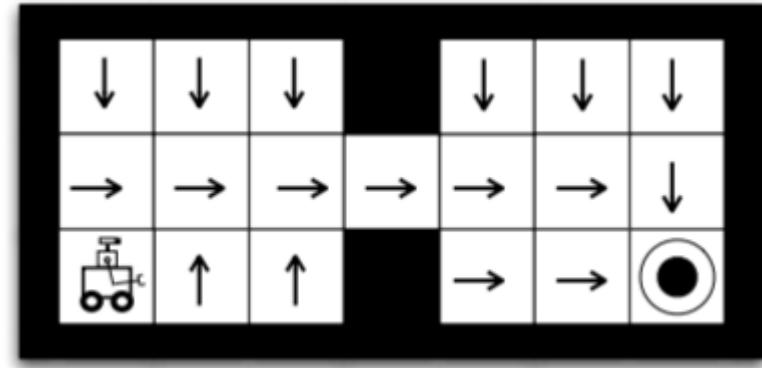
S : location (x, y) if the maze is a 2D grid
 A : move up, down, left, or right
 R : how good was the chosen action?
 T : where is the robot's new location?
 γ : how much does future reward worth?

- $0 \leq \gamma \leq 1$, [$\gamma \approx 0$: future reward is near 0 (immediate action is preferred)]

- Policy
 - $\pi: S \rightarrow A$
 - Maps states to actions
 - Gives an action for every state

- Return

- $$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$



Our goal:
Find π that maximizes expected return!

• State Value Function (V)

$$V^\pi(s) = E_\pi(R_t | s_t = s) = E_\pi\left(\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s\right)$$

- Expected return of starting at state s and following policy π
- How much return do I expect starting from state s ?

• Action Value Function (Q)

$$Q^\pi(s, a) = E_\pi(R_t | s_t = s, a_t = a) = E_\pi\left(\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s, a_t = a\right)$$

- Expected return of starting at state s , taking action a , and then following policy π
- How much return do I expect starting from state s and taking action a ?

- Our goal is to find the **optimal policy**

$$\pi^*(s) = \max_{\pi} R^{\pi}(s)$$

- If $T(s'|s, a)$ and $R(s, a, s')$ are known, this is a **planning** problem.
- We can use **dynamic programming** to find the optimal policy.

- Notes

- Bellman Equation
(Value Iteration)

$$\forall s \in S: V^*(s) = \max_a \sum_{s'} \{R(s, a, s') \cdot T(s, a, s') + \gamma V^*(s')\}$$



Artificial Intelligence and Mobility Lab

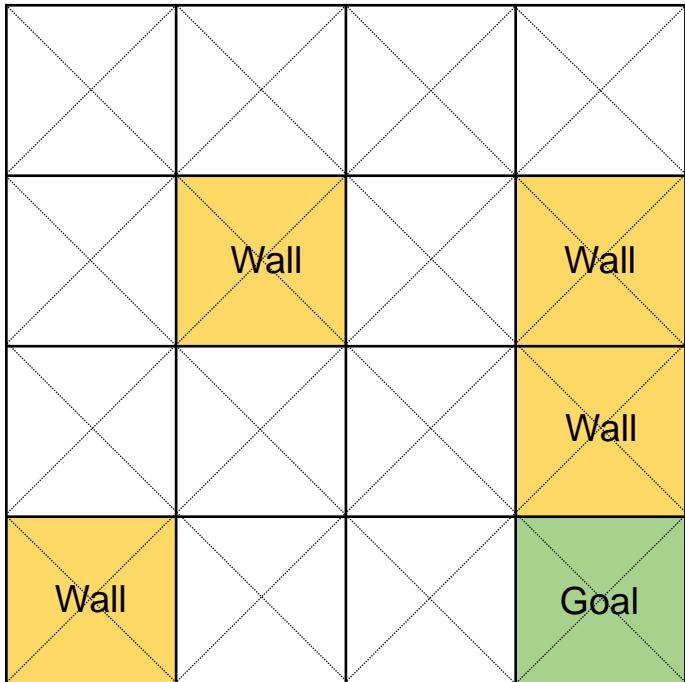
GAN and Reinforcement Learning

Reinforcement Learning (RL)

