

# EE2213 Introduction to AI

## Lecture 7: Introduction to Machine Learning

Dr. WANG Si

[si.wang@nus.edu.sg](mailto:si.wang@nus.edu.sg)

Electrical and Computer Engineering Department  
National University of Singapore

# OVERVIEW OF COURSE CONTENTS



- **Introduction (Shaojing)**

- What is AI
- Applications of AI
- AI agent

- **Search (Shaojing)**

- Uninformed search algorithms: breadth-first, depth-first, uniform-cost(Dijkstra's algorithm)
- Informed search algorithms: greedy best-first, A\*

- Applications

- **Optimisation (Shaojing)**

- Linear programming
- Convex problems
- Applications

- **Machine learning (Wang Si) (Weeks 6-9)**

- Supervised and unsupervised learning: regression, classification, clustering
- Neural networks and deep learning
- Applications

(Week 10, Monday)

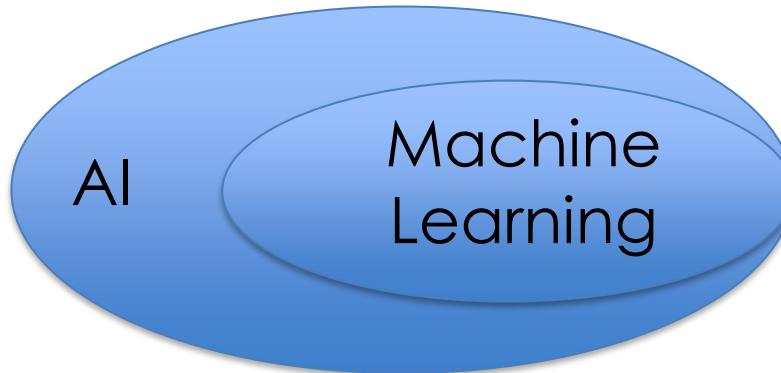
- **Knowledge representation (Wang Si) (Weeks 10-11)**

- Knowledge Representation and Reasoning
- Propositional Logic
- Applications

No Class on 20 Oct.! A recording will be uploaded on Canvas.

- **Ethical considerations (Shaojing)**

- Bias in AI
- Privacy concerns
- Societal impact



# What is Machine Learning?

AI Techniques that give computers **the ability to learn** without being explicitly programmed to do so

- Arthur Samuel



# What is Machine Learning?



A computer program is said to learn if its performance at task **T**, as measured by **P**, improves with **experience E**.

- Tom Mitchell

e.g., Google Maps finding the fastest route based on current traffic?

Google Maps predicting traffic jams based on historical traffic patterns?

millions of past trips from  
other users

Experience  
(E)

traffic jam  
prediction

Task  
(T)

Prediction accuracy

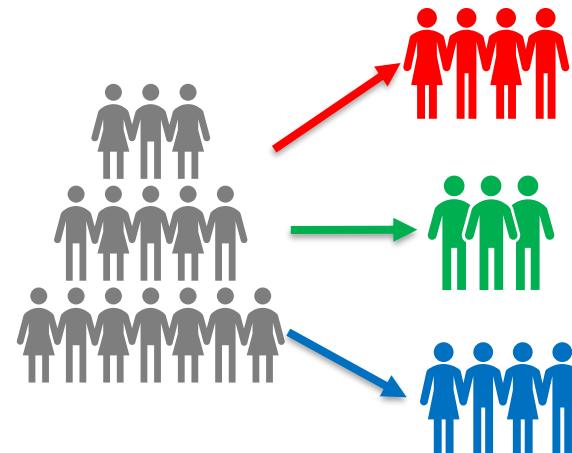
Performance  
Measure  
(P)

# What is Machine Learning?

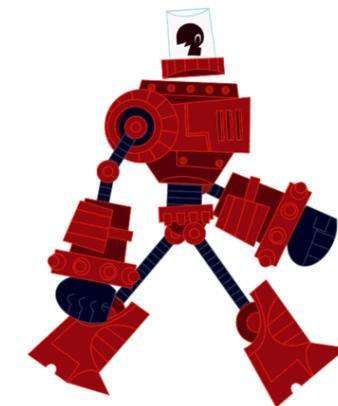
Examples: Task (**T**) Experience (**E**) Performance (**P**)



**T:** Spam Email Detection  
**E:** Email Data  
**P:** Detection Accuracy



**T:** Customer Segmentation  
**E:** Purchase data  
**P:** Silhouette Scores



**T:** Robot Walking  
**E:** Past records of Walking  
**P:** Walking time

*performance metric*

# Types of Machine Learning

-Supervised Learning

regression and classification tasks  
(training samples with desired output.  
Train a model to input -> function -> label)

-Unsupervised Learning

clustering tasks (use samples with no expected labels, and cluster to identify underlying patterns)

-Reinforcement Learning

seq of states, actions and delayed rewards

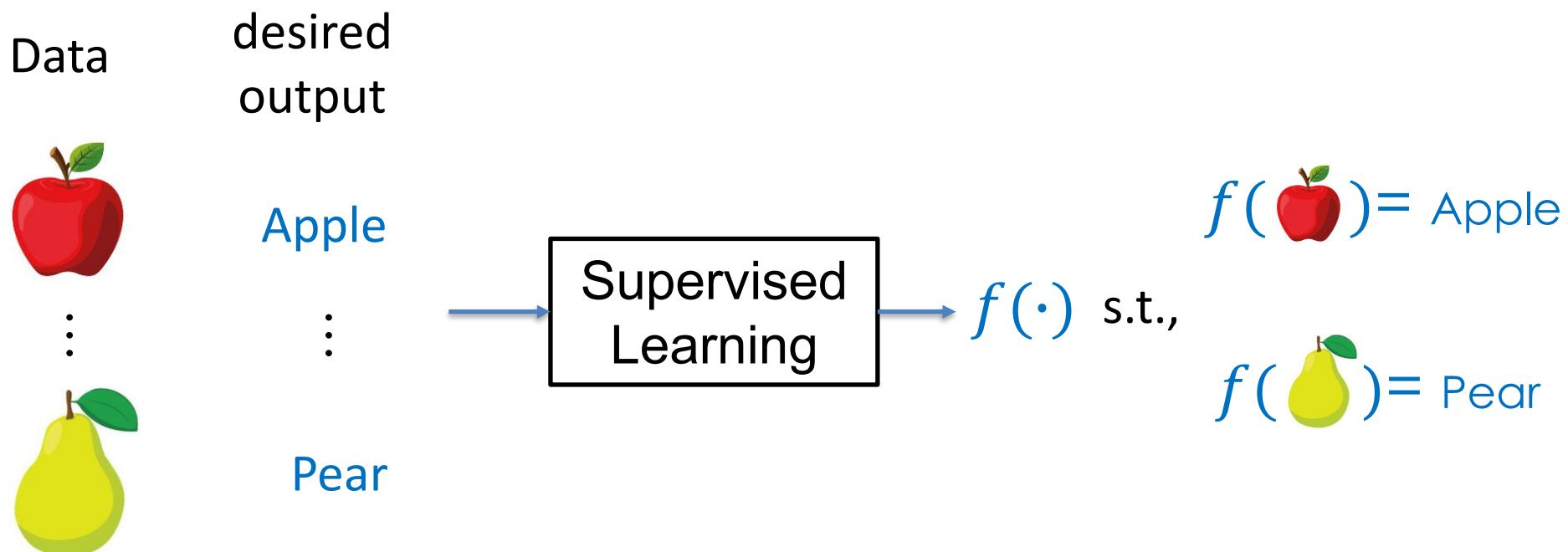
# Types of Machine Learning

## ➤ Supervised Learning

Input: data with desired output

by the output of the training  
data fed in.

Output: A function(i.e., model) that maps data to the desired output



# Types of Machine Learning

## ➤ Supervised Learning

Input: data with **desired output**

Output: A function(i.e., model) that maps data to the desired output

works with  
continuous labels

continuous

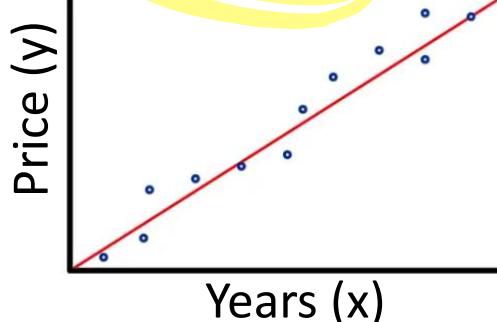
Numerical

Desired output

works with  
discrete labels  
(can be  
continuous too)

Regression

Lecture 8



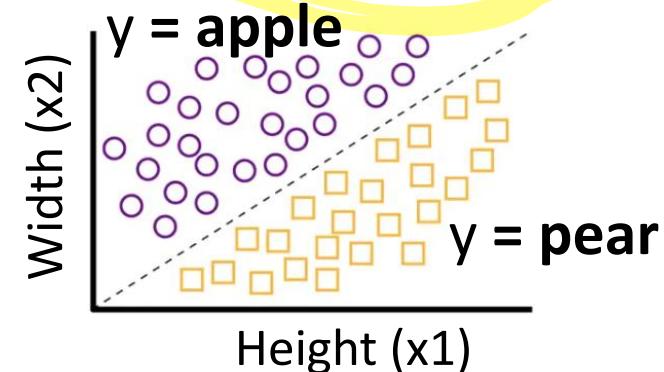
find the line that **best aligns**  
with samples

Lecture 9

Classification

predict  
category  
class  
given  
in test %.

Categorical



find the line that **separates**  
two classes

# Types of Machine Learning

## ➤ Unsupervised Learning

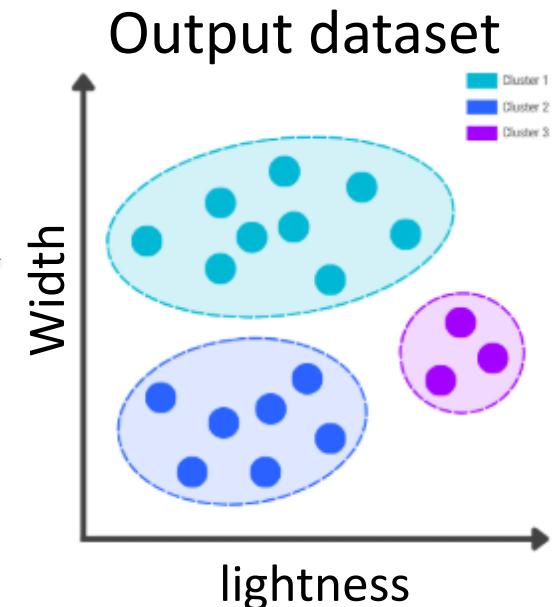
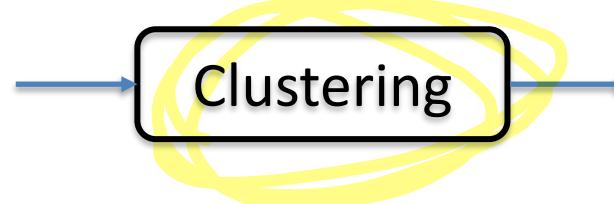
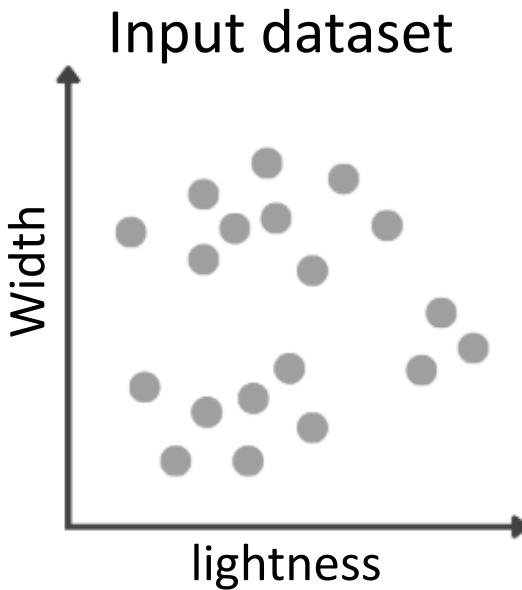
Input: data (**No desired output**)

No label given  
No expected output given

Output: Underlying patterns in data

Clustering

Lecture 10



# Types of Machine Learning

## ➤ Reinforcement Learning

Learn from experience to know what to do next time. (after each action)

Input: Sequence of States, Actions, and Delayed

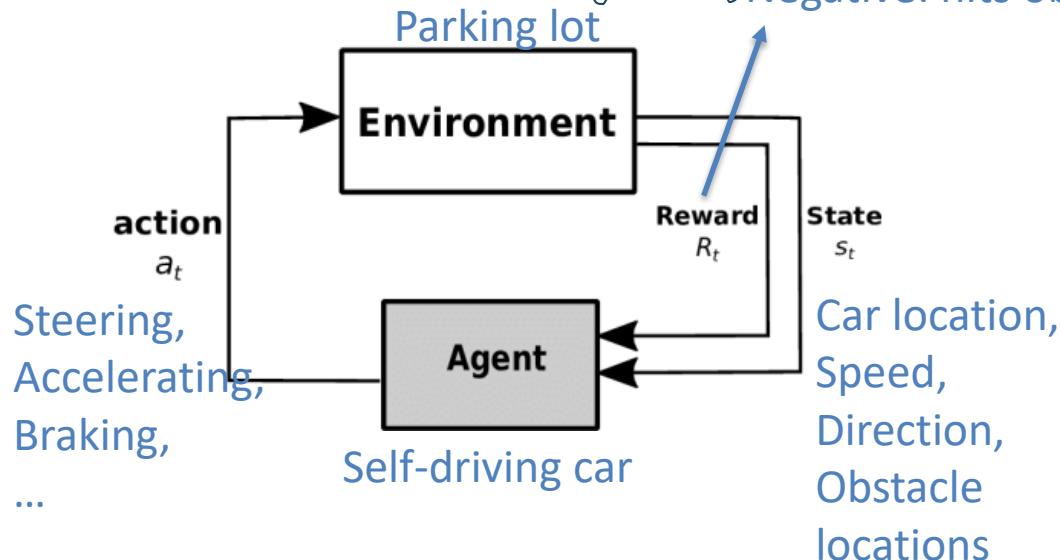
Rewards

Given seq. of  $S$  and actions  $A$  with delayed rewards  $R$ , output a policy (prob  $(a, s)$ ) to guide what action  $a$  to take on state  $s$ .

Output: Action Strategy (a rule that maps the environment to action)

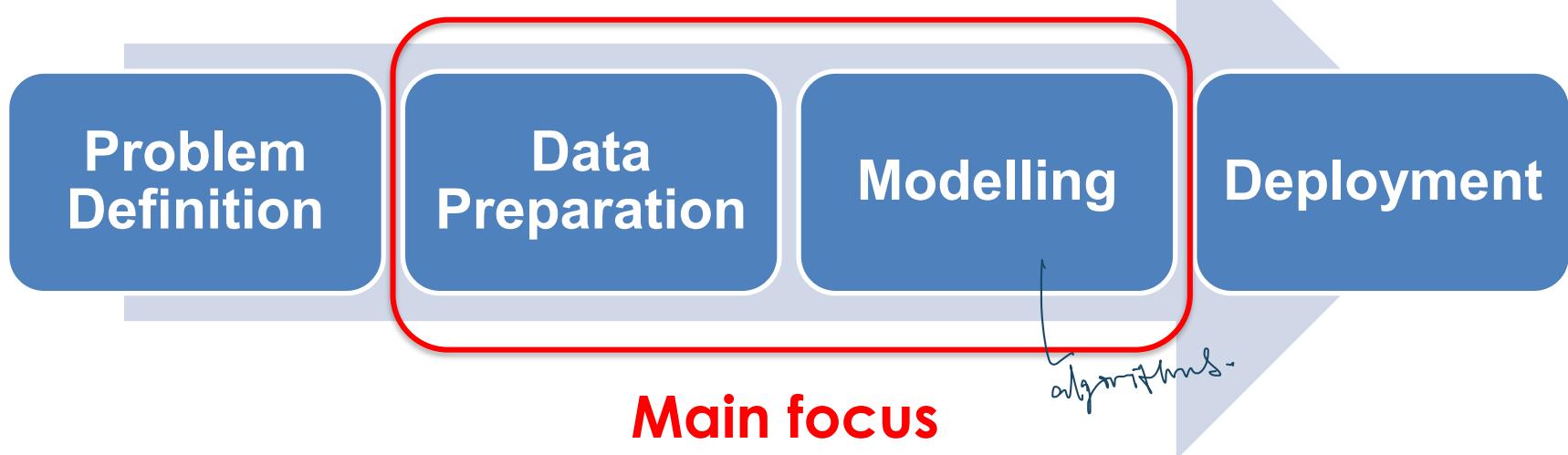
(demonstrates) Positive: closer to the correct spot

(punishes) Negative: hits obstacles or parks badly



Self-driving car learn to park  
example .

# Workflow of Machine Learning



# Workflow of Machine Learning



Problem Definition

Data Preparation

Modelling

Deployment

Example

## 1. Clarify the objective

What problem are you trying to solve?

## 2. Determine if ML is appropriate

Can this problem be learned from data?

## 3. Specify inputs and outputs

What input features will the model use?

What is the model trying to predict?

## 4. Choose the task type

Classification? Regression? Clustering?...

## 5. Decide on performance measure

Accuracy? Mean Squared Error?...

Predict House Prices

containing  $\Rightarrow$  determinants / price relation  
area of triangle / rectangle  $\Rightarrow$  all have formula  
or if outcome is purely random  
e.g. dice with only one roll.

Yes

can learn from experience (data)?  
 $\Rightarrow$  have some correlation.

Size, location, no. rooms, etc.

Sale price in SGD

Regression

Mean Squared Error

# Workflow of Machine Learning

Problem Definition

Data Preparation

Modelling

Deployment

## 1. Data Collection

Gather raw data from sensors,  
databases, files, etc.

See EERLL for types  
nominal : named  
ordinal : named + ordered  
interval : named + ordered + equal intervals  
ratio :  
↳ has true zero.

e.g., Predict House Prices

No.	Size	No. rooms	Location	Sold Price
1				
...				
100				

↳ attributes of samples  
for model to decide  
Labels to classify data  
INFO →

## 2. Exploratory Data Analysis

Flag any anomalies and inconsistencies.

Missing value

No.	Size	No. rooms	Location	Sold Price
1	112	2	"Pioneer"	600000
2	112	4	"Clementi"	801256
3	100	3	4	650K
4	116	3	"City Hall"	2300000
...				
100	112	2	"Pioneer"	600000

Inconsistent data values

Outlier

Duplicate data

# Workflow of Machine Learning



Problem Definition

Data Preparation

Modelling

Deployment

## 3. Data Preprocessing

- Address the problems flagged in the previous step.

e.g., handling missing value through dropping or imputation

If the number of samples is large.

Missing value

replace with zeros

No.	Size	No. rooms	Location	Sold Price
1	112	2	"Pioneer"	600000
2	112	4	"Clementi"	801256
3	100	3	4	"650K"
4	116	3	"City Hall"	1300000
...				
100	112	2	"Pioneer"	600000

Option 1: take the average value of the corresponding feature.

Option 2: take a value outside the normal range. e.g. -1 for size

Which of the following preprocessing steps would be the most suitable to address missing values in this dataset before applying linear regression?

- A. Remove all samples with missing values (CORRECT)
- B. Fill the missing values with zeros
- C. Remove all attributes (columns) with missing values
- D. Ignore the missing values as they will not impact model performance

# Workflow of Machine Learning

Problem Definition

Data Preparation

Modelling

Deployment

## 3. Data Preprocessing

- Format the data

- Convert non-numerical data into numerical form.

e.g., Use postal code for location

Convert categories into a binary vector through one-hot encoding.

*can do data wrangling  $\Rightarrow$  raw data to a more useful format before ML.*

After prediction

Categories: ['Cat', 'Dog', 'Bird']

one-hot encoding

Cat: [1,0,0]

Dog: [0,1,0]

Bird: [0,0,1]

$\hat{y} = [0.11\ldots, 1.2\ldots, 1.2\ldots]$

0

1

2

$\leftrightarrow$

cannot believe it  
except an order  
or numerical relationship

$0 \rightarrow 1$   
 $1 \rightarrow 0$

$1 \rightarrow 1$   
 $0 \rightarrow 0$

- Normalization (Ensure Features Contribute Fairly)

Across Samples for Each Feature

Min-max normalization:  $x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$

linear regression is not sensitive to scale of input features  
but it is important for iterative algos like gradient descent/MLP

Z-score standardization:  $x_{scaled} = \frac{x - \mu}{\sigma}$   $\rightarrow$  Mean  $\rightarrow$  Standard deviation

Use only  
training  
data  
Statistics!

No.	Feature 1	Feature 2
1	3 $\rightarrow$	-1 $\cancel{1}$
2	2 min 0	0 $\cancel{20}$
3	4 $\cancel{1}$	1 $\cancel{1}$

Normal distribution: mean is zero and standard deviation is one. But  $\mu=0$  and  $\sigma=1$  does not always work

# Workflow of Machine Learning



Problem Definition

Data Preparation

Modelling

Deployment

## 4. Feature Selection or Feature Extraction

Manipulating existing features to obtain **more relevant** features that improves performance

Feature Selection	Feature Extraction
Chooses from <b>existing features</b>	Creates <b>new features</b> from original ones
e.g., drop “Favorite color”	e.g., combine “height” and “weight” into “ <u>BMI</u> ”

Example task: Predict health risk *→ categorical w/ A ↗ use classification (supervised)*

	Height (m)	Weight (kg)	Age	Gender	Favorite color	Risk
1	1.75	59	33	Male	Red	High
...	...	...	...	...	...	...
100	1.65	50	31	Female	Blue	Medium

# Workflow of Machine Learning

Problem Definition

Data Preparation

Modelling

Deployment

## 4. Feature Selection and Feature Extraction

Manipulating existing features to obtain **more relevant** features that improves performance

-Feature Selection with **Pearson's correlation coefficient ( $r$ )**

(Measures linear relationship between two numerical variables)

Given  $N$  pairs of data points  $\{(a_1, b_1), \dots, (a_N, b_N)\}$

$$r = \frac{\frac{1}{N} \sum_{i=1}^N (a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\frac{1}{N} \sum_{i=1}^N (a_i - \bar{a})^2} \sqrt{\frac{1}{N} \sum_{i=1}^N (b_i - \bar{b})^2}}$$

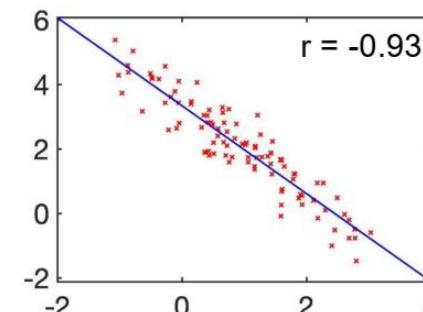
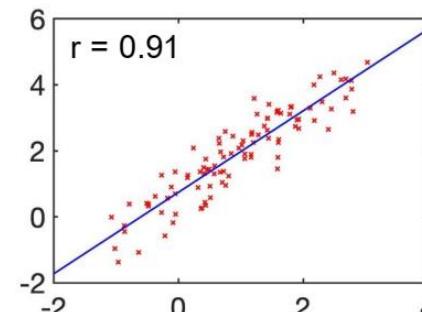
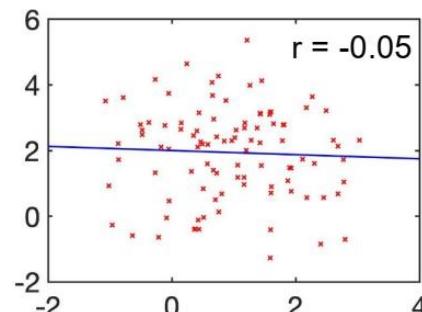
point  
into slope

$$\text{where } \bar{a} = \frac{1}{N} \sum_{i=1}^N a_i, \bar{b} = \frac{1}{N} \sum_{i=1}^N b_i$$

Use only  
training data!

recall 8 of n  
 $\Rightarrow$  -ve / +ve  
correlation  
one  $\uparrow$ , other  $\downarrow$  /  $\uparrow$

magnitude  
 $\Rightarrow$  strength  
of correlation.



# Workflow of Machine Learning



Problem Definition

Data Preparation

Modelling

Deployment

## 4. Feature Selection and Feature Extraction

Manipulating existing features to obtain **more relevant** features that improves performance

-Feature Selection with **Pearson's correlation coefficient (r)**

(Measures linear relationship between two numerical variables)

### Way 1:

- Compute r between each input feature and the target output in the **training set**.  
*→ more predictive, for supervised learning as regression*
- Select features with high absolute r.

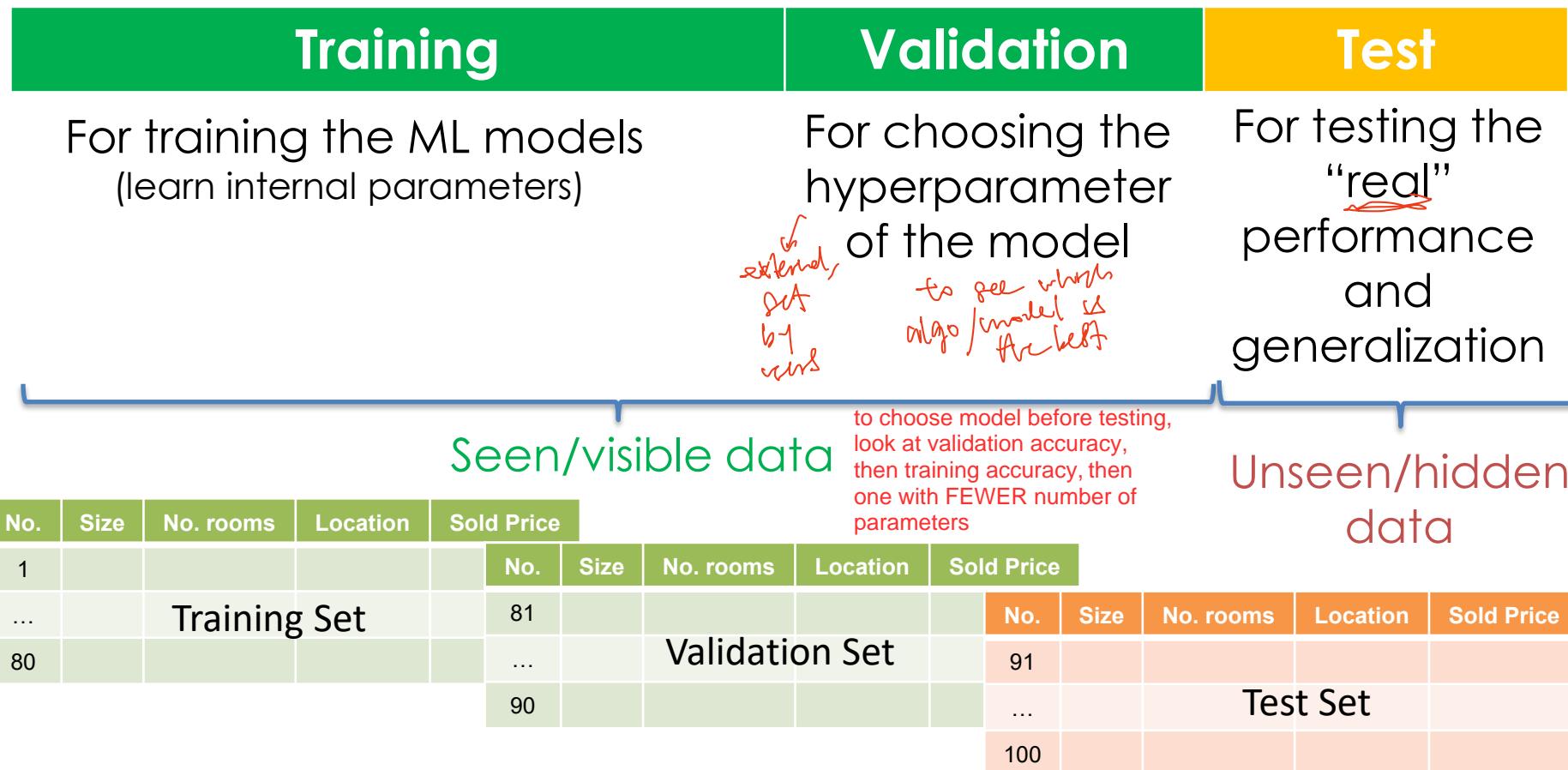
### Way 2:

- Compute r between every pair of input features in the **training set**  
*→ can view as redundant feature instead to reduce feature space.*
- Remove features with high absolute r.

# Workflow of Machine Learning



## 5. Split data into training, validation and test sets



# Workflow of Machine Learning

Problem Definition

Data Preparation

Modelling

Deployment

## 1. Train model on the training set

	No.	Size	No. rooms	Location	Sold Price
$x_1$	1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$y_1$
	...				
$x_{80}$	80	$x_{80,1}$	$x_{80,2}$	$x_{80,3}$	$y_{80}$

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_{80}, y_{80})\}$$

Learning Algorithm

Learned Model  
 $f(\cdot): X \rightarrow y$   
obtain this

## 2. Evaluate the model on the validation set

	No.	Size	No. rooms	Location	Sold Price
$x_{81}$	81	$x_{81,1}$	$x_{81,2}$	$x_{81,3}$	$y_{81}$
	...				
$x_{90}$	90	$x_{90,1}$	$x_{90,2}$	$x_{90,3}$	$y_{90}$

Learned Model  $f(\cdot)$

Performance measure P  
decide this

Predicted Prices:  
 $y'_{81}, y'_{82}, \dots, y'_{90}$

# Workflow of Machine Learning

Problem Definition

Data Preparation

Modelling

Deployment

## 3. Hyperparameter Tuning

Change the hyperparameters and repeat steps 1&2 to find the best hyperparameters that have the **best validation performance**.

	Hyperparameters	Parameters
<b>Definition</b>	External settings chosen before training	Internal variables learned from the data
<b>Learned during the training?</b>	No	Yes
<b>Set by user?</b>	Yes	No
<b>Examples</b>	-Learning rate -Training steps	<i>L.S. in gradient descent</i> -Regression Coefficients in Linear Regression

## 4. Test model on the test set

Unbiased evaluation of the final model

# Workflow of Machine Learning



Problem Definition

Data Preparation

Modelling

Deployment

- **Model Deployment (real-world environment)**

Web API : Model exposed as an endpoint that APPs can call

Cloud: AWS, GCP, or Azure provide managed ML serving platforms.

Edge: Model runs locally on devices like smartphones, IoT devices, etc.

# Summary by Quick Quiz



- ✓ Three Components in ML Definition

Task, experience, performance measure

- ✓ Three Types of ML

Supervised (regression for actual y/numerical, classification for label/class/categorical)

Unsupervised (no expected output)

Reinforcement learning (from seq of states and actions and delayed rewards/positive/negative feedback, to determine best actions)

- ✓ Three ways of Handling Missing Value

Dropping, average value, out of normal range value

- ✓ Data Formatting

Categorical data? Different Scales for Features?

One hot encoding

Normalisation (min-max, z score)

- ✓ One Criterion for Feature Selection

pearson correlation coefficient  $r$   
(only high  $r$  or remove high  $r$ )

- ✓ Dataset Partition

Training set to train the model, validation set to choose best hyperparameter for model., test set to evaluate model performance

# Review of Fundamentals

# Review of Fundamentals



## ❖ Scalars, Vectors, Matrices

- A **scalar** is a simple numerical value (e.g., 10, -4.2.).  
Denoted by an *italic* letter (e.g.,  $x, y$ )
- A **vector** is an ordered list of scalar values.  
Denoted by a **bold** letter (e.g.,  $\mathbf{x}, \mathbf{y}$ ) 
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$
- A **matrix** is a rectangular array of numbers arranged in rows and columns  
Denoted by a **bold** CAPITAL letter (e.g.,  $\mathbf{X}, \mathbf{Y}$ )

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & x_{1,3} \\ x_{2,1} & x_{2,2} & x_{2,3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1.5 \\ 2 & 9 & -3.2 \end{bmatrix}$$

# Review of Fundamentals

metabol!

## ❖ Operations on Vectors and Matrices

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{x}^T = [x_1 \ x_2]$$

$$\mathbf{x} - \mathbf{y} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 - y_1 \\ x_2 - y_2 \end{bmatrix}$$

$$a \mathbf{x} = a \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} ax_1 \\ ax_2 \end{bmatrix}$$

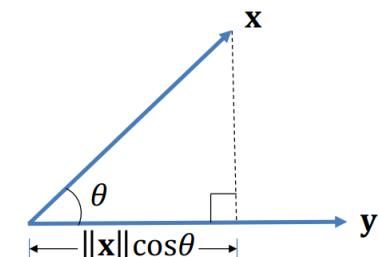
$$\begin{aligned} \mathbf{x} \cdot \mathbf{y} &= \mathbf{x}^T \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos\theta \\ &= [x_1 \ x_2] \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = x_1 y_1 + x_2 y_2 \end{aligned}$$

$\mathbf{x} \mathbf{w} = \mathbf{y}$

If  $(\mathbf{w}^T \mathbf{x})^T = (\mathbf{y}^T)^T \Rightarrow \mathbf{x}^T \mathbf{w} = \mathbf{y}$

$$\begin{aligned} \mathbf{Wx} &= \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \\ &= \begin{bmatrix} w_{1,1}x_1 + w_{1,2}x_2 + w_{1,3}x_3 \\ w_{2,1}x_1 + w_{2,2}x_2 + w_{2,3}x_3 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \mathbf{x}^T \mathbf{W} &= [x_1 \ x_2] \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \end{bmatrix} \\ &= [(x_1 w_{1,1} + x_2 w_{2,1}) \quad (x_1 w_{1,2} + x_2 w_{2,2}) \quad (x_1 w_{1,3} + x_2 w_{2,3})] \end{aligned}$$



**Output:**  
same

$$A_{3 \times 3} B_{3 \times 3} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} a_{11} \times b_{11} + a_{12} \times b_{21} + a_{13} \times b_{31} & a_{11} \times b_{12} + a_{12} \times b_{22} + a_{13} \times b_{32} & a_{11} \times b_{13} + a_{12} \times b_{23} + a_{13} \times b_{33} \\ a_{21} \times b_{11} + a_{22} \times b_{21} + a_{23} \times b_{31} & a_{21} \times b_{12} + a_{22} \times b_{22} + a_{23} \times b_{32} & a_{21} \times b_{13} + a_{22} \times b_{23} + a_{23} \times b_{33} \\ a_{31} \times b_{11} + a_{32} \times b_{21} + a_{33} \times b_{31} & a_{31} \times b_{12} + a_{32} \times b_{22} + a_{33} \times b_{32} & a_{31} \times b_{13} + a_{32} \times b_{23} + a_{33} \times b_{33} \end{bmatrix}$$

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

**-Matrix inverse:**  $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$

- $\det(A)$ : determinant of  $A$

$$\det(A) = \sum_{j=1}^k a_{ij} C_{ij}, \text{ where cofactor } C_{ij} = (-1)^{i+j} M_{ij}$$

$M_{ij}$ : minor of  $a_{ij}$  is the determinant obtained by deleting the row and column in which  $a_{ij}$  lies.

cofactor expansion along any row/col gives the same result.  $\Rightarrow$  most convenient result.

e.g.,  $A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$

$$M_{11} = \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} = a_{22} \times a_{33} - a_{23} \times a_{32}$$

$$C_{11} = (-1)^{1+1} M_{11} \rightarrow \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} \rightarrow a \begin{vmatrix} d & e \\ g & h \end{vmatrix} + c \begin{vmatrix} b & e \\ g & h \end{vmatrix} = a(dh - eg) + c(bh - eg)$$

$$\det(A) = \sum_{j=1}^3 a_{ij} C_{ij} = a_{i1} C_{i1} + a_{i2} C_{i2} + a_{i3} C_{i3}$$

or matlab

```
Python
import numpy as np
# Calculate inverse
A_inv = np.linalg.inv(A)
```

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

-**Matrix inverse:**  $\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A})$

- **adj(A):** adjugate/adjoint of  $\mathbf{A}$

$\text{adj}(\mathbf{A}) = \mathbf{C}^T$ , where cofactor  $\mathbf{C}_{ij} = (-1)^{i+j} \mathbf{M}_{ij}$

$\mathbf{M}_{ij}$ : minor of  $a_{ij}$  is the determinant obtained by deleting the row and column in which  $a_{ij}$  lies.

e.g.,  $\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$      $\mathbf{M}_{11} = \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} = a_{22} \times a_{33} - a_{23} \times a_{32}$

$\mathbf{C}_{11} = (-1)^{1+1} \mathbf{M}_{11}$

$$\text{adj}(\mathbf{A}) = \mathbf{C}^T = \begin{bmatrix} \mathbf{M}_{11} & -\mathbf{M}_{21} & \mathbf{M}_{31} \\ -\mathbf{M}_{12} & \mathbf{M}_{22} & \mathbf{M}_{32} \\ \mathbf{M}_{13} & -\mathbf{M}_{23} & \mathbf{M}_{33} \end{bmatrix}$$

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

-**Matrix inverse:**  $\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A})$

- $\det(\mathbf{A}) = \sum_{j=1}^k a_{ij} \mathbf{C}_{ij}$ , where cofactor  $\mathbf{C}_{ij} = (-1)^{i+j} \mathbf{M}_{ij}$
- $\text{adj}(\mathbf{A}) = \mathbf{C}^T$

$\mathbf{M}_{ij}$ : minor of  $a_{ij}$  is the determinant obtained by deleting the row and column in which  $a_{ij}$  lies.

e.g., What is the cofactor matrix of  $\mathbf{A} = \begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix}$        $\mathbf{C} = \begin{bmatrix} -3 & 6 & -3 \\ 6 & -12 & 6 \\ -3 & 6 & -3 \end{bmatrix}$

$$\mathbf{C}_{11} = (-1)^{1+1} \begin{vmatrix} 4 & 5 \\ 7 & 8 \end{vmatrix} = -3, \quad \mathbf{C}_{12} = (-1)^{1+2} \begin{vmatrix} 3 & 5 \\ 6 & 8 \end{vmatrix} = 6, \quad \mathbf{C}_{13} = (-1)^{1+3} \begin{vmatrix} 3 & 4 \\ 6 & 7 \end{vmatrix} = -3$$

$$\mathbf{C}_{21} = (-1)^{2+1} \begin{vmatrix} 1 & 2 \\ 7 & 8 \end{vmatrix} = 6, \quad \mathbf{C}_{22} = (-1)^{2+2} \begin{vmatrix} 0 & 2 \\ 6 & 8 \end{vmatrix} = -12, \quad \mathbf{C}_{23} = (-1)^{2+3} \begin{vmatrix} 0 & 1 \\ 6 & 7 \end{vmatrix} = 6,$$

$$\mathbf{C}_{31} = (-1)^{3+1} \begin{vmatrix} 1 & 2 \\ 4 & 5 \end{vmatrix} = -3, \quad \mathbf{C}_{32} = (-1)^{3+2} \begin{vmatrix} 0 & 2 \\ 3 & 5 \end{vmatrix} = 6, \quad \mathbf{C}_{33} = (-1)^{3+3} \begin{vmatrix} 0 & 1 \\ 3 & 4 \end{vmatrix} = -3,$$

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

-**Matrix inverse:**  $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$

Matrix inverse  $A^{-1}$  exists if and only if the matrix  $A$  is **square and nonsingular** (i.e., **invertible**)

*inv(A) in matlab-*

In other words, there exists a square matrix  $B$ , such that  **$AB=BA=I$**  (Identity matrix)

$$\mathbf{I}_{d \times d} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & 0 & \ddots & 0 & \vdots \\ 0 & \vdots & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

-**Matrix inverse:**  $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$

see 14 equivalent statements of invertibility -

Matrix inverse  $A^{-1}$  exists if and only if the matrix  $A$  is **square and nonsingular** (i.e., **invertible**)

$\det(A) \neq 0$

2 methods to check invertibility

$\text{rank}(A) = d$

$\det(A) = 0?$

full rank  
 $\text{rank}(A_{d \times d}) = d?$

$\text{rank}(A)$ : maximal number of linearly independent columns/rows of  $A$

$$\beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \cdots + \beta_m \mathbf{a}_m = 0 \text{ only holds for } \beta_1 = \cdots = \beta_m = 0$$

Compute  $\text{rank}(A)$ : Step 1. Transform into row echelon form using elementary row operations

Step 2. The number of non-zero rows in the row echelon form.

<https://stattrek.com/matrix-algebra/echelon-transform#MatrixA>

# Review of Fundamentals

## ❖ Operations on Vectors and Matrices

-**Matrix inverse:**  $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$

just materials

Matrix inverse  $A^{-1}$  exists if and only if the matrix  $A$  is **square and nonsingular** (i.e., **invertible**)

e.g.,  $\begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}$

$$\begin{bmatrix} 1 & 2 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 1 \\ 9 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 0 & 9 \\ 3 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \\ 3 & 4 & 5 \end{bmatrix}$$

$$2R_2 + R_1 = R_3$$

# Review of Fundamentals

## ❖ Some properties on Matrix Operations

$$\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A}$$

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

$$(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C})$$

$$(r\mathbf{A})^{-1} = r^{-1}\mathbf{A}^{-1}$$

$$(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$$

$$(\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC}$$

$$\mathbf{I}^{-1} = \mathbf{I}$$

$$\mathbf{I}_m \mathbf{A}_{m \times n} = \mathbf{A} = \mathbf{A}_{m \times n} \mathbf{I}_n$$

$$(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T$$

$$(\mathbf{A}^T)^T = \mathbf{A}$$

$$(\mathbf{A}^{-1})^{-1} = \mathbf{A}$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

$$(r\mathbf{A})^T = r\mathbf{A}^T$$

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

*very important for  
 $XW = y$ .*

# Review of Fundamentals



## ❖ Sets, Functions

-A **Set** is an unordered collection of unique elements

Denoted by a calligraphic capital letter (e.g.,  $\mathcal{X}, \mathcal{Y}$ )

$x \in \mathcal{X}$ : the element  $x$  belongs to the set  $\mathcal{X}$

✓ Finite set: a fixed number of elements

Denoted by accolades (e.g.,  $\{1, 3.5, 7, -1\}$ )

✓ Infinite set: include all values in some interval

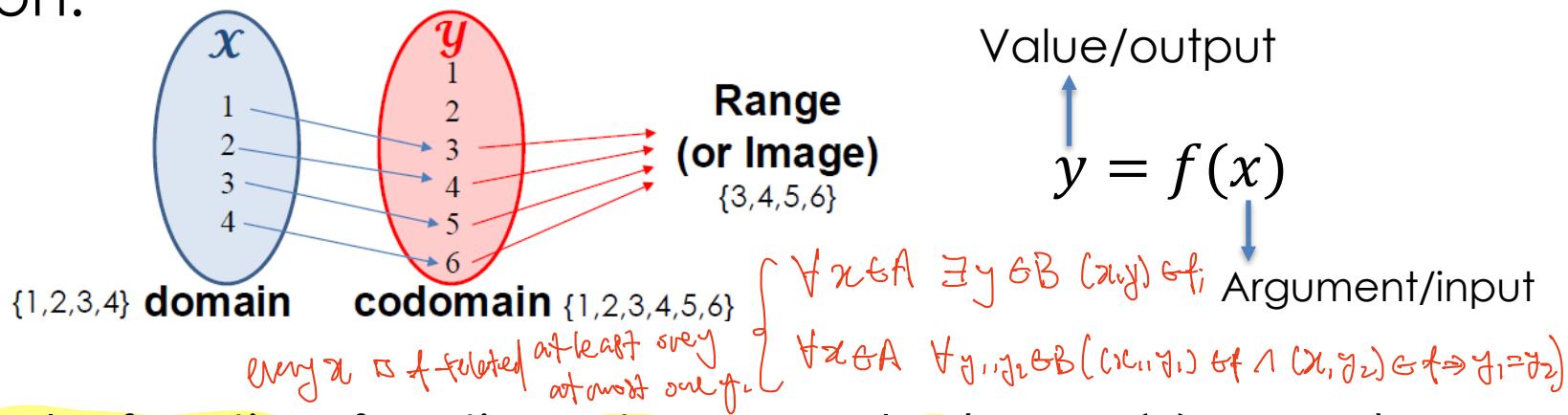
Denoted by square brackets as  $[a, b]$  if including  $a$  and  $b$   
(e.g.,  $[1,2]$ )

Denoted by parentheses as  $(a, b)$  if does not include  $a$  and  $b$   
(e.g.,  $(1,2)$ )

# Review of Fundamentals

## ❖ Sets, Functions

-A **Function** is a relation that associates **each element**  $x$  of a set  $X$ , the domain of the function, to a single element  $y$  of another set  $Y$ , the codomain of the function.