RL Implementation

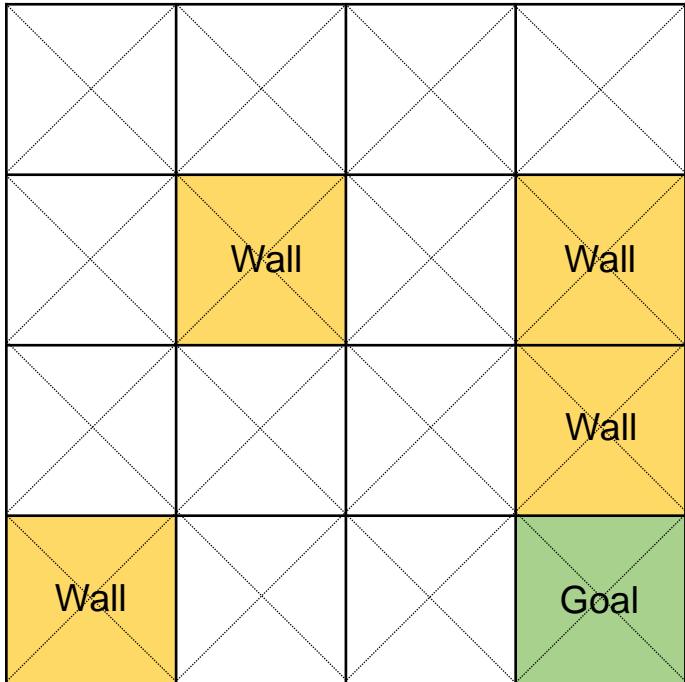
- RL Theory
- **RL Implementation**

Q-Learning (Basics)



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import gym
4 from gym.envs.registration import register
5 import random
6
7 """
8     Q-Table
9     | action | L | D | R | U |
10    -----
11 state: 0 |   |   |   |   |   |
12 -----
13 state: 1 |   |   |   |   |   |
14 -----
15 state: 2 |   |   |   |   |   |
16 -----
17 state: ... |   |   |   |   |   |
18 -----
19 """
20
21 register(
22     id='FrozenLake-v3',
23     entry_point='gym.envs.toy_text:FrozenLakeEnv',
24     kwargs={
25         'map_name': '4x4',
26         'is_slippery': False
27     }
28 )
29
30 env = gym.make("FrozenLake-v3")
```

Q-Learning (Basics)



```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import gym
4 from gym.envs.registration import register
5 import random
6
7 """
8     Q-Table
9     | action | L | D | R | U |
10    -----
11 state: 0 |   |   |   |   |   |
12 -----
13 state: 1 |   |   |   |   |   |
14 -----
15 state: 2 |   |   |   |   |   |
16 -----
17 state: ... |   |   |   |   |   |
18 ...
19
20
21 register(
22     id='FrozenLake-v3',
23     entry_point='gym.envs.toy_text:FrozenLakeEnv',
24     kwargs={
25         'map_name': '4x4',
26         'is_slippery': False
27     }
28 )
29
30 env = gym.make("FrozenLake-v3")
```

- Environment setting

```
32 # Initialization with 0 in Q-table
33 Q = np.zeros([env.observation_space.n, env.action_space.n]) # (16,4) where 16: 4*4 map, 4: actions
34 num_episodes = 1000 # Number of iterations
35
36 rList = []
37 successRate = []
38
39 def rargmax(vector):
40     m = np.amax(vector) # Return the maximum of an array or maximum along an axis (0 or 1)
41     indices = np.nonzero(vector == m)[0] # np.nonzero(True/False vector) => find the maximum
42     return random.choice(indices) # Random selection
43
44 for i in range(num_episodes): # Updates with num_episodes iterations
45     state = env.reset() # Reset
46     total_reward = 0 # Reward graph (1: success, 0: failure)
47     done = None
48
49     while not done: # The agent is not in the goal yet
50         action = rargmax(Q[state, :]) # Find maximum reward among 4 actions, find next action
51         new_state, reward, done, _ = env.step(action) # Result of the chosen action
52
53         Q[state, action] = reward + np.max(Q[new_state, :]) # Q-update
54         total_reward += reward
55         state = new_state
56
57         rList.append(total_reward) # Reward appending
58         successRate.append(sum(rList)/(i+1)) # Success rate appending
```

```

32 # Initialization with 0 in Q-table
33 Q = np.zeros([env.observation_space.n, env.action_space.n]) # (16,4) where 16: 4*4 map, 4: actions
34 num_episodes = 1000 # Number of iterations
35
36 rList = []
37 successRate = []
38
39 def rargmax(vector):
40     m = np.amax(vector) # Return the maximum of an array or maximum along an axis (0 or 1)
41     indices = np.nonzero(vector == m)[0] # np.nonzero(True/False vector) => find the maximum
42     return random.choice(indices) # Random selection
43
44 for i in range(num_episodes): # Updates with num_episodes iterations
45     state = env.reset() # Reset
46     total_reward = 0 # Reward graph (1: success, 0: failure)
47     done = None
48
49     while not done: # The agent is not in the goal yet
50         action = rargmax(Q[state, :]) # Find maximum reward among 4 actions, find next action
51         new_state, reward, done, _ = env.step(action) # Result of the chosen action
52
53         Q[state, action] = reward + np.max(Q[new_state, :]) # Q-update
54         total_reward += reward
55         state = new_state
56
57         rList.append(total_reward) # Reward appending
58         successRate.append(sum(rList)/(i+1)) # Success rate appending
    
```

- Randomly pick one when multiple argmax values exist

```

32 # Initialization with 0 in Q-table
33 Q = np.zeros([env.observation_space.n, env.action_space.n]) # (16,4) where 16: 4*4 map, 4: actions
34 num_episodes = 1000 # Number of iterations
35
36 rList = []
37 successRate = []
38
39 def rargmax(vector):
40     m = np.amax(vector) # Return the max
41     indices = np.nonzero(vector == m)[0]
42     return random.choice(indices) # Rand
43
44 for i in range(num_episodes): # Updates
45     state = env.reset() # Reset
46     total_reward = 0 # Reward graph (1:
47     done = None
48
49     while not done: # The agent is not in the goal yet
50         action = rargmax(Q[state, :]) # Find maximum reward among 4 actions, find next action
51         new_state, reward, done, _ = env.step(action) # Result of the chosen action
52
53         Q[state, action] = reward + np.max(Q[new_state, :]) # Q-update
54         total_reward += reward
55         state = new_state
56
57         rList.append(total_reward) # Reward appending
58         successRate.append(sum(rList)/(i+1)) # Success rate appending

```

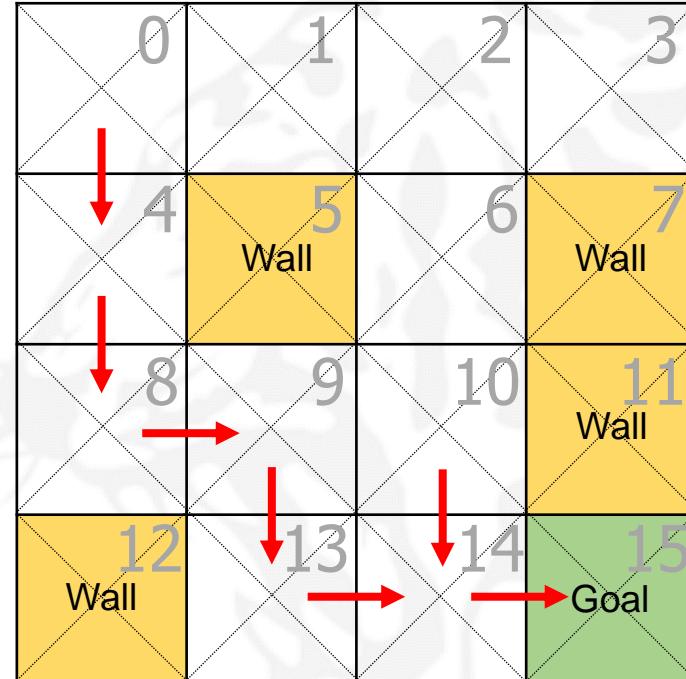
- Iteration until the agent arrives at the goal or it cannot move anymore.
- (line 50) find the action which returns max Q value.
- (line 51) take the action which is the result of (line 50).
 - done: if the agent is at goal or cannot move anymore, done → True
- (line 53) Q-update
- (line 54) reward value accumulation
- (line 55) state value update for next iteration