- ✓ **Scalar function:** function returns a scalar (e.g.,  $f(x)$  or  $f(\mathbf{x})$ )
- ✓ **Vector function:** function returns a vector (e.g.,  $\mathbf{f}(x)$  or  $\mathbf{f}(\mathbf{x})$ )

# Review of Fundamentals

## ❖ Operations on Functions

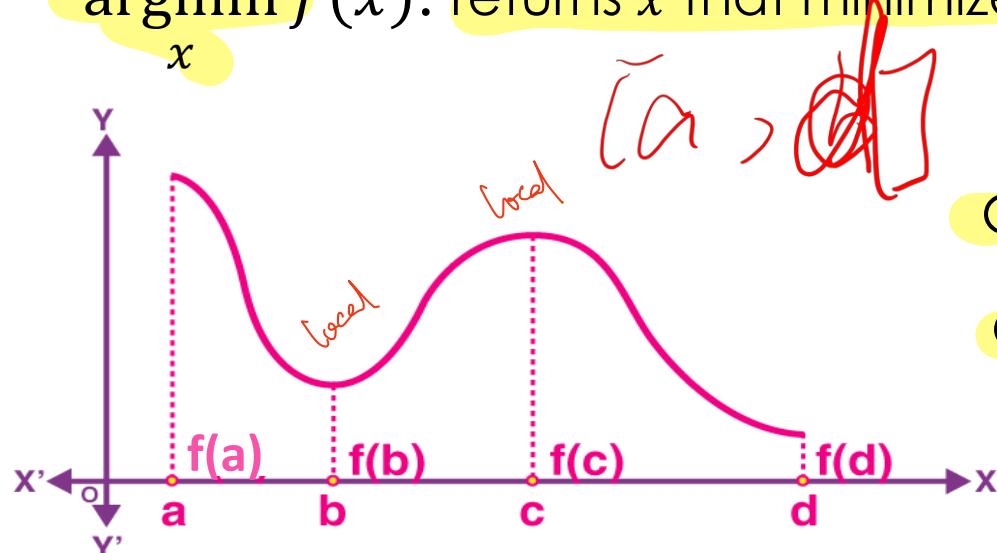
### -Maximum and Minimum

$\max_x f(x)$ : returns the highest value of  $f(x)$  across its domain

$\operatorname{argmax}_x f(x)$ : returns  $x$  that maximizes  $f(x)$  → a

$\min_x f(x)$ : returns the lowest value of  $f(x)$  across its domain

$\operatorname{argmin}_x f(x)$ : returns  $x$  that minimizes  $f(x)$



Global maximum =  $\operatorname{argmax}_x f(x)$

Global minimum =  $\operatorname{argmin}_x f(x)$

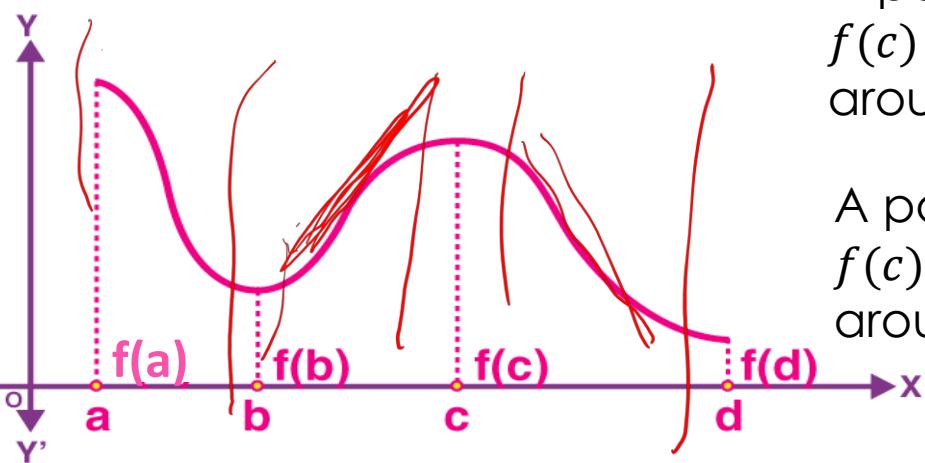
# Review of Fundamentals

## ❖ Operations on Functions

### -Derivative (single-input)

The derivative  $f'$  of a function  $f$  is a function that describes how fast  $f$  grows (or decreases)

- If  $f'$  is positive at some  $x$ ,  $f$  grows at this  $x$
- If  $f'$  is negative at some  $x$ ,  $f$  decreases at this  $x$
- If  $f'$  is zero at some  $x$ ,  $f$ 's slope at  $x$  is horizontal



A point  $x = c$  is a **local maximum** of  $f(x)$  if  $f(c) \geq f(x)$  for all  $x$  in some **open** interval around  $c$

A point  $x = c$  is a **local minimum** of  $f(x)$  if  $f(c) \leq f(x)$  for all  $x$  in some **open** interval around  $c$

# Review of Fundamentals

## ❖ Operations on Functions

### -Gradient (**multiple-input scalar function**)

Generalization of derivative for scalar functions with multiple inputs.

Denoted by  $\nabla_{\mathbf{x}} f$

diff of a  
scalar function  
 $f(\mathbf{x})$  wrt a vector  $\mathbf{x}$

A vector of **partial derivatives**

$$\nabla_{\mathbf{x}} f = \frac{df(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

Differentiation of  $f(\mathbf{x})$  w.r.t.  $\mathbf{x}$

if  $f(\mathbf{x})$  has  $d$  variables,  $\mathbf{x}$  is a  $d \times 1$  vector.

result  $\nabla_{\mathbf{x}} f$  is a  $d \times 1$  vector

e.g.,  $f(\mathbf{x}) = 2x_1 + 3x_2$ , what is  $\nabla_{\mathbf{x}} f$ ?

A matrix of **partial derivatives**

$$\nabla_{\mathbf{x}} f = \frac{df(\mathbf{X})}{d\mathbf{X}} = \begin{bmatrix} \frac{\partial f}{\partial x_{1,1}} & \cdots & \frac{\partial f}{\partial x_{1,K}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{d,1}} & \cdots & \frac{\partial f}{\partial x_{d,K}} \end{bmatrix}$$

same size as  $\mathbf{X}$

2  
[ 3 ]

# Review of Fundamentals

## ❖ Operations on Functions

### ~~-Jacobian (multiple-input vector function)~~

Generalization of derivative for vector functions with multiple inputs.

If  $f(\underline{x})$  is a vector function of size  $h \times 1$  and  $\underline{x}$  is a  $d \times 1$  vector,  
the Jacobian results in an  $h \times d$  matrix

$$\begin{bmatrix} 2 \\ 5x_1 \\ 3x_2 \end{bmatrix}$$

A matrix of **partial derivatives**

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_h}{\partial x_1} & \dots & \frac{\partial f_h}{\partial x_d} \end{bmatrix}$$

$$\overbrace{\partial(2x_1 + 3x_2)}^{2x_1 + 3x_2}$$

$$\overbrace{\partial(5x_1x_2)}^{5x_1x_2}$$

Differentiation of  $\mathbf{f}(\mathbf{x})$  of size  $h \times 1$  w.r.t.  $\mathbf{x}$  of size  $d \times 1$ .

e.g.,  $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} 2x_1 + 3x_2 \\ 5x_1x_2 \end{bmatrix}$ , what is the Jacobian matrix  $\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}}$ ?

$$\begin{bmatrix} 2 & 3 \\ 5x_2 & 5x_1 \end{bmatrix}$$

# Review of Fundamentals



## ❖ Common Differentiation Rules

- Constant Rule: the derivative of a constant is 0
- Sum & Difference Rule:  $(f(x) + g(x))' = f'(x) + g'(x)$   
 $(f(x) - g(x))' = f'(x) - g'(x)$
- Product Rule:  $(f(x)g(x))' = f(x)g'(x) + f'(x)g(x)$
- Quotient Rule:  $\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$
- Chain Rule:  $(f(g(x)))' = f'(g(x))g'(x)$  or  $\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$

# EE2213 Introduction to AI

## Lecture 8: Regression

Dr. WANG Si

[si.wang@nus.edu.sg](mailto:si.wang@nus.edu.sg)

Electrical and Computer Engineering Department  
National University of Singapore

# OVERVIEW OF COURSE CONTENTS



- **Introduction (Shaojing)**

- What is AI
- Applications of AI
- AI agent

- **Search (Shaojing)**

- Uninformed search algorithms: breadth-first, depth-first, uniform-cost(Dijkstra's algorithm)
- Informed search algorithms: greedy best-first, A\*

- Applications

- **Optimisation (Shaojing)**

- Linear programming
- Convex problems
- Applications

- **Machine learning (Wang Si) (Weeks 6-9)**

- Supervised and unsupervised learning: regression, classification, clustering
- Neural networks and deep learning
- Applications

- **Knowledge representation (Wang Si) (Weeks 10-11)**

- Knowledge Representation and Reasoning
  - Propositional Logic
  - Applications
- **Ethical considerations (Shaojing)**
- Bias in AI
  - Privacy concerns
  - Societal impact

(Week 10, Monday)

No Class on 20 Oct.! A recording will be uploaded on Canvas.

# AGENDA

➤ We will discuss:

- EE221L*
- Linear Regression
  - Polynomial Regression
  - Ridge Regression (for regularization)
  - ✖ ➤ Evaluation Metrics

At the end of this lecture, you should be able to:

- ✓ Explain the working principles
- ✓ Apply the techniques to solve regression problems.
- ✓ Evaluate the performance of a regression task

# Recap: Regression

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Output:  $f(\mathbf{x})$  that predict numerical  $y$  given  $\mathbf{x}$

map feature  $\mathbf{x}$  to a  $y$  value.

Applications:

House Price Prediction



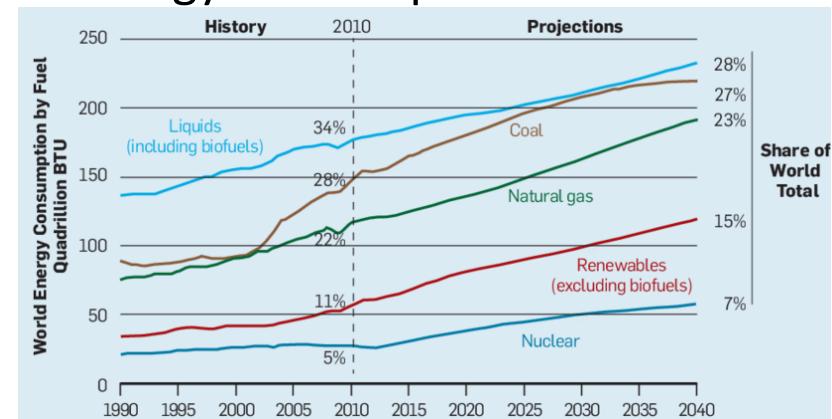
Student Performance Prediction

	Attendance (%)	Assignment	Quiz	Study Hours/Week	Final Exam
1	90	85	80	10	88
2	75	70	65	5	72
...	...	...	...	...	...

Sales Forecasting



Energy Consumption Estimation



# Linear Regression

If over determined, need left-inverse  
under determined, need right-inverse } see FE2211

# Linear Regression (Scalar Function)

- Linear Regression learns a model which is a **linear combination of features** of the input example.

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,

where  $\mathbf{x}$  is a **d-dimensional feature vector**

input  
feature  
vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Output:  $f_{\mathbf{w}}(\mathbf{x}) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d = \mathbf{x}^T \mathbf{w}$

where  $\mathbf{w}$  is a **d-dimensional vector of parameters**

weight  
vector

$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

e.g., House Price Prediction

No.	Size	No. rooms	Location	Sold Price
1	112	2	648310	600000
...				
90	100	4	129580	801256



# Linear Regression (Scalar Function)

## ➤ Learning (Training)

Find the optimal value  $\mathbf{w}^*$  which **minimizes**:

$$J(\mathbf{w}) = \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

↑ objective function  
↑ predicted  
↑ actual  
↑ cost function.

↓  
Squared error loss

↳ 2 slides later.

### Least Squares Regression

- **Objective function:** the expression to be minimized or maximized
- **Loss function:** a measure of the difference between **predicted output**  $f_{\mathbf{w}}(\mathbf{x}_i)$  and **actual/desired/target output**  $y_i$ .

*X X T is matrix way of writing sum of squares*

$$\sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \quad \text{in matrix form}$$

$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$

where  $\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,d} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \dots & x_{N,d} \end{bmatrix}$ ,  $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$

*↳ feature design matrix. ↳ rows → no. of samples ↳ columns → no. of features.*

# Linear Regression (Scalar Function)

## ➤ Learning (Training)

$$\underset{\mathbf{w}}{\operatorname{argmin}} \underbrace{(\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})}_{\hookrightarrow \text{cost function to minimise}}$$

returns the  $\mathbf{x}$  that gives the min  $J(\mathbf{x})$ .

in this case,  $\mathbf{w}$  that reduced our cost function (squared error) the most.

$$J(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = (\mathbf{w}^T \mathbf{X}^T - \mathbf{y}^T)(\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$= (\mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y}) = (\mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y})$$

The optimal value  $\mathbf{w}^*$  should have **zero gradient**:

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = 0$$

Why?

$$\frac{d(\mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y})}{d\mathbf{w}} = 0$$

using  
two rules

$$\begin{aligned} \frac{d(\mathbf{w}^T \mathbf{A}\mathbf{w})}{d\mathbf{w}} &= (\mathbf{A} + \mathbf{A}^T)\mathbf{w} \\ \frac{d(\mathbf{y}^T \mathbf{A}\mathbf{w})}{d\mathbf{w}} &= \mathbf{A}^T \mathbf{y} \end{aligned}$$

to find local minimum  
sqrt of squared errors is a convex function, any local min is the global minimum!

$$2\mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0$$

but multiply by  
 $(\mathbf{X}^T \mathbf{X})^{-1}$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$

If  $\mathbf{X}^T \mathbf{X}$  is invertible,

$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$



# Linear Regression (Scalar Function)

## ➤ Learning (Training)

- Geometric Interpretation of Least Squares Solution

$$\underbrace{\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})}_{\text{error term}} = \mathbf{0} = \mathbf{X}^T \mathbf{e}$$

$\mathbf{e}$   
derived from  $\frac{d J(\mathbf{w})}{d \mathbf{w}}$

$$\begin{bmatrix} x_{1,1} & \dots & x_{N,1} \\ \vdots & \ddots & \vdots \\ x_{1,d} & \dots & x_{N,d} \end{bmatrix} \mathbf{e} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

to find  $\mathbf{w}^*$ , iterated  
to zero to minimize  
 $J(\mathbf{w})$

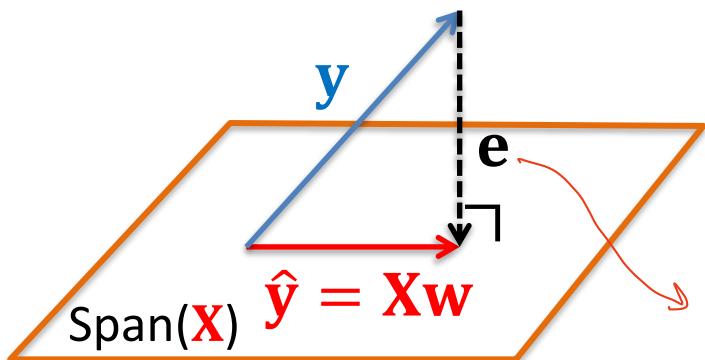
$\mathbf{X}\mathbf{w}^*$  is the projection of  $\mathbf{y}$  on the column space  $\text{Col}(\mathbf{X})$

↳ linear combination of col of  $\mathbf{X}$ .  $x_1 w_1 + x_2 w_2 + \dots$

$\mathbf{e}$  is perpendicular to each row of  $\mathbf{X}^T$   
(i.e., each column of  $\mathbf{X}$ )



$\mathbf{e}$  is perpendicular to the column space of  $\mathbf{X}$   
(i.e.,  $\text{Span}(\mathbf{X})$ )



Least Square Approximation (LSQ(Arb))  $\rightarrow$  from (SOP)

$y$ : actual correct outputs.

$\hat{y}$ : estimated outputs using given feature  $X$  and trained parameters/weights  $w$ .

$e$ : error between predicted and actual  $\Rightarrow \hat{y} - y \Rightarrow \mathbf{X}\mathbf{w} - \mathbf{y}$

We are minimizing the error  $e$  by finding  $w^*$  such that  $\|\mathbf{X}\mathbf{w}^* - \mathbf{y}\| \leq \|\mathbf{X}\mathbf{v} - \mathbf{y}\|$  for any other  $\mathbf{v} \in \mathbb{R}^n$

$\Rightarrow$  can find  $w^*$  the lsq solution  
such that  $\hat{y}$  is predicted best.

# Linear Regression (Scalar Function)

## ➤ Prediction (Testing)

$$f_{\mathbf{w}^*}(\mathbf{X}_{new}) = \mathbf{X}_{new}\mathbf{w}^* = \begin{bmatrix} x_{1,1}^{new} & \dots & x_{1,d}^{new} \\ \vdots & \ddots & \vdots \\ x_{m,1}^{new} & \dots & x_{m,d}^{new} \end{bmatrix} \begin{bmatrix} w_1^* \\ \vdots \\ w_d^* \end{bmatrix}$$

$$= \begin{bmatrix} (\mathbf{x}_1^{new})^T \mathbf{w}^* \\ \vdots \\ (\mathbf{x}_m^{new})^T \mathbf{w}^* \end{bmatrix} \Rightarrow$$

$y$  the predicted output vector.  
if  $w^*$  is best,  $y$  should be close to  $y$  by minimum squared errors.

How to ensure solution  $\mathbf{w}^*$  is the global minimum  
i.e. the true minimum best method  
to get best  $y = \mathbf{X}\mathbf{w}^*$

Is the obtained  $\mathbf{w}^*$  a global minimum?

Yes! sum of squared errors is convex  
solution  $\Rightarrow$  any local min is global min.

# Linear Regression (Scalar Function)



Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using least squares regression.

Step 1: Learning (work out  $\mathbf{w}^*$ )

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Since  $\mathbf{X}^T \mathbf{X} = [242]$  is invertible

$$= [242]^{-1} [-10, -8, -3, -1, 2, 8]$$

$$= -0.3512$$

$$\begin{bmatrix} 5 \\ 5 \\ 4 \\ 3 \\ 2 \\ 2 \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} -10 \\ -8 \\ -3 \\ -1 \\ 2 \\ 8 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 5 \\ 5 \\ 4 \\ 3 \\ 2 \\ 2 \end{bmatrix}$$

Step 2: Prediction

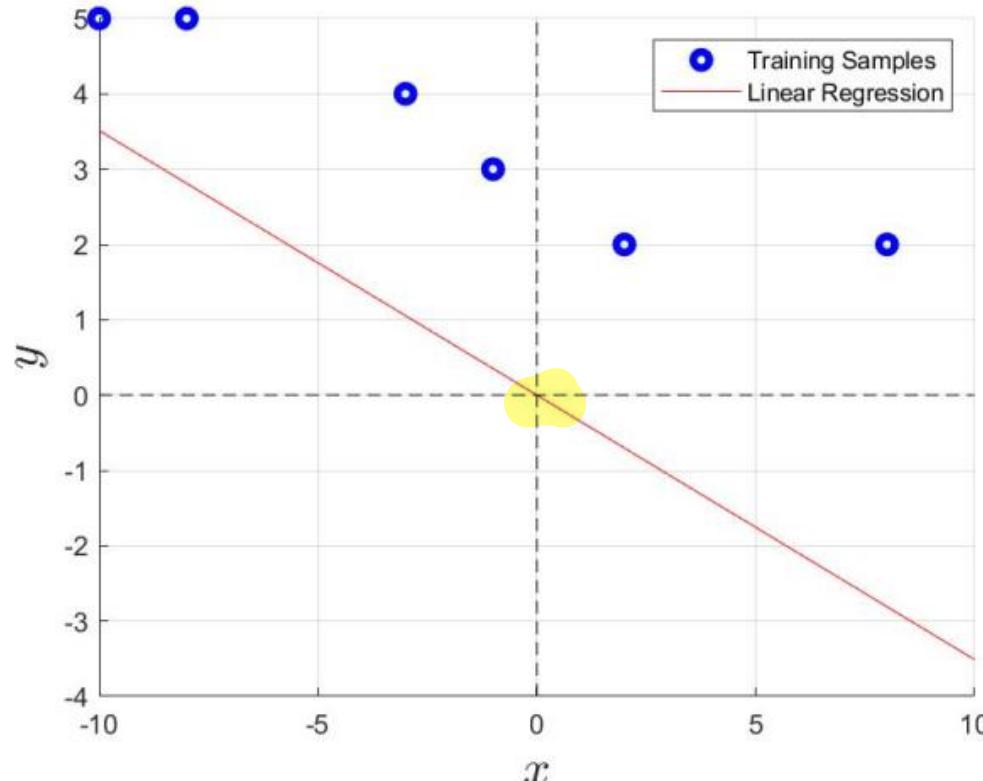
$$\mathbf{y}_t = \mathbf{f}_{\mathbf{w}^*}(\mathbf{X}_t) = \mathbf{x}_t \mathbf{w}^* = \begin{bmatrix} 7 \\ 9 \end{bmatrix} [-0.3512]$$

$$= \begin{bmatrix} -2.4587 \\ -3.1612 \end{bmatrix}$$

$$\mathbf{x}_t = \begin{bmatrix} 7 \\ 9 \end{bmatrix}$$

# Linear Regression (Scalar Function)

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using least squares regression.



$$f_{\mathbf{w}}(\mathbf{x}) = w_1 x_1 = -0.3512 x_1$$

The line always passes through the origin.

because no offset?

# Linear Regression (Scalar Function)

## ➤ Linear Regression with bias/offset

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,

Output:  $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = \mathbf{x}^T \mathbf{w}$

where  $\mathbf{x}$  is a  $(d+1)$ -dimensional feature vector

$\mathbf{w}$  is a  $(d+1)$ -dimensional vector of parameters

If  $\mathbf{x}$  is a column vector and not matrix  $\Rightarrow$  column vector or vector.

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

extra constant/intercept to allow the model to make non-zero predictions when all inputs are zero.

$$\hat{y} = \mathbf{x} \mathbf{w}^* + b$$

The bias/offset term  $w_0$  is responsible for translating the line/plane/hyperplane away from the origin

# Linear Regression (Scalar Function)

## ➤ Linear Regression with bias/offset

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,

Output:  $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = \mathbf{x}^T \mathbf{w}$

-Learning:  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

↙ see prev. slide  
on why  
written as  
vector of  
weights instead  
of matrix

where  $\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \ddots & & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$

-Prediction:  $\mathbf{f}_{\mathbf{w}^*}(\mathbf{X}_t) = \mathbf{X}_t \mathbf{w}^*$

where  $\mathbf{X}_t = \begin{bmatrix} 1 & x_{1,1}^t & \dots & x_{1,d}^t \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m,1}^t & \dots & x_{m,d}^t \end{bmatrix}, \mathbf{w}^* = \begin{bmatrix} w_0^* \\ \vdots \\ w_d^* \end{bmatrix}$

# Linear Regression (Scalar Function)

## ➤ Linear Regression with bias/offset

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using least squares regression.

$$\mathbf{X} = \begin{bmatrix} 1 & -10 \\ 1 & -8 \\ 1 & -3 \\ 1 & -1 \\ 1 & 2 \\ 1 & 8 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 5 \\ 5 \\ 4 \\ 3 \\ 2 \\ 2 \end{bmatrix}$$

$$\mathbf{x}_t = \begin{bmatrix} 1 & 7 \\ 1 & 9 \end{bmatrix}$$

### Step 1: Learning (work out $\mathbf{w}^*$ )

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Since  $\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 6 & -12 \\ -12 & 242 \end{bmatrix}$  is invertible

$$= \begin{bmatrix} 6 & -12 \\ -12 & 242 \end{bmatrix}^{-1} \begin{bmatrix} 1 & -10 \\ 1 & -8 \\ 1 & -3 \\ 1 & -1 \\ 1 & 2 \\ 1 & 8 \end{bmatrix}^T \begin{bmatrix} 5 \\ 5 \\ 4 \\ 3 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 3.1055 \\ -0.1972 \end{bmatrix}$$

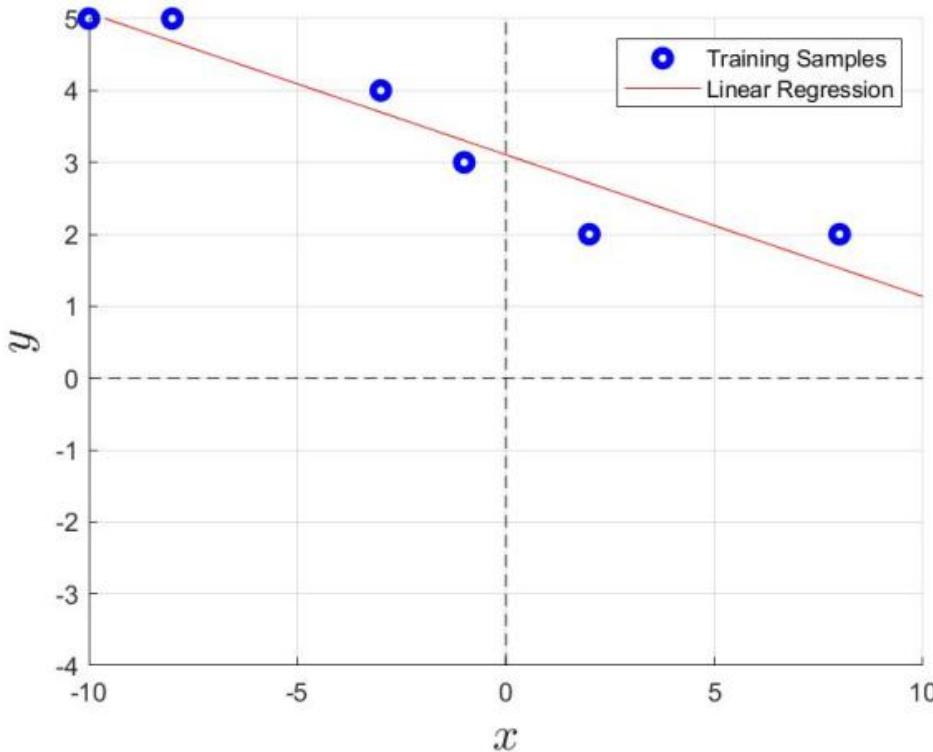
### Step 2: Prediction

$$\begin{aligned} \mathbf{y}_t &= \mathbf{f}_{\mathbf{w}^*}(\mathbf{x}_t) = \mathbf{x}_t \mathbf{w}^* = \begin{bmatrix} 1 & 7 \\ 1 & 9 \end{bmatrix} \begin{bmatrix} 3.1055 \\ -0.1972 \end{bmatrix} \\ &= \begin{bmatrix} 1.7248 \\ 1.3303 \end{bmatrix} \end{aligned}$$

# Linear Regression (Scalar Function)

## ➤ Linear Regression with bias/offset

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using least squares regression.

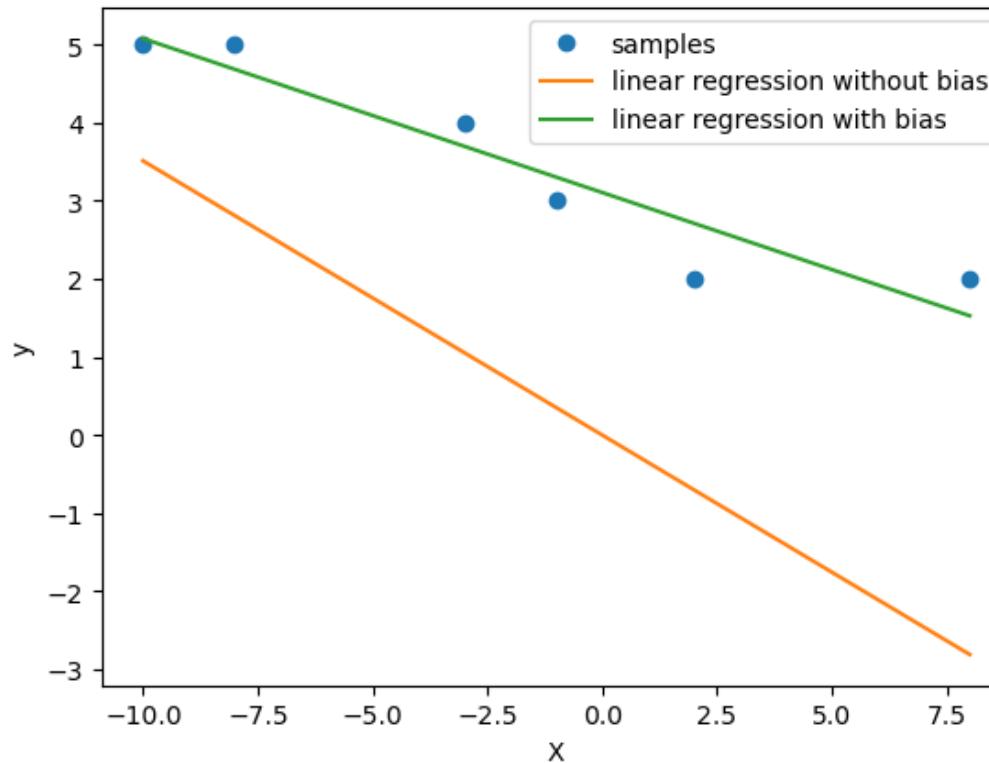


$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 \\ = 3.1055 - 0.1972 x_1$$

# Linear Regression (Scalar Function)

## ➤ Comparison

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using least squares regression.



# Linear Regression (Vector Function)

## ➤ Extend to Multiple Outputs

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ,

where  $\mathbf{x}$  is a  $(d+1)$ -dimensional feature vector  
 $\mathbf{y}$  is a  $h$ -dimensional target vector

e.g., predict multiple physiological signals

	Age	Weight	Physical activity level	Oxygen saturation	Heart rate	Blood pressure	Respiration rate
1							
...							

Output:  $\mathbf{f}_{\mathbf{W}}(\mathbf{x}) = \mathbf{x}^T \mathbf{W}$ , where  $\mathbf{W}$  is a matrix of size  $(d+1) \times h$

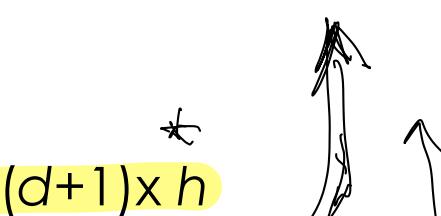
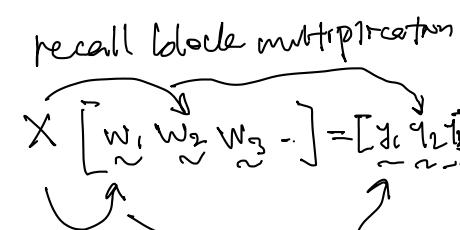
$$\mathbf{W} = \begin{bmatrix} w_{0,1} & \dots & w_{0,h} \\ \vdots & \ddots & \vdots \\ w_{d,1} & \dots & w_{d,h} \end{bmatrix}$$

weight vector  $\mathbf{w}_1$  for  $y_1$

weight vector  $\mathbf{w}_h$  for  $y_h$

Equivalent to performing Linear Regression for each target output variable ( $h$  times)

$$\mathbf{f}_{\mathbf{W}}(\mathbf{x}) = [f_{\mathbf{w}_1}(\mathbf{x}), \dots, f_{\mathbf{w}_h}(\mathbf{x})]$$



# Linear Regression (Vector Function)

## ➤ Learning

Find the optimal value  $\mathbf{w}^*$  which minimizes:

$$\begin{aligned} J(\mathbf{W}) &= \sum_{k=1}^h \sum_{i=1}^N (f_{\mathbf{w}_k}(\mathbf{x}_i) - y_{i,k})^2 \\ &= \sum_{k=1}^h (\mathbf{X}\mathbf{w}_k - \mathbf{y}_k)^T (\mathbf{X}\mathbf{w}_k - \mathbf{y}_k) \end{aligned}$$

*view as matrix - multiplication  
⇒ matrix requires.  
any concept remains the same*

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix}$$

The optimal value  $\mathbf{W}^*$  should have **zero gradient**:

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial J}{\partial w_{0,1}} \\ \vdots \\ \frac{\partial J}{\partial w_{d,1}} \end{bmatrix} \quad \dots \quad \begin{bmatrix} \frac{\partial J}{\partial w_{0,h}} \\ \vdots \\ \frac{\partial J}{\partial w_{d,h}} \end{bmatrix} = \mathbf{0} \quad \rightarrow \quad \frac{\partial J(\mathbf{W})}{\partial \mathbf{w}_k} = \mathbf{0}, k = 1, 2, \dots, h$$

*w.grad = w - αt y \* Appt.*

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{w}_1} \quad \dots \quad \frac{\partial J(\mathbf{W})}{\partial \mathbf{w}_h}$$

$$\mathbf{w}_k^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_k$$

# Linear Regression (Vector Function)

## ➤ Learning

Find the optimal value for  $\mathbf{W}^*$  which minimizes:

$$\begin{aligned} J(\mathbf{W}) &= \sum_{k=1}^h \sum_{i=1}^N (f_{\mathbf{w}_k}(\mathbf{x}_i) - y_{i,k})^2 \\ &= \sum_{k=1}^h (\mathbf{X}\mathbf{w}_k - \mathbf{y}_k)^T (\mathbf{X}\mathbf{w}_k - \mathbf{y}_k) \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \cdots & x_{N,d} \end{bmatrix}$$

$$\begin{aligned} \mathbf{W}^* &= [\mathbf{w}_1^*, \dots, \mathbf{w}_h^*] = [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_1, \dots, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_h] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \end{aligned}$$

$$\mathbf{Y} = \begin{bmatrix} y_{1,1} & \cdots & y_{1,h} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \cdots & y_{N,h} \end{bmatrix}$$

Is the obtained  $\mathbf{W}^*$  a global minimum?

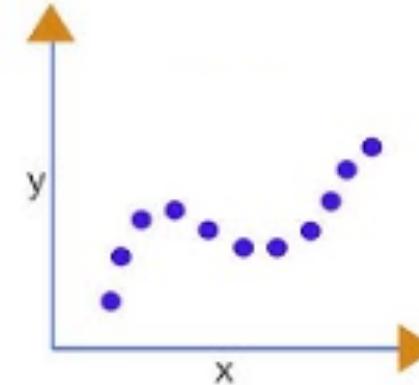
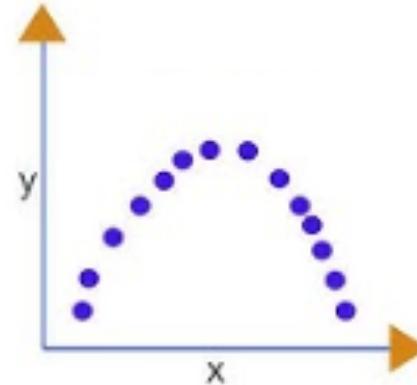
same thing nothing changed.

# Linear Regression (Vector Function)

## ➤ Prediction

$$\mathbf{F}_{\mathbf{W}^*}(\mathbf{X}_t) = \mathbf{X}_t \mathbf{W}^* = \begin{bmatrix} \mathbf{1} & x_{1,1}^t & \dots & x_{1,d}^t \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{1} & x_{m,1}^t & \dots & x_{m,d}^t \end{bmatrix} \begin{bmatrix} w_{0,1}^* & \dots & w_{0,h}^* \\ \vdots & \ddots & \vdots \\ w_{d,1}^* & \dots & w_{d,h}^* \end{bmatrix}$$

$$= \begin{bmatrix} (\mathbf{x}_1^t)^T \mathbf{w}_1^* & \dots & (\mathbf{x}_1^t)^T \mathbf{w}_h^* \\ \vdots & \ddots & \vdots \\ (\mathbf{x}_m^t)^T \mathbf{w}_1^* & \dots & (\mathbf{x}_m^t)^T \mathbf{w}_h^* \end{bmatrix}$$



# Polynomial Regression

non-linear relationship between input and output

# Polynomial Regression

## ➤ Polynomial Model

$$f_w(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k + \dots$$

**1<sup>st</sup> order (Linear)**

Inherent representation of PVA  
as a special form of polynomial regression, when  $n=1$

**2<sup>nd</sup> order (Quadratic)**

**3<sup>rd</sup> order (Cubic)**

e.g.,  $d=1$ , 3<sup>rd</sup> order:  $f_w(\mathbf{x}) = w_0 + w_1 x_1 + w_{11} x_1^2 + w_{111} x_1^3$

$d=2$ , 2<sup>nd</sup> order:  $f_w(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{11} x_1^2 + w_{12} x_1 x_2 + w_{22} x_2^2$

**Order number:  
Hyperparameter!**

(determined,  
size of  
model, set  
by user not  
model)

see  
total

No. parameters for  $n$ -th order polynomial  
model with  $d$  input variables?

→ w-d features

$\{x_1, x_1^2, x_1 x_2, x_2, x_2^2\}$   
 $\{x_1^2, x_1 x_2, x_2^2\}$   
 $\{x_1, x_2\}$   
 $\{x_1, x_2, x_1^2, x_1 x_2, x_2^2\}$   
 $\{x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^3, x_1^2 x_2, x_1 x_2^2, x_2^3\}$

No. parameters for *degree n terms* with  
 $d$  input variables?

total power =  $n$  only

$$\binom{n+d}{d} = \frac{(n+d)!}{n! d!}$$

$$\binom{n+d-1}{d-1} = \frac{(n+d-1)!}{n! (d-1)!}$$

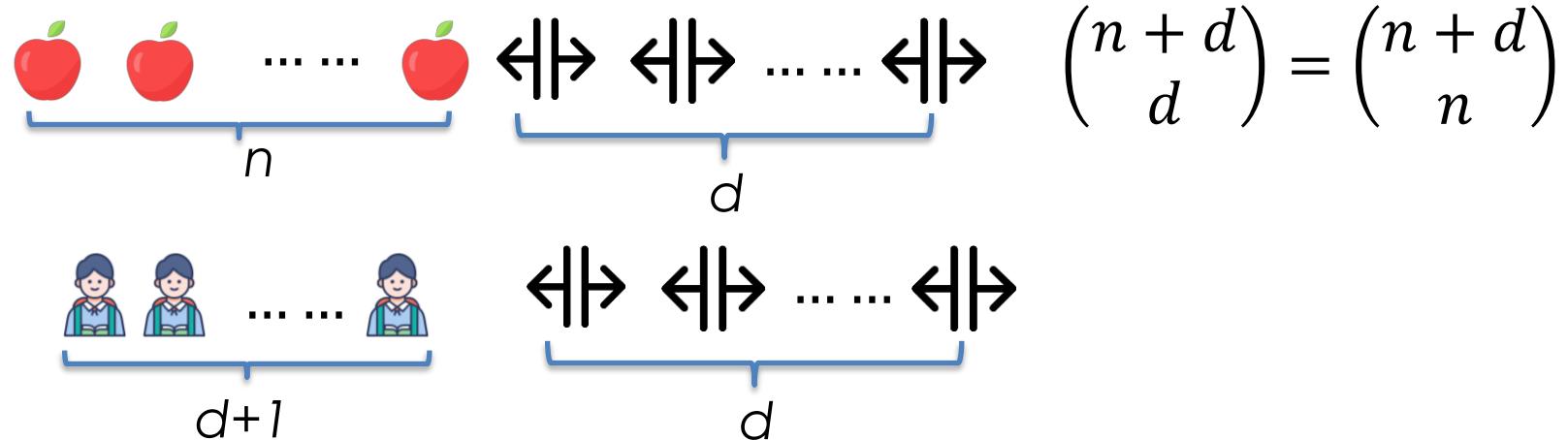
- Proof: total No. paras for  $n$ -th order polynomial model with  $d$  input variables is  $\binom{n+d}{d} = \frac{(n+d)!}{n!d!}$

$$x_1^{k_1} x_2^{k_2} \dots x_d^{k_d}: k_1 + k_2 + \dots + k_d \leq n$$



$$k_1 + k_2 + \dots + k_d + k_{d+1} = n$$

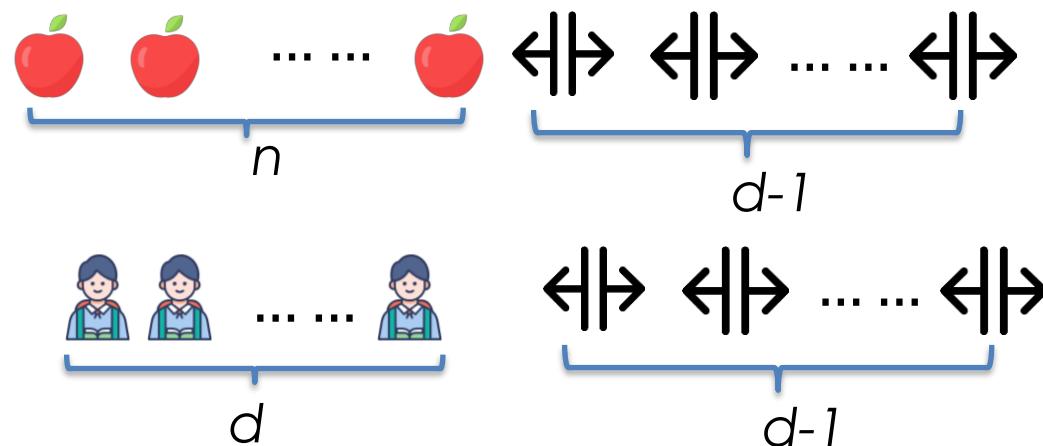
Counting: equivalent to allocating  $n$  apples to  $d+1$  students.



- Proof: No. paras for degree  $n$  terms with  $d$  input variables is  $\binom{n+d-1}{d-1} = \frac{(n+d-1)!}{n!(d-1)!}$

$$x_1^{k_1} x_2^{k_2} \dots x_d^{k_d}: k_1 + k_2 + \dots + k_d = n$$

Counting: equivalent to allocating  $n$  apples to  $d$  students.



$$\binom{n+d-1}{d-1} = \binom{n+d-1}{n}$$

# Polynomial Regression

## ➤ Polynomial Model

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k + \dots$$

$$= f_{\mathbf{w}}(\mathbf{p}(\mathbf{x})) = \mathbf{p}^T \mathbf{w}$$

$$\mathbf{p} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ \vdots \\ x_i x_j \\ \vdots \\ x_i x_j x_k \\ \vdots \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \\ \vdots \\ w_{ij} \\ \vdots \\ w_{ijk} \\ \vdots \end{bmatrix}$$

Annotations:

- Grouping braces:
  - Brace from 1 to  $x_d$ : "0<sup>th</sup> degree linear term"
  - Brace from  $x_i x_j$  to  $x_i x_j x_k$ : "2<sup>nd</sup> degree 3<sup>rd</sup> degree"

Equivalent to:  
polynomial transformation  
+  
linear regression

# Polynomial Regression

## ➤ Learning

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,

Output:  $f_{\mathbf{w}}(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{p}(\mathbf{x})) = \mathbf{p}^T \mathbf{w}$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

Find the optimal value  $\mathbf{w}^*$  which minimizes:

$$J(\mathbf{w}) = \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{p}(\mathbf{x}_i)) - y_i)^2 = (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y})$$

*Some cost functions*  
*↑ just replace X with P which has all degree terms.*

$$\mathbf{P} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} & \dots & x_{1,i}x_{1,j} & \dots & x_{1,i}x_{1,j}x_{1,k} & \dots \\ \vdots & \vdots & & \vdots & & \vdots & & \vdots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} & \dots & x_{N,i}x_{N,j} & \dots & x_{N,i}x_{N,j}x_{N,k} & \dots \end{bmatrix}$$

Solution:  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{y}$ , If  $\mathbf{P}^T \mathbf{P}$  is invertible

yes,  
global  
minimum?