```

68 """
69 Final Q-Table
70 [[0. 1. 0. 0.] 0 (D)
71 [0. 0. 0. 0.] 1
72 [0. 0. 0. 0.] 2
73 [0. 0. 0. 0.] 3
74 [0. 1. 0. 0.] 4 (D)
75 [0. 0. 0. 0.] 5
76 [0. 0. 0. 0.] 6
77 [0. 0. 0. 0.] 7
78 [0. 0. 1. 0.] 8 (R)
79 [0. 1. 0. 0.] 9 (D)
80 [0. 1. 0. 0.] 10 (D)
81 [0. 0. 0. 0.] 11
82 [0. 0. 0. 0.] 12
83 [0. 0. 1. 0.] 13 (R)
84 [0. 0. 1. 0.] 14 (R)
85 [0. 0. 0. 0.] 15
86 Success Rate : 0.903
87 """

```

[L, D, R, U]



Outline

- Introduction
 - Theory
 - TensorFlow
- Generative Adversarial Network
- Reinforcement Learning
- **Deep Reinforcement Learning**





Artificial Intelligence and Mobility Lab

GAN and Reinforcement Learning

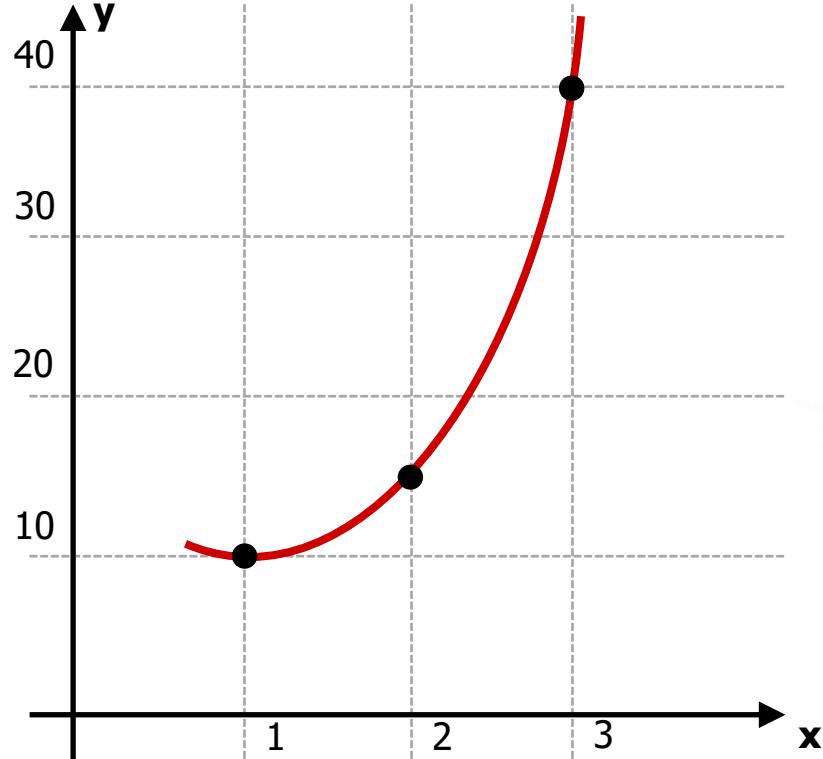
Deep Reinforcement Learning (DRL)

DRL Theory

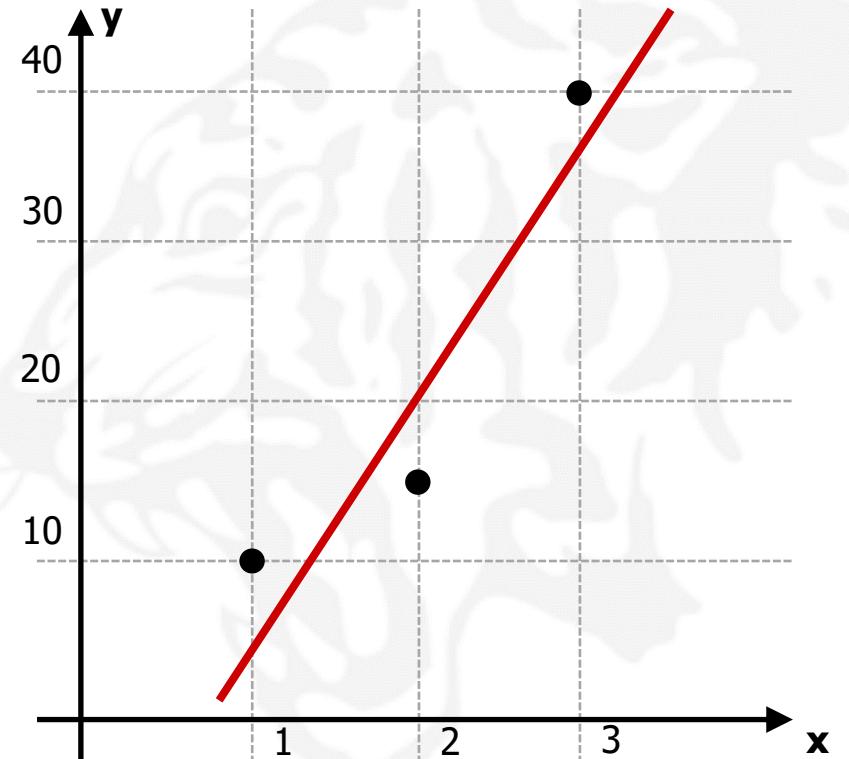
- **DRL Theory**

Interpolation vs. Linear Regression