# Polynomial Regression

## ➤ Prediction

$$f_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{w}^*$$

$$= \begin{bmatrix} 1 & x_{1,1}^t & \cdots & x_{1,d}^t & \cdots & x_{1,i}^t x_{1,j}^t & \cdots & x_{1,i}^t x_{1,j}^t x_{1,k}^t & \cdots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{m,1}^t & \cdots & x_{m,d}^t & \cdots & x_{m,i}^t x_{m,j}^t & \cdots & x_{m,i}^t x_{m,j}^t x_{m,k}^t & \cdots \end{bmatrix} \begin{bmatrix} w_0^* \\ w_1^* \\ \vdots \\ w_d^* \\ \vdots \\ w_{ij}^* \\ \vdots \\ w_{ijk}^* \\ \vdots \end{bmatrix}$$

$$= \begin{bmatrix} (\mathbf{p}_1^t)^T \mathbf{w}^* \\ \vdots \\ (\mathbf{p}_m^t)^T \mathbf{w}^* \end{bmatrix}$$

# Polynomial Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using 3<sup>rd</sup> order polynomial regression.

## Step 0: Polynomial Transformation ( $X \rightarrow P$ )

3<sup>rd</sup> order polynomial model with 1 input variable:

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_{11} x_1^2 + w_{111} x_1^3 = [1 \quad x_1 \quad x_1^2 \quad x_1^3] \begin{bmatrix} w_0 \\ w_1 \\ w_{11} \\ w_{111} \end{bmatrix}$$

$$\mathbf{P} = \begin{bmatrix} 1 & x_{1,1} & x_{1,1}^2 & x_{1,1}^3 \\ 1 & x_{2,1} & x_{2,1}^2 & x_{2,1}^3 \\ 1 & x_{3,1} & x_{3,1}^2 & x_{3,1}^3 \\ 1 & x_{4,1} & x_{4,1}^2 & x_{4,1}^3 \\ 1 & x_{5,1} & x_{5,1}^2 & x_{5,1}^3 \\ 1 & x_{6,1} & x_{6,1}^2 & x_{6,1}^3 \end{bmatrix} = \begin{bmatrix} 1 & -10 & 100 & -1000 \\ 1 & -8 & 64 & -512 \\ 1 & -3 & 9 & -27 \\ 1 & -1 & 1 & -1 \\ 1 & 2 & 4 & 8 \\ 1 & 8 & 64 & 512 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 5 \\ 5 \\ 4 \\ 3 \\ 2 \\ 2 \end{bmatrix}$$

*The only  
additional  
step.  
Steps are all the same.*

# Polynomial Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using 3<sup>rd</sup> order polynomial regression.

Step 1: Learning (work out  $\mathbf{w}^*$ )

Since  $\mathbf{P}^T \mathbf{P}$  is invertible,  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{y} = \begin{bmatrix} 2.6894 \\ -0.3772 \\ 0.0134 \\ 0.0029 \end{bmatrix}$

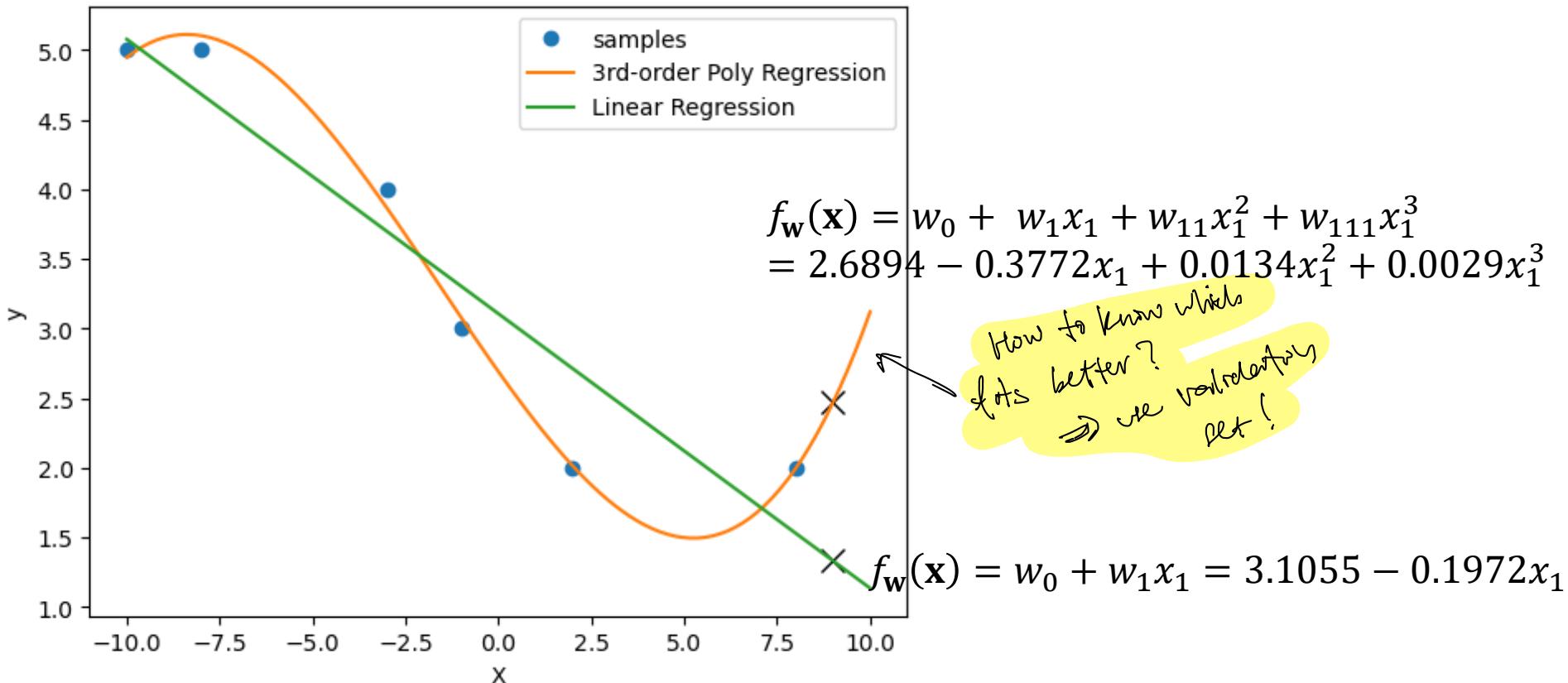
Step 2: Prediction

$$\mathbf{P}_t = \begin{bmatrix} 1 & x_{1,1}^t & (x_{1,1}^t)^2 & (x_{1,1}^t)^3 \\ 1 & x_{2,1}^t & (x_{2,1}^t)^2 & (x_{2,1}^t)^3 \end{bmatrix} = \begin{bmatrix} 1 & 7 & 49 & 343 \\ 1 & 9 & 81 & 729 \end{bmatrix}$$

$$\begin{aligned} \mathbf{y}_t &= \mathbf{f}_{\mathbf{w}^*}(\mathbf{P}_t) = \mathbf{P}_t \mathbf{w}^* = \begin{bmatrix} 1 & 7 & 49 & 343 \\ 1 & 9 & 81 & 729 \end{bmatrix} \begin{bmatrix} 2.6894 \\ -0.3772 \\ 0.0134 \\ 0.0029 \end{bmatrix} \\ &= \begin{bmatrix} 1.6874 \\ 2.4661 \end{bmatrix} \end{aligned}$$

# Polynomial Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^6 = \{(-10, 5), (-8, 5), (-3, 4), (-1, 3), (2, 2), (8, 2)\}$ , predict the test dataset  $\{7, 9\}$  using 3<sup>rd</sup> order polynomial regression.



# Polynomial Regression

## ➤ Extend to Multiple Outputs

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ,

where  $\mathbf{x}$  is a  $d$ -dimensional feature vector  
 $\mathbf{y}$  is a  $h$ -dimensional target vector

Output:  $\mathbf{f}_{\mathbf{W}}(\mathbf{x}) = \mathbf{f}_{\mathbf{W}}(\mathbf{p}(\mathbf{x})) = \mathbf{p}^T \mathbf{W}$ ,

*w become matrix*

-Learning:  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Y}$

-Prediction:  $\mathbf{F}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{W}^*$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_h \end{bmatrix}$$

# Ridge Regression

If  $\mathbf{X}^T \mathbf{X}$  or  $\mathbf{P}^T \mathbf{P}$  is not invertible

# Ridge Regression

➤ Ridge Regression = Regression + L2 regularization

Shrinks the regression coefficients  $w$  by imposing a penalty to their size.

recall EE2211, reduce weights to prevent overfitting.

go to slide 42.

Find the optimal value  $\mathbf{w}^*$  which minimizes:

$$J(\mathbf{w}) = \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{p}(\mathbf{x}_i)) - y_i)^2 + \lambda \sum_i w_i^2$$

L2 regularization:  
 $\sum_i w_i^2 = \|\mathbf{w}\|_2^2$ ,  
squared L2 norm of  $\mathbf{w}$

$$= (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

lambda places more weight on Regularisation term. Means minimisation on actual loss function is lesser. Increasing lambda results in greater losses, but smaller bias and variance ( $\mathbf{W}$  is lesser, but this also means its less accurate)

$\uparrow \lambda \Rightarrow \downarrow \mathbf{w}$  magnitude

(Hyperparameter)

$\lambda > 0$  is a complexity/regularization parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage.

# Ridge Regression

## ➤ Learning

$$\operatorname{argmin}_{\mathbf{w}} (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

Solution:  $\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \mathbf{0}$

$$\frac{d ((\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w})}{d\mathbf{w}} = \mathbf{0}$$

$$2\mathbf{P}^T \mathbf{P}\mathbf{w} - 2\mathbf{P}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0}$$

$$\mathbf{P}^T \mathbf{P}\mathbf{w} + \lambda \mathbf{w} = \mathbf{P}^T \mathbf{y}$$

$$(\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})\mathbf{w} = \mathbf{P}^T \mathbf{y}$$

Since  $\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}$  is always invertible,  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{y}$

# Ridge Regression

**Proof:**  $\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}$  is always invertible when  $\lambda > 0$

*In contradiction* Suppose  $\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}$  is not invertible.

Columns/rows in  $\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}$  are linearly dependent

There exists a vector  $\mathbf{a} \neq \mathbf{0}$ , s.t.,  $(\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})\mathbf{a} = \mathbf{0}$

$$\mathbf{a}^T (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}) \mathbf{a} = 0$$

$$\mathbf{a}^T \mathbf{P}^T \mathbf{P} \mathbf{a} + \lambda \mathbf{a}^T \mathbf{a} = 0$$

$$\|\mathbf{P}\mathbf{a}\|^2 + \lambda \|\mathbf{a}\|^2 = 0$$

However,  $\|\mathbf{P}\mathbf{a}\|^2 + \lambda \|\mathbf{a}\|^2 = 0$  only holds if  $\|\mathbf{P}\mathbf{a}\|^2 = \|\mathbf{a}\|^2 = 0$

$$\mathbf{a} = \mathbf{0}$$

Columns/rows in  $\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I}$  are linearly independent

$\Rightarrow$  full rank  $\Rightarrow$  invertible

# Ridge Regression

## ➤ Learning

### Geometric Interpretation of Ridge Regression Solution

$$\underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

Lagrangian Duality  $\parallel$

for constrained optimization problems  
in this form,

$$\underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}) + \lambda g(\mathbf{x})$$

$f(x)$   $g(x)$

$$\underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y}), \text{ subject to } \mathbf{w}^T \mathbf{w} \leq c$$

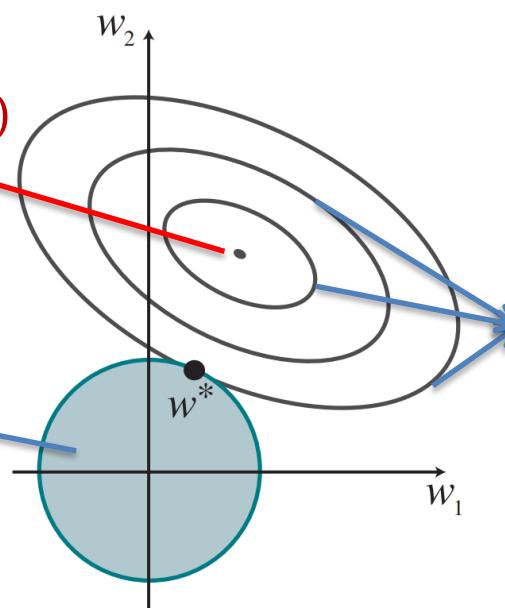
Primal form

$$\underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}), \text{ s.t. } g(\mathbf{x}) \leq 0$$

$$\min_{\mathbf{w}} (\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y})$$

Set of points  $\mathbf{w}$  in 2D weight space that have  $\mathbf{w}^T \mathbf{w} \leq c$

$$\text{radius} = \sqrt{c}$$



Contours of  $(\mathbf{P}\mathbf{w} - \mathbf{y})^T (\mathbf{P}\mathbf{w} - \mathbf{y})$

# Ridge Regression

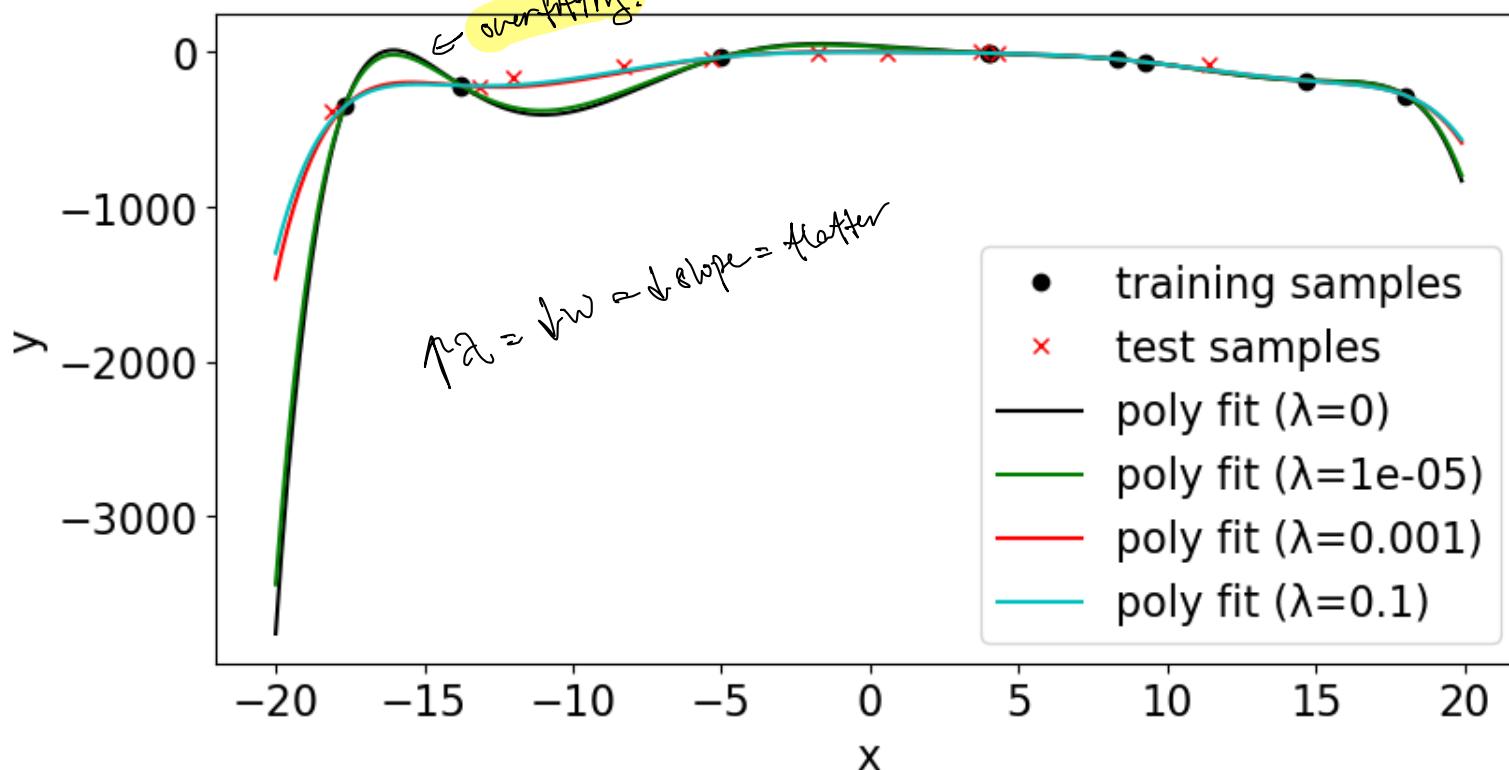
## ➤ Prediction

$$f_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{w}^*$$

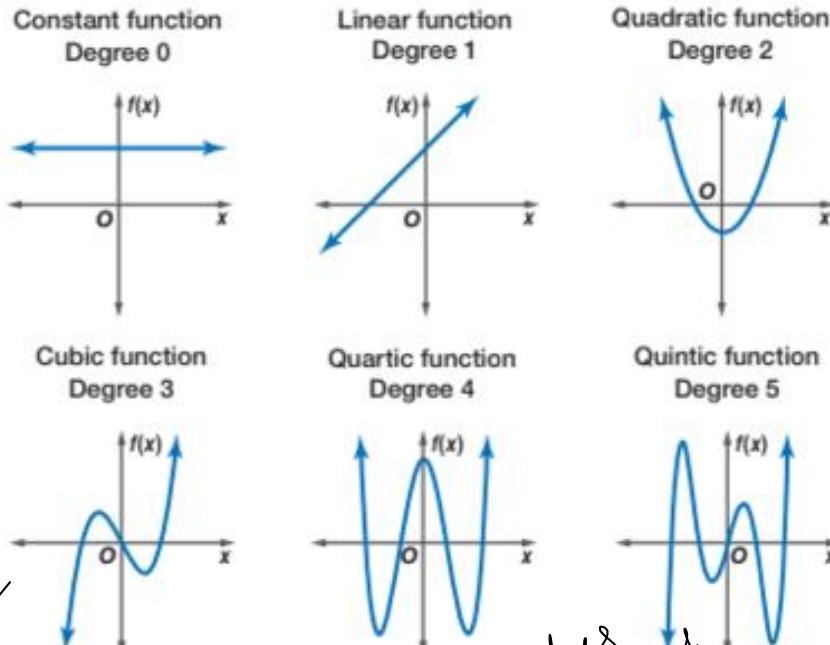
How is weight shrinkage connected with  $\mathbf{P}^T \mathbf{P}$  being not invertible?

Example:

Order = 8, no. of input variable = 1



## Polynomial model with different degrees



When is  $\mathbf{P}^T \mathbf{P}$  not invertible?

$$\textcircled{1} \quad \text{Rank}(\mathbf{P}) = \text{Rank}(\mathbf{P}^T \mathbf{P})$$

Can be proved by **rank-nullity theorem**

$$\text{Rank}(\mathbf{A}) + \text{nullity}(\mathbf{A}) = \text{No. cols}$$

\textcircled{2} more no. of rows / cols that are linearly independent  $\Rightarrow$  dimension of column space , and each column represents a feature.

\textcircled{3}

↓W.

Regularisation has similar effect of reducing the polynomial degree/features (flatter).

**Flattening**  
(Fewer ups and downs)

**Reducing degree**  
(Reduce no. polynomial features)

From \textcircled{1}-\textcircled{6}: rank(\mathbf{P}) rep. no. of unique features that are just enough to describe the data

polynomial features: actual features we use to describe the data

If \textcircled{7}, we used features that are combination/dependent on other features.

\textcircled{8} highly correlated features (redundancy)  
↳ regularisation doesn't fix features, but it suppresses the effective contribution of dependent feature  $\Rightarrow$  smaller effect as removes redundant poly features

\textcircled{9} rank(\mathbf{P}) < no. polynomial features

\textcircled{10} unique features of data

\textcircled{11}

**Intrinsic dimension of the feature space**

\textcircled{12}

\textcircled{13}

# Prove $\text{Rank}(P) = \text{Rank}(P^T P)$ by rank-nullity theorem



Prove  $\text{Nullity}(P) = \text{Nullity}(P^T P)$



Prove

- if  $x$  satisfies  $Px = 0, P^T Px = 0$
- and
- if  $x$  satisfies  $P^T Px = 0, Px = 0$

$$x^T P^T P x = 0$$



$$\|Px\|^2 = 0$$



$$Px = 0$$

# Ridge Regression

## ➤ Extend to Multiple Outputs

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ,

where  $\mathbf{x}$  is a  $d$ -dimensional feature vector  
 $\mathbf{y}$  is a  $h$ -dimensional target vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_h \end{bmatrix}$$

Output:  $\mathbf{f}_{\mathbf{W}}(\mathbf{x}) = \mathbf{f}_{\mathbf{W}}(\mathbf{p}(\mathbf{x})) = \mathbf{p}^T \mathbf{W}$ ,

-Learning:

Find the optimal values for  $\mathbf{W}^*$  which minimizes:

$$J(\mathbf{W}) = \sum_{k=1}^h \left( \sum_{i=1}^N (f_{\mathbf{w}_k}(\mathbf{p}(\mathbf{x}_i)) - y_{i,k})^2 + \lambda \sum_i w_{i,k}^2 \right)$$

$$= \sum_{k=1}^h ((\mathbf{P}\mathbf{w}_k - \mathbf{y}_k)^T (\mathbf{P}\mathbf{w}_k - \mathbf{y}_k) + \lambda \mathbf{w}_k^T \mathbf{w}_k)$$

$$\mathbf{Y} = \begin{bmatrix} y_{1,1} & \cdots & y_{1,h} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \cdots & y_{N,h} \end{bmatrix}$$

**Solution in Compact Form:**  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{Y}$

↓ now a matrix      ↓ now a matrix

# Ridge Regression

## ➤ Extend to Multiple Outputs

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ,

where  $\mathbf{x}$  is a  $d$ -dimensional feature vector  
 $\mathbf{y}$  is a  $h$ -dimensional target vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_h \end{bmatrix}$$

Output:  $\mathbf{f}_{\mathbf{W}}(\mathbf{x}) = \mathbf{f}_{\mathbf{W}}(\mathbf{p}(\mathbf{x})) = \mathbf{p}^T \mathbf{W}$ ,

-Learning:  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{Y}$

-Prediction:  $\mathbf{F}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{W}^*$



# Evaluation Metrics for Regression

# Evaluation Metrics for Regression

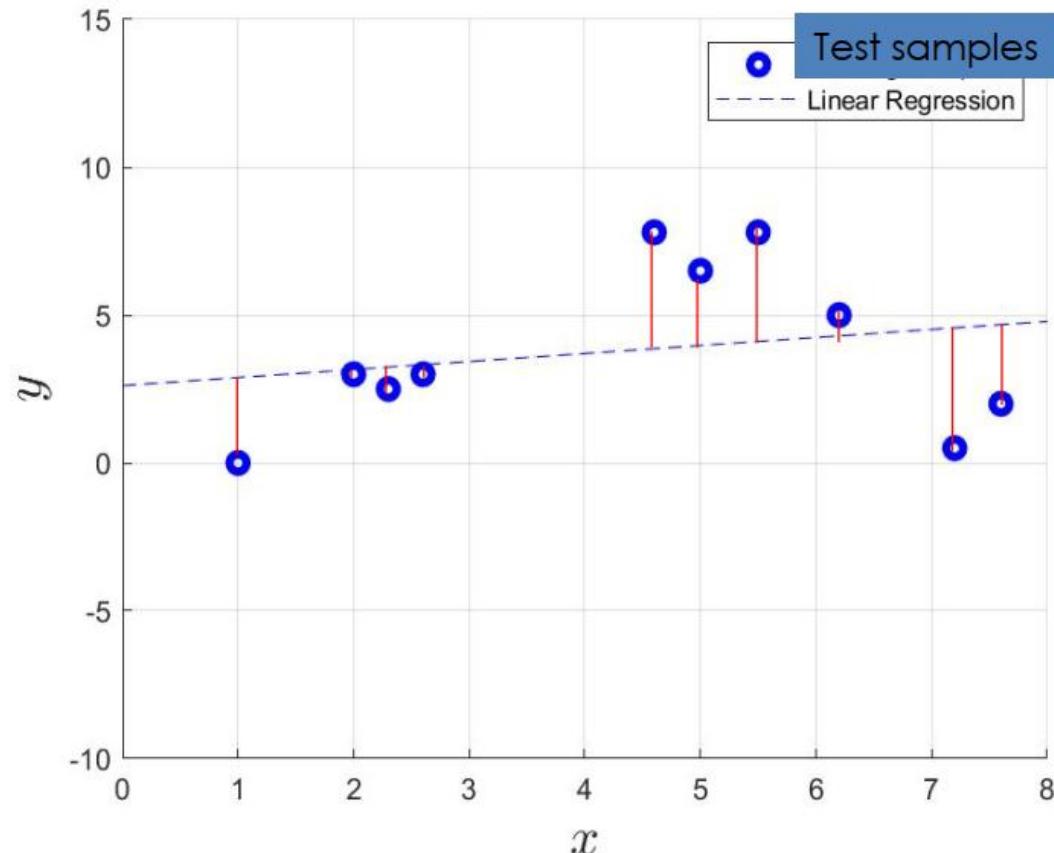
## Mean Square Error

$$(MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n})$$

## Mean Absolute Error

$$(MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n})$$

where  $y_i$  denotes the target output and  $\hat{y}_i$  denotes the predicted output for sample  $i$ .



```
from sklearn.metrics import mean_absolute_error, mean_squared_error

mae = mean_absolute_error(y_true, y_pred)
mse = mean_squared_error(y_true, y_pred)
```

# Summary

## ✓ Linear Regression

Scalar output for each sample

-Learning:  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

-Prediction:  $\mathbf{f}_{\mathbf{w}^*}(\mathbf{X}_t) = \mathbf{X}_t \mathbf{w}^*$

Vector output for each sample

-Learning:  $\mathbf{W}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

-Prediction:  $\mathbf{F}_{\mathbf{W}^*}(\mathbf{X}_t) = \mathbf{X}_t \mathbf{W}^*$

## ✓ Polynomial Regression

Scalar output for each sample

-Learning:  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{y}$

-Prediction:  $\mathbf{f}_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{w}^*$

Vector output for each sample

-Learning:  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Y}$

-Prediction:  $\mathbf{F}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{W}^*$

## ✓ Ridge Regression

Scalar output for each sample

-Learning:  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{y}$

-Prediction:  $\mathbf{f}_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{w}^*$

Vector output for each sample

-Learning:  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P} + \lambda \mathbf{I})^{-1} \mathbf{P}^T \mathbf{Y}$

-Prediction:  $\mathbf{F}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \mathbf{P}_t \mathbf{W}^*$

## ✓ 2 Evaluation Metrics for Regression Task

mean squared error, mean absolute error

# EE2213 Introduction to AI

## Lecture 9: Classification

Dr. WANG Si

[si.wang@nus.edu.sg](mailto:si.wang@nus.edu.sg)

Electrical and Computer Engineering Department  
National University of Singapore

# OVERVIEW OF COURSE CONTENTS

- **Introduction (Shaojing)**

- What is AI
- Applications of AI
- AI agent

- **Search (Shaojing)**

- Uninformed search algorithms: breadth-first, depth-first, uniform-cost(Dijkstra's algorithm)
- Informed search algorithms: greedy best-first, A\*

- Applications

- **Optimisation (Shaojing)**

- Linear programming
- Convex problems
- Applications

- **Machine learning (Wang Si) (Weeks 6-9)**

- Supervised and unsupervised learning: regression, classification, clustering
- Neural networks and deep learning
- Applications

- **Knowledge representation (Wang Si) (Weeks 10-11)**

- Knowledge Representation and Reasoning
  - Propositional Logic
  - Applications
- **Ethical considerations (Shaojing)**
- Bias in AI
  - Privacy concerns
  - Societal impact

(Week 10, Monday)

No Class on 20 Oct.! A recording will be uploaded on Canvas.

# AGENDA

- We will discuss:
    - Polynomial Regression for Classification
    - Logistic Regression
    - Evaluation Metrics for Classification
- not recommended -*
- use this*

At the end of this lecture, you should be able to:

- ✓ Explain the underlying principles of logistic regression
- ✓ Apply the two algorithms for classification
- ✓ Evaluate the classification performance

# Recap: Classification Task

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Output:  $f(\mathbf{x})$  that predict **categorical**  $y$  given  $\mathbf{x}$

Applications:

Spam Detection



Sentiment Analysis

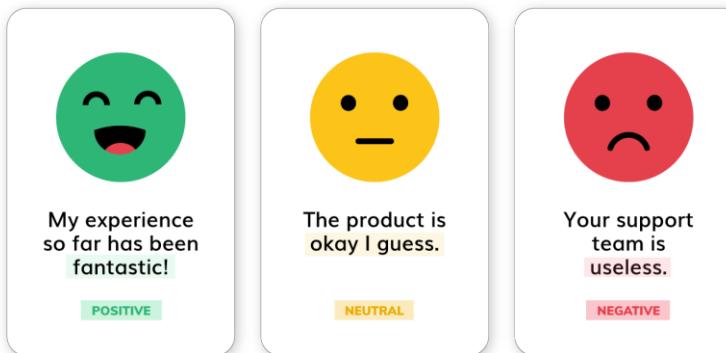
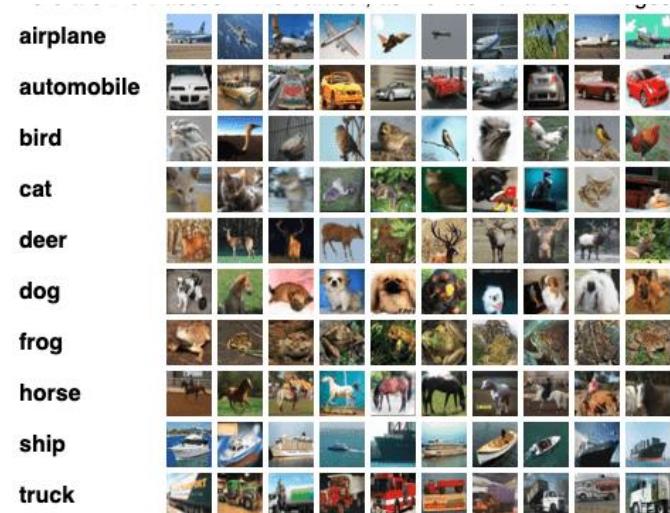


Image Classification



That is life	English
C'est la vie	Portuguese
Yahee jeevan hai	Hindi
Bubomi obo	isiXhosa
それが人生です	Japanese

Language  
Identification

(Linear regression included)

# **Polynomial Regression for**

## **Classification?**

# Polynomial Regression for Classification?

## ➤ Binary Classification

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Output:  $f_{\mathbf{w}}(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{p}(\mathbf{x})) = \text{sign}(\mathbf{p}^T \mathbf{w})$

Threshold function

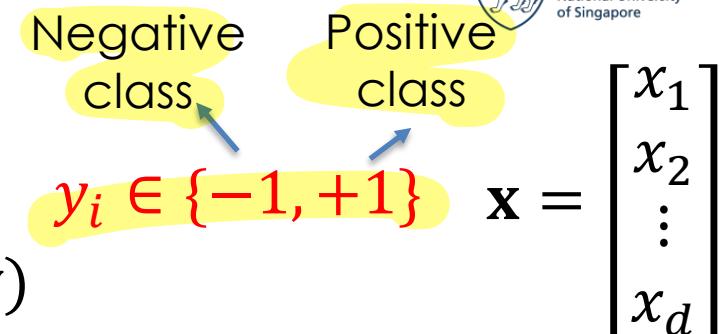
If  $\mathbf{P}^T \mathbf{P}$  is invertible, then

-Learning:  $\mathbf{w}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{y}$  CHM the same

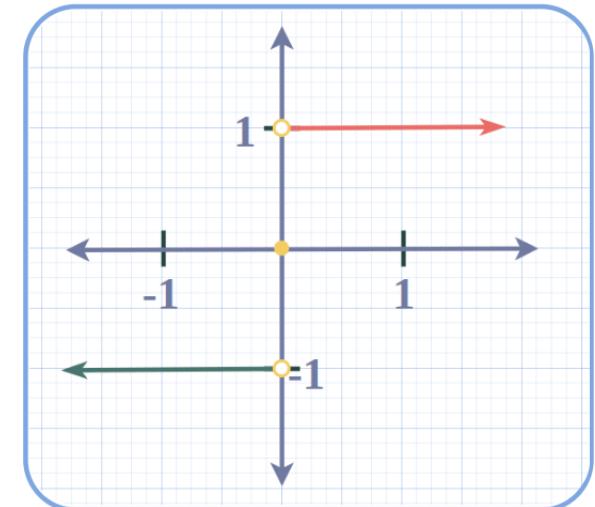
-Prediction:  $\mathbf{f}_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \text{sign}(\mathbf{P}_t \mathbf{w}^*)$

input defined?  
 $(\mathbf{p}_i^t)^T \mathbf{w}^* = 0?$

$$= \begin{bmatrix} \text{sign}((\mathbf{p}_1^t)^T \mathbf{w}^*) \\ \vdots \\ \text{sign}((\mathbf{p}_m^t)^T \mathbf{w}^*) \end{bmatrix}$$



Graph of Signum Function



$$\text{sign}(z) = \begin{cases} -1, z < 0 \\ 0, z = 0 \\ +1, z > 0 \end{cases}$$

# Polynomial Regression for Classification?

## ➤ Binary Classification

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Output:  $f_{\mathbf{w}}(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{p}(\mathbf{x})) = \text{sign}(\mathbf{p}^T \mathbf{w})$

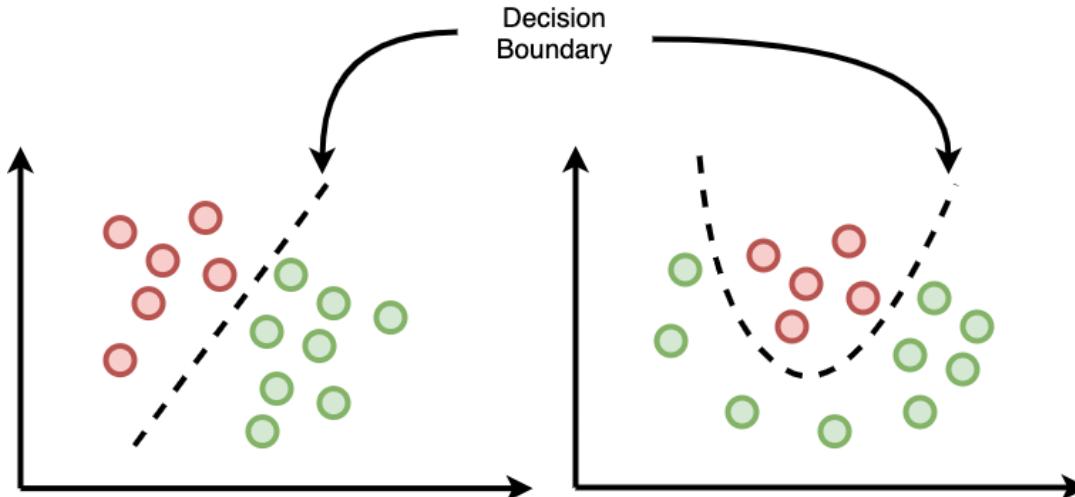
Negative class  
Positive class

$$y_i \in \{-1, +1\}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

### ■ Decision Boundary: $\mathbf{p}^T \mathbf{w} = 0$

(a line or a surface separates the two classes)



$$\text{sign}(z) = \begin{cases} -1, z < 0 \\ 0, z = 0 \\ +1, z > 0 \end{cases}$$

What if threshold function

$$th(z) = \begin{cases} -1, z < 0.3 \\ 0, z = 0.3 \\ +1, z > 0.3 \end{cases}$$

$\mathbf{p}^T \mathbf{w} = 0.3$  vs the decision boundary

# Polynomial Regression for Classification?

## ➤ Multi-category Classification

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$

$\mathbf{y}_i$  is K-Dimensional one-hot encoded vector

e.g.,  $K=3$  categories: dog, cat, rabbit.

One-hot encoding

dog=[1,0,0], cat=[0,1,0], rabbit=[0,0,1]

Output:  $f_{\mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}}(\mathbf{p}(\mathbf{x})) = \underset{k=1, \dots, K}{\operatorname{argmax}}(\mathbf{p}^T \mathbf{W})$

look at which  
k gave the largest  
value. e.g. [0.1, -0.4, 3]  
class 3

If  $\mathbf{P}^T \mathbf{P}$  is invertible, then

-Learning:  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Y}$

-Prediction:

$$\mathbf{f}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \begin{bmatrix} \underset{k=1, \dots, K}{\operatorname{argmax}}((\mathbf{p}_1^t)^T \mathbf{W}^*) \\ \vdots \\ \underset{k=1, \dots, K}{\operatorname{argmax}}((\mathbf{p}_m^t)^T \mathbf{W}^*) \end{bmatrix}$$

How does decision  
boundary look like?

# Polynomial Regression for Classification?

## ➤ Multi-category Classification

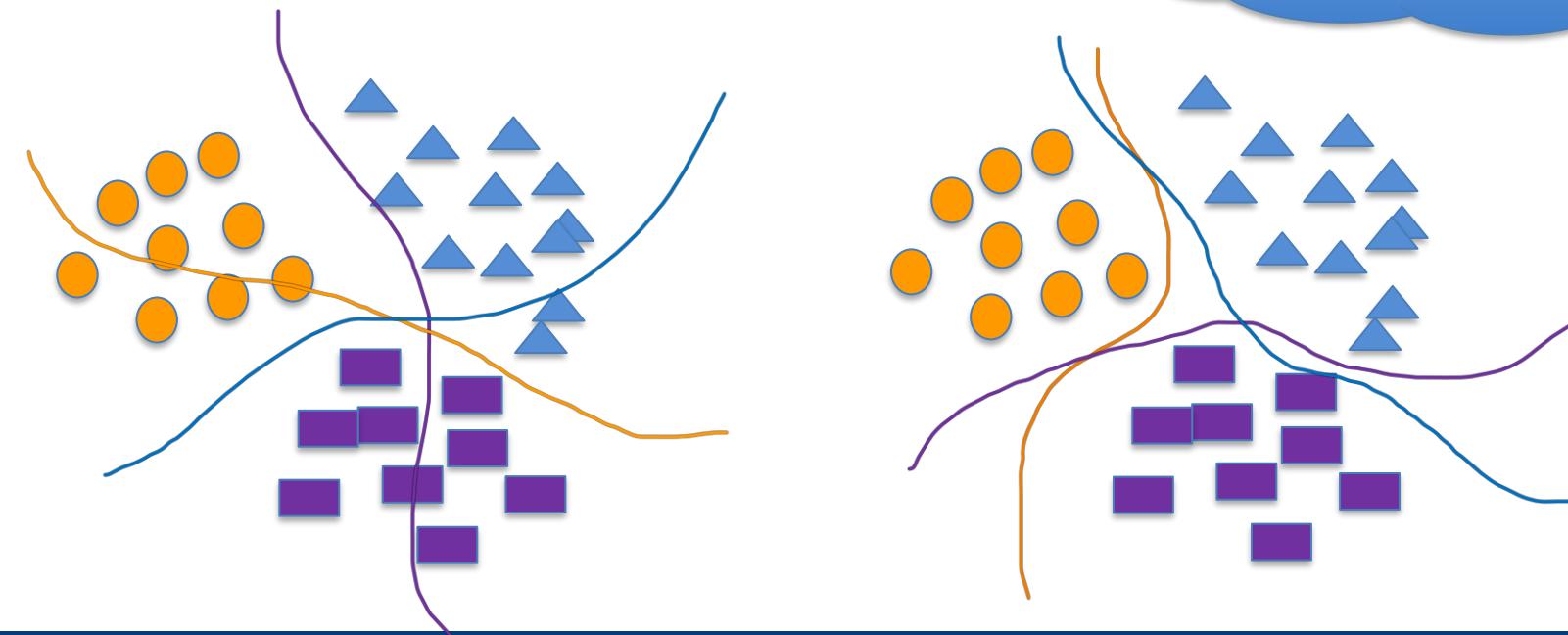
Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$

$\mathbf{y}_i$  is  $K$ -Dimensional one-hot encoded vector

Output:  $f_{\mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}}(\mathbf{p}(\mathbf{x})) = \underset{k=1, \dots, K}{\text{argmax}}(\mathbf{p}^T \mathbf{W})$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

How does decision boundary look like?



# Polynomial Regression for Classification?

## ➤ Multi-category Classification

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the class label for the test dataset  $\{-0.1, 0.4\}$  using 1<sup>st</sup> order polynomial regression.

↳ Linear regression

2

# Polynomial Regression for Classification?

## ➤ Multi-category Classification

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the class label for the test dataset  $\{-0.1, 0.4\}$  using 1<sup>st</sup> order polynomial regression.

### Step 0: Polynomial Transformation ( $\mathbf{X} \rightarrow \mathbf{P}$ ) and One-hot Encoding

1<sup>st</sup> order polynomial model:

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 = [1 \quad x_1] \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\mathbf{P} = \begin{bmatrix} 1 & -1 \\ 1 & 0 \\ 1 & 0.5 \\ 1 & 0.3 \\ 1 & 0.8 \end{bmatrix}$$

1  
class1: [1,0,0]  
class2: [0,1,0]  
class3: [0,0,1]

$$\mathbf{Y} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

# Polynomial Regression for Classification?

## ➤ Multi-category Classification

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the class label for the test dataset  $\{-0.1, 0.4\}$  using 1<sup>st</sup> order polynomial regression.

### Step 1: Learning (work out $\mathbf{W}^*$ )

Since  $\mathbf{P}^T \mathbf{P}$  is invertible,  $\mathbf{W}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Y} = \begin{bmatrix} 0.4780 & 0.3333 & 0.1887 \\ -0.6499 & 0.5556 & 0.0943 \end{bmatrix}$

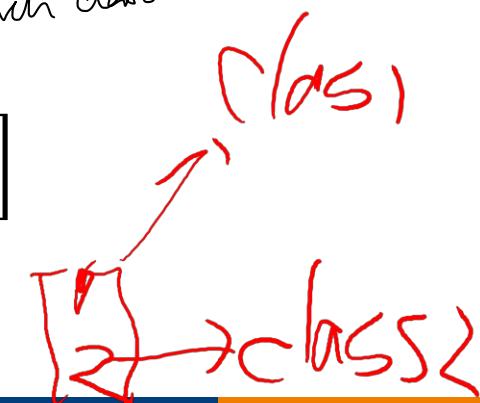
### Step 2: Prediction

$$\mathbf{P}_t = \begin{bmatrix} 1 & x_{1,1}^t \\ 1 & x_{2,1}^t \end{bmatrix} = \begin{bmatrix} 1 & -0.1 \\ 1 & 0.4 \end{bmatrix}$$

$$\mathbf{Y}_t^* = \mathbf{P}_t \mathbf{W}^* = \begin{bmatrix} 1 & -0.1 \\ 1 & 0.4 \end{bmatrix} \begin{bmatrix} 0.4780 & 0.3333 & 0.1887 \\ -0.6499 & 0.5556 & 0.0943 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5430 & 0.2778 & 0.1792 \\ 0.2180 & 0.5556 & 0.2264 \end{bmatrix}$$

argmax



# Polynomial Regression for Classification?

## ➤ Not Recommended for Classification

1. Hard Threshold Function, e.g., "sign()"

Completely confident prediction

$$\mathbf{P}_t \mathbf{w}^* = \begin{bmatrix} -5 \\ 10 \\ 0.01 \end{bmatrix} \xrightarrow{\text{sign}()} \hat{\mathbf{y}}_t = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

Near decision boundary

2.  $\mathbf{p}^T \mathbf{w}$  is a number between  $-\infty$  and  $+\infty$ .

The regression model output values are not interpretable

$$\mathbf{P}_t \mathbf{w}^* = \begin{bmatrix} -4 \\ 1 \\ 4.2 \\ -0.8 \\ 2.5 \\ 0.2 \end{bmatrix}$$

(With **soft** threshold function)

# Logistic Regression

-Used for classification task, not regression task

# Binomial Logistic Regression

## ➤ Model (Binary)

Input:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ ,  $y_i \in \{0, 1\}$

Output:  $f_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w})$ , where  $\sigma(z) = \frac{1}{1+e^{-z}}$

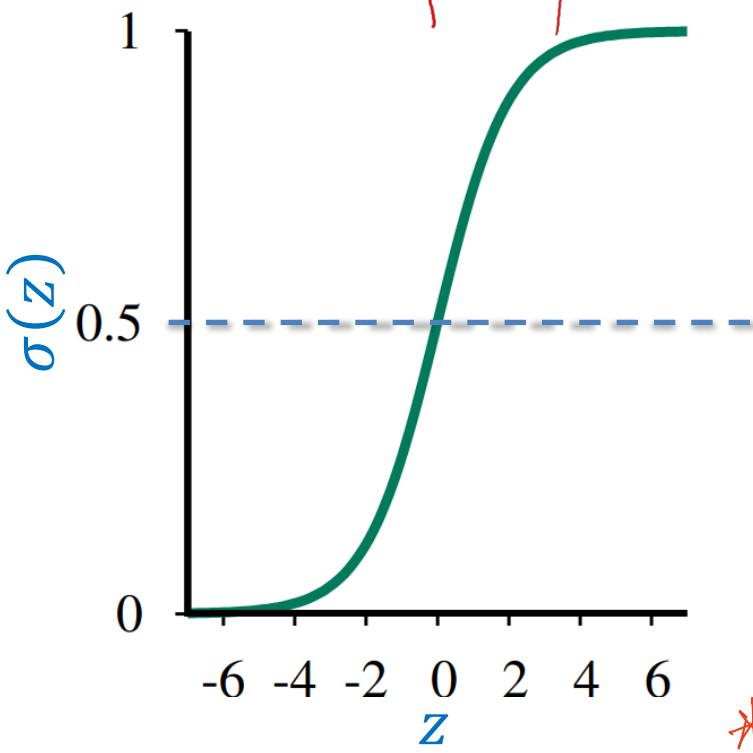
Negative class      Positive class

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

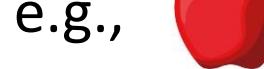
Sigmoid (logistic) function

- Value range:  $(0, 1)$  → Probability

asymptotic  
cannot exactly 0 or 1



Label ( $y_i$ ):



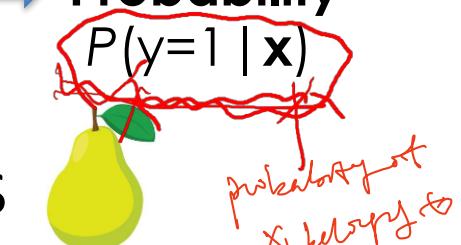
VS



1

If  $f_{\mathbf{w}}(\mathbf{x}_{new}) = 0.2$ ,  $P(?) = 0.2$   
If  $f_{\mathbf{w}}(\mathbf{x}_{new}) = 0.6$ ,  $P(?) = 0.6$

\* cannot assert probability < 1 depends on z threshold most likely on pear.



probability of pear is 0.2 ⇒ higher probability of it being apple

# Binomial Logistic Regression

## ➤ Learning (Binary)

Find the optimal values for  $\mathbf{w}^*$  which minimizes:

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \text{Loss}(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$$

*new w/ annotations*

recall from EE2211 how entropy is calculated (using log/ln), and it represents how pure/impure a prediction is (how many in right class, how many in wrong class)

where  $\text{Loss}(f_{\mathbf{w}}(\mathbf{x}_i), y_i) = -y_i \ln(f_{\mathbf{w}}(\mathbf{x}_i)) - (1 - y_i) \ln(1 - f_{\mathbf{w}}(\mathbf{x}_i))$

-/- **Binary Cross-Entropy Loss**

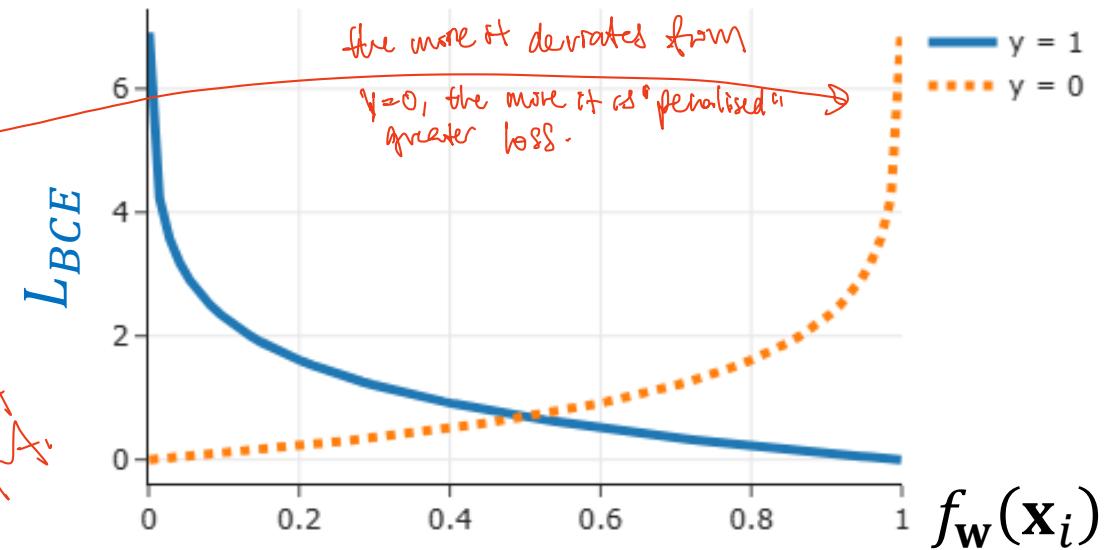
$$L_{BCE}(f_{\mathbf{w}}(\mathbf{x}_i), y_i) =$$

*goes this*

$$\begin{cases} -\ln(1 - f_{\mathbf{w}}(\mathbf{x}_i)), & \text{if } y_i = 0 \\ -\ln(f_{\mathbf{w}}(\mathbf{x}_i)), & \text{if } y_i = 1 \end{cases}$$

*goes this*

$y_i = \text{target output}$



# Binomial Logistic Regression

## ➤ Learning (Binary)

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N L_{BCE}(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$$

$$\operatorname{argmax}_{\mathbf{w}} F = \operatorname{argmin}_{\mathbf{w}} -F$$



$$\ln \left( \prod_{i=1}^N a_i \right) = \sum_{i=1}^N \ln(a_i)$$

Entropy loss  
optimized  
function.

Maximum Likelihood Estimation

$$\operatorname{argmax}_{\mathbf{w}} \ln \left( \prod_{i=1}^N f_{\mathbf{w}}(\mathbf{x}_i)^{y_i} (1 - f_{\mathbf{w}}(\mathbf{x}_i))^{(1-y_i)} \right)$$

Why maximize?  
relevant predictions  
to be as close  
to target classes  
as possible.

(Assumption: data samples  
are **independent**.)

↳ or else will not  
be product of all  
probabilities.

Independence  $\Rightarrow$  outcome  
of one does  
not affect outcome another.

$$\text{if } y_i = 0, \quad 1 - f_{\mathbf{w}}(\mathbf{x}_i)$$

belongs to

$$P(y_i = 0 | \mathbf{x}_i)$$

$$\text{if } y_i = 1, \quad f_{\mathbf{w}}(\mathbf{x}_i)$$

belongs to

$$P(y_i = 1 | \mathbf{x}_i)$$

(Assumption: data samples  
are **identically distributed**.)

↳ must be same distribution

Estimated  $\mathbf{w}$  by maximization over  
the ln of likelihood of training data  
↳ probability of training data  
belonging to their corresponding  
target class.

# Binomial Logistic Regression

## ➤ Learning (Binary)

$$\operatorname{argmin}_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L_{BCE}(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$$

polynomial regression can just use matrix multiplication to train the weights, but logistic regression need to use gradient descent because we cannot solve the derivative (equate to 0 to get min), ie. no closed form solution. So we need to iteratively move towards the min using steps of the complex gradient function

no closed-form solution!

Recall: **Gradient Descent Algorithm!**

Goal:  $\operatorname{argmin}_{\mathbf{w}} J(\mathbf{w})$

Initialize  $\mathbf{w}_0$  and learning rate  $\eta$   
**while** true do

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \eta \nabla_{\mathbf{w}} J(\mathbf{w}_k)$$

**if** converge **then**

**return**  $\mathbf{w}_{k+1}$

**end**

**end**

learning rate ·  
differentiate again ·

If you take the derivative and set it to zero, you get:

$$X^T(\sigma(X\mathbf{w}) - y) = 0$$

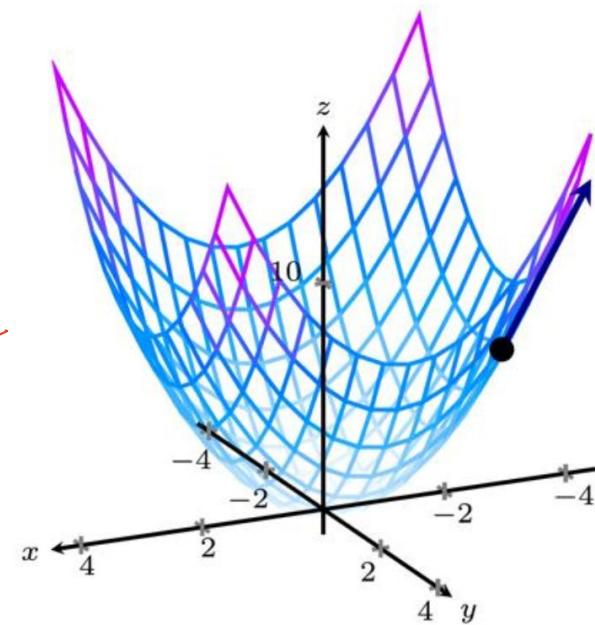
This equation is:

✗ nonlinear in  $\mathbf{w}$

because  $\sigma(\cdot)$  wraps a non-linear sigmoid around  $X\mathbf{w}$ .

There is no algebraic way to isolate  $\mathbf{w}$ .

No matrix inverse trick works here.



# Binomial Logistic Regression

## ➤ Learning (Binary)

$$\begin{aligned}
 \nabla_w J(w) &= \nabla_w \frac{1}{N} \sum_{i=1}^N L_{BCE}(f_w(\mathbf{x}_i), y_i) \\
 &= \frac{1}{N} \sum_{i=1}^N \nabla_w L_{BCE}(f_w(\mathbf{x}_i), y_i) \\
 &= \frac{1}{N} \sum_{i=1}^N \nabla_w (-y_i \ln(f_w(\mathbf{x}_i)) - (1-y_i) \ln(1-f_w(\mathbf{x}_i))) \\
 &\stackrel{\text{negate}}{=} -\frac{1}{N} \sum_{i=1}^N (y_i \nabla_w \ln(f_w(\mathbf{x}_i)) + (1-y_i) \nabla_w \ln(1-f_w(\mathbf{x}_i)))
 \end{aligned}$$

cost function

write out the terms of forward cross entropy loss

By chain rule:

$$\begin{aligned}
 \frac{\partial y(u(x))}{\partial x} &= \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}, \quad \text{chain rule} \\
 &\quad \frac{d \ln(f_w(\mathbf{x}_i))}{d f_w(\mathbf{x}_i)} \frac{d f_w(\mathbf{x}_i)}{d w} \\
 &\quad \leftarrow \text{key is to compute the two gradients} \quad -\frac{1}{1-f_w(\mathbf{x}_i)} \nabla_w f_w(\mathbf{x}_i) \\
 &\quad \leftarrow \text{chain rule} \\
 &= -\frac{1}{N} \sum_{i=1}^N \left( \frac{y_i}{f_w(\mathbf{x}_i)} \boxed{\nabla_w f_w(\mathbf{x}_i)} - \frac{1-y_i}{1-f_w(\mathbf{x}_i)} \boxed{\nabla_w f_w(\mathbf{x}_i)} \right)
 \end{aligned}$$

At this point, we need to find gradient of  $f$ , the logistic regression model wrt.  $w$ .

Continue on next page...

# Binomial Logistic Regression

## ➤ Learning (Binary)

$$\begin{aligned}\nabla_{\mathbf{w}} J(\mathbf{w}) &= \nabla_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N L_{BCE}(f_{\mathbf{w}}(\mathbf{x}_i), y_i) \\ &= -\frac{1}{N} \sum_{i=1}^N \left( \frac{y_i}{f_{\mathbf{w}}(\mathbf{x}_i)} \nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}_i) - \frac{1-y_i}{1-f_{\mathbf{w}}(\mathbf{x}_i)} \nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}_i) \right)\end{aligned}$$

*left*  
*whether this*

$$\nabla_{\mathbf{w}} f_{\mathbf{w}}(\mathbf{x}_i) = \nabla_{\mathbf{w}} \sigma(\mathbf{x}_i^T \mathbf{w})$$

$$= \sigma'(z_i) \mathbf{x}_i$$

By chain rule and  $\frac{d(\mathbf{b}^T \mathbf{w})}{d\mathbf{w}} = \mathbf{b}$

$$= \sigma(z_i)(1 - \sigma(z_i)) \mathbf{x}_i$$

$$= f_{\mathbf{w}}(\mathbf{x}_i)(1 - f_{\mathbf{w}}(\mathbf{x}_i)) \mathbf{x}_i$$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (f_{\mathbf{w}}(\mathbf{x}_i) - y_i) = \frac{1}{N} \mathbf{X}^T (\mathbf{f}_{\mathbf{w}}(\mathbf{X}) - \mathbf{y})$$

one again all  
matrix to condense  
the sum of  
portion.

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix}$$

# Binomial Logistic Regression

Proof.

$$\begin{aligned}\frac{\partial \sigma(a)}{\partial a} &= \frac{\partial}{\partial a} \left( \frac{1}{1 + \exp(-\beta a)} \right) \\&= -\frac{1}{(1 + e^{-\beta a})^2} \frac{\partial(1 + e^{-\beta a})}{\partial a} \\&= \frac{\beta}{(1 + e^{-\beta a})^2} e^{-\beta a} \\&= \frac{\beta}{(1 + e^{-\beta a})^2} (1 + e^{-\beta a} - 1) \\&= \beta \left( \frac{1}{1 + e^{-\beta a}} - \frac{1}{(1 + e^{-\beta a})^2} \right) \\&= \beta \left( \sigma(a) - \sigma^2(a) \right) \\&= \beta \sigma(a)(1 - \sigma(a))\end{aligned}$$

Take  $\beta = 1$  for our case

# Binomial Logistic Regression

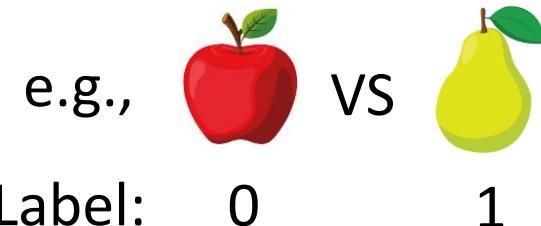
## ► Prediction (Binary)

Now can we gradient descent to update  $w$  each iteration

Still need this

$$f_{\mathbf{w}^*}(\mathbf{X}_t) = \sigma(\mathbf{X}_t \mathbf{w}^*) = \sigma \left( \begin{bmatrix} 1 & x_{1,1}^t & \dots & x_{1,d}^t \\ \vdots & \ddots & \vdots & \vdots \\ 1 & x_{m,1}^t & \dots & x_{m,d}^t \end{bmatrix} \begin{bmatrix} w_0^* \\ \vdots \\ w_d^* \end{bmatrix} \right) = \begin{bmatrix} \sigma((\mathbf{x}_1^t)^T \mathbf{w}^*) \\ \vdots \\ \sigma((\mathbf{x}_m^t)^T \mathbf{w}^*) \end{bmatrix}$$

Each entry is the probability of each test sample belonging to the binary label 1.



$$f_{\mathbf{w}}(\mathbf{x}_{new}) = 0.2$$

Probability of  $\mathbf{x}_{new}$  belonging to pear class (label 1) is 0.2

$$f_{\mathbf{w}}(\mathbf{x}_{new}) = 0.6$$

Probability of  $\mathbf{x}_{new}$  belonging to pear class (label 1) is 0.6

# Binomial Logistic Regression



Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class1}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to class 1 using binomial logistic regression.

# Binomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class1}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to class 1 using binomial logistic regression.

## Step 0: Prepare the training data

### Binary encoding of the target output

$$X = \begin{bmatrix} 1 & -1 \\ 1 & 0 \\ 1 & 0.5 \\ 1 & 0.3 \\ 1 & 0.8 \end{bmatrix}, y = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

→ up to us what to assign.  
but since qn wants to class 1, we put label 1 to class 1.  
output will directly be the probability of  $x = -0.1$  belonging to class 1.

# Binomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class1}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to class 1.

## Step 1: Learning (work out $\mathbf{w}^*$ )

Gradient Descent:

-Gradient  $\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{N} \mathbf{X}^T (\mathbf{f}_{\mathbf{w}}(\mathbf{X}) - \mathbf{y})$

-Initialize  $\mathbf{w}_0 = \begin{bmatrix} 0.1 \\ -1 \end{bmatrix}$ , learning rate  $\eta = 0.1$

-At each iteration,  $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} J(\mathbf{w}_k)$

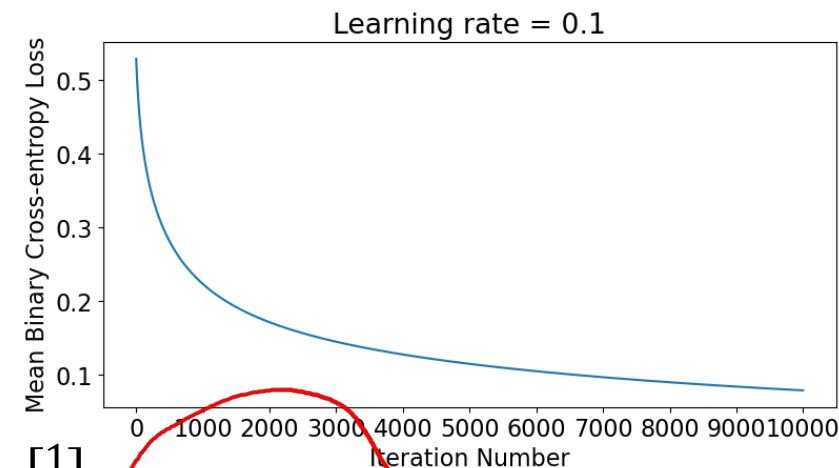
$$\mathbf{w}_0 = \begin{bmatrix} 0.1 \\ -1 \end{bmatrix}$$

$$\nabla_{\mathbf{w}} J(\mathbf{w}_0) = \frac{1}{5} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -1 & 0 & 0.5 & 0.3 & 0.8 \end{bmatrix} \left( \begin{bmatrix} 1 \\ \frac{1}{1 + e^{-\mathbf{x}_1^T \mathbf{w}_0}} \\ \vdots \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix} \right)$$

$$\mathbf{w}_1 = \mathbf{w}_0 - \eta \nabla_{\mathbf{w}} J(\mathbf{w}_0) = \begin{bmatrix} 0.1108 \\ -1.0110 \end{bmatrix}$$

$\vdots$

$$\mathbf{w}_{10000} = \begin{bmatrix} 6.0786 \\ -15.3051 \end{bmatrix}$$



$$\begin{bmatrix} 1 \\ \frac{1}{1 + e^{-\mathbf{x}_1^T \mathbf{w}_0}} \\ \vdots \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -0.1083 \\ 0.1102 \end{bmatrix}$$

# Binomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class1}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to class 1.

## Step 2: Prediction

Augment 1 for the feature vectors

$$\mathbf{x}_t = \begin{bmatrix} 1 \\ -0.1 \end{bmatrix}$$

$$y_t^* = \sigma(\mathbf{x}_t^T \mathbf{w}^*)$$

$$= \sigma([1 \quad -0.1] \begin{bmatrix} 6.0786 \\ -15.3051 \end{bmatrix})$$

$$= \frac{1}{1 + e^{-7.6091}}$$

$$= 0.9995$$

after 10000 iterations

$$\mathbf{w}_{10000} = \begin{bmatrix} 6.0786 \\ -15.3051 \end{bmatrix}$$

NEED TO KNOW THRESHOLD TO CLASSIFY, or else it is just % confident that test data belongs to class

If the threshold is set to 0.5,

$$\text{i.e., } \sigma(z) = \frac{1}{1 + e^{-z}} \begin{cases} > 0.5, \text{class1} \\ < 0.5, \text{class2} \end{cases}$$

$$e^{-z} = 1 \Rightarrow z = \mathbf{x}^T \mathbf{w} = 0$$

Then  $\frac{1}{1 + e^{-(\mathbf{x}^T \mathbf{w}^*)}} = 0.5$  is the decision boundary.

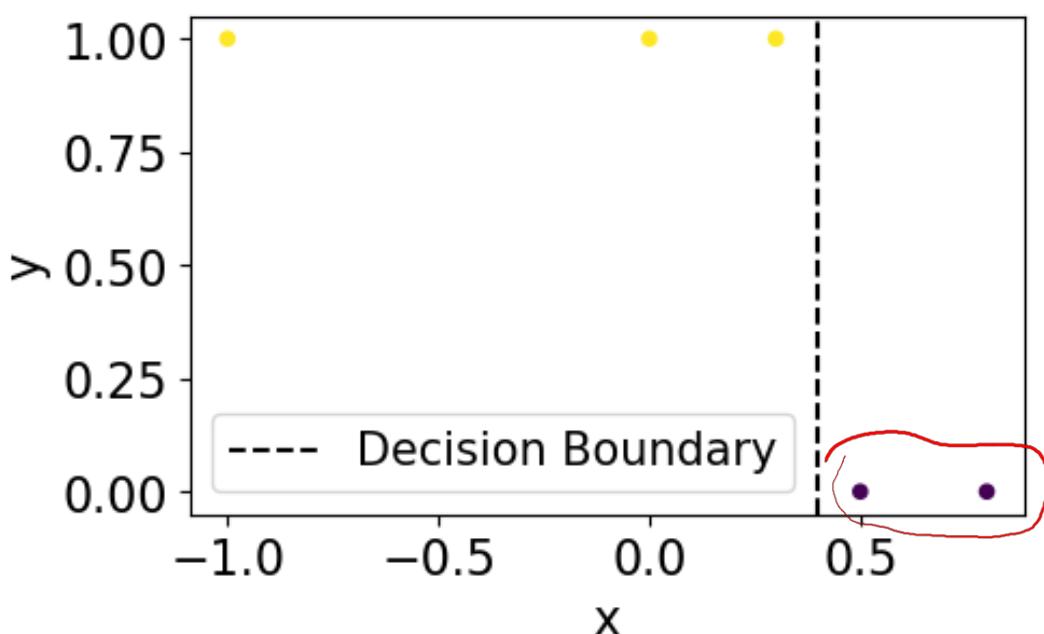
$$\mathbf{x}^T \mathbf{w}^* = 0$$

$$[1 \quad x_1] \begin{bmatrix} 6.0786 \\ -15.3051 \end{bmatrix} = 0$$

derives  
boundary  
 $x_1 = 0.3971$

# Binomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class1}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to class 1.



If the threshold is set to 0.5,

$$\text{i.e., } \sigma(z) = \frac{1}{1 + e^{-z}} \begin{cases} > 0.5, \text{class1} \\ < 0.5, \text{class2} \end{cases}$$

Then,  $\frac{1}{1+e^{-(x^T w^*)}} = 0.5$  is the decision boundary.

$$x^T w^* = 0$$

$$\begin{bmatrix} 1 & x_1 \end{bmatrix} \begin{bmatrix} 6.0786 \\ -15.3051 \end{bmatrix} = 0$$

$$x_1 = 0.3971$$

# Multinomial Logistic Regression

## ➤ Model (Multi-class)

Input:  $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ,

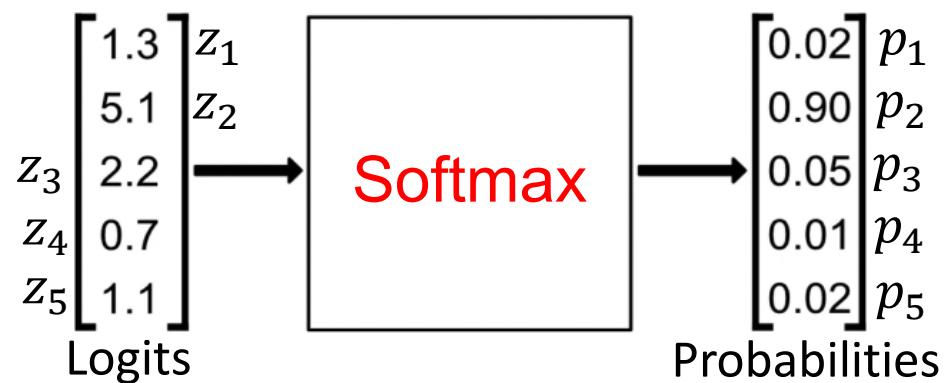
Target output  $\mathbf{y}_i$  is one-hot encoded vector

$$\begin{aligned} \text{Output: } \mathbf{f}_{\mathbf{W}}(\mathbf{x}) &= \sigma(\mathbf{x}^T \mathbf{W}) = \sigma([1 \quad x_1 \quad \cdots \quad x_d]) \\ &= \sigma([\mathbf{x}^T \mathbf{w}_1 \quad \cdots \quad \mathbf{x}^T \mathbf{w}_K]) \\ &= \sigma([z_1 \quad \cdots \quad z_K]) \\ &= \left[ \frac{e^{z_1}}{\sum_{j=1}^K e^{z_j}} \quad \cdots \quad \frac{e^{z_K}}{\sum_{j=1}^K e^{z_j}} \right] \end{aligned}$$

*K classes*

where  $\sigma(z)_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}, k = 1, 2, \dots, K$

is **Softmax** Function



# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

Find the optimal values for  $\mathbf{W}^*$  which minimizes:

$$J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)$$

probability of sample  $\mathbf{x}_i$  belongs to class  $K$

where  $L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i) = -\sum_{k=1}^K y_{i,k} \ln(p_{i,k})$ .

Categorical cross-entropy loss

$p_{i,k}$ :  $k$ -th element of output vector  $\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i)$

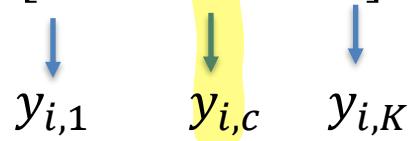
$y_{i,k}$ :  $k$ -th element of target vector  $\mathbf{y}_i$

↳ only one entry is 1 (one-hot)

~~$L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i) = -\ln(p_{i,c})$ , where  $y_{i,c} = 1$~~

Only the target/ground-truth class contributes

$\mathbf{x}_i$  from class  $c$ ,  $\mathbf{y}_i = [0, \dots, 0, 1, 0, \dots 0]$



# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\operatorname{argmin}_{\mathbf{W}} \frac{1}{N} \sum_{i=1}^N L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)$$

**no closed-form solution!**

the sigmoid function makes the loss non linear, derivative leads to equations that cannot be solved algebraically, so cannot equate to zero to get min,

we thus use gradient descent which iteratively goes towards minimum. Still convex

## Use Gradient Descent Algorithm!

$$\frac{dJ(\mathbf{W})}{d\mathbf{W}} = \frac{d}{d\mathbf{W}} \left( \frac{1}{N} \sum_{i=1}^N L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i) \right) = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\mathbf{W}} L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)$$

$$= \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{0,1}} & \dots & \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{0,K}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{d,1}} & \dots & \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{d,K}} \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} w_{0,1} & \dots & w_{0,K} \\ w_{1,1} & \dots & w_{1,K} \\ \vdots & \ddots & \vdots \\ w_{d,1} & \dots & w_{d,K} \end{bmatrix}$$

Start with  
work out one  
and the rest  
are the same.

Continue on next page...

# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\diamond \frac{\partial L_{CCE}(\mathbf{f}_W(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{q,g}} = - \frac{\partial \sum_{k=1}^K y_{i,k} \ln(p_{i,k})}{\partial w_{q,g}}$$

$$p_{i,k} = \sigma(\mathbf{z}_i)_k \\ = \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}}$$

$$= - \sum_{k=1}^K y_{i,k} \frac{\partial \ln(p_{i,k})}{\partial w_{q,g}}$$

By chain rule

$$= - \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} \frac{\partial}{\partial w_{q,g}} \left( \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \right)$$

By chain rule

$$= - \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} \frac{\partial}{\partial z_{i,g}} \left( \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \right) \frac{\partial z_{i,g}}{\partial w_{q,g}}$$

Continue on next page...

# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\frac{\partial}{\partial z_{i,g}} \left( \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \right)$$

*at other outputs*

$$\begin{cases} k = g: e^{z_{i,g}} \\ k \neq g: 0 \end{cases}$$

$$e^{z_{i,g}}$$

By quotient rule:

$$\frac{d}{dx} \left[ \frac{f(x)}{g(x)} \right] = \frac{g(x)f'(x) - f(x)g'(x)}{(g(x))^2}$$

$$= \frac{\left( \sum_{j=1}^K e^{z_{i,j}} \right) \frac{\partial e^{z_{i,k}}}{\partial z_{i,g}} - e^{z_{i,k}} \frac{\partial \left( \sum_{j=1}^K e^{z_{i,j}} \right)}{\partial z_{i,g}}}{\left( \sum_{j=1}^K e^{z_{i,j}} \right)^2}$$

$$= \frac{\left( \sum_{j=1}^K e^{z_{i,j}} \right) \delta_{kg} e^{z_{i,g}} - e^{z_{i,k}} e^{z_{i,g}}}{\left( \sum_{j=1}^K e^{z_{i,j}} \right)^2}$$

$$= \frac{\delta_{kg} e^{z_{i,g}}}{\sum_{j=1}^K e^{z_{i,j}}} - \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \frac{e^{z_{i,g}}}{\sum_{j=1}^K e^{z_{i,j}}}$$

$$= \delta_{kg} p_{i,g} - p_{i,k} p_{i,g}$$

↓ output of softmax function

$$= (\delta_{kg} - p_{i,k}) p_{i,g}$$

$$\delta_{kg} = \begin{cases} 1, & \text{if } k = g \\ 0, & \text{if } k \neq g \end{cases}$$

# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\diamond \frac{\partial L_{CCE}(\mathbf{f}_W(\mathbf{x}_i), y_i)}{\partial w_{q,g}} = - \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} \frac{\partial}{\partial z_{i,g}} \left( \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \right) \boxed{\frac{\partial z_{i,g}}{\partial w_{q,g}}}$$

$$\frac{\partial z_{i,g}}{\partial w_{q,g}} = \frac{\partial (\mathbf{x}_i^T \mathbf{w}_g)}{\partial w_{q,g}}$$

$$= \frac{\partial}{\partial w_{q,g}} (x_{i,0}w_{0,g} + \dots + x_{i,d}w_{d,g})$$

$$= x_{i,q}$$

$$\mathbf{W} = \begin{bmatrix} w_{0,1} & \cdots & w_{0,K} \\ w_{1,1} & \cdots & w_{1,K} \\ \vdots & \ddots & \vdots \\ w_{d,1} & \cdots & w_{d,K} \end{bmatrix}$$

$$= [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_K]$$

# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\diamond \frac{\partial L_{CCE}(\mathbf{f}_W(\mathbf{x}_i), y_i)}{\partial w_{q,g}} = - \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} \frac{\partial}{\partial z_{i,g}} \left( \frac{e^{z_{i,k}}}{\sum_{j=1}^K e^{z_{i,j}}} \right) \frac{\partial z_{i,g}}{\partial w_{q,g}}$$

$$= - \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} (\delta_{kg} - p_{i,k}) p_{i,g} x_{i,q}$$

$$\delta_{kg} = \begin{cases} 1, & \text{if } k = g \\ 0, & \text{if } k \neq g \end{cases}$$

$$= p_{i,g} x_{i,q} \sum_{k=1}^K \frac{y_{i,k}}{p_{i,k}} (p_{i,k} - \delta_{kg})$$

$$= p_{i,g} x_{i,q} \left( \sum_{k=1}^K y_{i,k} - \frac{y_{i,g}}{p_{i,g}} \right)$$

$$= p_{i,g} x_{i,q} \left( 1 - \frac{y_{i,g}}{p_{i,g}} \right)$$

$$= x_{i,q} (p_{i,g} - y_{i,g})$$

only one term  
we select zero  
among all entries of matrix  
or just one,

# Multinomial Logistic Regression

## ➤ Learning (Multi-class)

$$\diamond \frac{dJ(\mathbf{W})}{d\mathbf{W}} = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\mathbf{W}} L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)$$

$$= \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{0,1}} & \dots & \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{0,K}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{d,1}} & \dots & \frac{\partial L_{CCE}(\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i), \mathbf{y}_i)}{\partial w_{d,K}} \end{bmatrix}$$

$$= \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} x_{i,0}(p_{i,1} - y_{i,1}) & \dots & x_{i,0}(p_{i,K} - y_{i,K}) \\ \vdots & \ddots & \vdots \\ x_{i,d}(p_{i,1} - y_{i,1}) & \dots & x_{i,d}(p_{i,K} - y_{i,K}) \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix}$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (\mathbf{f}_{\mathbf{W}}(\mathbf{x}_i) - \mathbf{y}_i) = \frac{1}{N} \mathbf{X}^T (\mathbf{F}_{\mathbf{W}}(\mathbf{X}) - \mathbf{Y})$$

$\mathbf{Y}$  has a size of  $N \times K$

↳ same form! just in matrix and cost vector.

Exactly the same as binomial Logistic Regression! Why?

# Multinomial Logistic Regression

should hold!

- **Softmax() = Logistic(), when K=2**

Categorical.

binomial

VS Binomial (K=2)

- ✓ 1. Gradient in matrix vs vector
- 2. Loss function - categorical vs binary  
Cross entropy loss

✓ 3. softmax vs logistic function

Suppose class 1 is assigned with label 1

$$\sigma(\mathbf{z})_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

$$= \frac{1}{1 + e^{z_2 - z_1}}$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-\mathbf{x}^T \mathbf{w}}}$$

$$= P(y = 1 | \mathbf{x})$$

Let  $a = z_1 - z_2 = \mathbf{x}^T (\underbrace{\mathbf{w}_1 - \mathbf{w}_2}_{\mathbf{w}})$

$$\sigma(\mathbf{z})_1 = \frac{1}{1 + e^{-a}} = \frac{1}{1 + e^{-\mathbf{x}^T \mathbf{w}}} \leftarrow \text{same}$$

$$= P(y = 1 | \mathbf{x})$$

# Multinomial Logistic Regression

- $L_{CCE} = L_{BCE}$ , when  $K=2$

Suppose class 1 is assigned with label 1

$$L_{CCE} = - \sum_{k=1}^2 y_k \ln(p_k)$$

$$= -y_1 \ln(p_1) - y_2 \ln(p_2)$$

$$= -y_1 \ln(p_1) - (1 - y_1) \ln(1 - p_1)$$

✓ 1. Gradient in matrix vs vector  
✓ 2. Loss function ~ categorical vs binary  
Cross entropy loss  
✓ 3. softmax vs logistic function

# Multinomial Logistic Regression

## ➤ Prediction (Multi-class)

$$\begin{aligned} \mathbf{F}_{\mathbf{W}^*}(\mathbf{X}_t) &= \sigma(\mathbf{X}_t \mathbf{W}^*) = \sigma \left( \begin{bmatrix} 1 & x_{1,1}^t & \dots & x_{1,d}^t \\ \vdots & \ddots & \vdots & \vdots \\ 1 & x_{m,1}^t & \dots & x_{m,d}^t \end{bmatrix} \begin{bmatrix} w_{0,1}^* & \dots & w_{0,K}^* \\ w_{1,1}^* & \dots & w_{1,K}^* \\ \vdots & \ddots & \vdots \\ w_{d,1}^* & \dots & w_{d,K}^* \end{bmatrix} \right) \\ &= \begin{bmatrix} \sigma((\mathbf{x}_1^t)^T \mathbf{W}^*) \\ \vdots \\ \sigma((\mathbf{x}_m^t)^T \mathbf{W}^*) \end{bmatrix} \end{aligned}$$

Each row is the probability distribution for each test sample

To predict the class label, output column index that has the largest probability for each test sample.

e.g., 2 test samples, 4 classes

Sample 1 [0.1 0.7, 0, 0.2] output  
Sample 2 [0.3, 0.3 0.4, 0] arg max  
class 2  
class 3

# Multinomial Logistic Regression



Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to each class using multinomial logistic regression.

# Multinomial Logistic Regression



Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to each class using multinomial logistic regression.

## Step 0: Prepare the training data

One-hot encoding of the target output

*class1*: [1,0,0]    *class2*: [0,1,0]    *class3*: [0,0,1]

$$\mathbf{X} = \begin{bmatrix} 1 & -1 \\ 1 & 0 \\ 1 & 0.5 \\ 1 & 0.3 \\ 1 & 0.8 \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

# Multinomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to each class using multinomial logistic regression.

## Step 1: Learning (work out $\mathbf{W}^*$ )

Gradient Descent:

-Gradient  $\nabla_{\mathbf{W}} J(\mathbf{W}) = \frac{1}{N} \mathbf{X}^T (\mathbf{F}_{\mathbf{W}}(\mathbf{X}) - \mathbf{Y})$

-Initialize  $\mathbf{W}_0 = \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ -1 & 0.5 & 0.1 \end{bmatrix}$ ,  
learning rate  $\eta = 0.5$

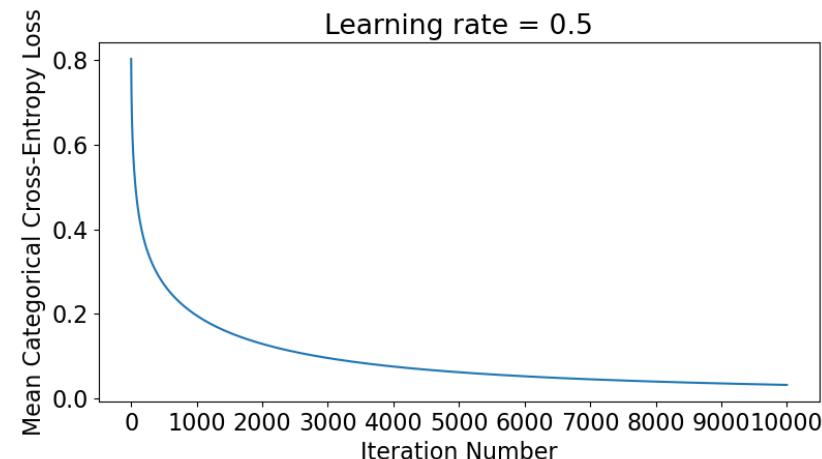
-At each iteration,  $\mathbf{W}_{k+1} = \mathbf{W}_k - \eta \nabla_{\mathbf{W}} J(\mathbf{W}_k)$

$$\mathbf{W}_0 = \begin{bmatrix} 0.5 & 0.3 & 0.2 \\ -1 & 0.5 & 0.1 \end{bmatrix}$$

$$\nabla_{\mathbf{W}} J(\mathbf{W}_0) = \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -1 & 0 & 0.5 & 0.3 & 0.8 \end{bmatrix} \left( \begin{bmatrix} \text{Softmax}(\mathbf{x}_1^T \mathbf{W}_0) \\ \vdots \\ \text{Softmax}(\mathbf{x}_5^T \mathbf{W}_0) \end{bmatrix} - \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \right)$$

$$\therefore = \begin{bmatrix} -0.0355 & -0.0479 & 0.0834 \\ 0.1321 & -0.1402 & 0.0080 \end{bmatrix}$$

$$\mathbf{W}_{10000} = \begin{bmatrix} 6.5540 & -8.5683 & 3.0143 \\ -25.6024 & 27.1366 & -1.9341 \end{bmatrix}$$



# Multinomial Logistic Regression

Example: Given a training set  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{class1}), (0, \text{class1}), (0.5, \text{class2}), (0.3, \text{class3}), (0.8, \text{class2})\}$ , predict the probability of  $\{x=-0.1\}$  belonging to each class using multinomial logistic regression.

## Step 2: Prediction

Augment 1 for the feature vector

$$\mathbf{x}_t = \begin{bmatrix} 1 \\ -0.1 \end{bmatrix} \quad \mathbf{w}_{10000} = \begin{bmatrix} 6.5540 & -8.5683 & 3.0143 \\ -25.6024 & 27.1366 & -1.9341 \end{bmatrix}$$

$$y_t^* = \sigma(\mathbf{x}_t^T \mathbf{w}^*)$$

$$= \sigma\left([1 \quad -0.1] \begin{bmatrix} 6.5540 & -8.5683 & 3.0143 \\ -25.6024 & 27.1366 & -1.9341 \end{bmatrix}\right)$$

$$= \left[ \frac{e^{9.1143}}{e^{9.1143} + e^{-11.2820} + e^{3.2077}}, \frac{e^{-11.2820}}{e^{9.1143} + e^{-11.2820} + e^{3.2077}}, \frac{e^{3.2077}}{e^{9.1143} + e^{-11.2820} + e^{3.2077}} \right]$$

$$= [9.9729 \times 10^{-1}, 1.3830 \times 10^{-9}, 2.7142 \times 10^{-3}]$$

*class 1*      *class 2*      *class 3*  
*largest probability.*

Please feel free to provide your thoughts under the Section of “**Discussions**” in the thread: “**Ask Questions Here for Week 7**” on Canvas.

# Why Cross-Entropy Loss,

entropy! recall EE2211 Gini, Entropy, it is to test pureness of each classification, and we want to reduce that.

# not Squared Error Loss ?

in logistic regression

# Evaluation Metrics for Classification

# Evaluation Metrics for Classification

## ➤ Binary Classification

Class1 (label +1 or 1): Positive Class

Class2 (label -1 or 0): Negative Class

correctly predicted to positive  
wrongly predicted to negative class

- TP:** True Positive
- FN:** False Negative  
(i.e., Type II Error)
- FP:** False Positive  
(i.e., Type I Error)
- TN:** True Negative

## Confusion Matrix for Binary Classification

		$\hat{P}$ (predicted)	$\hat{N}$ (predicted)
P (actual)	TP	FN	
N (actual)	FP	TN	



Precision  
 $TP/(TP+FP)$



Recall  
 $TP/(TP+FN)$



Accuracy  
 $(TP+TN)/(TP+TN+FP+FN)$

# Evaluation Metrics for Classification

## ➤ Binary Classification

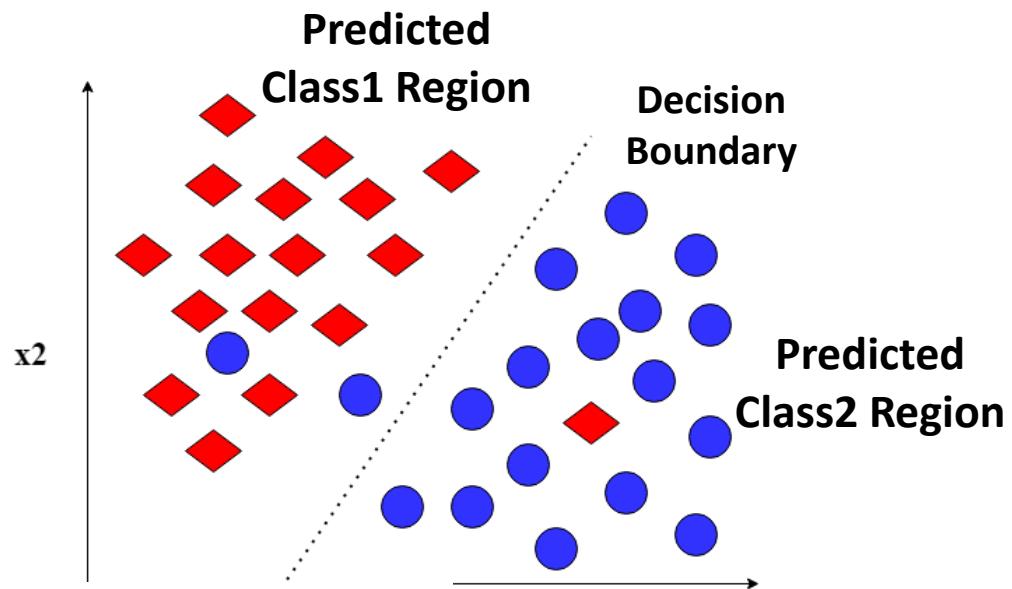
e.g.,

- ◆ Class1: Positive Class
- Class2: Negative Class



Confusion Matrix for Binary Classification

	$\hat{P}$ (predicted)	$\hat{N}$ (predicted)
P (actual)	16 TP	1 FN
N (actual)	2 FP	16 TN



only binary has these:

**TP**: True Positive

**FN**: False Negative  
(i.e., Type II Error)

**FP**: False Positive  
(i.e., Type I Error)

**TN**: True Negative

Precision  
 $TP/(TP+FP)$

Recall  
 $TP/(TP+FN)$

Accuracy  
 $(TP+TN)/(TP+TN+FP+FN)$

# Evaluation Metrics for Classification

## ➤ Multicategory Classification

### Confusion Matrix for Multicategory Classification

	$P_{\widehat{1}}$ (predicted)	$P_{\widehat{2}}$ (predicted)		$P_{\widehat{C}}$ (predicted)
$P_1$ (actual)	$P_{1,\widehat{1}}$	$P_{1,\widehat{2}}$	...	$P_{1,\widehat{C}}$
$P_2$ (actual)	$P_{2,\widehat{1}}$	$P_{2,\widehat{2}}$	...	$P_{2,\widehat{C}}$
⋮	⋮	⋮	⋮	⋮
$P_C$ (actual)	$P_{C,\widehat{1}}$	$P_{C,\widehat{2}}$		$P_{C,\widehat{C}}$

```
from sklearn.metrics import confusion_matrix, accuracy_score
cm = confusion_matrix(y_true, y_pred)
acc = accuracy_score(y_true, y_pred)
```

only accuracy.

diagonal entries are the correct ones.

= sum of diagonal entries / total sample count

# Summary

## ✓ Polynomial Regression for Classification

Binary Classification

-Learning:

$$\mathbf{w}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{y} \quad y_i \in \{-1, +1\}$$

-Prediction:

$$\mathbf{f}_{\mathbf{w}^*}(\mathbf{P}(\mathbf{X}_t)) = \text{sign}(\mathbf{P}_t \mathbf{w}^*)$$

Multi-category Classification

-Learning:

$$\mathbf{W}^* = (\mathbf{P}^T \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Y}$$

-Prediction:

$$\mathbf{F}_{\mathbf{W}^*}(\mathbf{P}(\mathbf{X}_t)) = \begin{bmatrix} \underset{k=1, \dots, K}{\text{argmax}} ((\mathbf{p}_1^t)^T \mathbf{W}^*) \\ \vdots \\ \underset{k=1, \dots, K}{\text{argmax}} ((\mathbf{p}_m^t)^T \mathbf{W}^*) \end{bmatrix}$$

## ✓ Logistic Regression

Binary Classification

-Learning: Gradient Descent

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{N} \mathbf{X}^T (\mathbf{f}_{\mathbf{w}}(\mathbf{X}) - \mathbf{y})$$

-Prediction:

$$\mathbf{f}_{\mathbf{w}^*}(\mathbf{X}_t) = \sigma(\mathbf{X}_t \mathbf{w}^*) = \begin{bmatrix} \sigma((\mathbf{x}_1^t)^T \mathbf{w}^*) \\ \vdots \\ \sigma((\mathbf{x}_m^t)^T \mathbf{w}^*) \end{bmatrix}$$

Multi-category Classification

-Learning: Gradient Descent

$$\nabla_{\mathbf{W}} J(\mathbf{W}) = \frac{1}{N} \mathbf{X}^T (\mathbf{F}_{\mathbf{w}}(\mathbf{X}) - \mathbf{Y})$$

-Prediction:

$$\mathbf{F}_{\mathbf{W}^*}(\mathbf{X}_t) = \sigma(\mathbf{X}_t \mathbf{W}^*) = \begin{bmatrix} \sigma((\mathbf{x}_1^t)^T \mathbf{W}^*) \\ \vdots \\ \sigma((\mathbf{x}_m^t)^T \mathbf{W}^*) \end{bmatrix}$$

## ✓ Evaluation Metrics for Classification Task

binary  
1. Precision  
2. Recall  
3. Accuracy

multinomial  
only this

# EE2213 Introduction to AI

## Lecture 10: Clustering

Dr. WANG Si

[si.wang@nus.edu.sg](mailto:si.wang@nus.edu.sg)

Electrical and Computer Engineering Department  
National University of Singapore

# OVERVIEW OF COURSE CONTENTS



- **Introduction (Shaojing)**

- What is AI
- Applications of AI
- AI agent

- **Search (Shaojing)**

- Uninformed search algorithms: breadth-first, depth-first, uniform-cost(Dijkstra's algorithm)
- Informed search algorithms: greedy best-first, A\*

- Applications

- **Optimisation (Shaojing)**

- Linear programming
- Convex problems
- Applications

- **Machine learning (Wang Si) (Weeks 6-9)**

- Supervised and unsupervised learning: regression, classification, clustering
- Neural networks and deep learning
- Applications

- **Knowledge representation (Wang Si) (Weeks 10-11)**

- Knowledge Representation and Reasoning
  - Propositional Logic
  - Applications
- **Ethical considerations (Shaojing)**
- Bias in AI
  - Privacy concerns
  - Societal impact

(Week 10, Monday)

No Class on 20 Oct.! A recording  
will be uploaded on Canvas.

# AGENDA

- We will discuss:
  - types of clusters and types of clustering
  - K-means clustering
  - Fuzzy C-means clustering

At the end of this lecture, you should be able to:

- ✓ Differentiate between types of clusters and clustering
- ✓ Apply K-means and Fuzzy C-means
- ✓ Describe the limitations of K-means clustering

# Recap: Clustering



Groups a set of objects in such a way that objects in the same group (called a **cluster**) are **more similar** (in some sense) to each other than to those in other groups (clusters).

## Applications:

A screenshot of the Netflix mobile application. At the top, there's a navigation bar with "NETFLIX" and links for "Home", "TV Shows", "Movies", "Recently Added", and "My List". Below this, under the heading "Western Movies", there are two movie thumbnails: "THE RED SEA DIVING RESORT" and "COCO". At the bottom, under the heading "Because you watched Cars", there are three movie thumbnails: "Cars 3", "Tayo The Little Bus", and "Cars 2".

## Recommender systems

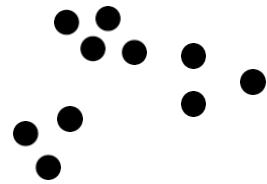
## Image Segmentation



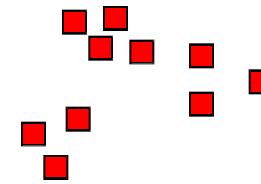
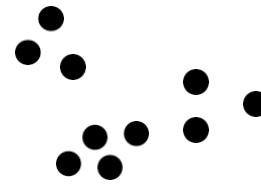
## Social Network Analysis



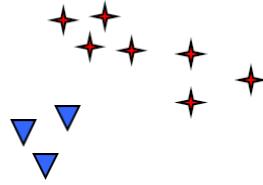
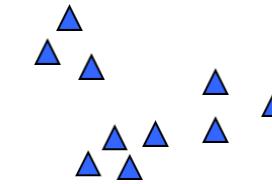
# Notion of a Cluster can be Ambiguous



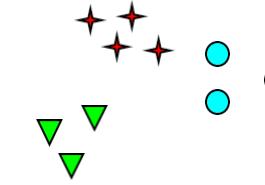
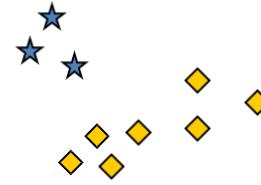
How many clusters?