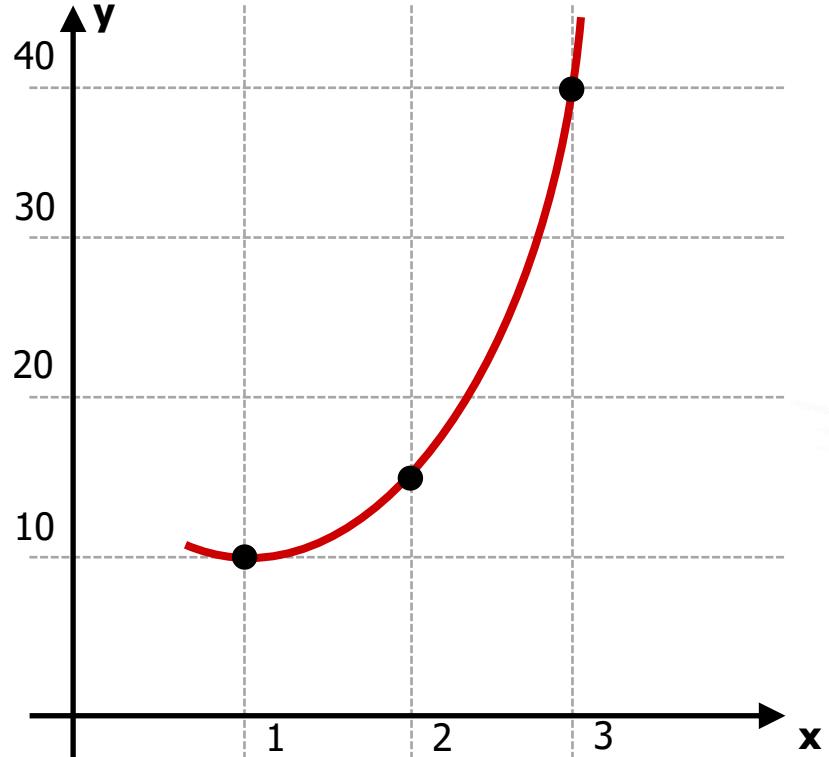
Interpolation



Linear Regression



Interpolation

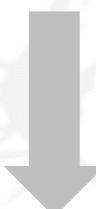


Interpolation with Polynomials

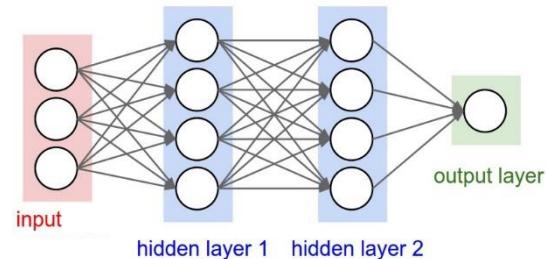
$$y = a_2 x^2 + a_1 x^1 + a_0$$

where three points are given.

→ Unique coefficients (a_0, a_1, a_2) can be calculated.



Is this related to
Neural Network Training?

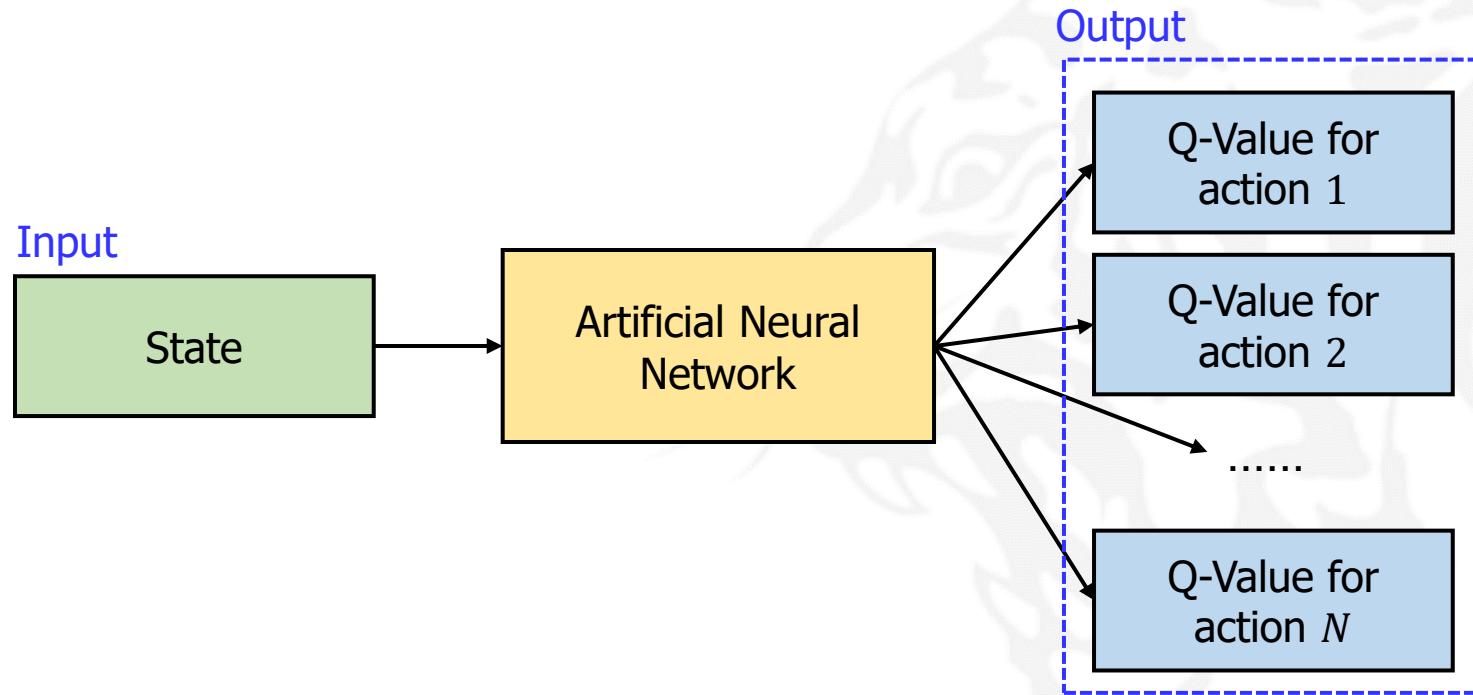


$$Y = a(a(a(X \cdot W_1 + b_1) \cdot W_2 + b_2) \cdot W_o + b_o)$$

where training data/labels (X : data, Y : labels) are given.

- Find $W_1, b_1, W_2, b_2, W_o, b_o$
- This is the mathematical meaning of neural network training.
- **Function Approximation**
- The most well-known function approximation with neural network:
Deep Reinforcement Learning

- It is inefficient to make the Q-table for each state-action pair.