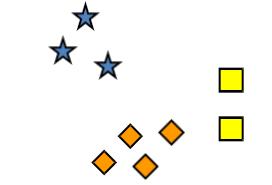
Two Clusters



Four Clusters



Six Clusters



# Types of Clusters

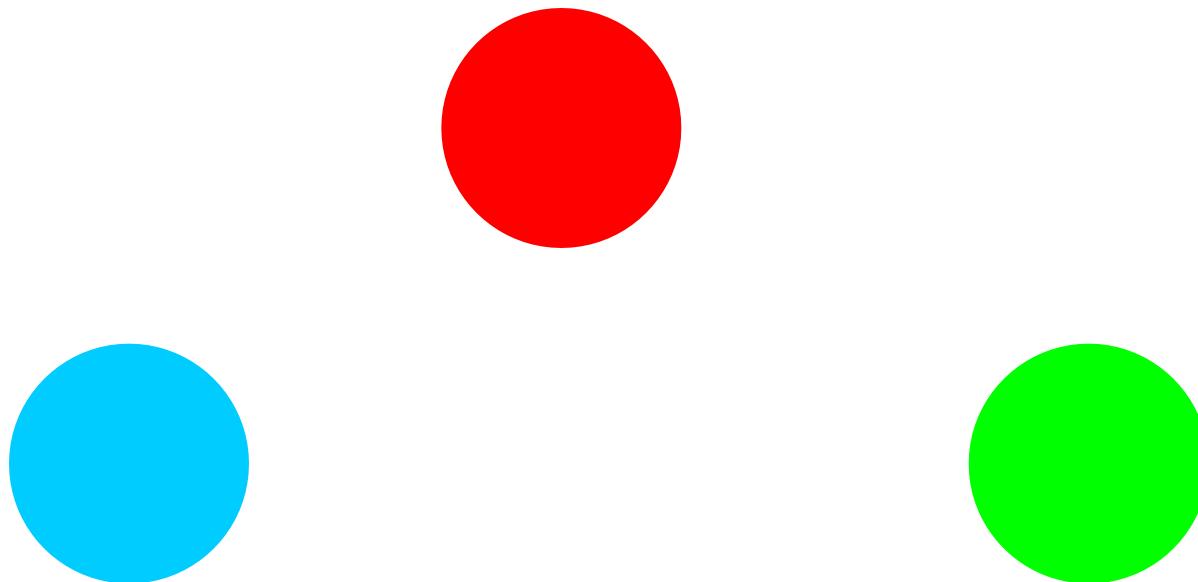
- Well-separated Clusters
- Center-based Clusters
- Contiguous Clusters
- Density-based Clusters

# Types of Clusters

- Well-Separated Cluster:

- A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.

*by distance between points*



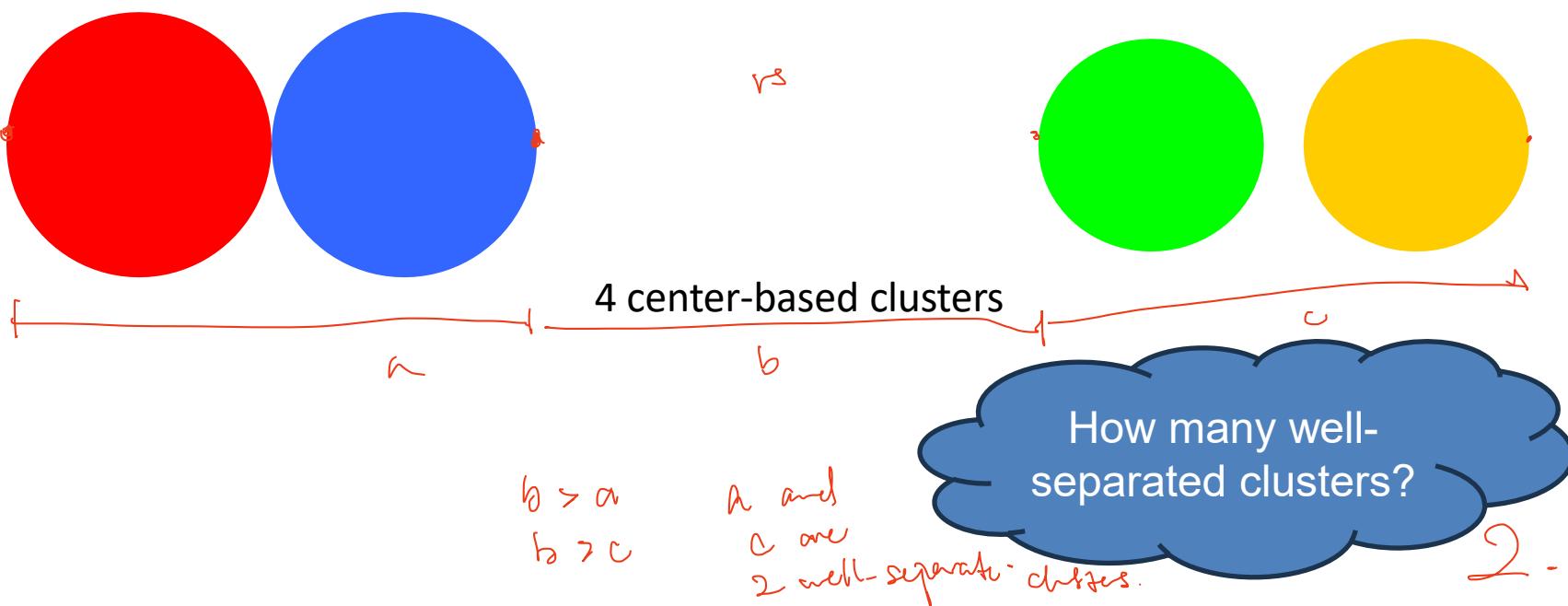
3 well-separated clusters

# Types of Clusters

## Center-based Cluster:

as more said:  
distance to center of cluster.

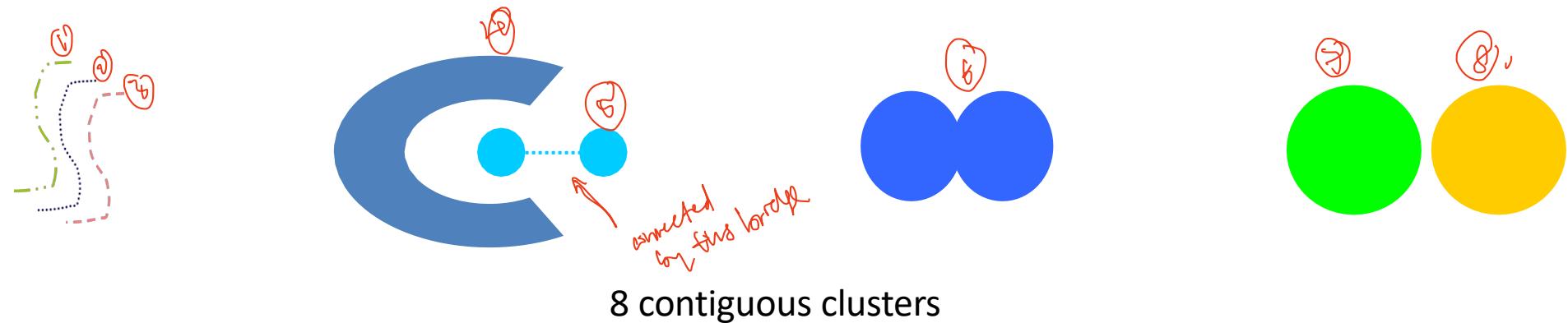
- A cluster is a set of objects such that a point in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster
- The center of a cluster is often a centroid, the average of all the points in the cluster, or the most “representative” point of a cluster



# Types of Clusters

- Contiguous Cluster (Nearest neighbor)

- A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.



How many well-separated clusters?

4.

How many center-based clusters?

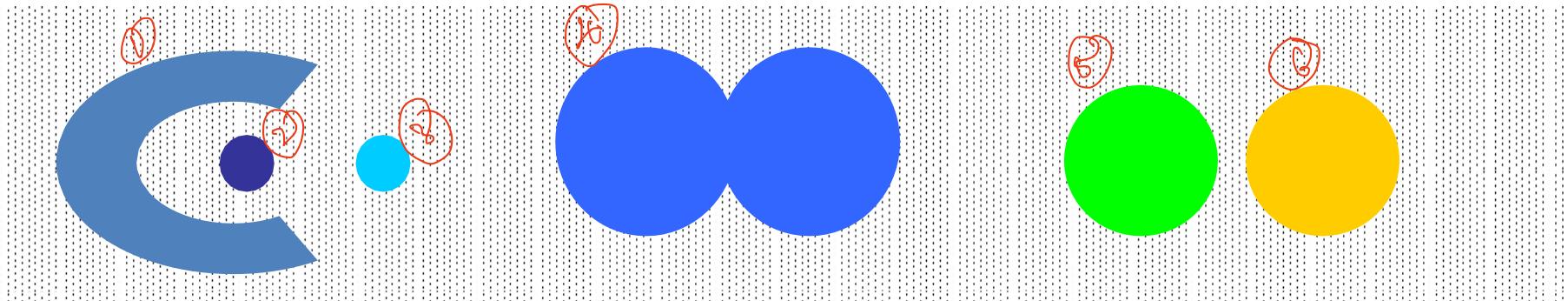
5

# Types of Clusters

- Density-based Cluster

- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
- Used when noise and outliers are present.

↳ in low-density region.



6 density-based clusters

# Types of Clustering

*algorithm -*

- **Hard Clustering**

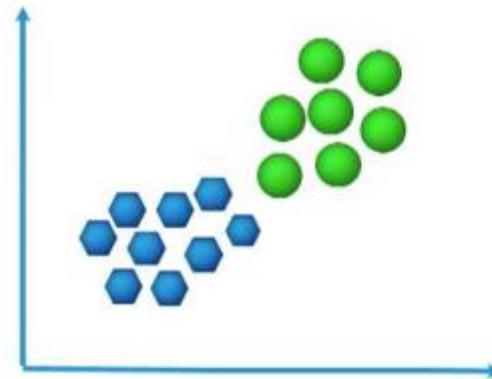
- Each data point belongs to a cluster completely or not.
- e.g., **K-means**, Hierarchical, etc. *one and only one -*

- **Soft Clustering**

- Each data point belongs to each cluster with a probability or degree of membership.
- e.g., **Fuzzy C-means**, Gaussian mixture model, etc.

Data Points	Hard	Soft	
	Clusters	P(C1)	P(C2)
A	C1	0.91	0.09
B	C2	0.3	0.7
C	C2	0.17	0.83
D	C1	1	0

*Probability  
more not  
sum to 1,  
depends  
on alg.*



# K-means Clustering

-Hard Clustering

What type of clusters?

# K-means Clustering

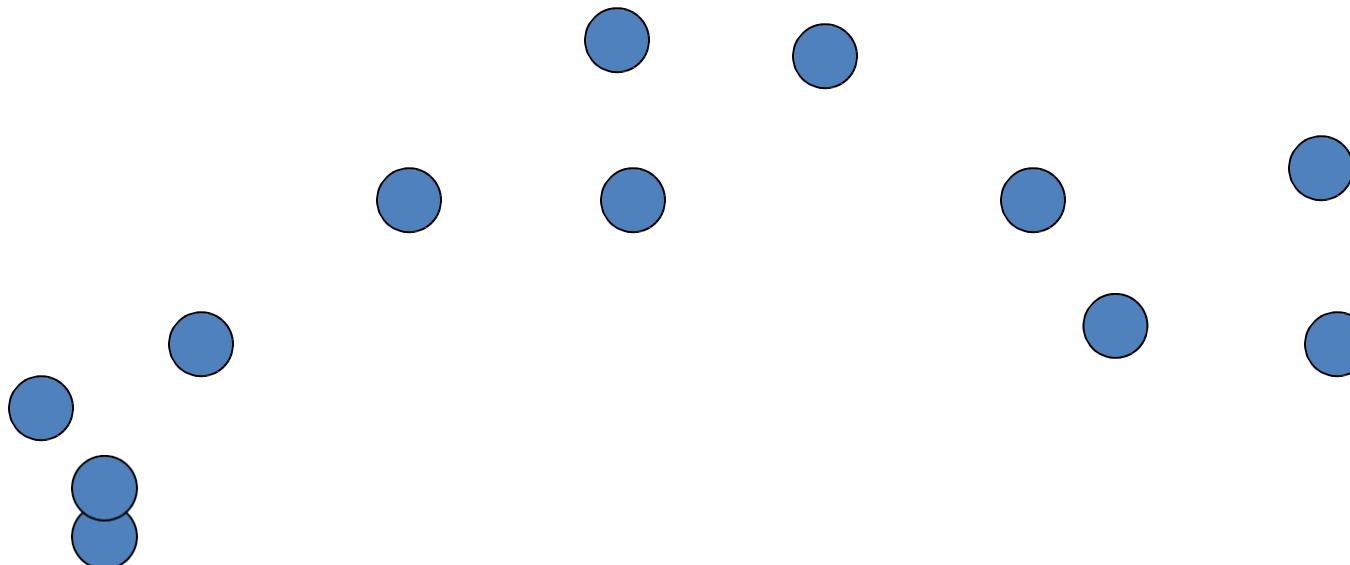
center-based clustering



- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

*center-based  
clusters  
will be produced.*

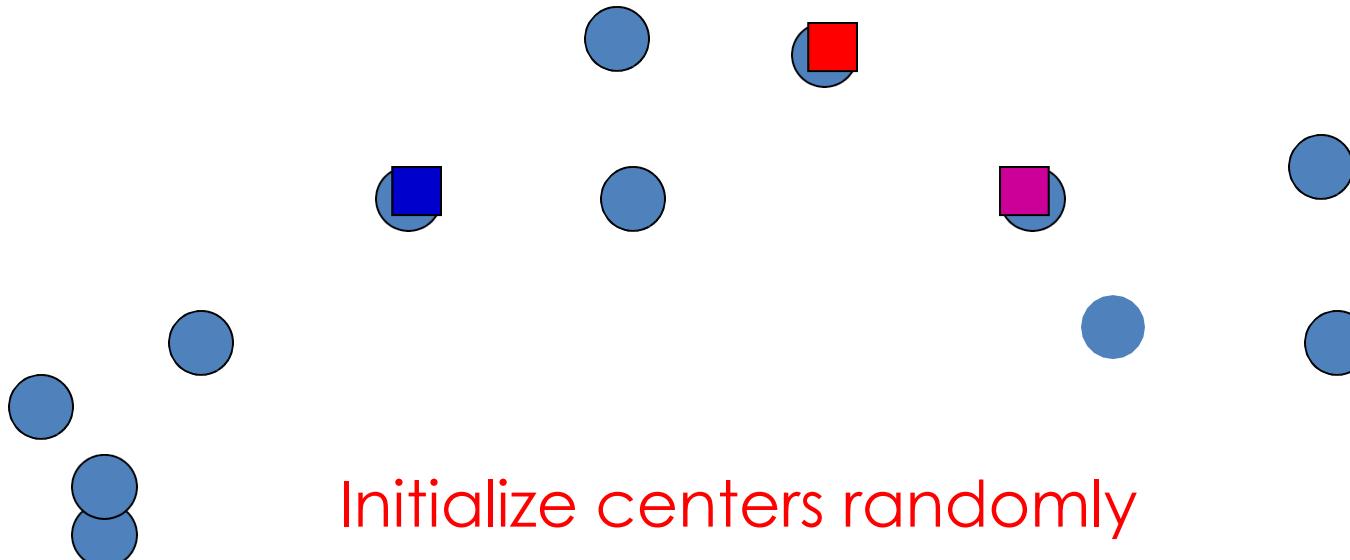


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 1: Randomly choose 3 points as the centers of the 3 clusters

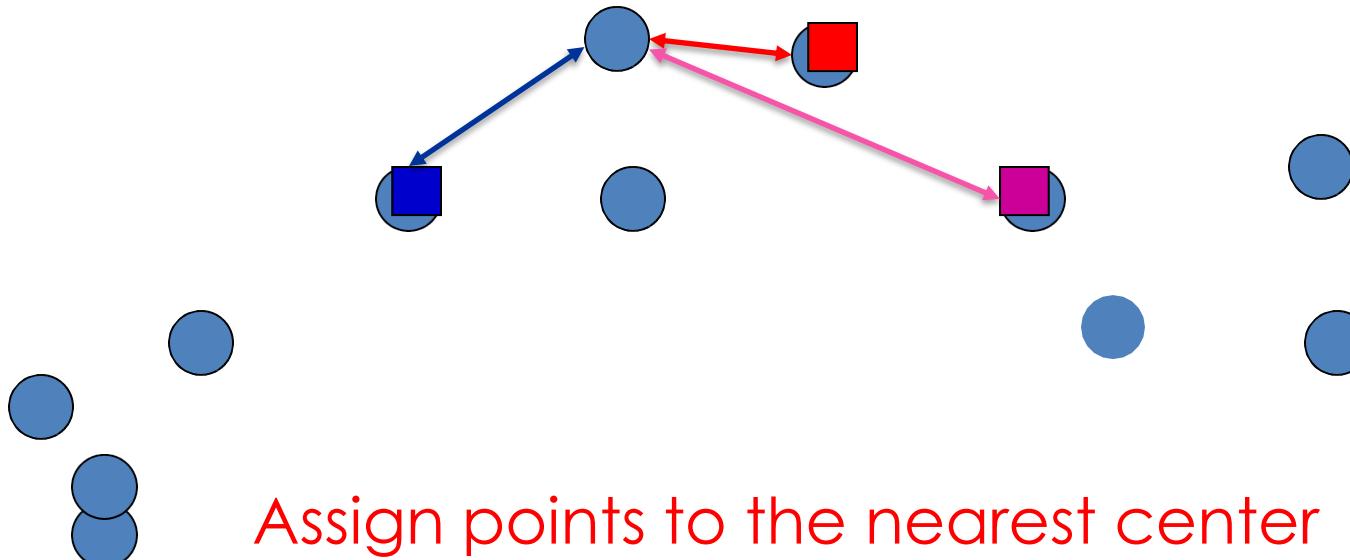


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 2: For each point, calculate the <sup>euclidean</sup> distance between the point and the 3 centers, and assign it to the nearest one.

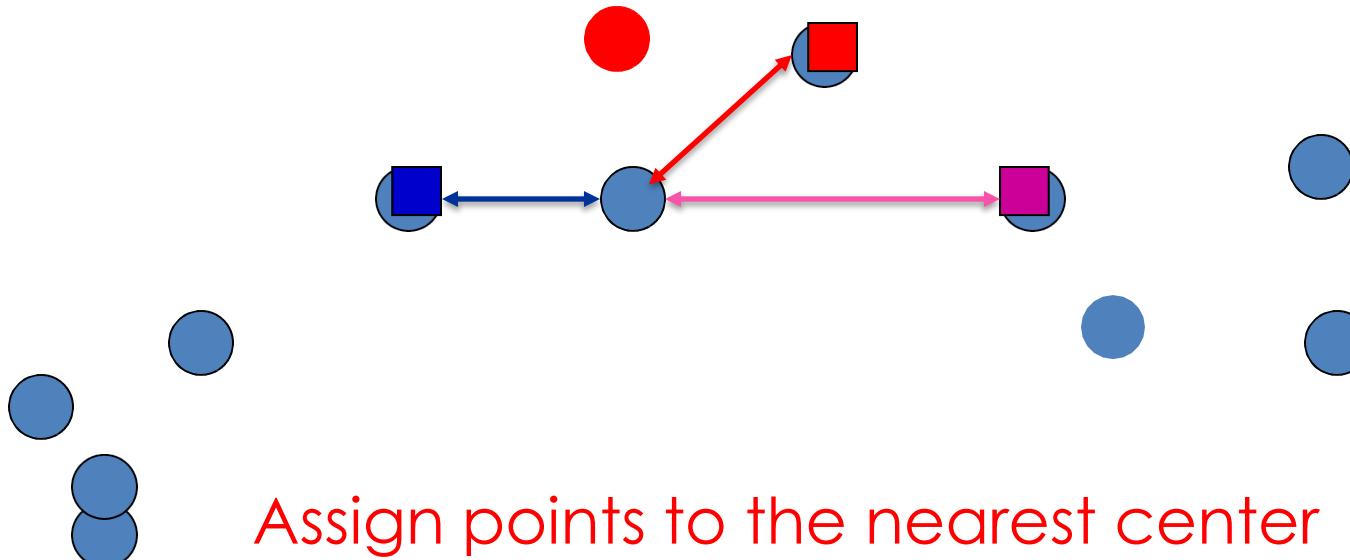


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 2: For each point, calculate the distance between the point and the 3 centers, and assign it to the nearest one.

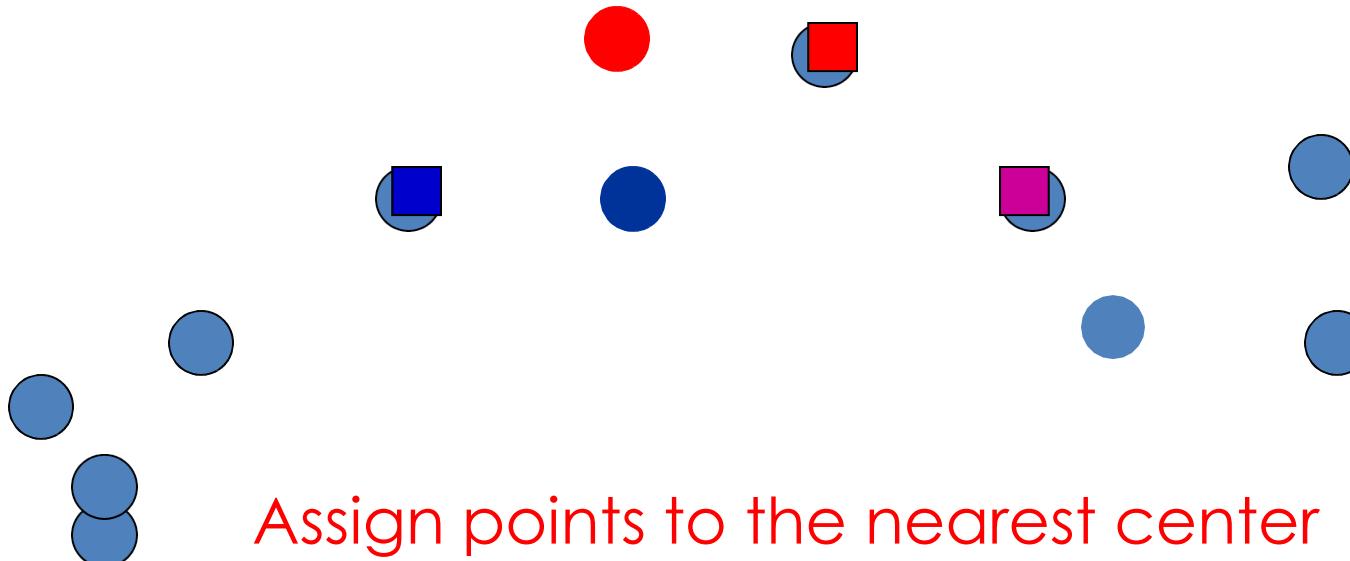


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 2: For each point, calculate the distance between the point and the 3 centers, and assign it to the nearest one.

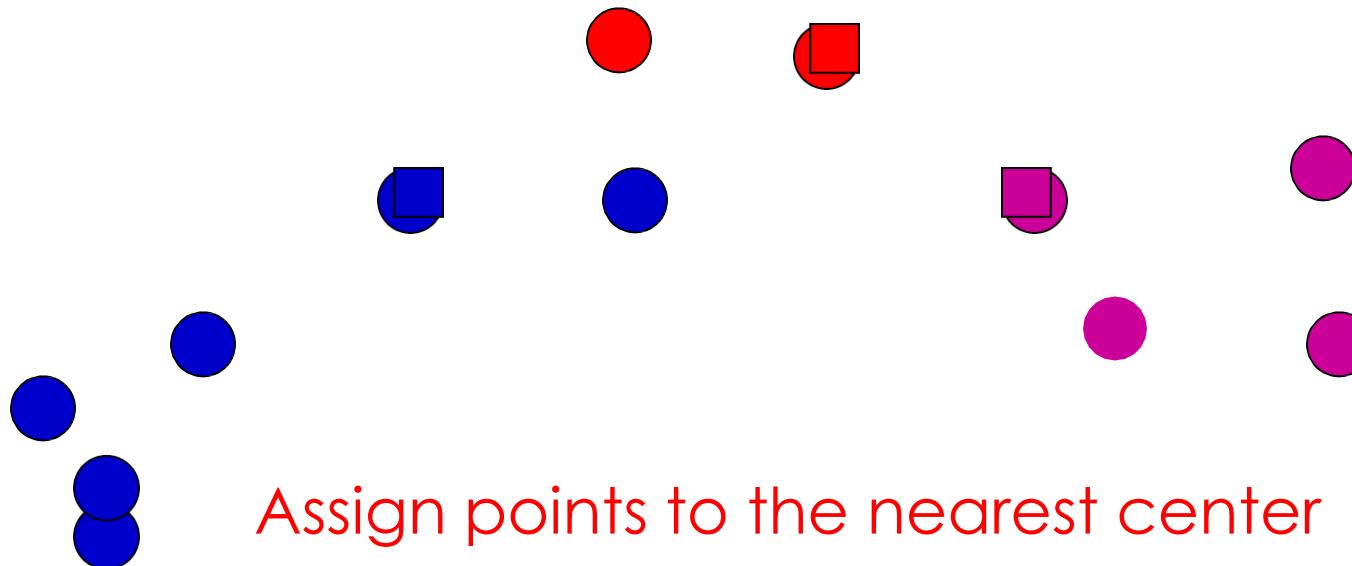


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 2: For each point, calculate the distance between the point and the 3 centers, and assign it to the nearest one.



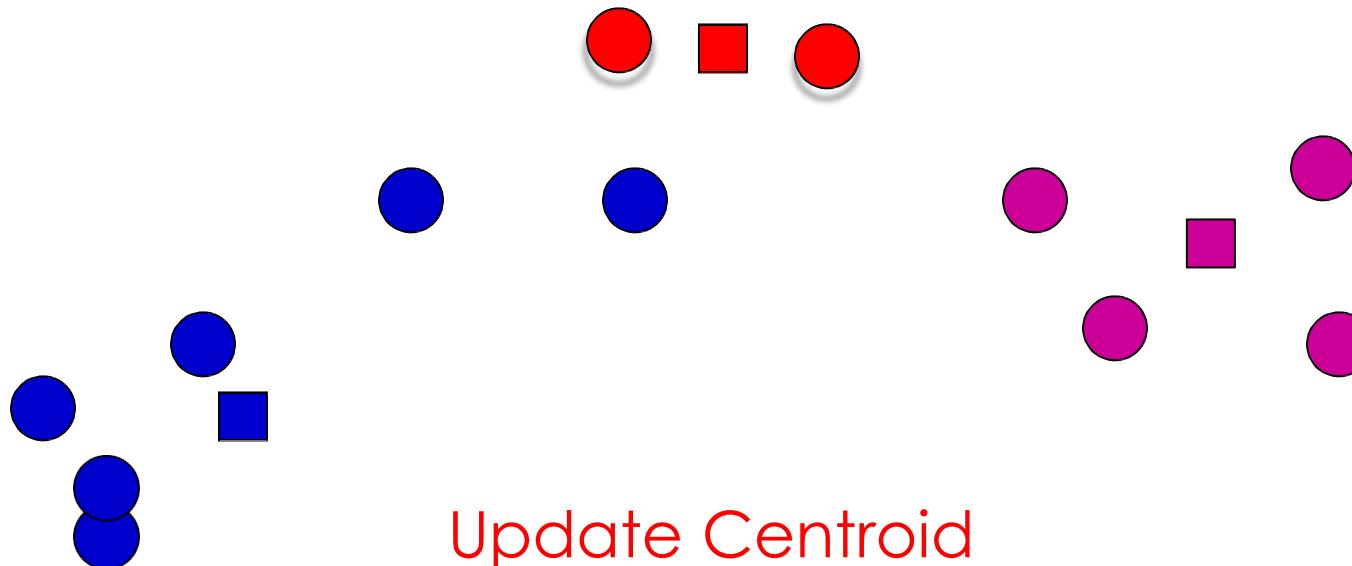
# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 3: Recompute the 3 centers given the current cluster assignments

↳ update each to the centroid of points in new cluster

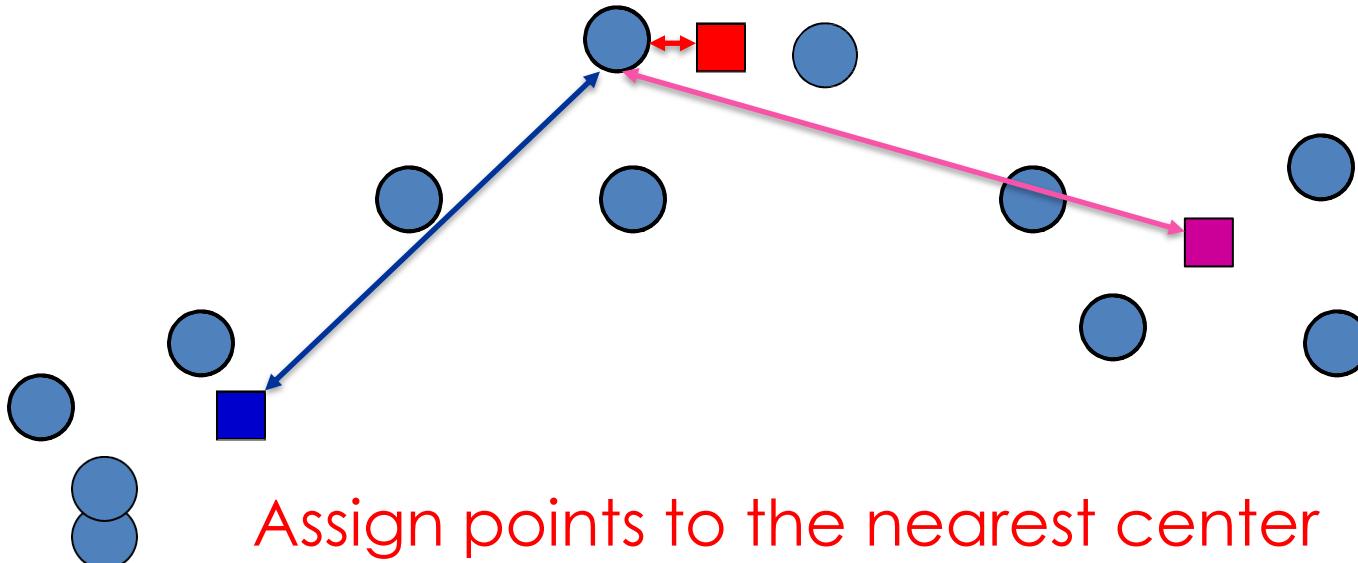


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

 Step 4 (repeat Step 2): For each point, calculate the distance between the point and the updated 3 centers, and assign it to the nearest one.

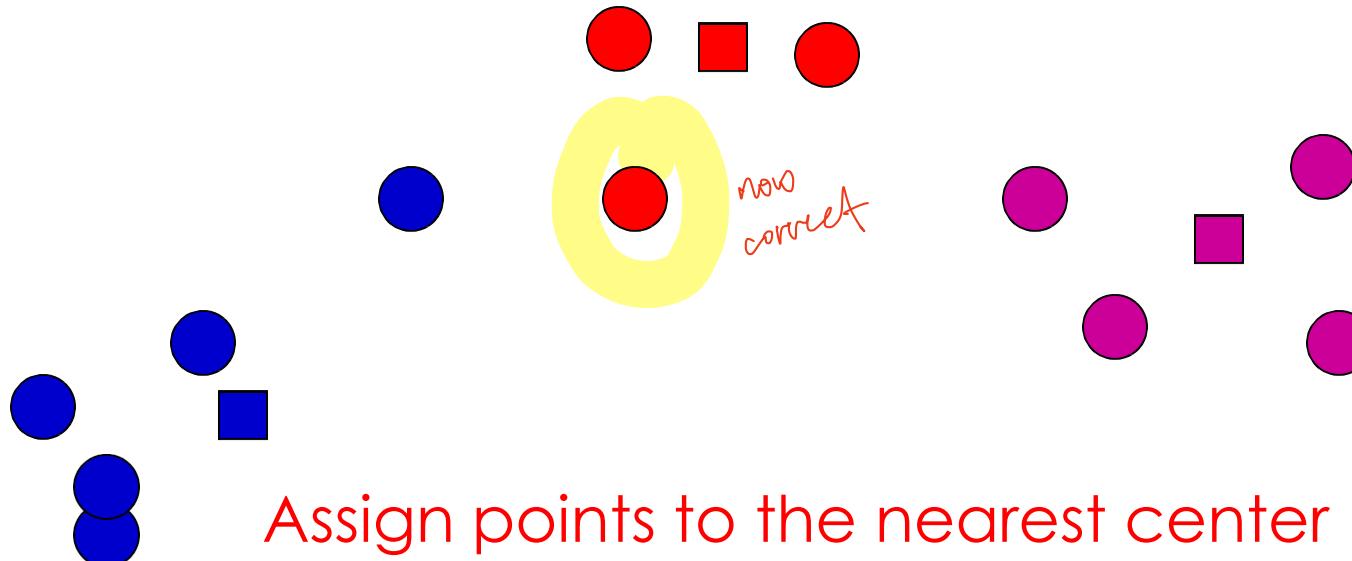


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 4 (repeat Step 2): For each point, calculate the distance between the point and the updated 3 centers, and assign it to the nearest one.

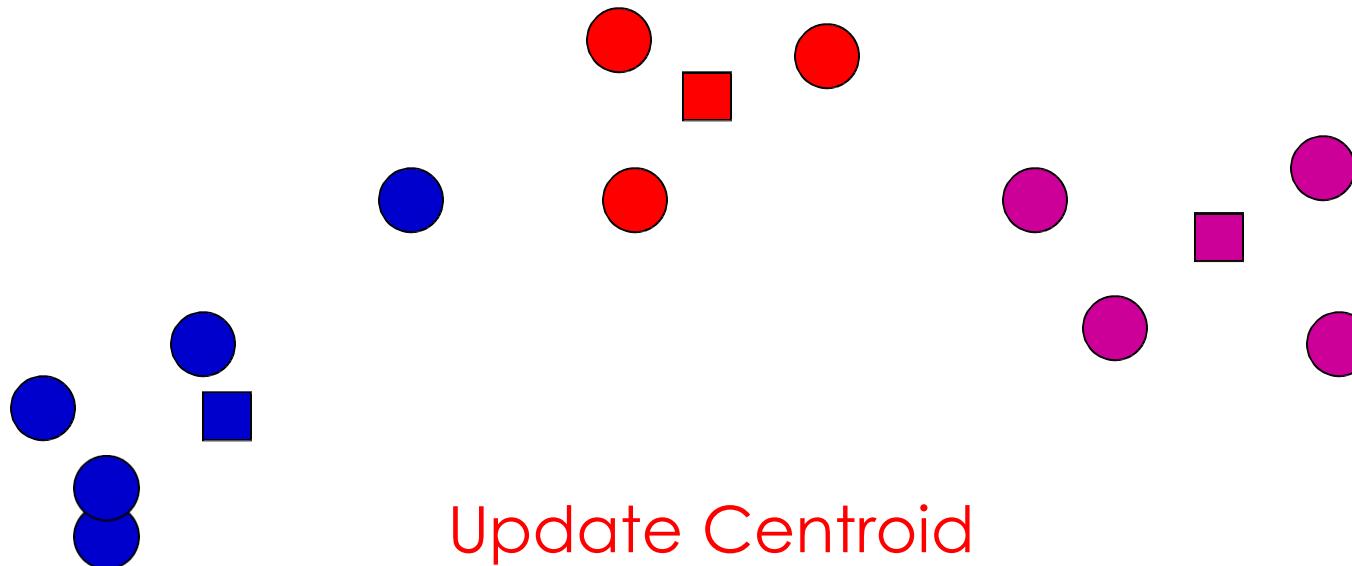


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

*again* Step 5 (repeat Step 3): Recompute the 3 centers given the current cluster assignments.

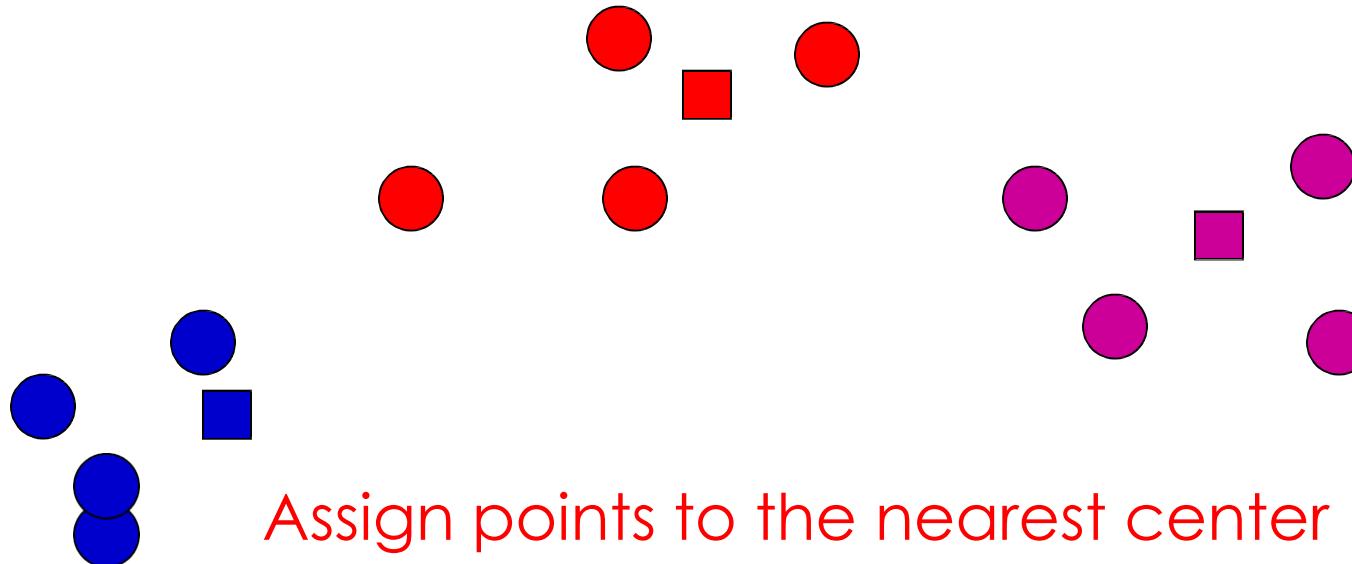


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 6 (repeat Step 2): For each point, calculate the distance between the point and the updated 3 centers, and assign it to the nearest one.

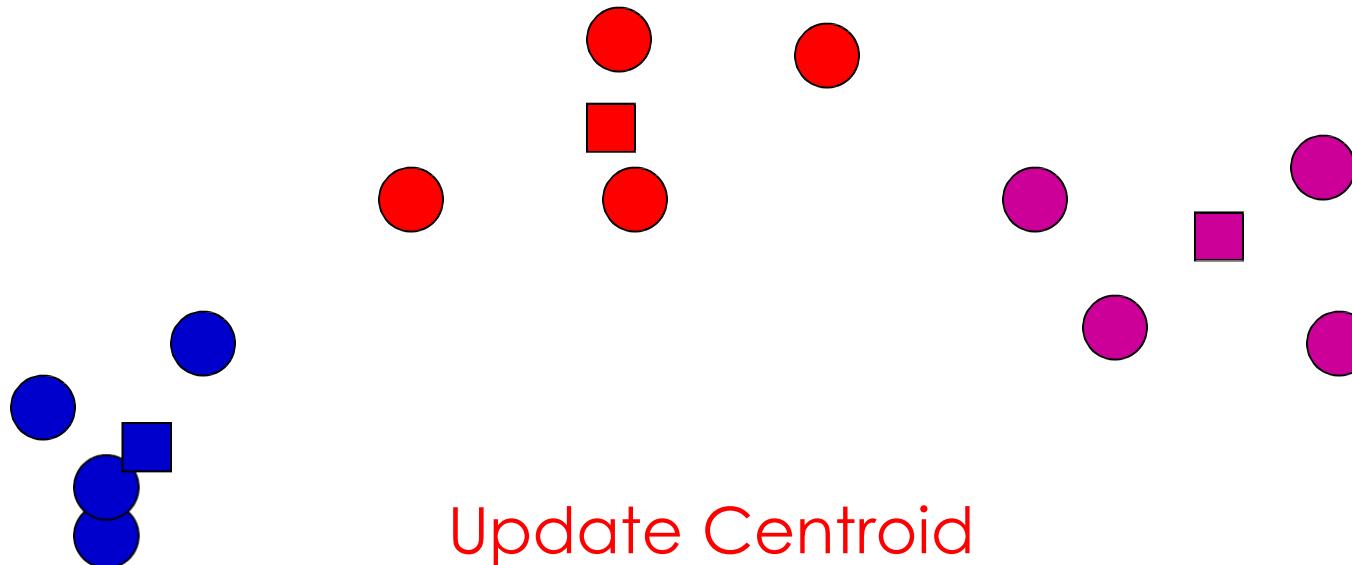


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 7 (repeat Step 3): Recompute the 3 centers given the current cluster assignments.

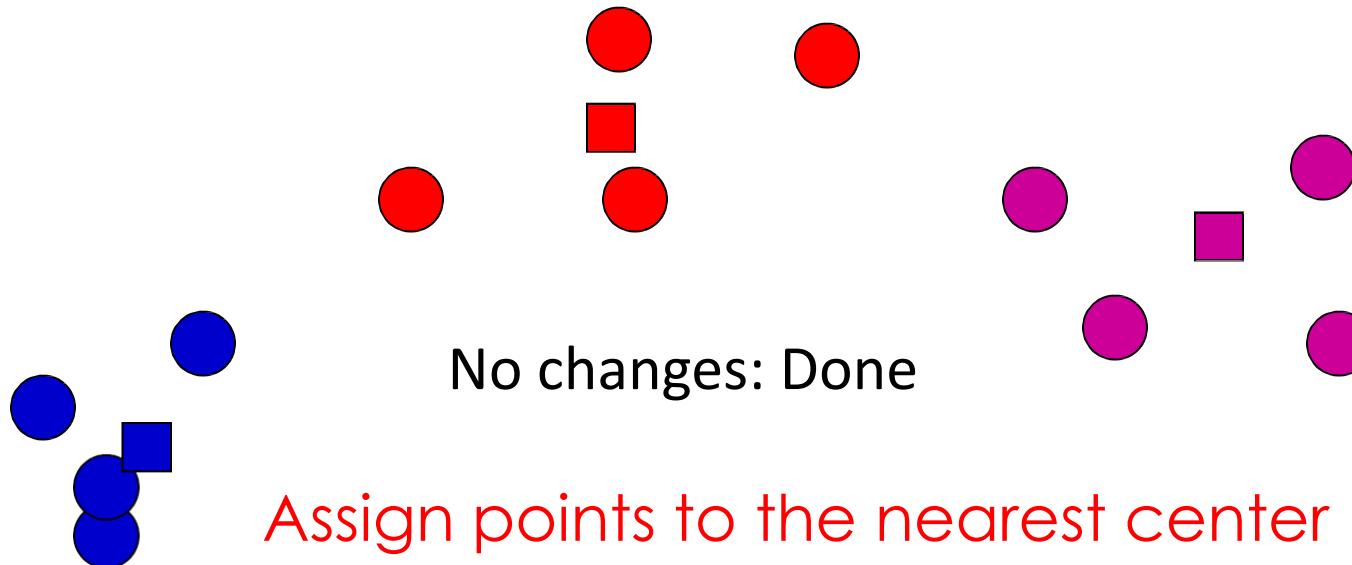


# K-means Clustering

- ✓ Each cluster is associated with a centroid (center point)
- ✓ Each point is assigned to the cluster with the closest centroid

An example: K=3 (we seek 3 clusters)

Step 8 (repeat Step 2): For each point, calculate the distance between the point and the updated 3 centers, and assign it to the nearest one.





# K-means Clustering

## ➤ Algorithm

Looping between Assignment and Centroid Update

- 
- 1: Select  $K$  points as the initial centroids.
  - 2: **repeat**
  - 3:   Form  $K$  clusters by assigning all points to the closest centroid.  
*euclidean distance*
  - 4:   Recompute the centroid of each cluster. *update*
  - 5: **until** The centroids don't change *stop condition*
-

# K-means Clustering

## ➤ Objective Function (within-cluster variance)

Find the optimal values for  $\mathbf{W}^*_{N \times K}$  and  $\mathbf{C}^*_{K \times d}$  which minimizes:

$$J(\mathbf{W}, \mathbf{C}) = \sum_{i=1}^N \sum_{k=1}^K w_{ik} \|\mathbf{x}_i - \mathbf{c}_k\|^2$$

*within cluster variance*

*samples* → *clusters*

*w* → *membership of sample + belongs to cluster k.*

Squared Euclidean distance

where  $w_{ik} \in \{0,1\}$ :

$$\begin{cases} w_{ik} = 1, & \text{if } \mathbf{x}_i \in \text{cluster } S_k \\ w_{ik} = 0, & \text{else} \end{cases}$$

0 means don't belong

C means clustering  $\mathbf{c}_k$  (d-dimensional vector) is the centroid of cluster  $S_k$

e.g.,  $\mathbf{x}_1$  ●  
 $\mathbf{x}_2$  ●  
 $c_1$  ◆

●  $\mathbf{x}_3$   
◆  $c_2$   
●  $\mathbf{x}_4$

$w_{11} = ?$  ○  $w_{12} = ?$  ○  $w_{21} = ?$  ○  $w_{22} = ?$  ○

$w_{31} = ?$  ○  $w_{32} = ?$  ○  $w_{41} = ?$  ○  $w_{42} = ?$  ○

$J(\mathbf{W}, \mathbf{C}) = ?$

W shows numbering.  
only 4 terms left  
works like variance for cluster  
 $(\|\mathbf{x}_1 - c_1\|^2 + (\|\mathbf{x}_3 - c_2\|^2 + (\|\mathbf{x}_2 - c_1\|^2 + (\|\mathbf{x}_4 - c_2\|^2))$

# K-means Clustering

## ► Formal Formulation

1. Assignment Step (Fix  $\mathbf{C}_{K \times d}$  and update  $\mathbf{W}_{N \times K}$ )

If  $\|\mathbf{x}_i - \mathbf{c}_k\|^2 < \|\mathbf{x}_i - \mathbf{c}_{\forall j \neq k}\|^2$ :  $w_{ik} = 1$   
 (else  $w_{ik} = 0$ )

$i = 1, \dots, N; j, k = 1, \dots, K$

update until

2. Update Step (Fix  $\mathbf{W}_{N \times K}$  and update  $\mathbf{C}_{K \times d}$ )

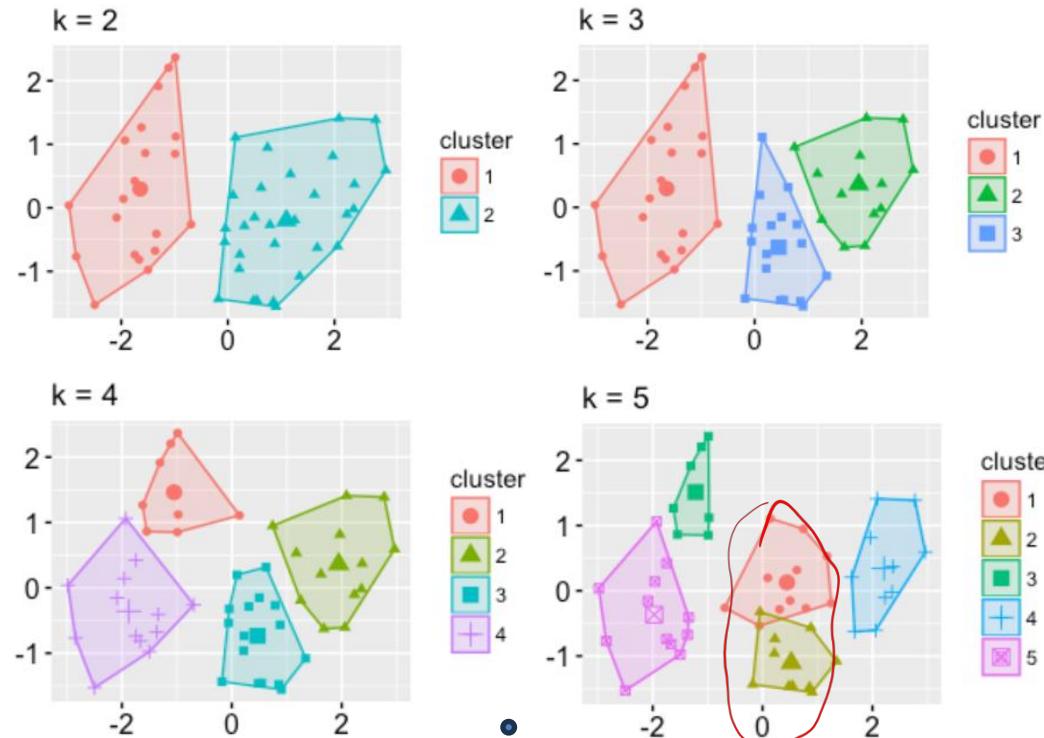
$$\frac{dJ(\mathbf{W}, \mathbf{C})}{d\mathbf{C}} = \begin{bmatrix} \frac{\partial \sum_{i=1}^N \sum_{k=1}^K w_{ik} (\mathbf{x}_i - \mathbf{c}_k)^T (\mathbf{x}_i - \mathbf{c}_k)}{\partial \mathbf{c}_1} \\ \vdots \\ \frac{\partial \sum_{i=1}^N \sum_{k=1}^K w_{ik} (\mathbf{x}_i - \mathbf{c}_k)^T (\mathbf{x}_i - \mathbf{c}_k)}{\partial \mathbf{c}_K} \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N w_{i1} (\mathbf{x}_i - \mathbf{c}_1) \\ \vdots \\ -2 \sum_{i=1}^N w_{iK} (\mathbf{x}_i - \mathbf{c}_K) \end{bmatrix} = 0$$

$$\mathbf{c}_k = \frac{\sum_{i=1}^N w_{ik} \mathbf{x}_i}{\sum_{i=1}^N w_{ik}}, k = 1, \dots, K$$

center for each cluster.

# K-means Clustering: Limitations

➤ Need to predefine  $K$  *→ hyperparameter*

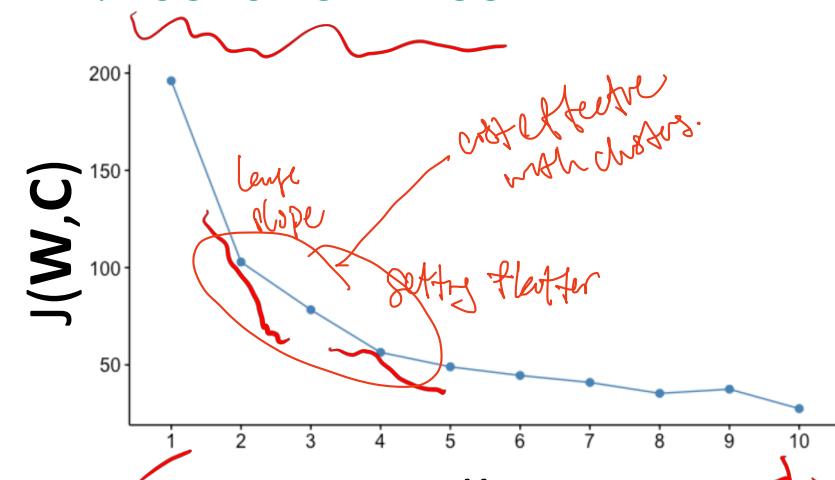


What is the optimal  $K$ ?

Use elbow method.

Elbow Method:

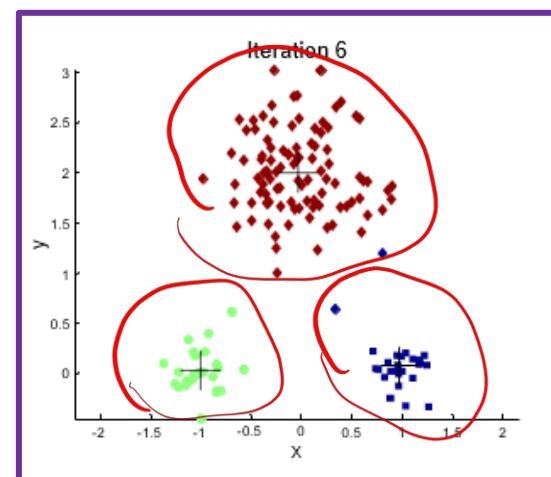
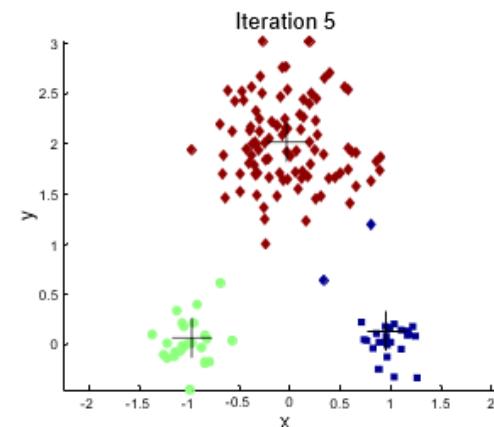
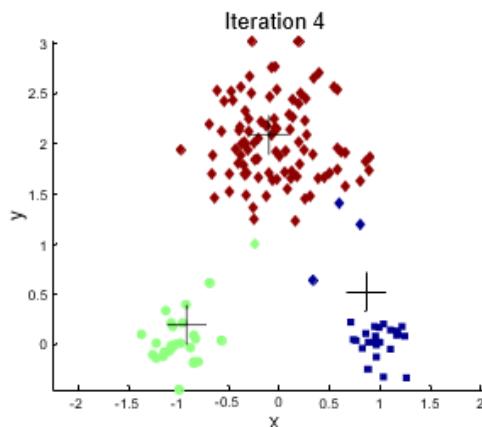
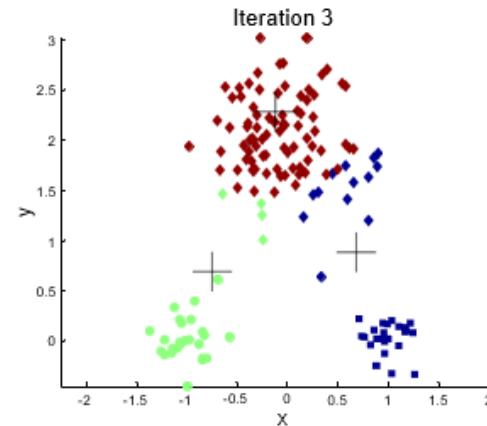
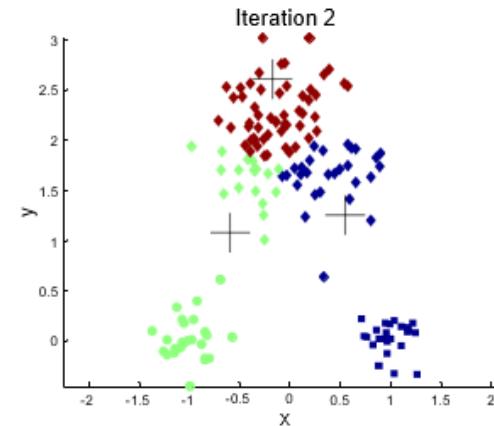
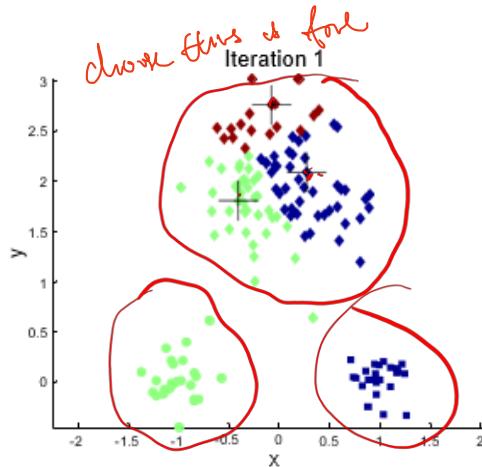
1. K-means for different  $K$
2. For each  $K$ , calculate  $J(\mathbf{W}, \mathbf{C})$
3. Plot the curve of  $J(\mathbf{W}, \mathbf{C})$  over  $K$
4. "bend" or "knee"



why not minimum cost function  $J(\mathbf{w}, \mathbf{c})$ ?  
then  $K$  will be number of samples themselves

# K-means Clustering: Limitations

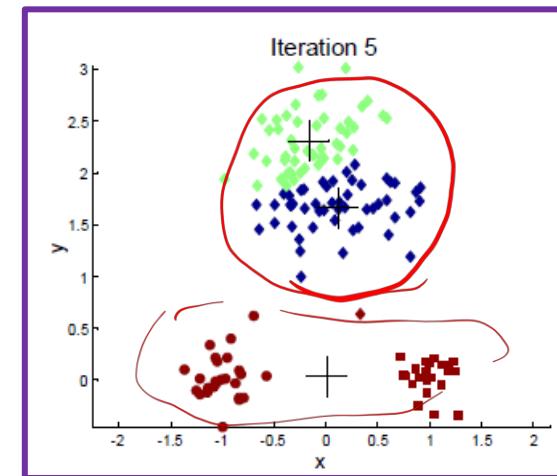
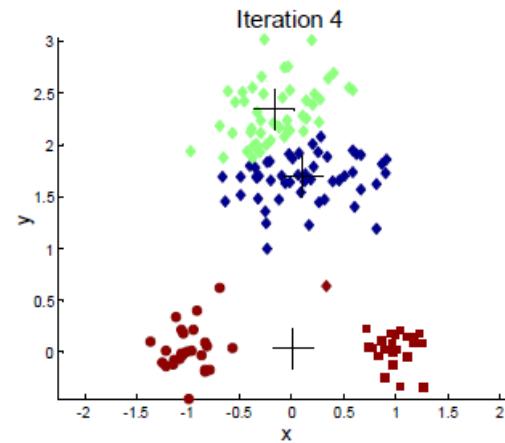
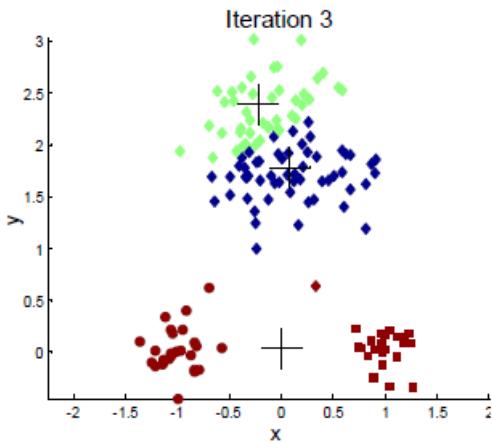
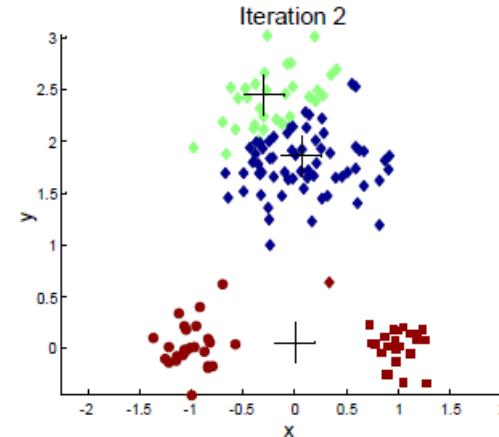
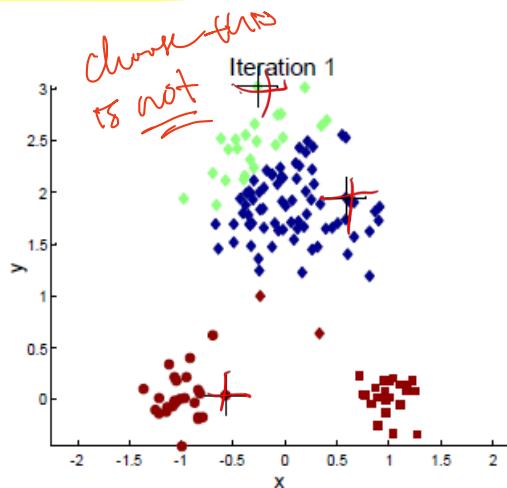
➤ **Sensitive to initial centroid selection**



Optimal

# K-means Clustering: Limitations

➤ Sensitive to initial centroid selection

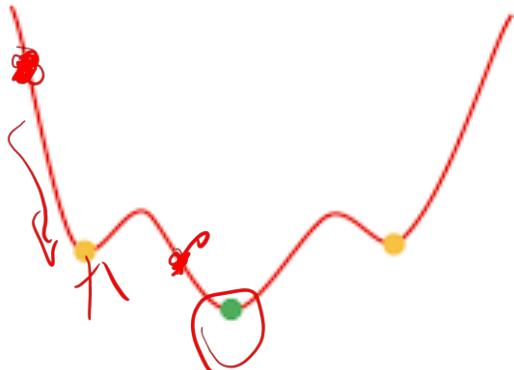


Suboptimal

# K-means Clustering: Limitations

➤ **Sensitive to initial centroid selection**

$J(W, C)$  is non-convex



The K-means algorithm always converges.

The K-means algorithm can be configured to use Manhattan Distance.

K-means is not guaranteed to find a global minimum; it finds only local minimum.

↳ which depends on your initial centroid chosen

How to guarantee a global minimum?

Given the number of clusters and data points, compute the cost function for every set possible, minimum resulting class of errors will be the global minimum.

Computational impractical, so can try the centroids in first place, but not guaranteed, only possible to find.

This is an  
NP-hard problem

# K-means Clustering: Limitations

## ➤ **Sensitive to initial centroid selection**

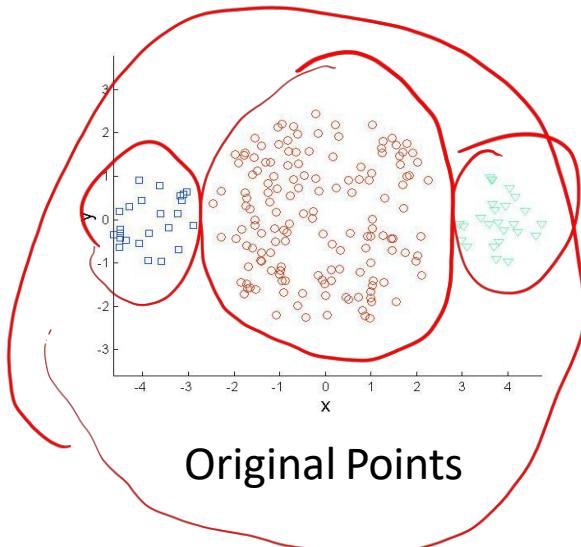
### Practically Possible Trials:

- Multiple runs with different random initializations
  - Helps, but probability is not on your side
- Select more than  $K$  initial centroids and then select among these initial centroids
  - Select most widely separated
- Post-processing
  - Eliminate small clusters that may represent outliers
  - Split ‘loose’ clusters, i.e., clusters with relatively high  $J(W, C)$
  - Merge clusters that are ‘close’ and that have relatively low  $J(W, C)$

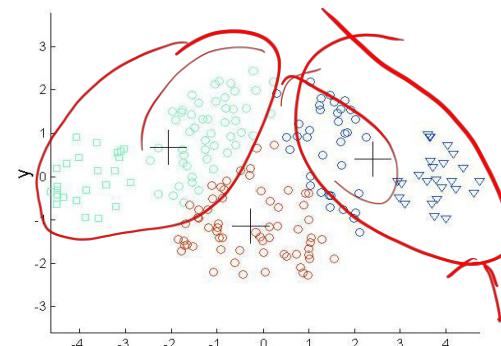
### K-Means++: Improved initialization for K-Means

# K-means Clustering: Limitations

➤ **Assume spherical, equal-sized, similar-density clusters**

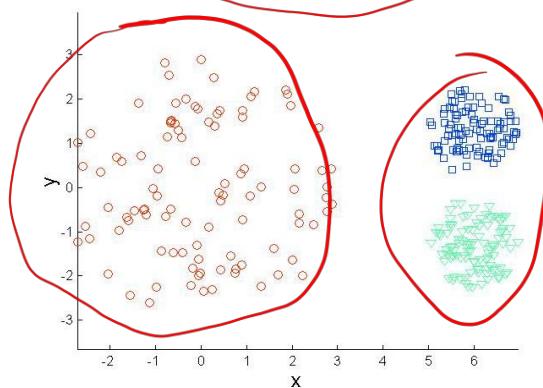


Original Points

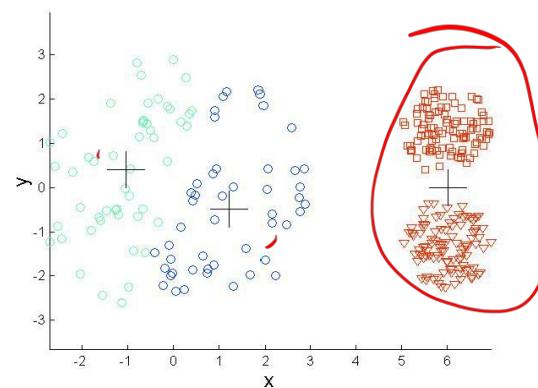


K-means (3 Clusters)

Differ in size!



Original Points



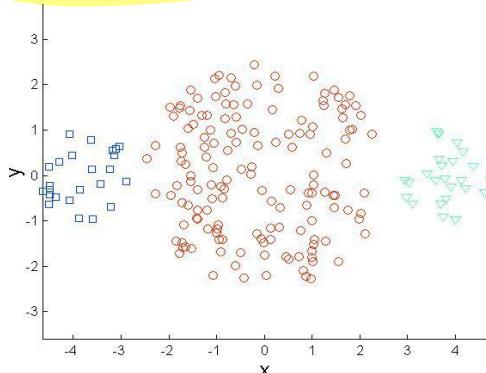
K-means (3 Clusters)

Differ in Density!

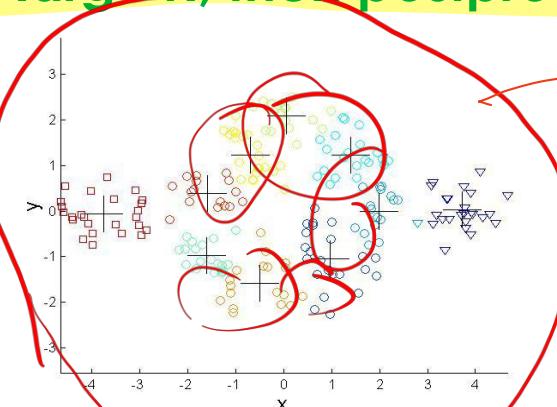
# K-means Clustering: Limitations

➤ Assume spherical, equal-sized, similar-density clusters

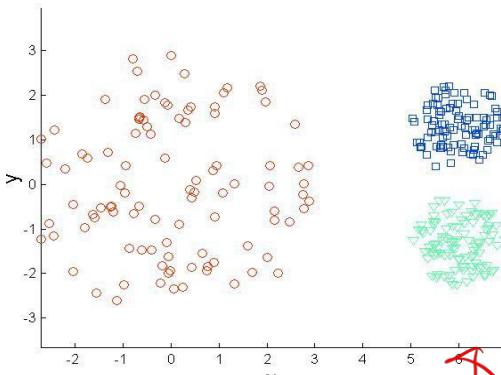
Possible Solution: Use large K, then postprocess



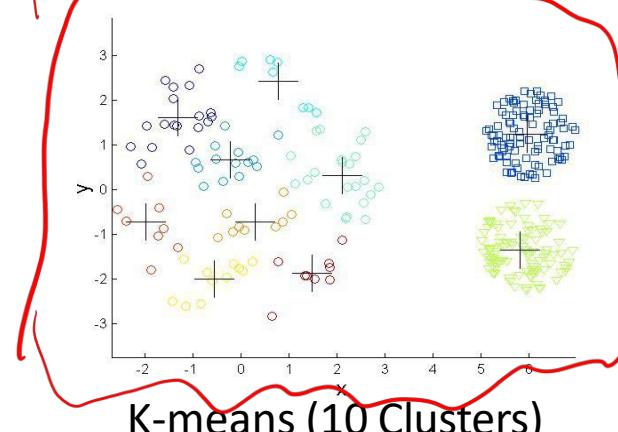
Original Points



Differ in size



Original Points



Differ in Density

# K-means Clustering: Example

Consider the following data points:

Sample ID	1	2	3	4	5	6
Sample	1	2	3	10	11	12

We'd like to group them into two groups by K-means according to Euclidean distance. Assume the initial centroid is Sample 1 for Group A, and Sample 4 for Group B. Upon convergence, what will be the two centres?

Which of the following statements is/are correct?

- (a) K-means clustering is suitable for datasets with spiral-shaped clusters.
- (b) Given the same set of data points, the number of well-separated clusters that can be formed is always more than the number of center-based clusters that can be formed.
- (c) Running K-means clustering multiple times with different random initializations can guarantee reaching the global minimum solution.
- (d) Each data point can belong to multiple clusters for K-means clustering.
- (e) None

Answer: (e)

- (a) is wrong because K-means clustering is suitable for spherical clusters.
- (b) is wrong because the number of well-separated clusters that can be formed is always no more than the number of center-based clusters that can be formed, given the same set of data points.
- (c) is wrong because it cannot guarantee reaching the global minimum solution. It only increases the likelihood.
- (d) is wrong because each data point can belong to only one cluster for K-means clustering. K-means clustering belongs to hard clustering which assigns each data point to exact one cluster.

# K-means Clustering: Example

Consider the following data points:

Sample ID	1	2	3	4	5	6
Sample	1	2	3	10	11	12

We'd like to group them into two groups by K-means according to Euclidean distance. Assume the initial centroid is Sample 1 for Group A, and Sample 4 for Group B. Upon convergence, what will be the two centres?

Initialization:  $c_A = 1, c_B = 10$

1<sup>st</sup> iteration:

Assignment:  $D(x_1, c_A) = 0, D(x_1, c_B) = 9$

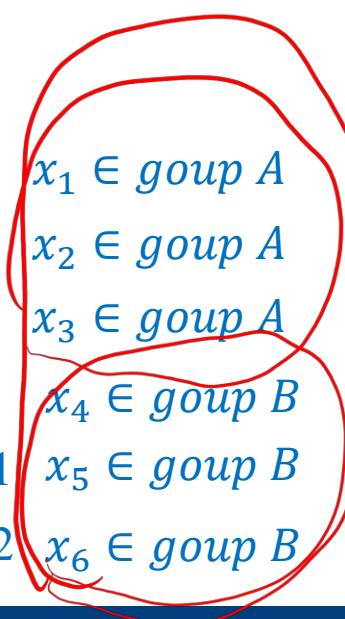
$\nwarrow D(x_2, c_A) = 1, D(x_2, c_B) = 8$

$\nearrow D(x_3, c_A) = 2, D(x_3, c_B) = 7$

$\nexists D(x_4, c_A) = 9, D(x_4, c_B) = 0$

$\nearrow D(x_5, c_A) = 10, D(x_5, c_B) = 1$

$\nearrow D(x_6, c_A) = 11, D(x_6, c_B) = 2$



Update: *Centroid based on center & current points AB*

$$c_A = \frac{x_1 + x_2 + x_3}{3} = 2$$

$$c_B = \frac{x_4 + x_5 + x_6}{3} = 11$$

# K-means Clustering: Example

Consider the following data points:

Sample ID	1	2	3	4	5	6
Sample	1	2	3	10	11	12

We'd like to group them into two groups by K-means according to Euclidean distance. Assume the initial centroid is Sample 1 for Group A, and Sample 4 for Group B. Upon convergence, what will be the two centres?

$$c_A = 2, c_B = 11$$

2<sup>nd</sup> iteration:

Assignment:  $D(x_1, c_A) = 1, D(x_1, c_B) = 10$

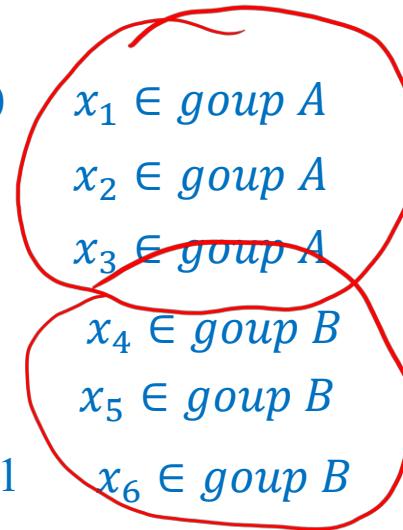
$D(x_2, c_A) = 0, D(x_2, c_B) = 9$

$D(x_3, c_A) = 1, D(x_3, c_B) = 8$

$D(x_4, c_A) = 8, D(x_4, c_B) = 1$

$D(x_5, c_A) = 9, D(x_5, c_B) = 0$

$D(x_6, c_A) = 10, D(x_6, c_B) = 1$



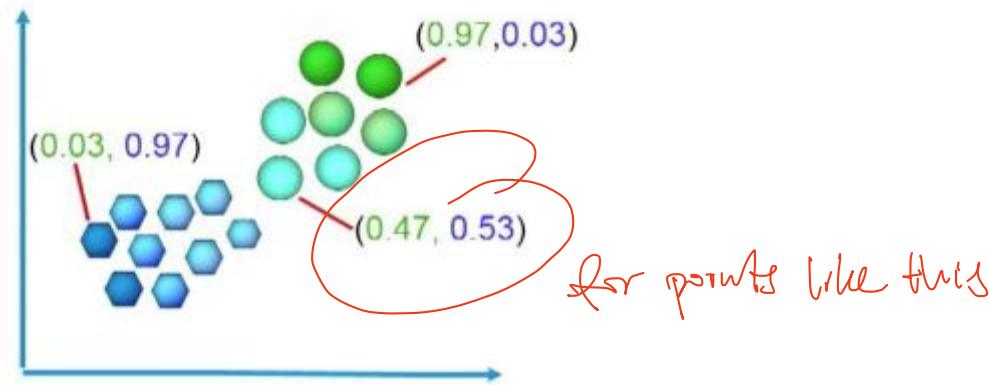
Update:

$$c_A = \frac{x_1 + x_2 + x_3}{3} = 2$$

$$c_B = \frac{x_4 + x_5 + x_6}{3} = 11$$

No changes!

converged,  
stop?



# Fuzzy C-Means Clustering

## -Soft Clustering

# Fuzzy C-Means Clustering

## ➤ Objective Function

Find the optimal values for  $\mathbf{W}^*_{N \times C}$  and  $\mathbf{C}^*_{C \times d}$  which minimizes:

$$J(\mathbf{W}, \mathbf{C}) = \sum_{i=1}^N \sum_{k=1}^C (w_{ik})^r \| \mathbf{x}_i - \mathbf{c}_k \|^2$$

Similar to k-means -  
Degree of membership -  
*(W) is no longer a discrete value  
of 0 and 1. If is between 0-1.*  
Squared Euclidean distance

where  $w_{ik} \in [0,1]$ :  $w_{ik} = \frac{1}{\sum_{j=1}^C \left( \frac{\|\mathbf{x}_i - \mathbf{c}_k\|}{\|\mathbf{x}_i - \mathbf{c}_j\|} \right)^{\frac{2}{r-1}}}$  tells the degree of  $\mathbf{x}_i$  belonging to cluster  $\mathbf{c}_k$

*level of a point overlap or uncertainty for multiple classes*  
 $r > 1$  is the fuzzifier: determines the level of cluster fuzziness.  
(usually,  $1.25 \leq r \leq 2$ )

In Fuzzy C-Means clustering, the membership degrees of a data point across all clusters always sum to 1.

# Fuzzy C-Means Clustering

## Algorithm Formulation

Save algorithm,

### 1. Assignment Step (Fix $\mathbf{c}_{C \times d}$ and update $\mathbf{W}_{N \times c}$ )

$$w_{ik} = \frac{1}{\sum_{j=1}^C \left( \frac{\|\mathbf{x}_i - \mathbf{c}_k\|}{\|\mathbf{x}_i - \mathbf{c}_j\|} \right)^{r-1}}$$

### 2. Update Step (Fix $\mathbf{W}_{N \times c}$ and update $\mathbf{c}_{C \times d}$ )

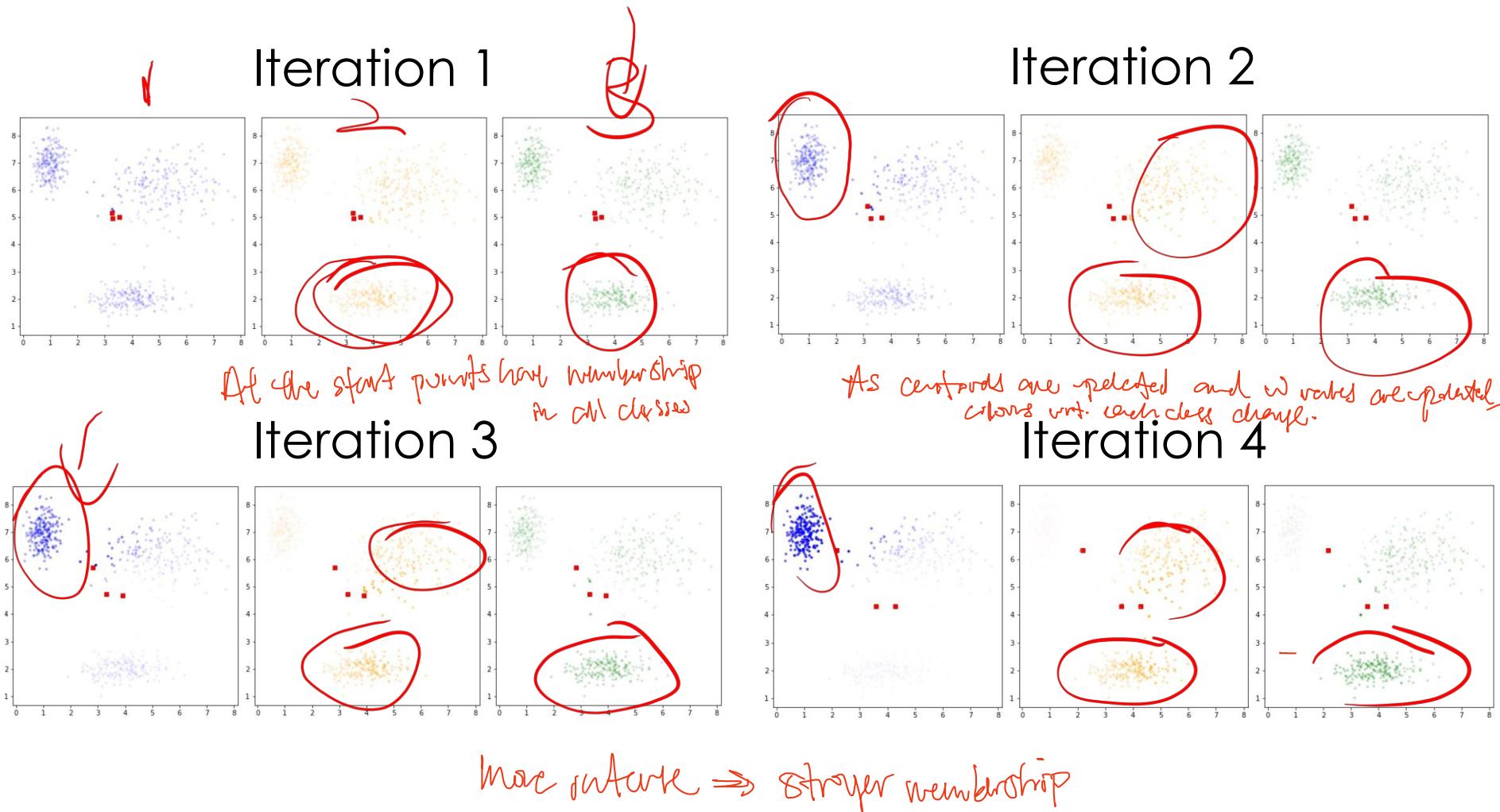
$$\frac{dJ(\mathbf{W}, \mathbf{C})}{d\mathbf{C}} = \begin{bmatrix} \frac{\partial \sum_{i=1}^N \sum_{k=1}^C (w_{ik})^r (\mathbf{x}_i - \mathbf{c}_k)^T (\mathbf{x}_i - \mathbf{c}_k)}{\partial \mathbf{c}_1} \\ \vdots \\ \frac{\partial \sum_{i=1}^N \sum_{k=1}^C (w_{ik})^r (\mathbf{x}_i - \mathbf{c}_k)^T (\mathbf{x}_i - \mathbf{c}_k)}{\partial \mathbf{c}_K} \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^N (w_{i1})^r (\mathbf{x}_i - \mathbf{c}_1) \\ \vdots \\ -2 \sum_{i=1}^N (w_{iC})^r (\mathbf{x}_i - \mathbf{c}_K) \end{bmatrix} = \mathbf{0}$$

$$\mathbf{c}_k = \frac{\sum_{i=1}^N (w_{ik})^r \mathbf{x}_i}{\sum_{i=1}^N (w_{ik})^r}, k = 1, \dots, C$$

only with fixed  $r$ , the cost function is convex

# Fuzzy C-Means Clustering

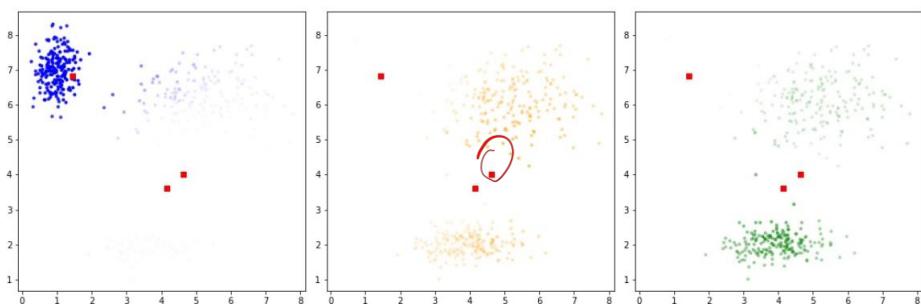
## ➤ Visualization



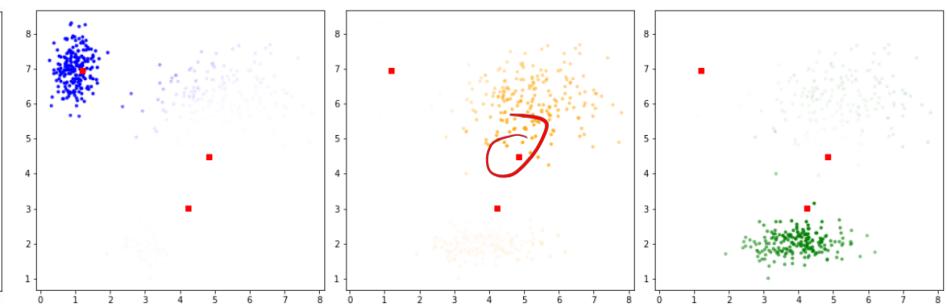
# Fuzzy C-Means Clustering

## ➤ Visualization

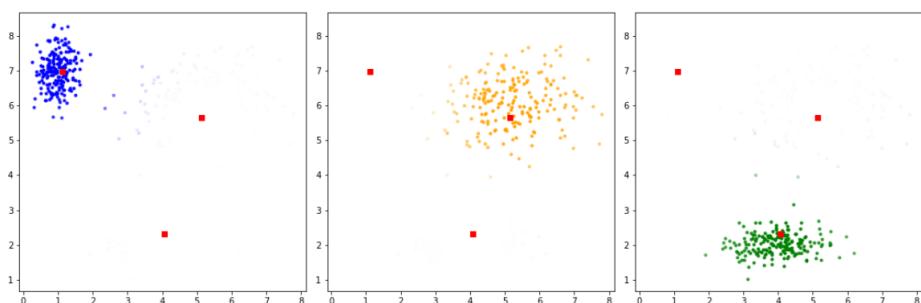
Iteration 5



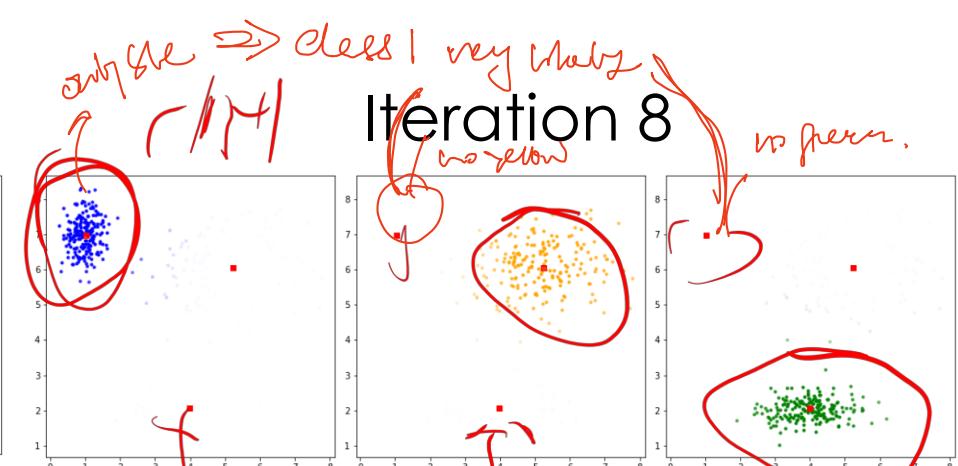
Iteration 6



Iteration 7



Iteration 8



# Summary By Quick Quiz

✓ Four Types of Clusters

allow method to select optimal

✓ Two Types of Clustering

✓ K-means Clustering

1. Centroid Initialization

2. Assignment (to nearest centroid)

3. Update

$$\mathbf{c}_k = \frac{\sum_{i=1}^N w_{ik} \mathbf{x}_i}{\sum_{i=1}^N w_{ik}}, k = 1, \dots, K$$

Stop when/centre new-centroid < threshold

✓ Fuzzy C-means Clustering

1. Centroid Initialization

2. Assignment

$$w_{ik} = \frac{1}{\sum_{j=1}^C \left( \frac{\|\mathbf{x}_i - \mathbf{c}_k\|}{\|\mathbf{x}_i - \mathbf{c}_j\|} \right)^{\frac{2}{r-1}}}$$

3. Update

$$\mathbf{c}_k = \frac{\sum_{i=1}^N (w_{ik})^r \mathbf{x}_i}{\sum_{i=1}^N (w_{ik})^r}, k = 1, \dots, C$$

# EE2213 Introduction to AI

## Lecture 11: Neural Networks

Dr. WANG Si

[si.wang@nus.edu.sg](mailto:si.wang@nus.edu.sg)

Electrical and Computer Engineering Department  
National University of Singapore

# OVERVIEW OF COURSE CONTENTS



- **Introduction (Shaojing)**

- What is AI
- Applications of AI
- AI agent

- **Search (Shaojing)**

- Uninformed search algorithms: breadth-first, depth-first, uniform-cost(Dijkstra's algorithm)
- Informed search algorithms: greedy best-first, A\*

- Applications

- **Optimisation (Shaojing)**

- Linear programming
- Convex problems
- Applications

- **Machine learning (Wang Si) (Weeks 6-9)**

- Supervised and unsupervised learning: regression, classification, clustering
- Neural networks and deep learning
- Applications

- **Knowledge representation (Wang Si) (Weeks 10-11)**

- Knowledge Representation and Reasoning
  - Propositional Logic
  - Applications
- **Ethical considerations (Shaojing)**
- Bias in AI
  - Privacy concerns
  - Societal impact

(Week 10, Monday)

No Class on 20 Oct.! A recording  
will be uploaded on Canvas.

# AGENDA

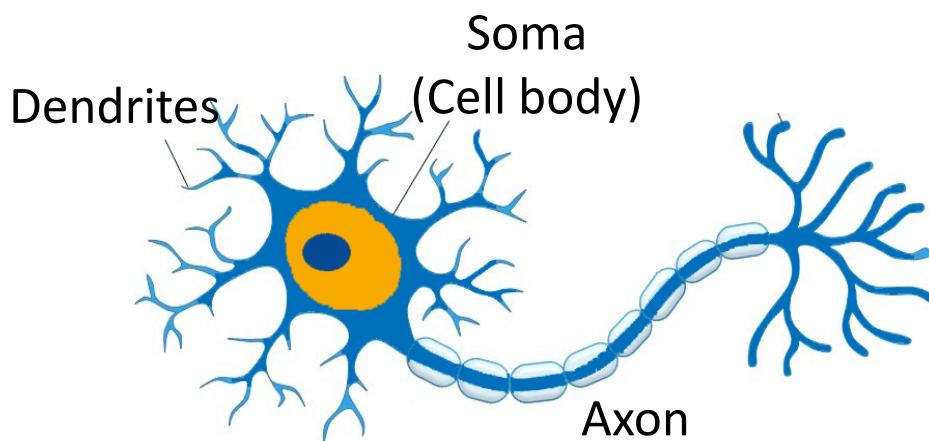
- We will discuss:
  - Perceptron
  - Multilayer Perceptron (MLP)
  - MLPs for Image Classification

At the end of this lecture, you should be able to:

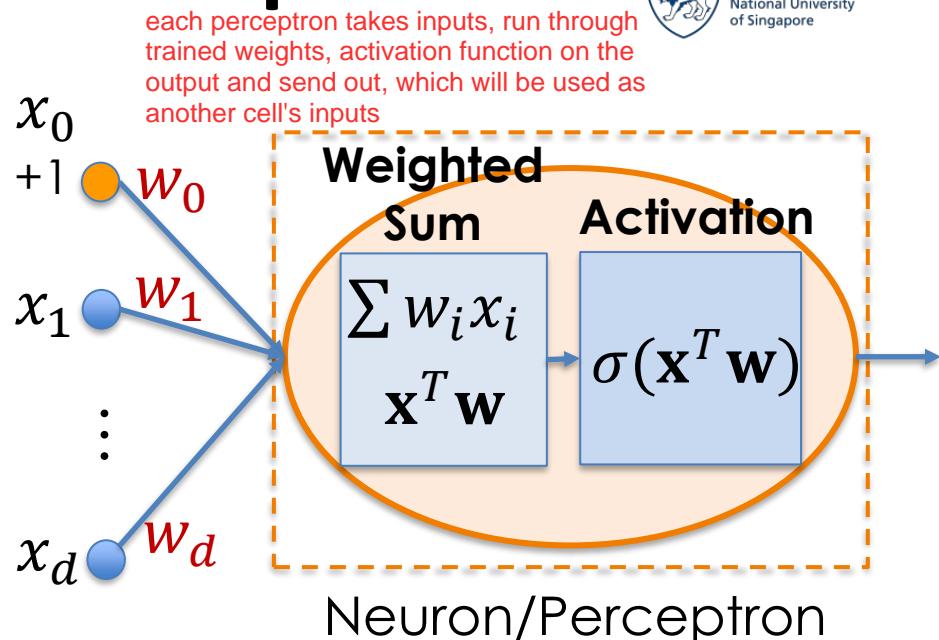
- ✓ Describe perceptron model
- ✓ Explain the architecture and learning process of MLP
- ✓ Apply MLP for image classification and identify the limitations

# Perceptron

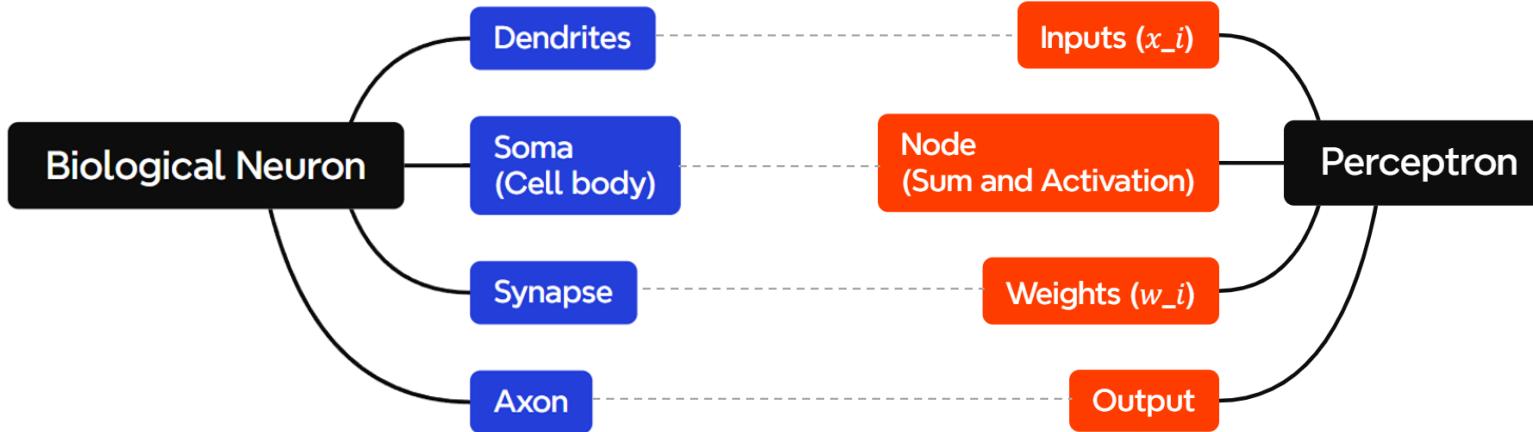
# Biological Neuron & Perceptron



The input to the activation function ( $xTw$ ) of the perceptron is a linear combination of features.  
relu means not probabilistic output



The activation function of the Perceptron is used to transform the predicted output, not every input feature individually.