→ ANN is used to **approximate the Q-function**.



- ICML 2018 Tutorial
 - <https://sites.google.com/view/icml2018-imitation-learning/>



Imitation Learning Tutorial ICML 2018

- ICML 2019 Tutorial
 - <https://slideslive.com/38917941/imitation-prediction-and-modelbased-reinforcement-learning-for-autonomous-driving>



Imitation, Prediction, and Model-Based Reinforcement Learning for Autonomous Driving

Sergey Levine

15th June 2019 - 10:50am



- Gameplay

Pro-Gamer



Trained Agent

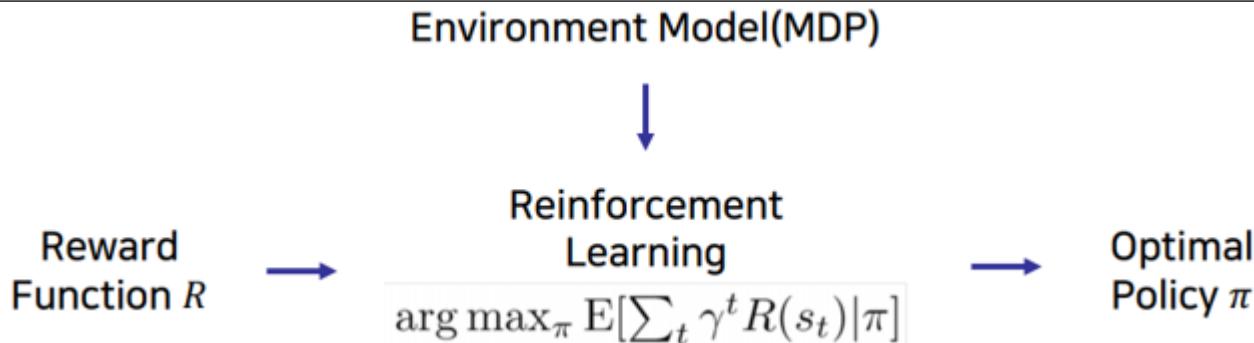


The goal of Imitation Learning is to train a policy to mimic
the expert's demonstrations

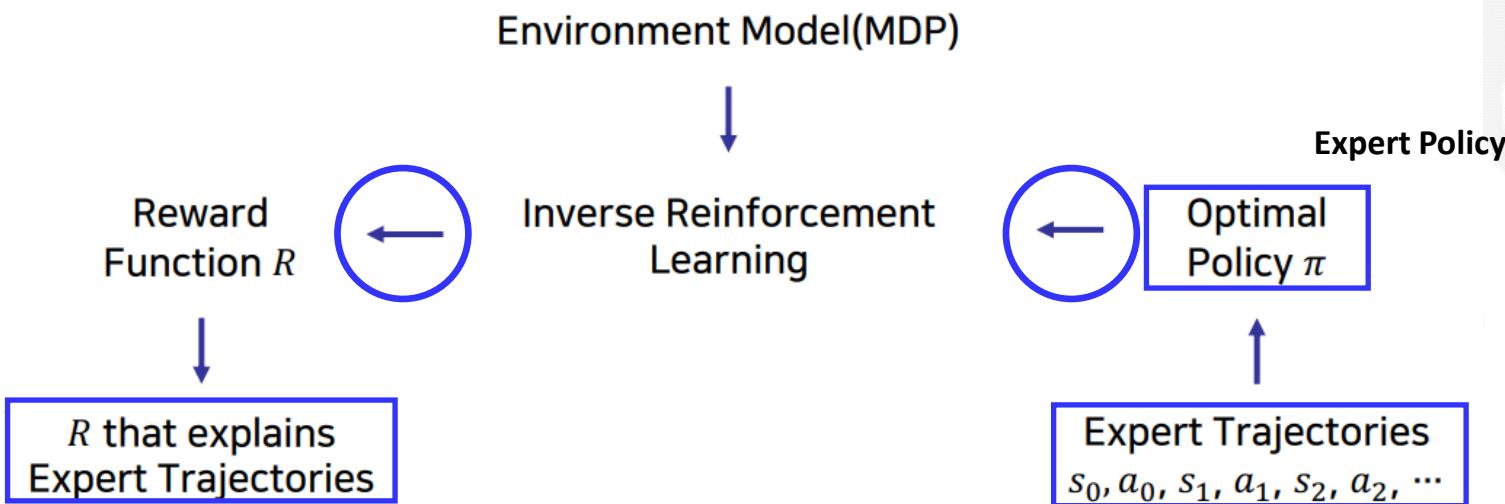
Inverse Reinforcement Learning (IRL)



RL



IRL



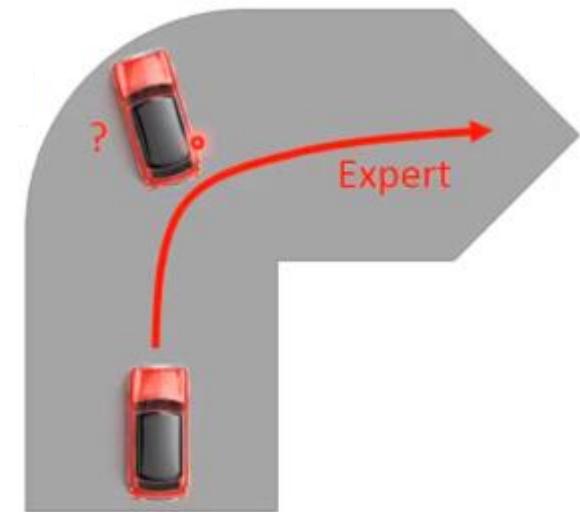
• Behavior Cloning

- Define $P^* = P(s|\pi^*)$ (distribution of states visited by **expert**)