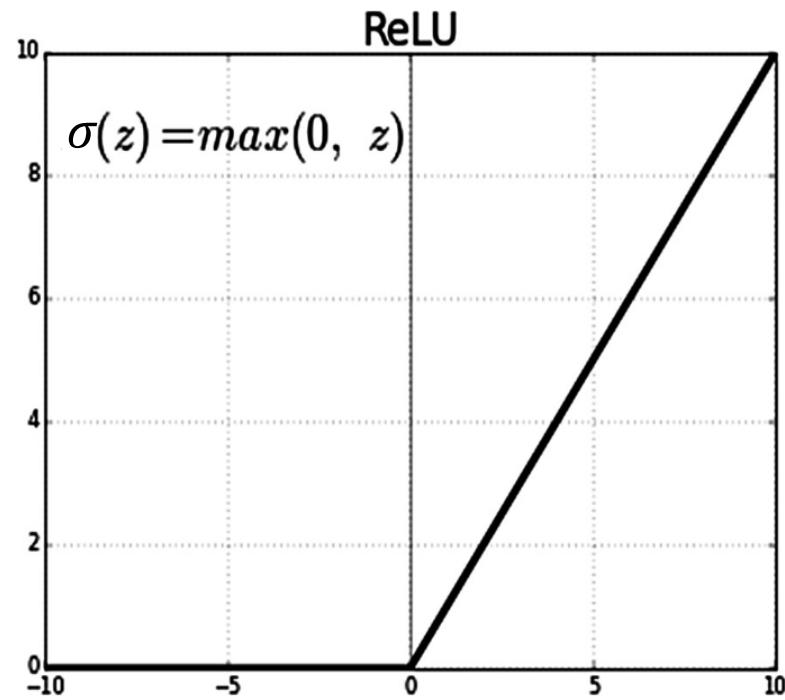
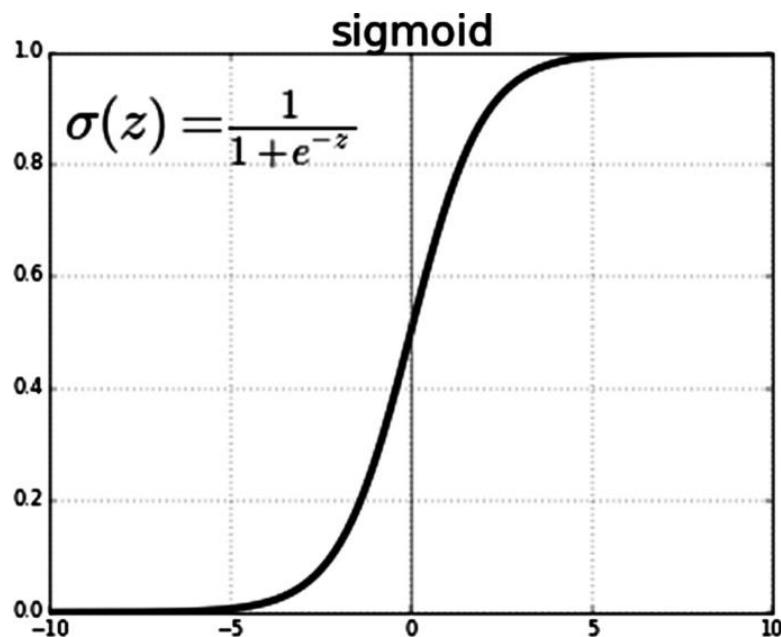


$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

# Perceptron

## ➤ Activation Function (Non-Linear)



rectified linear unit

if  $z < 0$ , function is  
0, else,  $z$   
not probabilistic

What might be the reasons for the difference in performance between the single-layer neural network and the multi-layer neural network?

**Solution:**

The difference in performance between the single-layer neural network and the multi-layer neural network can be attributed to the fact that the single-layer neural network is a linear classifier, while the multi-layer neural network can learn non-linear decision boundaries. In this case, the multi-layer neural network is able to capture more complex relationships between the input features and the output label, leading to better performance on the validation set.

How might you modify the single-layer neural network to improve its performance, and what are the advantages and disadvantages of doing so?

**Solution:**

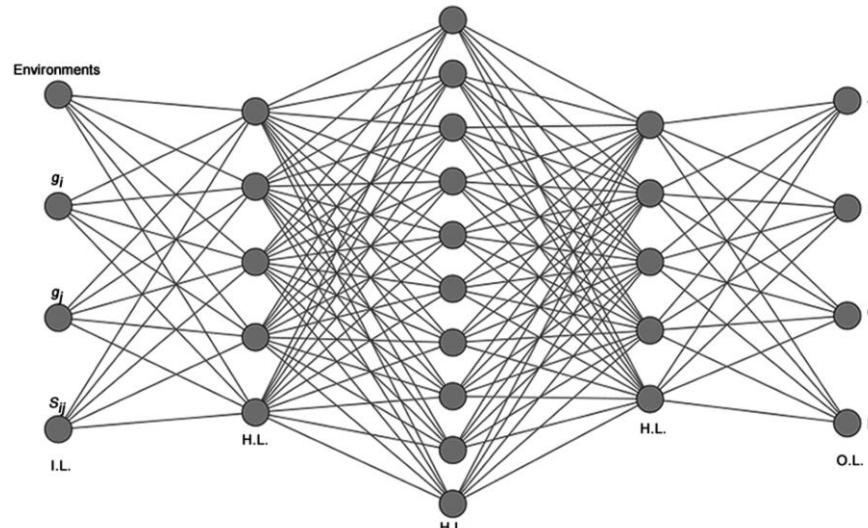
To improve the performance of the single-layer neural network, you could add more features, such as polynomial or interaction terms, to capture non-linear relationships between the input features and the output labels. However, this approach can quickly become computationally expensive and may lead to overfitting if the number of features is too high. Note that the final decision is still based on a linear combination of the transformed features.

What techniques could you use to improve the performance of the multi-layer neural network?

**Solution:**

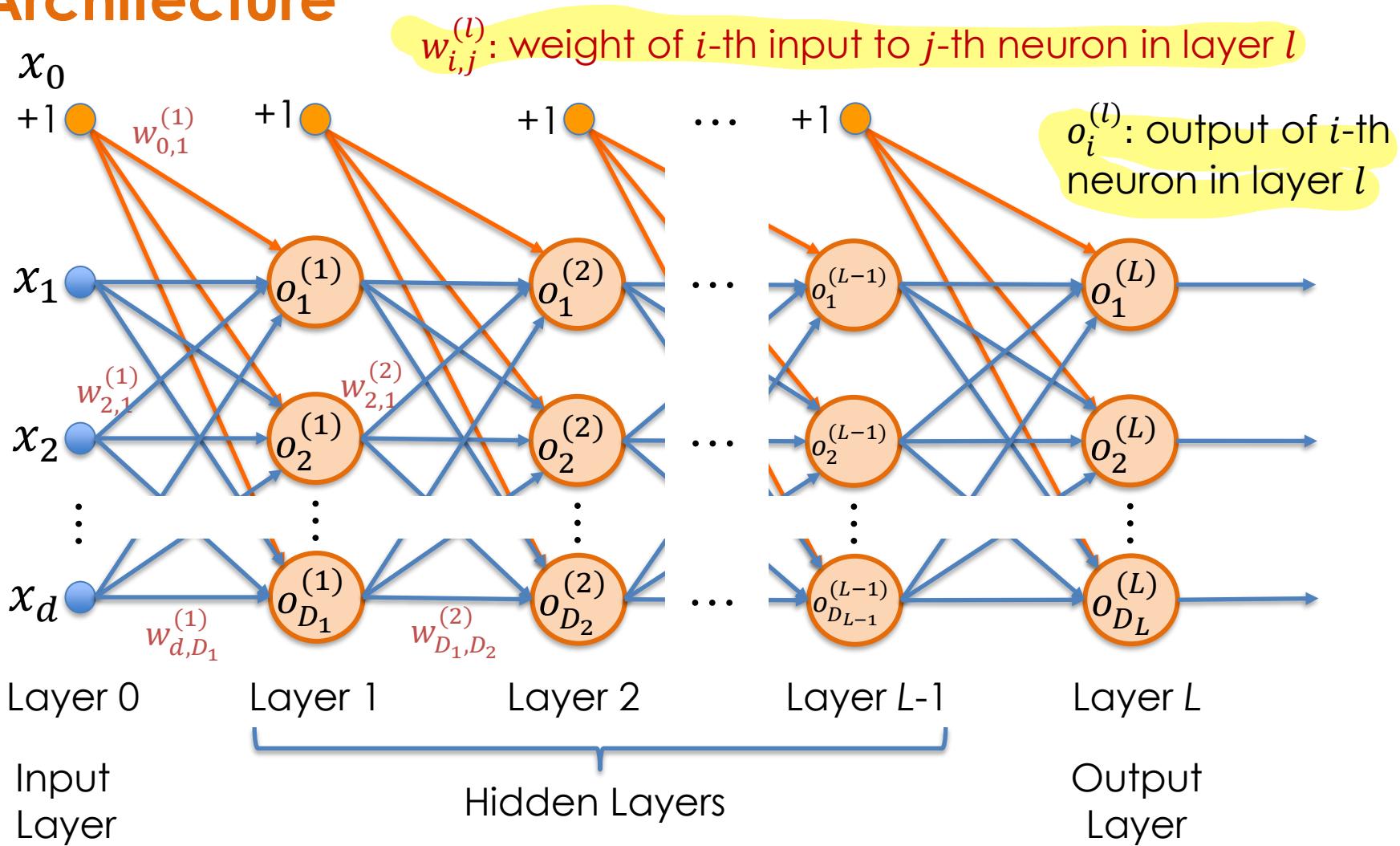
You could try adding more hidden layers (increasing the depth of the network) or increasing the number of neurons in the hidden layer (increasing the width). However, the complexity of the computed function grows exponentially with the depth of an architecture, and increases more slowly with the width. You can find out more in the paper by Eldan and Shamir (2016) on the expressive power of deep neural networks. For further details, see [this link](#). Hence, increasing the depth is generally more efficient. Additionally, while both approaches lead to the multi-layer neural network increasing its capacity to learn more complex patterns, it can also increase the risk of overfitting.

# Multilayer Perceptron (MLP)



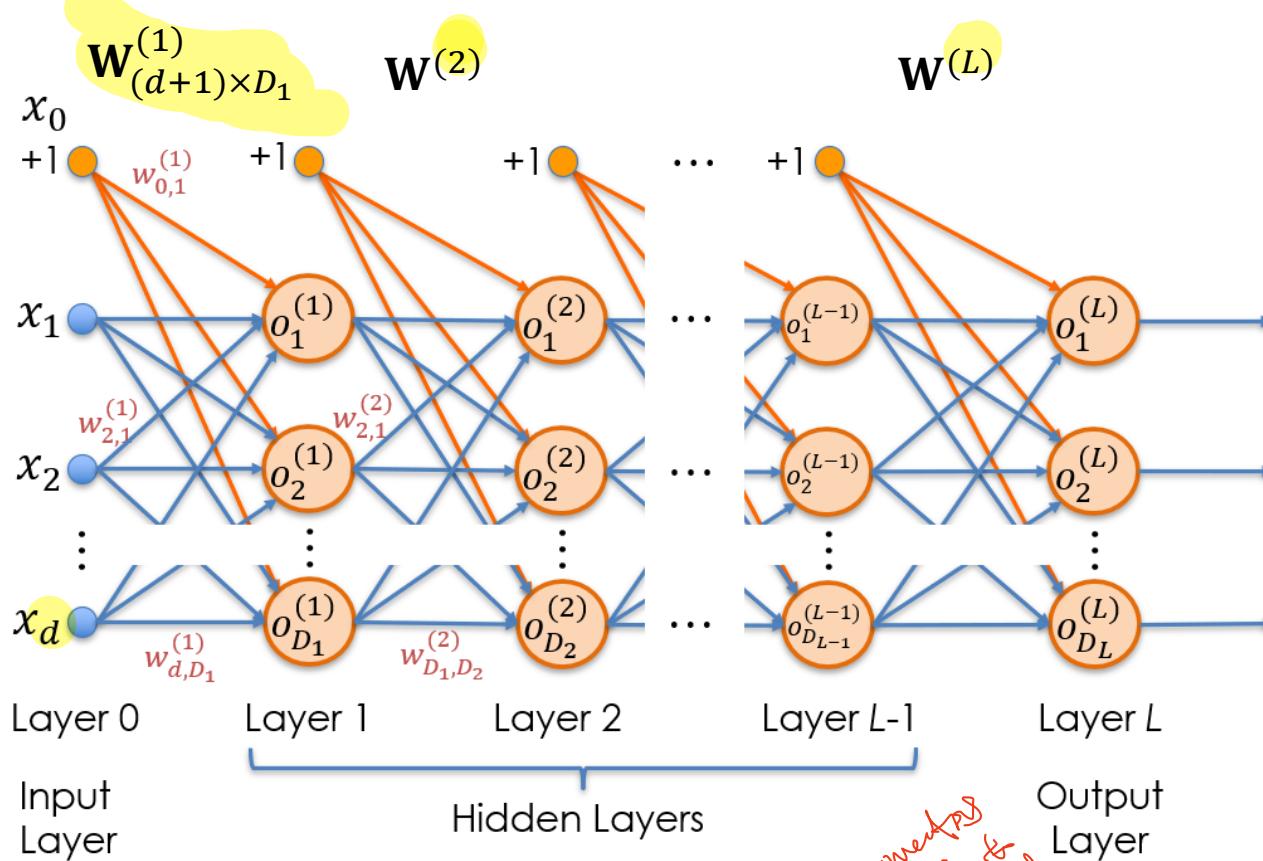
# Multilayer Perceptron (MLP)

## ➤ Architecture



# Multilayer Perceptron (MLP)

## ➤ Architecture



e.g., 3 inputs to layer 3,  
2 neurons in layer 3

$\mathbf{W}^L$   
input x neuron (next)

Inputs (combine forward)  
↓

$$\mathbf{W} = (d+1) \times D_l$$

↑ number of neurons.

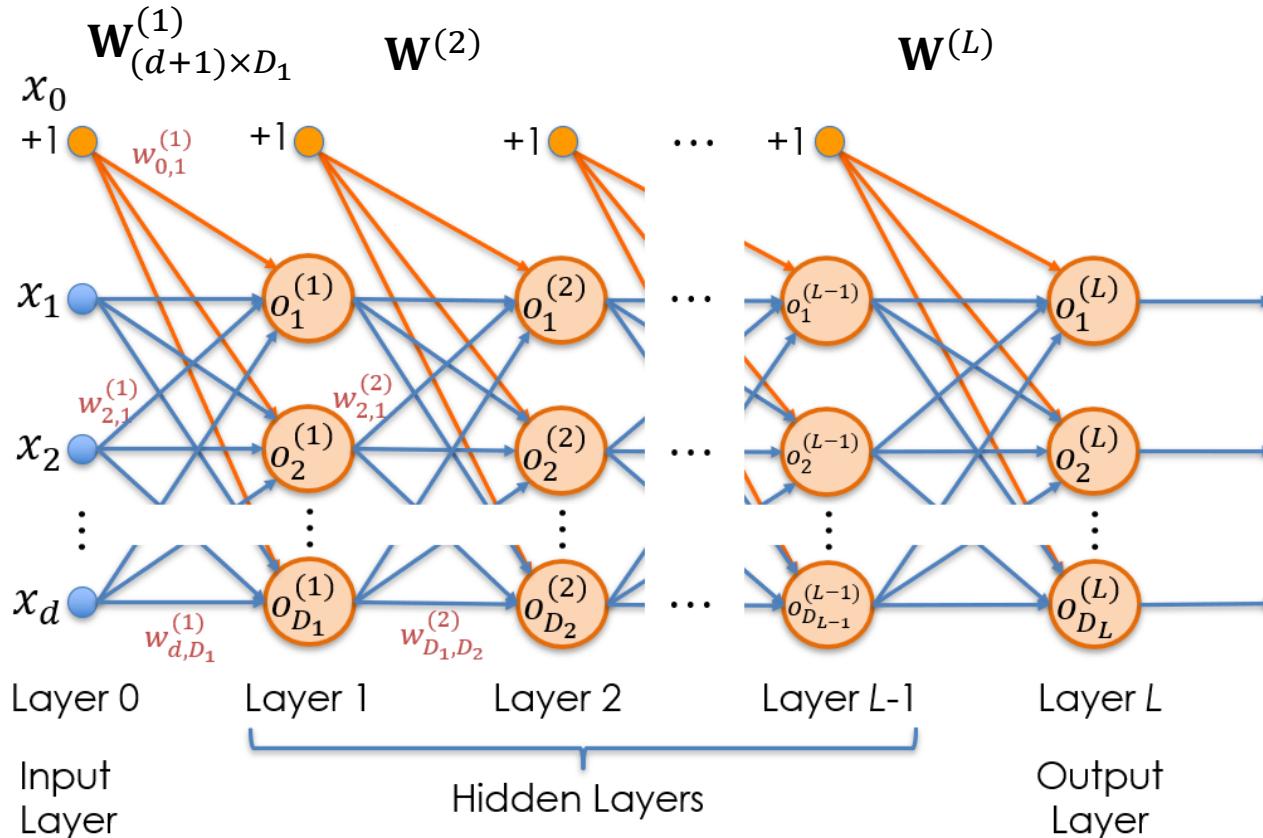
$$\mathbf{W}^{(l)} = \begin{bmatrix} w_{0,1}^{(l)} & \cdots & w_{0,D_l}^{(l)} \\ \vdots & \ddots & \vdots \\ w_{D_{l-1},1}^{(l)} & \cdots & w_{D_{l-1},D_l}^{(l)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^{(l)} & \cdots & \mathbf{w}_{D_l}^{(l)} \end{bmatrix}$$

weight vectors  
from +1 to  
neurons

Weight vector for first neuron in layer  $l$

# Multilayer Perceptron (MLP)

# ➤ Architecture



$$\mathbf{W}^{(l)} = \begin{bmatrix} w_{0,1}^{(l)} & \dots & w_{0,D_l}^{(l)} \\ \vdots & \ddots & \vdots \\ w_{D_{l-1},1}^{(l)} & \dots & w_{D_{l-1},D_l}^{(l)} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^{(l)} & \dots & \mathbf{w}_{D_l}^{(l)} \end{bmatrix}$$

Weight vector for first neuron in layer  $l$

$$\mathbf{f}_W(\mathbf{x}) = \text{Softmax}([\mathbf{1}, \text{ReLU}([\mathbf{1}, \text{ReLU}(\mathbf{x}^T \mathbf{W}^{(1)})] \mathbf{W}^{(2)})] \mathbf{W}^{(3)})$$

where  $\mathbf{x}$  is the input vector augmented with 1,  $\mathbf{W}^{(1)}$ ,  $\mathbf{W}^{(2)}$  and  $\mathbf{W}^{(3)}$  are weight matrices. Which of the following statements is/are correct?

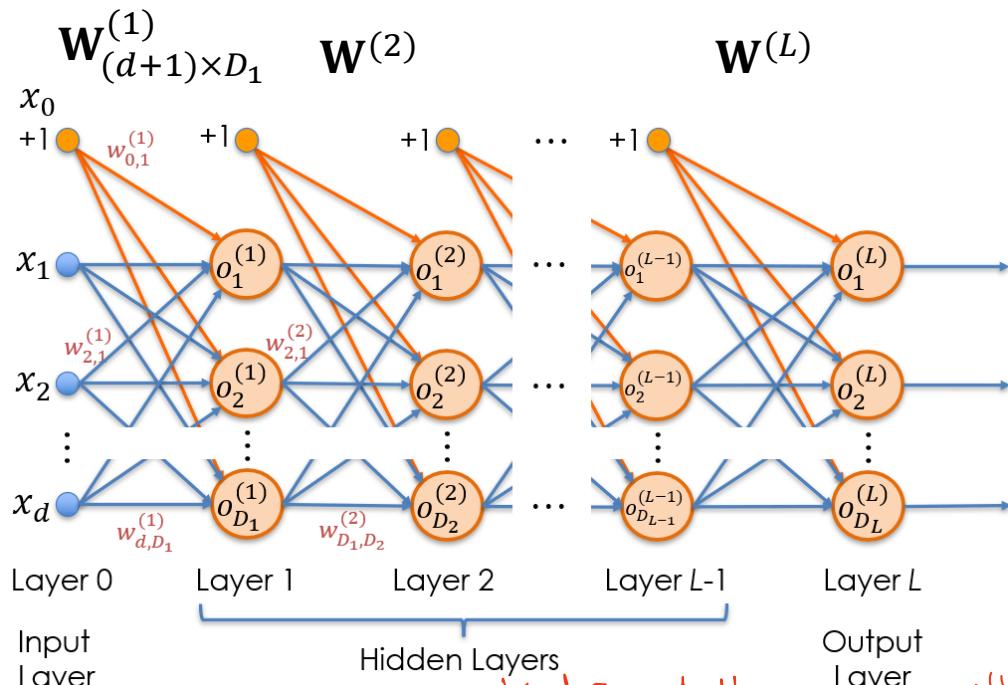
- (a) This MLP has 3 layers of neurons. *hidden layers*
  - (b) The number of neurons in each layer of this MLP is equal to the number of rows in the weight matrix for that layer.
  - (c) The number of image classes is equal to the number of columns in  $W^{(3)}$ .
  - (d) The dimension of  $x$  is equal to the dimension of the flattened image.
  - (e) None.

**Answer:** (a)(c)

- (b) is wrong because the number of neurons in each layer of this MLP should be equal to the number of columns in the weight matrix for that layer.
  - (d) is wrong because the dimension of  $x$  is equal to the dimension of the flattened image + 1 (for the bias)

# Multilayer Perceptron (MLP)

## ➤ Architecture



$$o_i^{(l)} = \sigma^{(l)} \left( \sum_{j=1}^{D_{l-1}} o_j^{(l-1)} w_{j,i}^{(l)} + w_{0,i}^{(l)} \right) = \sigma^{(l)} ([1, \mathbf{o}^{(l-1)}] \mathbf{w}_i^{(l)})$$

$$\mathbf{o}^{(l)} = \sigma^{(l)} ([1, \mathbf{o}^{(l-1)}] \mathbf{W}^{(l)})$$

$$\mathbf{W}^{(l)} = \begin{bmatrix} w_{0,1}^{(l)} & \dots & w_{0,D_l}^{(l)} \\ \vdots & \ddots & \vdots \\ w_{D_{l-1},1}^{(l)} & \dots & w_{D_{l-1},D_l}^{(l)} \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{w}_1^{(l)} & \dots & \mathbf{w}_{D_l}^{(l)} \end{bmatrix}$$

Weight vector for first neuron in layer  $l$

*weighted sum of all neurons connected to its current layer*

*start from input  $\rightarrow$  output*

*apply activation function*

A MLP is essentially a nested function

$$f_{\mathbf{W}}(\mathbf{x}) = \sigma^{(L)}([1, \dots \sigma^{(2)}([1, \sigma^{(1)}(\mathbf{x}^T \mathbf{W}^{(1)})] \mathbf{W}^{(2)}) \dots] \mathbf{W}^{(L)})$$

*nested by going through layers, of weight and activation function.*

# Multilayer Perceptron (MLP)



## ➤ Architecture

Example: A MLP of 2 layers has been constructed as:

$$\mathbf{f}_W(\mathbf{x}) = \sigma([1, \sigma(\mathbf{x}^T \mathbf{W}^{(1)})] \mathbf{W}^{(2)}),$$

where  $\mathbf{x} = [1, 1, 3]^T$ ,  $\mathbf{W}^{(1)} = \mathbf{W}^{(2)} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$ ,  $\sigma(z) = \max(0, z)$ .

Compute network output  $\mathbf{f}_W(\mathbf{x})$ .

# Multilayer Perceptron (MLP)

## ➤ Architecture

Example: A MLP of 2 layers has been constructed as:

$$f_W(x) = \sigma([1, \sigma(x^T W^{(1)})] W^{(2)}),$$

↗ input  
 ↗ sigmoid.  
 ↗ neurons

$$\text{where } x = [1, 1, 3]^T, W^{(1)} = W^{(2)} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}, \sigma(z) = \max(0, z).$$

Compute network output  $f_W(x)$ .

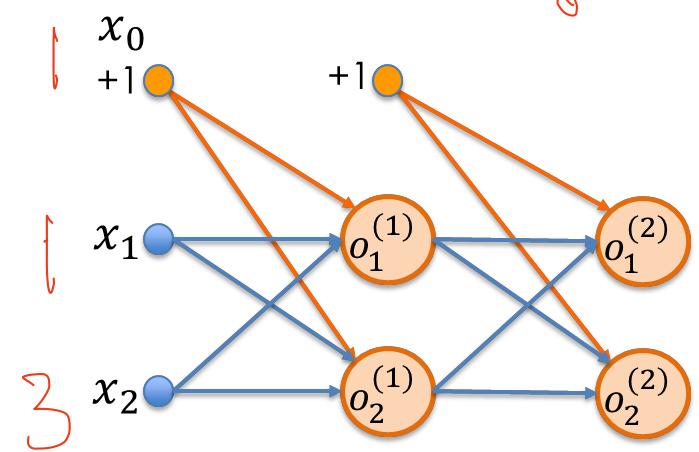
$b + e^{-x^T w}$   
 initial  
is random  
we gradient  
descent

**Solution:**

$$\begin{aligned} \sigma(x^T W^{(1)}) &= \sigma\left([1, 1, 3] \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}\right) \\ &= \sigma([2, -1]) = [2, 0] \end{aligned}$$

$$\begin{aligned} \sigma([1, \sigma(x^T W^{(1)})] W^{(2)}) &= \sigma\left([1, 2, 0] \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}\right) \\ &= \sigma([-1, -2]) = [0, 0] \end{aligned}$$

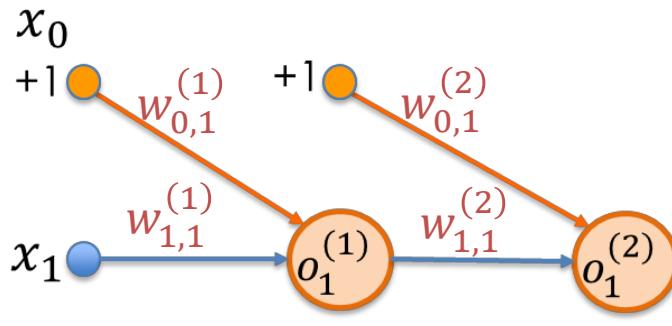
↗ add bias.  
 ↗



# Multilayer Perceptron (MLP)

## ➤ Architecture

Why non-linear activation function  $\sigma(z)$ ?



If linear/affine activation function:

$$\sigma(z) = az + b$$

$$\begin{aligned} o_1^{(1)} &= \sigma(w_{0,1}^{(1)} + w_{1,1}^{(1)}x_1) \\ &= a(w_{0,1}^{(1)} + w_{1,1}^{(1)}x_1) + b \\ &= (aw_{0,1}^{(1)} + b) + aw_{1,1}^{(1)}x_1 \end{aligned}$$

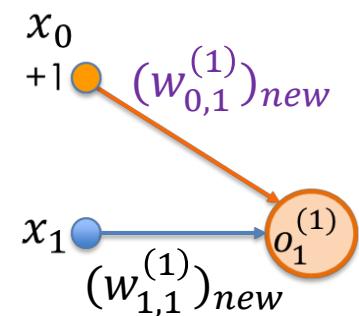
linear: output whatever it receives  
affine = add a bias term

$$\begin{aligned} o_1^{(2)} &= \sigma(w_{0,1}^{(2)} + w_{1,1}^{(2)}o_1^{(1)}) \\ &= a(w_{0,1}^{(2)} + w_{1,1}^{(2)}o_1^{(1)}) + b \\ &= (aw_{0,1}^{(2)} + b) + aw_{1,1}^{(2)}o_1^{(1)} \\ &= (aw_{0,1}^{(2)} + b) + aw_{1,1}^{(2)}(aw_{0,1}^{(1)} + b) + aw_{1,1}^{(2)}x_1 \\ &= (aw_{0,1}^{(2)} + a^2w_{1,1}^{(2)}w_{0,1}^{(1)} + abw_{1,1}^{(2)} + b) \\ &\quad + a^2w_{1,1}^{(2)}w_{1,1}^{(1)}x_1 \end{aligned}$$

$(w_{0,1}^{(1)})_{new}$   
 $(w_{1,1}^{(1)})_{new}$

then the whole network is just a  $\rightarrow$  linear regression

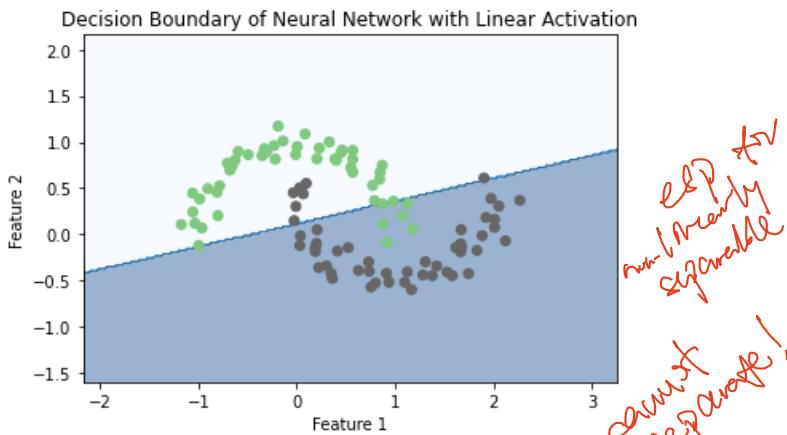
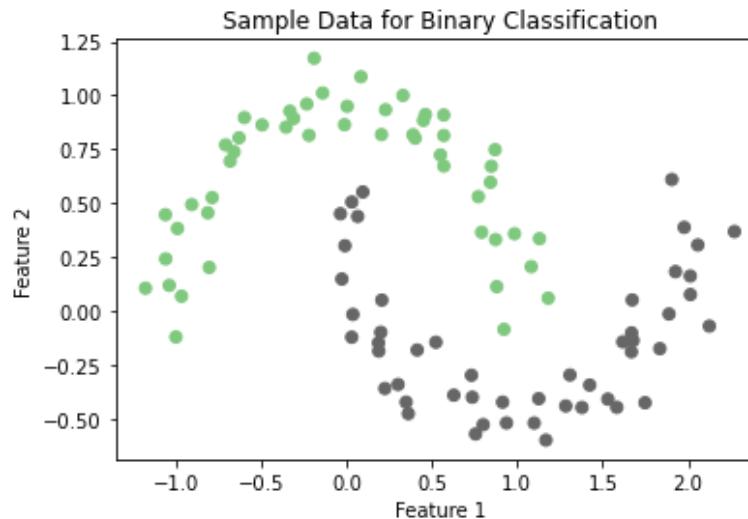
poor performance'



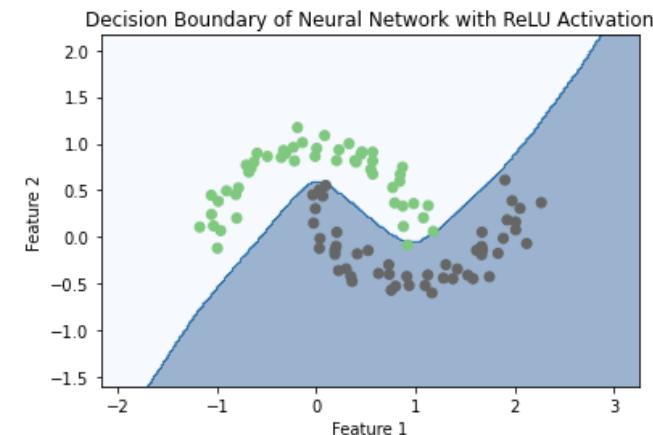
With Linear/affine  $\sigma(z)$ , MLP=Single-layer Perceptron!

# Multilayer Perceptron (MLP)

## ➤ Visualization of Decision Boundaries



Linear/Affine Activation



ReLU Activation

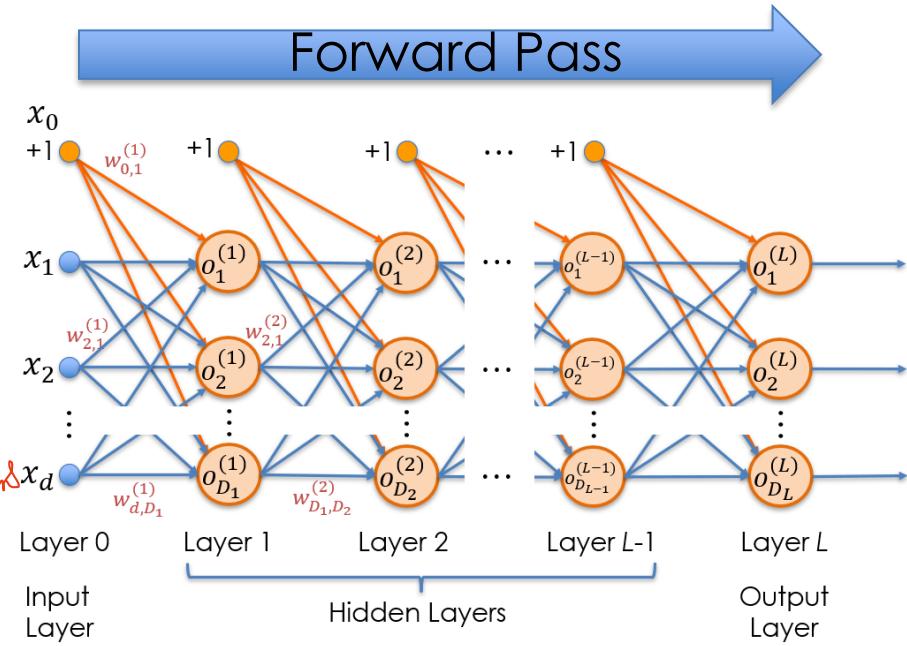
# Multilayer Perceptron (MLP)

## Learning ( $\mathbf{W}$ )

1. Random Initialization of  $\mathbf{W}$
2. Forward Pass (fixed  $\mathbf{W}$ )
  - i. Compute output responses
3. Backward Pass (update  $\mathbf{W}$ )  
i.e., **Backpropagation**
  - i. Propagate the errors backward
  - ii. Update  $\mathbf{W}$

$$\mathbf{W}_{new} = \mathbf{W}_{old} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}_{old})$$

In essence, Backpropagation= Gradient Decent!



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

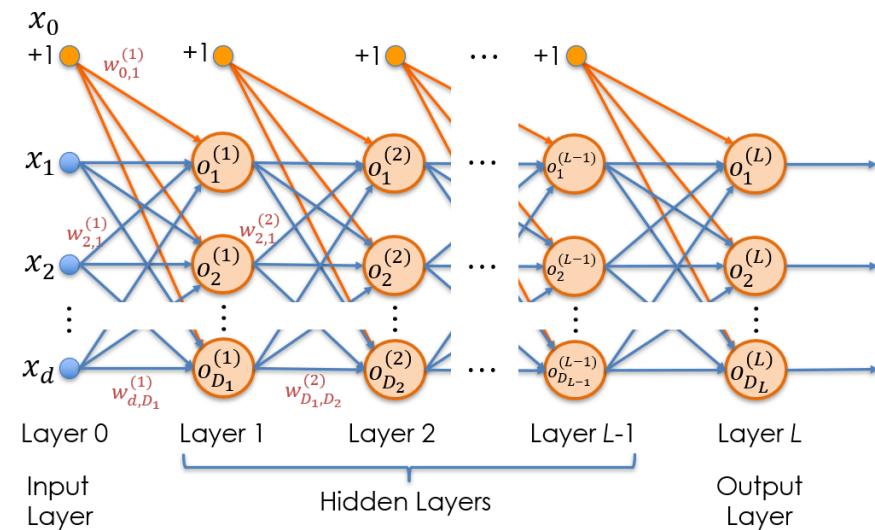
$$\mathbf{W}_{new} = \mathbf{W}_{old} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}_{old})$$

$$\nabla_{\mathbf{W}} J(\mathbf{W}_{old}) = \frac{1}{N} \sum \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}_{old})$$

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

where  $z_j^{(l)} = (\mathbf{a}^{(l)})^T \mathbf{w}_j^{(l)}$

$\mathbf{a}^l$ : input vector to layer  $l = \begin{cases} \mathbf{x}, & \text{if } l = 1 \\ [1, \mathbf{o}^{(l-1)}]^T, & \text{else} \end{cases}$



Backward Pass

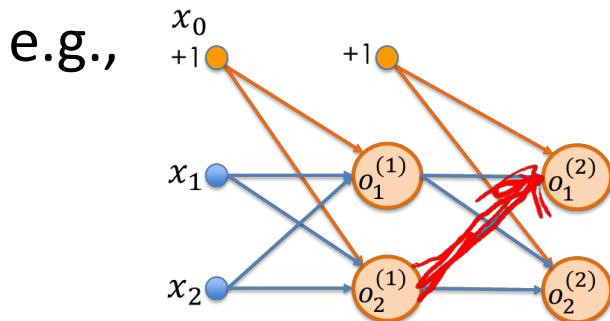
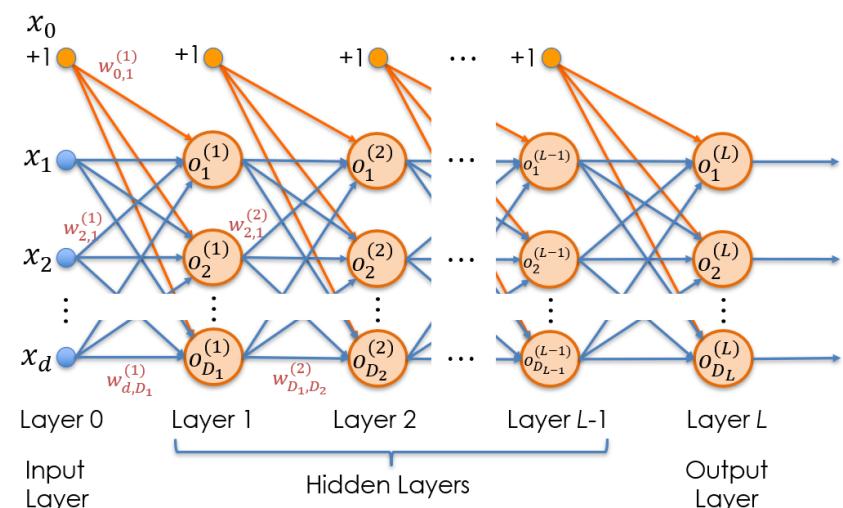
# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

$$\begin{aligned} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}} &= \frac{\partial}{\partial w_{i,j}^{(l)}} \left( \sum_m w_{m,j}^{(l)} a_m^{(l)} \right) \\ &= a_i^{(l)} \end{aligned}$$

$m=j$



$$\frac{\partial z_1^{(2)}}{\partial w_{2,1}^{(2)}} = ? \quad \frac{\partial z_1^{(1)}}{\partial w_{0,1}^{(1)}} = ?$$

# Multilayer Perceptron (MLP)

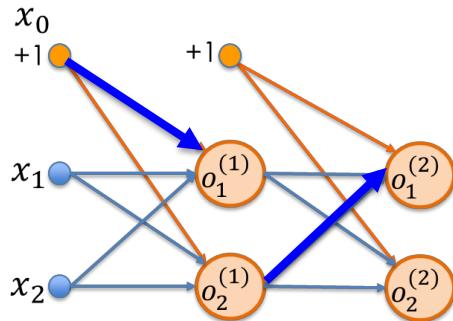
## ➤ Learning (W): Backpropagation

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

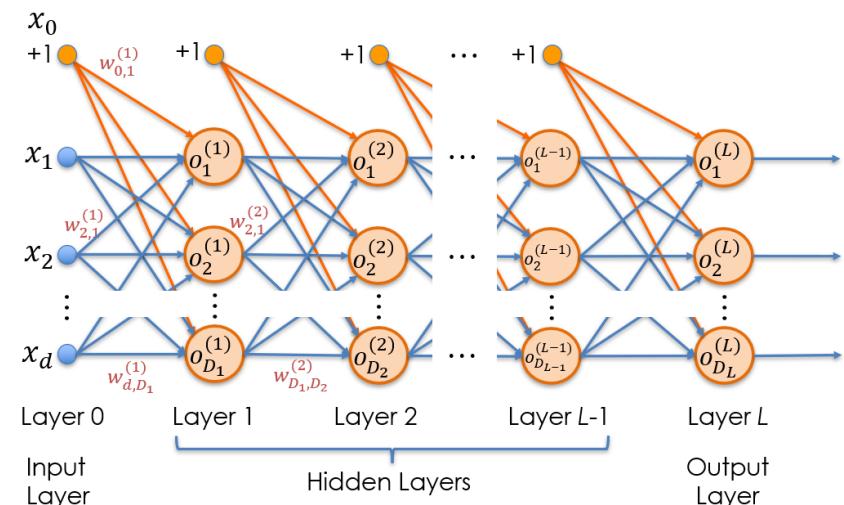
$$\frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}} = \frac{\partial}{\partial w_{i,j}^{(l)}} \left( \sum_m w_{m,j}^{(l)} a_m^{(l)} \right)$$

$$= a_i^{(l)}$$

e.g.,



$$\frac{\partial z_1^{(2)}}{\partial w_{2,1}^{(2)}} = o_2^{(1)} \quad \frac{\partial z_1^{(1)}}{\partial w_{0,1}^{(1)}} = x_0 = 1$$



Backward Pass

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

*Error term*  $\epsilon_j^{(l)}$

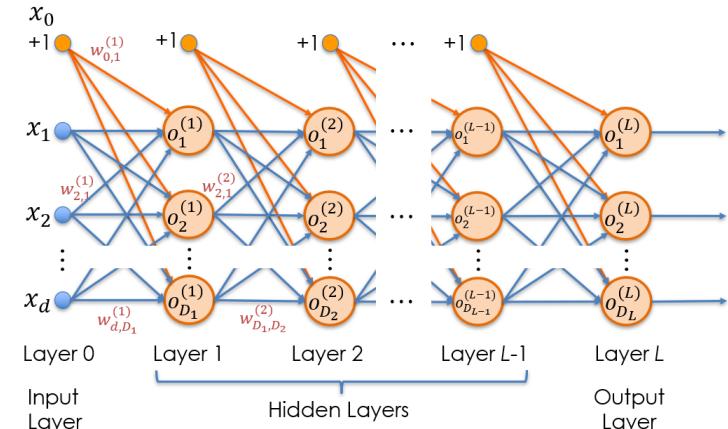
$a_i^{(l)}$

- If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial z_j^{(L)}} = \frac{\partial \mathcal{L}(o_1^{(L)}, \dots, o_{D_L}^{(L)})}{\partial z_j^{(L)}} = \sum_{m=1}^{D_L} \frac{\partial \mathcal{L}(o_m^{(L)})}{\partial z_j^{(L)}}$$

$$= \sum_{m=1}^{D_L} \frac{\partial \mathcal{L}(o_m^{(L)})}{\partial o_m^{(L)}} \frac{\partial o_m^{(L)}}{\partial z_j^{(L)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}}$$



Backward Pass

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

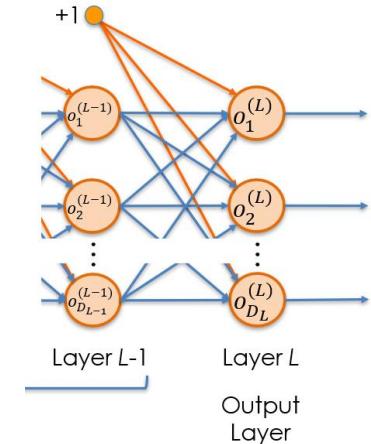
If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}$$

e.g., If  $\sigma^{(L)}(\mathbf{z})_k = \frac{e^{z_k^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}}$ ,  $\mathcal{L}(\mathbf{W}) = - \underbrace{\sum_{k=1}^{D_L} y_k \ln(p_k)}_{\text{Categorical Cross Entropy}}.$  What is  $\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(L)}}$ ?