- **Learning objective**

$$\operatorname{argmin}_{\theta} E_{(s,a_E) \sim P^*} L(a_E, \pi_{\theta}(s))$$
$$L(a_E, \pi_{\theta}(s)) = (a_E - \pi_{\theta}(s))^2$$

• Discussion

- Works well when P^* close to the distribution of states visited by π_{θ}
- **Minimize 1-step deviation error** along the expert trajectories



- Starcraft2

States: $s = \text{minimap, screen}$

Action: $a = \text{select, drag}$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_\theta(s) \rightarrow a$

States: s

Action: a

Policy: π_θ

- Policy maps states to actions : $\pi_\theta(s) \rightarrow a$
- Distributions over actions : $\pi_\theta(s) \rightarrow P(a)$

State Dynamics: $P(s'|s,a)$

- Typically not known to policy
- Essentially the simulator/environment

Rollout: sequentially execute $\pi_\theta(s_0)$ on initial state

- Produce trajectories τ

$P(\tau|\pi)$: distribution of trajectories induced by a policy

$P(s|\pi)$: distribution of states induced by a policy



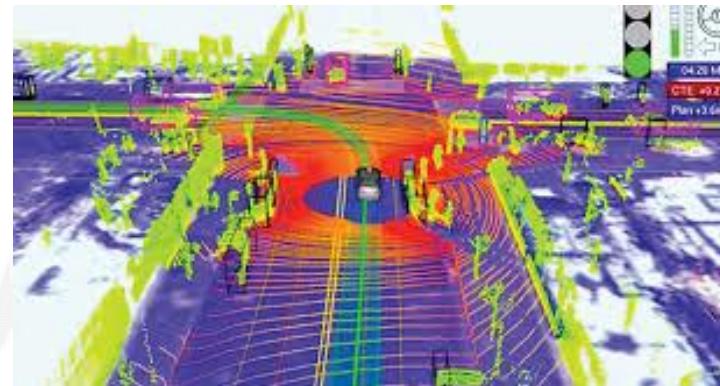
- Autonomous Driving Control

States: $s = \text{sensors}$

Action: $a = \text{steering wheel, brake, ...}$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_\theta(s) \rightarrow a$



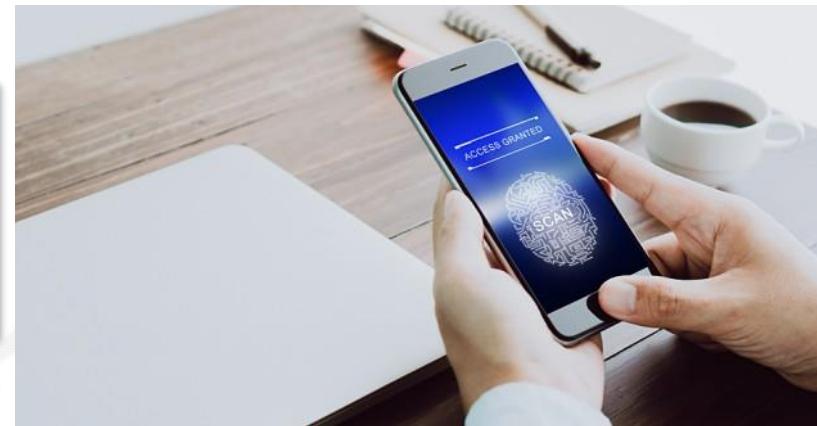
- Smartphone Security

States: $s = \text{apps}, \dots$

Action: $a = \text{use patterns}, \dots$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_\theta(s) \rightarrow a$



- PPF/RFTN Injection Control in Medicine

States: $s = \text{BIS, BP, ...}$

Action: $a = \text{PPF, RFTN, ...}$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_\theta(s) \rightarrow a$



Thank you for your attention!

- More questions?
 - joongheon@korea.ac.kr