Softmax

Categorical Cross Entropy



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}$$

e.g., If  $\sigma^{(L)}(\mathbf{z})_k = \frac{e^{z_k^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}}$ ,  $\mathcal{L}(\mathbf{W}) = - \sum_{k=1}^{D_L} y_k \ln(p_k)$ . What is  $\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(L)}}$ ?

$\underbrace{\sigma^{(L)}(\mathbf{z})_k = \frac{e^{z_k^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}}$  Softmax       $\underbrace{- \sum_{k=1}^{D_L} y_k \ln(p_k)}$  Categorical Cross Entropy

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(L)}} = \epsilon_j^{(L)} a_i^{(L)} = a_i^{(L)} (o_j^{(L)} - y_j)$$

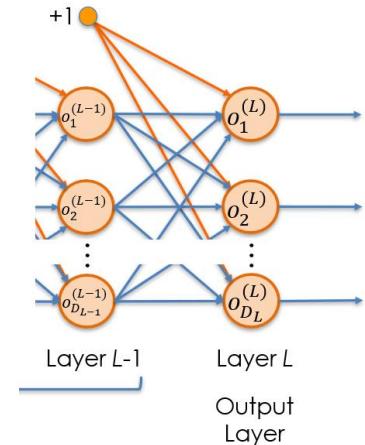
a is activation  
output, e is error  
term

$$\begin{aligned} \epsilon_j^{(L)} &= \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m \\ &= \sum_{m=1}^{D_L} -\frac{y_m}{o_m^{(L)}} (\delta_{mj} - o_m^{(L)}) o_j^{(L)} \\ &= o_j^{(L)} - y_j \end{aligned}$$

$$\mathcal{L}'(o_m^{(L)}) = - \sum_{k=1}^{D_L} y_k \frac{d(\ln(o_k^{(L)}))}{d(o_m^{(L)})} = -\frac{y_m}{o_m^{(L)}}$$

$$\sigma^{(L)'}(z_j^{(L)})_m = \frac{d}{d(z_j^{(L)})} \left( \frac{e^{z_m^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}} \right) = (\delta_{mj} - o_m^{(L)}) o_j^{(L)}$$

$$\delta_{mj} = \begin{cases} 1, & \text{if } m = j \\ 0, & \text{if } m \neq j \end{cases}$$



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

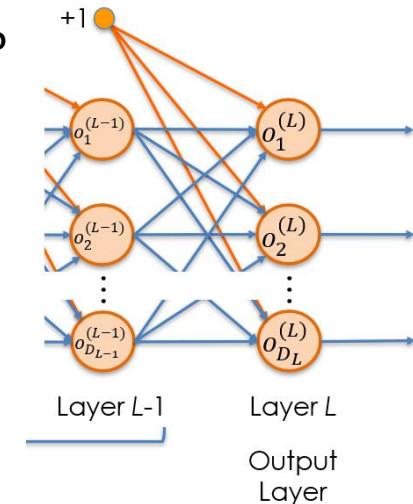
If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}$$

e.g., If  $\sigma^{(L)}(\mathbf{z})_k = \frac{e^{z_k^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}}$ ,  $\mathcal{L}(\mathbf{W}) = - \underbrace{\sum_{k=1}^{D_L} y_k \ln(p_k)}_{\text{Categorical Cross Entropy}}$ . What is  $\frac{d\mathcal{L}}{d\mathbf{W}^{(L)}}$ ?

Softmax

Categorical Cross Entropy



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}$$

e.g., If  $\sigma^{(L)}(\mathbf{z})_k = \frac{e^{z_k^{(L)}}}{\sum_{s=1}^{D_L} e^{z_s^{(L)}}}$ ,  $\mathcal{L}(\mathbf{W}) = - \sum_{k=1}^{D_L} y_k \ln(p_k)$ . What is  $\frac{d\mathcal{L}}{d\mathbf{W}^{(L)}}$ ?



Categorical Cross Entropy  
same size as no. of output neurons.

$$\frac{d\mathcal{L}}{d\mathbf{W}^{(L)}} = \mathbf{a}^{(L)} \epsilon^{(L)} = \mathbf{a}^{(L)} (\mathbf{o}^{(L)} - \mathbf{y}) = \mathbf{a}^{(L)} (\mathbf{f}_\mathbf{W}(\mathbf{x}) - \mathbf{y})$$

Size:  
 $(D_{L-1} + 1) \times D_L$

need to check if size of output is correct

Size:  
 $(D_{L-1} + 1) \times 1$

Size:  
 $1 \times D_L$

$\mathbf{w}_{d+1}^{(L)} \times D$

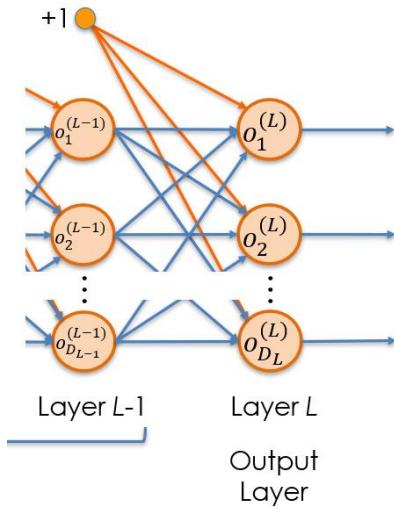
give direction of w.

$$\frac{dJ}{d\mathbf{W}^{(L)}} = \frac{1}{N} \sum \frac{d\mathcal{L}}{d\mathbf{W}^{(L)}} = \frac{1}{N} (\mathbf{A}^{(L)})^T (\mathbf{F}_\mathbf{W}(\mathbf{X}) - \mathbf{Y})$$

mean of loss function

$\mathbf{Y}$  has a size of  $N \times D_L$

$\mathbf{A}^{(L)}$  has a size of  $N \times (D_{L-1} + 1)$



$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix}$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

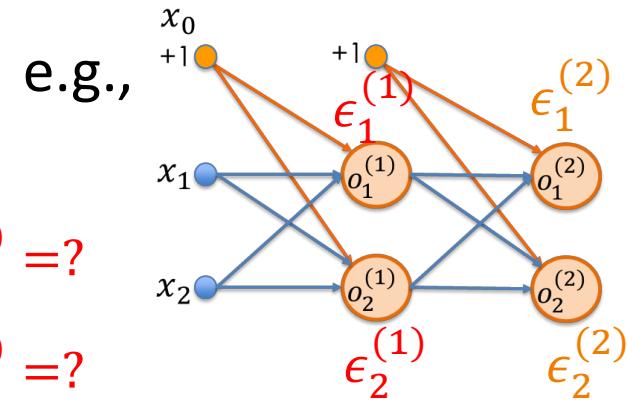
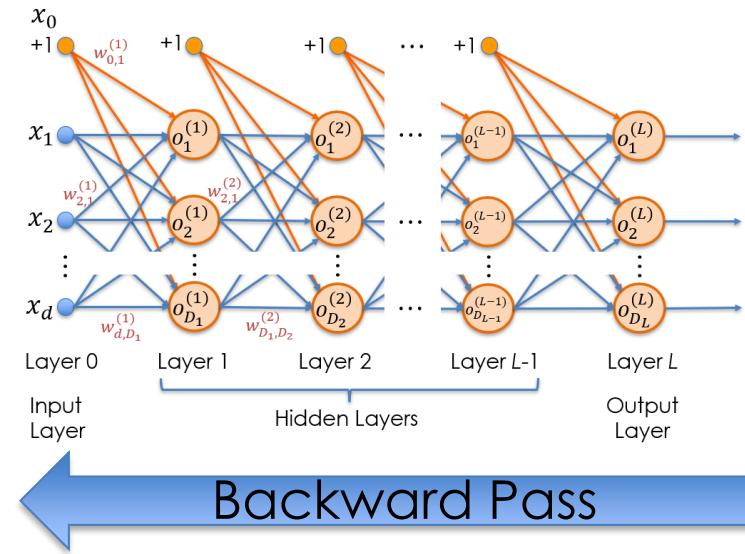
$\epsilon_j^{(l)}$

- If hidden layer ( $l < L$ ):

$$\frac{\partial \mathcal{L}}{\partial z_j^{(l)}} = \frac{\partial \mathcal{L}(z_1^{(l+1)}, \dots, z_{D_{l+1}}^{(l+1)})}{\partial z_j^{(l)}} = \sum_{m=1}^{D_{l+1}} \frac{\partial \mathcal{L}(z_m^{(l+1)})}{\partial z_j^{(l)}}$$

$$= \sum_{m=1}^{D_{l+1}} \frac{\partial \mathcal{L}(z_m^{(l+1)})}{\partial z_m^{(l+1)}} \frac{\partial z_m^{(l+1)}}{\partial z_j^{(l)}} = \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} \frac{\partial z_m^{(l+1)}}{\partial o_j^{(l)}} \frac{\partial o_j^{(l)}}{\partial z_j^{(l)}}$$

$$= \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)}$$



$$\epsilon_1^{(1)} = ?$$

$$\epsilon_2^{(1)} = ?$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{i,j}^{(l)}}$$

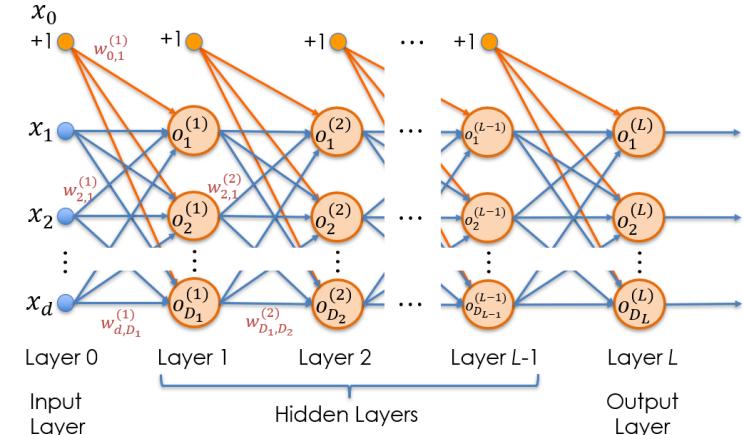
$\epsilon_j^{(l)}$

- If hidden layer ( $l < L$ ):

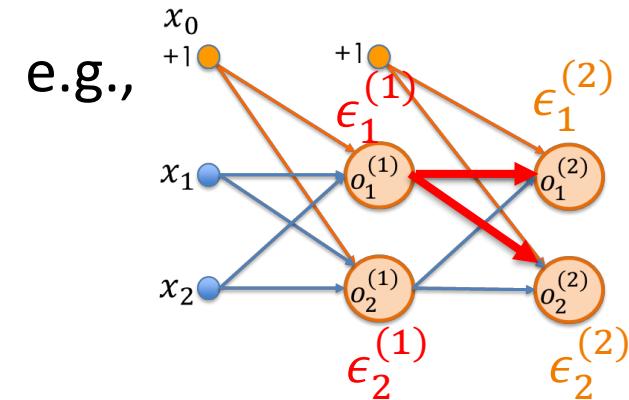
$$\frac{\partial \mathcal{L}}{\partial z_j^{(l)}} = \frac{\partial \mathcal{L}(z_1^{(l+1)}, \dots, z_{D_{l+1}}^{(l+1)})}{\partial z_j^{(l)}} = \sum_{m=1}^{D_{l+1}} \frac{\partial \mathcal{L}(z_m^{(l+1)})}{\partial z_j^{(l)}}$$

$$= \sum_{m=1}^{D_{l+1}} \frac{\partial \mathcal{L}(z_m^{(l+1)})}{\partial z_m^{(l+1)}} \frac{\partial z_m^{(l+1)}}{\partial z_j^{(l)}} = \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} \frac{\partial z_m^{(l+1)}}{\partial o_j^{(l)}} \frac{\partial o_j^{(l)}}{\partial z_j^{(l)}}$$

$$= \sigma^{(l)\prime}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)}$$



Backward Pass



$$\epsilon_1^{(1)} = \sigma^{(1)\prime} \times (w_{1,1}^{(2)} \epsilon_1^{(2)} + w_{1,2}^{(2)} \epsilon_2^{(2)})$$

# Multilayer Perceptron (MLP)

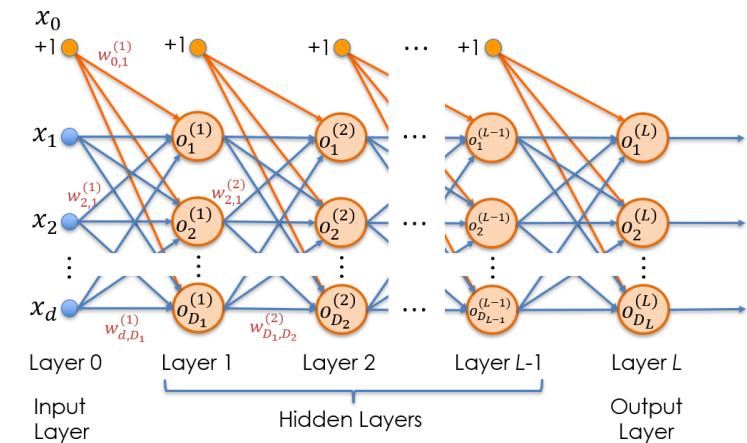
## ➤ Learning (W): Backpropagation

If hidden layer ( $l < L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)} a_i^{(l)}$$

e.g., If  $\sigma^{(l)}(z) = \max(0, z)$ , what is  $\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}}$ ?

ReLU



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

If hidden layer ( $l < L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)} a_i^{(l)}$$

e.g., If  $\sigma^{(l)}(z) = \max(0, z)$ , what is  $\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}}$ ?

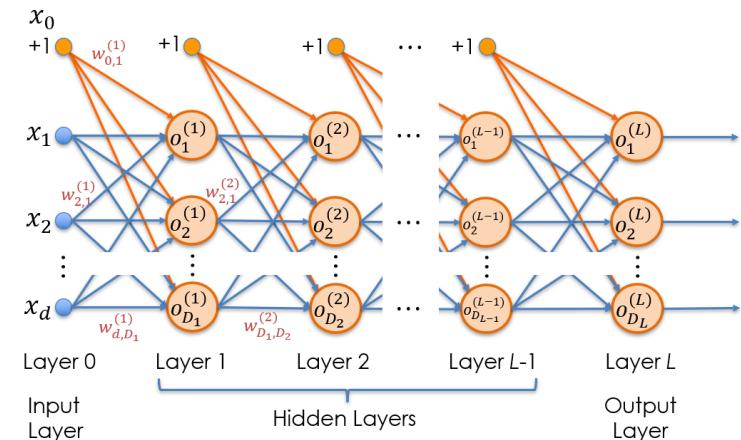
ReLU

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \epsilon_j^{(l)} a_i^{(l)} = \delta_{z_j^{(l)}} \epsilon^{(l+1)} (\mathbf{w}_{j,:}^{(l+1)})^T a_i^{(l)}$$

$$\epsilon_j^{(l)} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)}$$

$$= \delta_{z_j^{(l)}} \epsilon^{(l+1)} (\mathbf{w}_{j,:}^{(l+1)})^T$$

Size:  
 $1 \times D_{l+1}$       Size:  
 $D_{l+1} \times 1$



$$\sigma^{(l)'}(z_j^{(l)}) = \frac{d(\max(0, z_j^{(l)}))}{d(z_j^{(l)})} = \begin{cases} 0, & \text{if } z_j^{(l)} < 0 \\ 1, & \text{if } z_j^{(l)} > 0 \end{cases}$$

$$\delta_{z_j^{(l)}} = \begin{cases} 0, & \text{if } z_j^{(l)} < 0 \\ 1, & \text{if } z_j^{(l)} > 0 \end{cases}$$

# Multilayer Perceptron (MLP)

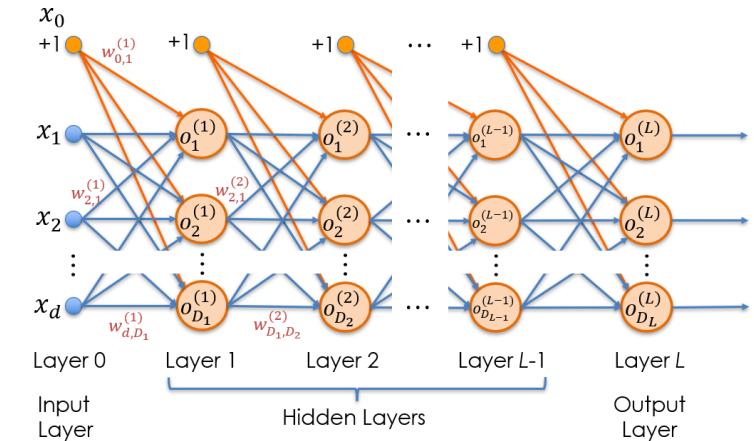
## ➤ Learning (W): Backpropagation

If hidden layer ( $l < L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)} a_i^{(l)}$$

e.g., If  $\sigma^{(l)}(z) = \max(0, z)$ , what is  $\frac{d\mathcal{L}}{d\mathbf{w}^{(l)}}$ ?

ReLU



# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

If hidden layer ( $l < L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)} a_i^{(l)}$$

e.g., If  $\sigma^{(l)}(z) = \max(0, z)$ , what is  $\frac{d\mathcal{L}}{d\mathbf{W}^{(l)}}$ ?

$\underbrace{\quad}_{\text{ReLU}}$

Size:  
 $(D_{l-1} + 1) \times D_l$

$\frac{d\mathcal{L}}{d\mathbf{W}^{(l)}} = \mathbf{a}^{(l)} \boldsymbol{\epsilon}^{(l)} = \mathbf{a}^{(l)} (\delta_{Z^{(l)}} * (\epsilon^{(l+1)} (\mathbf{W}_{1:D_l,:}^{(l+1)})^T))$

Size:  
 $(D_{l-1} + 1) \times 1$

Size:  
 $1 \times D_l$

Size:  
 $1 \times D_{l+1}$

$$\mathbf{W}^{(l+1)} = \begin{bmatrix} w_{0,1}^{(l+1)} & \dots & w_{0,D_{l+1}}^{(l+1)} \\ \vdots & \ddots & \vdots \\ w_{D_l,1}^{(l+1)} & \dots & w_{D_l,D_{l+1}}^{(l+1)} \end{bmatrix}$$

$\left\{ \begin{array}{l} \mathbf{A}^{(l)} \text{ has a size of } N \times (D_{l-1} + 1) \\ \boldsymbol{\delta}_{Z^{(l)}} \text{ has a size of } N \times D_l \\ \mathbf{E}^{(l+1)} \text{ has a size of } N \times D_{l+1} \end{array} \right.$

$$\frac{dJ}{d\mathbf{W}^{(l)}} = \frac{1}{N} \sum \frac{d\mathcal{L}}{d\mathbf{W}^{(l)}} = \frac{1}{N} (\mathbf{A}^{(l)})^T (\boldsymbol{\delta}_{Z^{(l)}} * (\mathbf{E}^{(l+1)} (\mathbf{W}_{1:D_l,:}^{(l+1)})^T))$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

If output layer ( $l = L$ ):

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)'}(z_j^{(L)})_m a_i^{(L)}$$

*Summary ↴*  
If softmax activation, categorical cross entropy loss:

$$\frac{d\mathcal{L}}{d\mathbf{W}^{(L)}} = \mathbf{a}^{(L)} \underbrace{(\mathbf{f}_W(\mathbf{x}) - \mathbf{y})}_{\epsilon^{(L)}}, \quad \frac{dJ}{d\mathbf{W}^{(L)}} = \frac{1}{N} (\mathbf{A}^{(L)})^T \underbrace{(\mathbf{F}_W(\mathbf{X}) - \mathbf{Y})}_{\mathbf{E}^{(L)}}$$

If hidden layer ( $l < L$ ):

Why? softmax is used for multi-class classification, takes raw scores and outputs into probabilities for class prediction, and softmax forces every hidden vector to sum to 1 and be positive, which is not what we want. Can use sigmoid here for binary classification.

Hidden layers use non-linear activations like relu to learn representations of the input, must use non-linear activation to avoid collapsing into just linear regression

$$\frac{\partial \mathcal{L}}{\partial w_{i,j}^{(l)}} = \sigma^{(l)'}(z_j^{(l)}) \sum_{m=1}^{D_{l+1}} \epsilon_m^{(l+1)} w_{j,m}^{(l+1)} a_i^{(l)}$$

If ReLU activation:

$$\frac{d\mathcal{L}}{d\mathbf{W}^{(l)}} = \mathbf{a}^{(l)} (\delta_{\mathbf{z}^{(l)}} * (\epsilon^{(l+1)} (\mathbf{W}_{1:D_l,:}^{(l+1)})^T))$$

$$\frac{dJ}{d\mathbf{W}^{(l)}} = \frac{1}{N} (\mathbf{A}^{(l)})^T (\delta_{\mathbf{z}^{(l)}} * (\mathbf{E}^{(l+1)} (\mathbf{W}_{1:D_l,:}^{(l+1)})^T))$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

Example: A MLP of 2 layers has been constructed as:

$$f_W(x) = \sigma^{(2)}([1, \sigma^{(1)}(x^T W^{(1)})]W^{(2)}),$$

where  $\sigma^{(1)}$  is ReLU, and  $\sigma^{(2)}$  is the Softmax.

Training data:  $X = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$ ,  $Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Loss function: Categorical Cross Entropy

Weight initialization:  $W^{(1)}(0) = W^{(2)}(0) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$

What are  $W^{(1)}(1)$  and  $W^{(2)}(1)$  after 1<sup>st</sup> iteration of weight updates if learning rate  $\eta=0.1$  ?

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

Example: A MLP of 2 layers has been constructed as:

$$f_W(x) = \sigma^{(2)}([1, \sigma^{(1)}(x^T W^{(1)})]W^{(2)}),$$

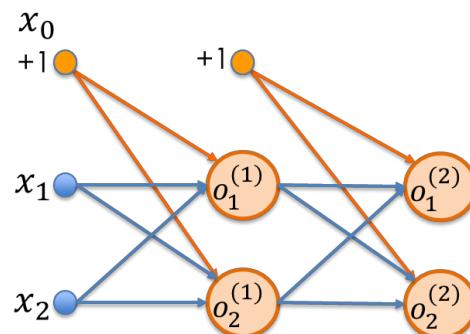
where  $\sigma^{(1)}$  is ReLU, and  $\sigma^{(2)}$  is the Softmax.

Training data:  $X = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$ ,  $Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Loss function: Categorical Cross Entropy

Weight initialization:  $W^{(1)}(0) = W^{(2)}(0) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$

What are  $W^{(1)}(1)$  and  $W^{(2)}(1)$  after 1<sup>st</sup> iteration of weight updates if learning rate  $\eta=0.1$  ?



### Step 1: Forward Pass

$$F_{W(0)}(X) = \sigma^{(2)}([1, \sigma^{(1)}(XW^{(1)}(0))]W^{(2)}(0))$$

z output  
 $\begin{bmatrix} 2 & 0 \\ 1.5 & 0 \end{bmatrix}$

activation output  
 $\begin{bmatrix} -1 & -2 \\ -1 & -1.5 \end{bmatrix}$

pre activation z output  
 $= \begin{bmatrix} 0.7311 & 0.2689 \\ 0.6225 & 0.3775 \end{bmatrix}$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

Example: A MLP of 2 layers has been constructed as:

$$f_W(x) = \sigma^{(2)}([1, \sigma^{(1)}(x^T W^{(1)})] W^{(2)}),$$

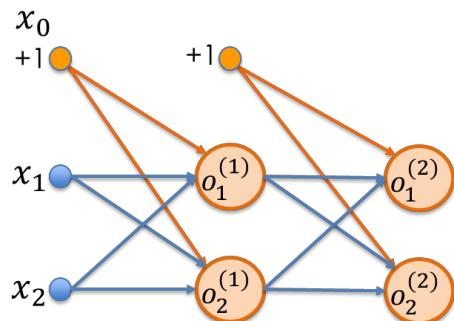
where  $\sigma^{(1)}$  is ReLU, and  $\sigma^{(2)}$  is the Softmax.

Training data:  $X = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$ ,  $Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Loss function: Categorical Cross Entropy

Weight initialization:  $W^{(1)}(0) = W^{(2)}(0) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$

What are  $W^{(1)}(1)$  and  $W^{(2)}(1)$  after 1<sup>st</sup> iteration of weight updates if learning rate  $\eta=0.1$  ?



$W^{(2)}$  update:

$$W^{(2)}(1) = W^{(2)}(0) - \eta \frac{dJ(W^{(2)}(0))}{dW^{(2)}} = \begin{bmatrix} -1.0177 & 0.0177 \\ -0.0198 & -0.9802 \\ 1 & 0 \end{bmatrix}$$

Step 2: Backward Pass

Output layer:

$$\frac{dJ(W^{(2)}(0))}{dW^{(2)}} = \frac{1}{N} (\mathbf{A}^{(2)})^T (\mathbf{E}_{W(0)}(\mathbf{X}) - \mathbf{Y})$$

A is activation output

E is the error term

$$\mathbf{E}^{(2)} = \begin{bmatrix} -0.2689 & 0.2689 \\ 0.6225 & -0.6225 \end{bmatrix}$$

$$\mathbf{A}^{(2)} = \begin{bmatrix} 1 & 2 & 0 \\ 1 & 1.5 & 0 \end{bmatrix}$$

$$\frac{dJ(W^{(2)}(0))}{dW^{(2)}} = \begin{bmatrix} 0.1768 & -0.1768 \\ 0.1979 & -0.1979 \\ 0 & 0 \end{bmatrix}$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

Example: A MLP of 2 layers has been constructed as:

$$f_W(x) = \sigma^{(2)}([1, \sigma^{(1)}(x^T W^{(1)})] W^{(2)}),$$

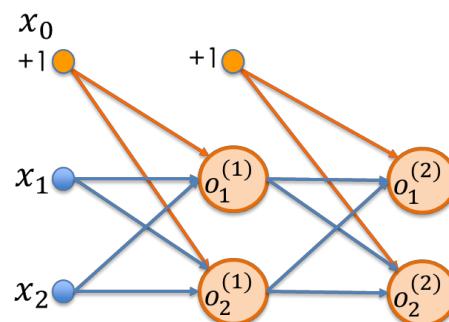
where  $\sigma^{(1)}$  is ReLU, and  $\sigma^{(2)}$  is the Softmax.

Training data:  $X = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$ ,  $Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Loss function: Categorical Cross Entropy

Weight initialization:  $W^{(1)}(0) = W^{(2)}(0) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$

What are  $W^{(1)}(1)$  and  $W^{(2)}(1)$  after 1<sup>st</sup> iteration of weight updates if learning rate  $\eta=0.1$  ?



$W^{(1)}$  update:

$$W^{(1)}(1) = W^{(1)}(0) - \eta \frac{dJ(W^{(1)}(0))}{dW^{(1)}} = \begin{bmatrix} -1.0177 & 0 \\ -0.0488 & -1 \\ 0.9625 & 0 \end{bmatrix}$$

Step 2: Backward Pass

Hidden layer:  $\frac{dJ(W^{(1)}(0))}{dW^{(1)}} = \frac{1}{N} (\mathbf{A}^{(1)})^T (\delta_{Z^{(1)}} * (\mathbf{E}^{(2)} (\mathbf{W}_{1:2,:}^{(2)}(0))^T))$

$$\mathbf{A}^{(1)} = \mathbf{X} = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$$

$$\begin{aligned} \mathbf{E}^{(1)} &= \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} * (\mathbf{E}^{(2)} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}) \\ &= \begin{bmatrix} -0.2689 & 0 \\ 0.6225 & 0 \end{bmatrix} \end{aligned}$$

$$\frac{dJ(W^{(1)}(0))}{dW^{(1)}} = \begin{bmatrix} 0.1768 & 0 \\ 0.4880 & 0 \\ 0.3747 & 0 \end{bmatrix}$$

# Multilayer Perceptron (MLP)

## ➤ Learning (W): Backpropagation

Example: A MLP of 2 layers has been constructed as:

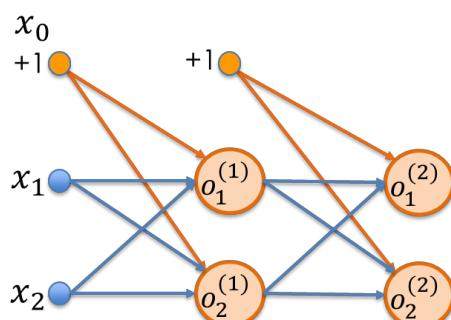
$$f_W(x) = \sigma^{(2)}([1, \sigma^{(1)}(x^T W^{(1)})] W^{(2)}),$$

where  $\sigma^{(1)}$  is ReLU, and  $\sigma^{(2)}$  is the Softmax.

Training data:  $X = \begin{bmatrix} 1 & 1 & 3.0 \\ 1 & 2 & 2.5 \end{bmatrix}$ ,  $Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Loss function: Categorical Cross Entropy

Weight initialization:  $W^{(1)}(0) = W^{(2)}(0) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \end{bmatrix}$



3<sup>rd</sup> Iteration...

4<sup>th</sup> Iteration...

⋮

2<sup>nd</sup> Iteration:

Step 1: Forward Pass

$$F_{W(1)}(X) = \sigma^{(2)}([1, \sigma^{(1)}(X W^{(1)}(1))] W^{(2)}(1))$$

Step 2: Backward Pass

✓ Output layer:

$$\frac{dJ(W^{(2)}(1))}{dW^{(2)}} = \frac{1}{N} (\mathbf{A}^{(2)})^T (F_{W(1)}(X) - Y)$$

$$W^{(2)}(2) = W^{(2)}(1) - \eta \frac{dJ(W^{(2)}(1))}{dW^{(2)}}$$

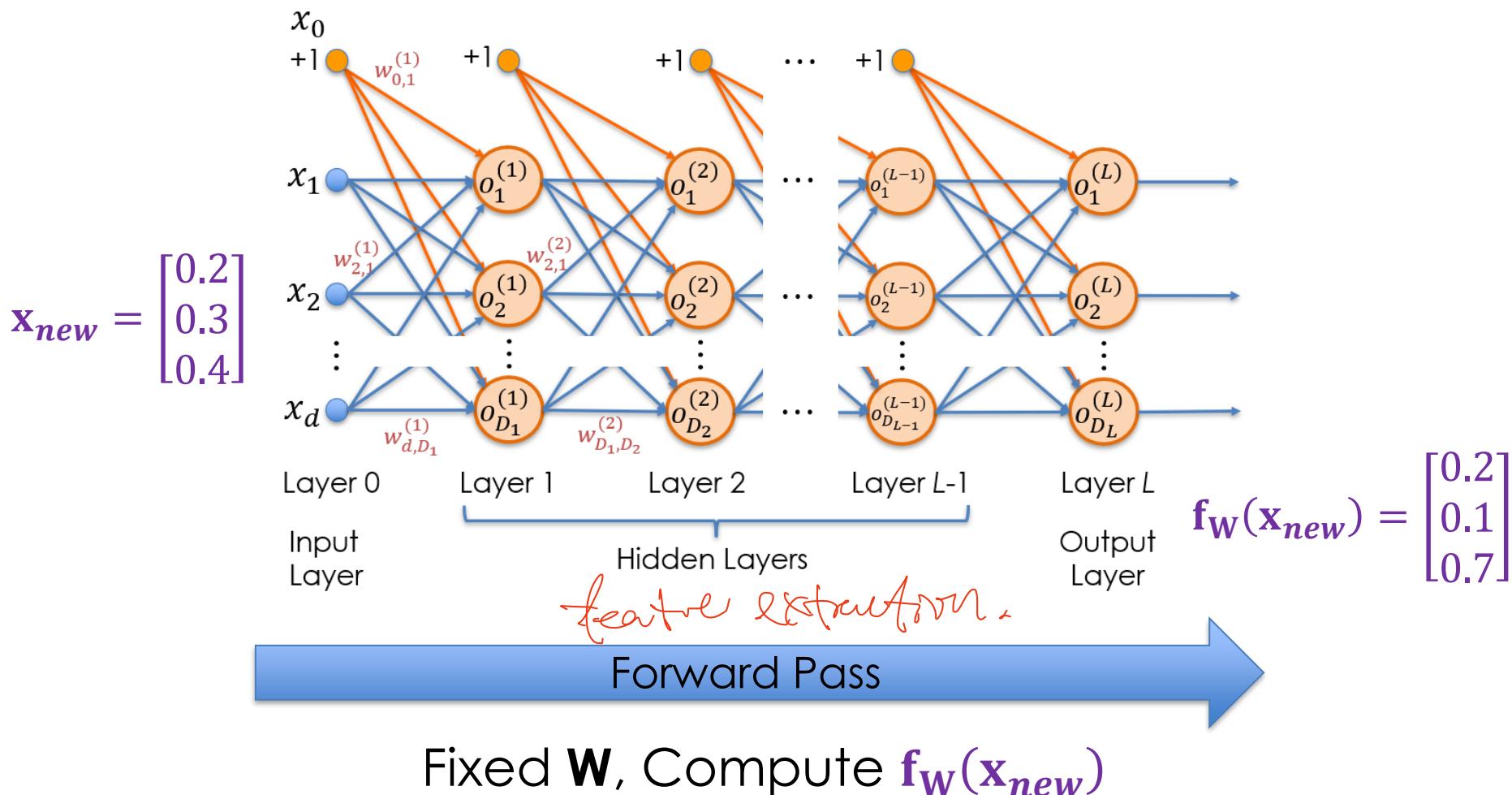
✓ Hidden layer:

$$\frac{dJ(W^{(1)}(1))}{dW^{(1)}} = \frac{1}{N} (\mathbf{A}^{(1)})^T (\delta_{Z(1)} * (E^{(2)}(W_{1:2,:}^{(2)}(1))^T))$$

$$W^{(1)}(2) = W^{(1)}(1) - \eta \frac{dJ(W^{(1)}(1))}{dW^{(1)}}$$

# Multilayer Perceptron (MLP)

## ➤ Prediction (Forward Pass)



# MLP for Image Classification

# MLP for Image Classification

➤ Image is stored as a matrix in a computer



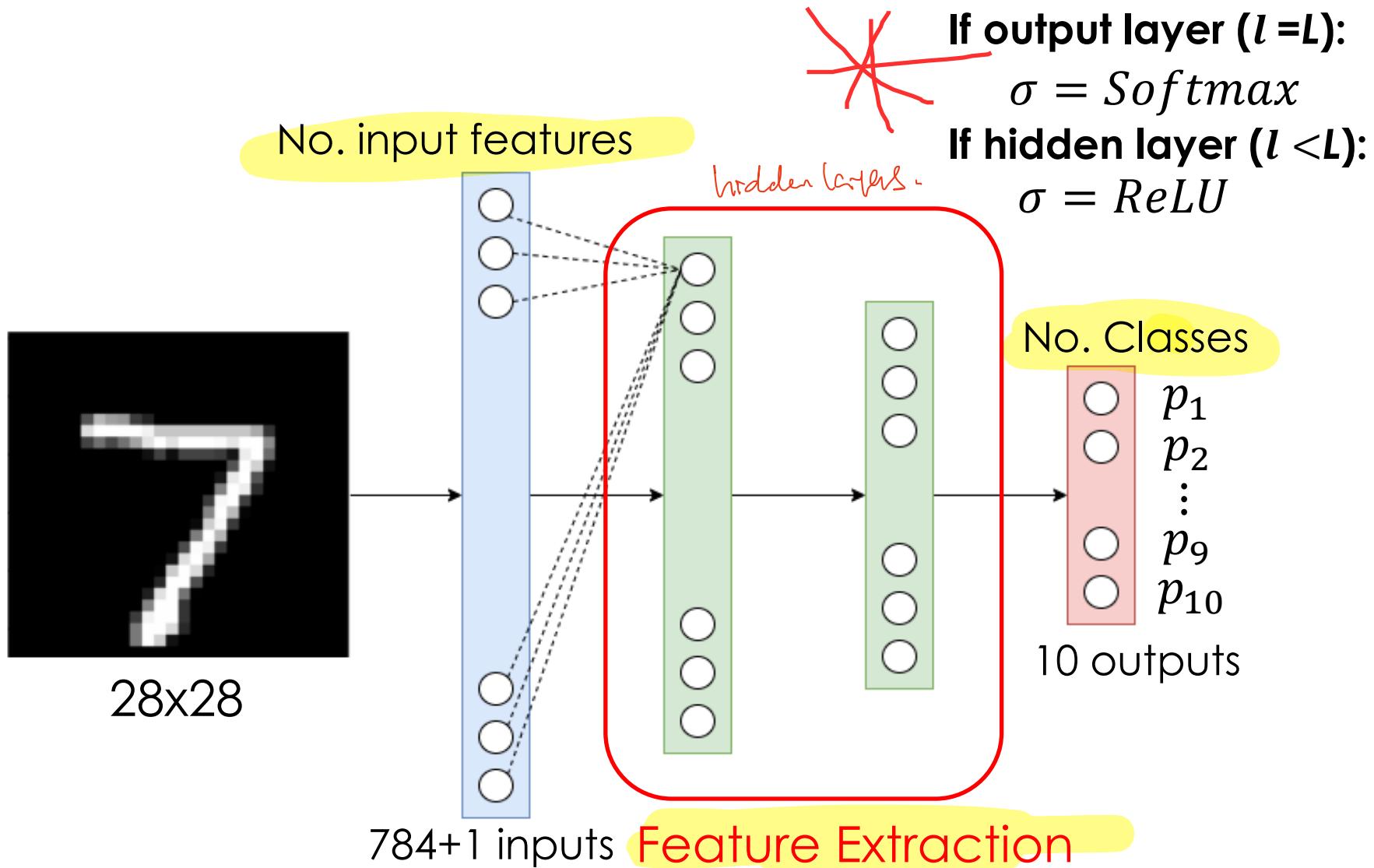
0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0	
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0	
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5		
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1	
0	0	5	5	0	0	0	0	14	1	0	6	6	6	0	0	



Pixel: smallest unit of a digital image

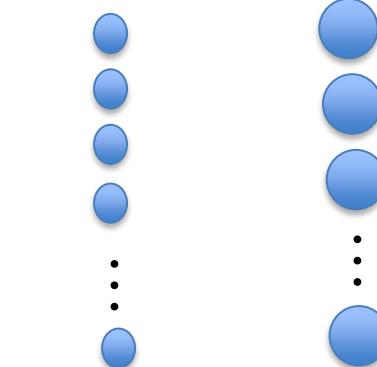
0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0	
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0	
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5		
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1	
0	0	5	5	0	0	0	0	14	1	0	6	6	6	0	0	

# MLP for Image Classification



# MLP for Image Classification

## ➤ Limitations



How many  
input features?



Flattening

Loss of Spatial Structure

How many parameters to  
learn for the first layer?

10K neurons in the  
first layer

$4000 \times 10000 \approx 400$  millions

Parameter explosion

»Convolutional Neural Network!

# Summary

- ✓ Perceptron
- ✓ Multilayer Perceptron (MLP)
- Learning ( $\mathbf{W}$ )**

1. Random Initialization of  $\mathbf{W}$
2. Forward Pass (fixed  $\mathbf{W}$ )
3. Backward Pass (update  $\mathbf{W}$ )

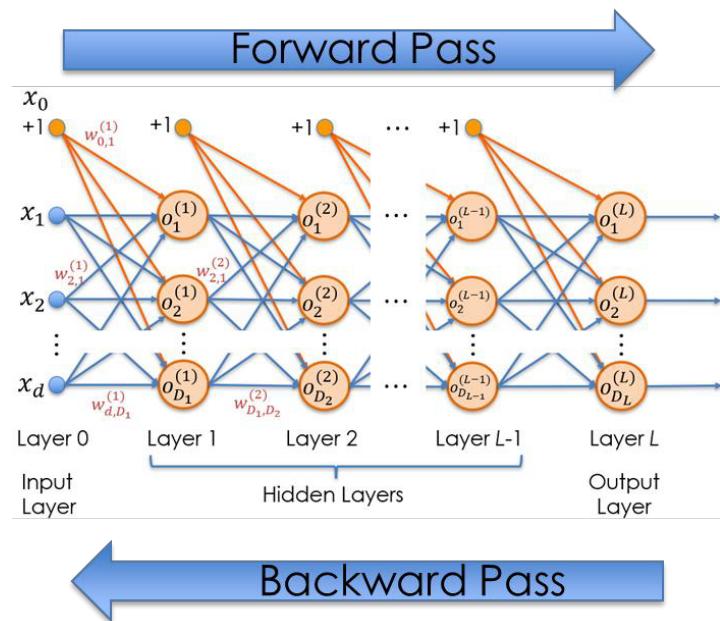
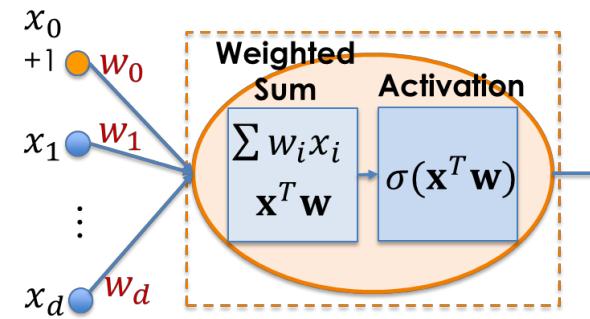
$$\mathbf{W}_{new} = \mathbf{W}_{old} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}_{old})$$

## **-Prediction (Forward Pass)**

$$f_{\mathbf{W}}(\mathbf{x}_{new})$$

- ✓ MLP for Image Classification

Limitations



# THE END

**NATIONAL UNIVERSITY OF SINGAPORE  
DEPARTMENT OF ELECTRICAL & COMPUTER ENGINEERING**

**ACADEMIC YEAR 2025-2026  
SEMESTER 1**

**EE2213: Introduction to AI**

**TUTORIAL 6: L7-L9**

---

**Question 1 (Lecture 7)**

What is the relationship between Machine Learning (ML) and Artificial Intelligence (AI)?

**Answer:** ML is a subset of AI

**Question 2 (Lecture 7)**

Suppose you are working on a machine learning algorithm to navigate a robot through a room from start to goal while avoiding obstacles within some time, according to previous runs through the room with sensor feedback (e.g., successful and failed attempts). What are the Task, Performance, and Experience involved according to the definition of machine learning.

**Answer:**

Task (T): Navigate the robot through a room from start to goal while avoiding obstacles

Performance (P): Success rate of reaching the goal without collision, time taken to reach the goal

Experience (E): Previous runs through the room with sensor feedback.

**Question 3 (Lecture 7)**

Suppose you are working on a machine learning algorithm that groups customers into market segments based on purchasing behavior.

- (i) What type of task is it?
  - (a) Regression
  - (b) Classification
  - (c) Clustering
  - (d) None of these

**Answer:** (c)

- (ii) If the data you have collected involve millions of attributes, what would you do?

**Answer:** Feature Selection or Feature Extraction

**Question 4 (Lecture 7)**

You are given a set of data for supervised learning. A sample block of data looks like this:

0.0012	20123	1.545	12.47	1
0.0034	29087	1.908	20.66	-1
0.0029	19456	2.109	27.09	1
0.0013	32098	1.009	18.76	1
0.0088	25908	1.398	30.08	-1

Each row corresponds to a sample data measurement with 4 input features and 1 response.

- (i) What kind of undesired effect can you anticipate if this set of raw data is used for

learning?

**Answer:** Those features with very large values may overshadow those with very small values.

(ii) How can the data be preprocessed to handle this issue?

**Answer:** We can either use min-max or z-score normalization to resolve the problem.

### Question 5 (Lecture 7)

Please find the synthetic data provided in the “synthetic-data.xlsx”. You may need to install openpyxl library to read .xlsx files. Read the dataset by “from pandas import read\_excel”. Use Python to perform the following tasks.

- (i) Print the summary statistics of this dataset. (You may use “.describe()” or “.info()”)
- (ii) Check and locate if any missing values or duplicates. (You may use “.isnull()” and “.duplicated()”)

**Answer:** See TUT6-Q5.ipynb

### Question 6 (Lecture 8)

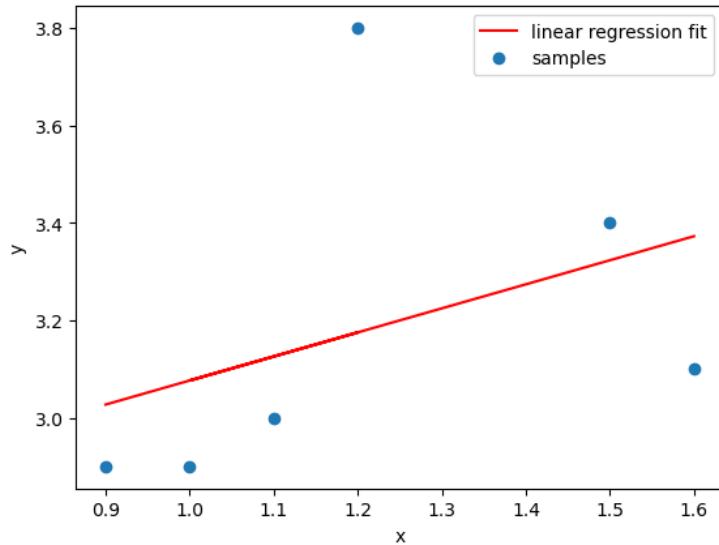
Given the following dataset for training:

Sample No.	x	y
1	1.6	3.1
2	1.5	3.4
3	1.0	2.9
4	1.1	3.0
5	1.2	3.8
6	0.9	2.9

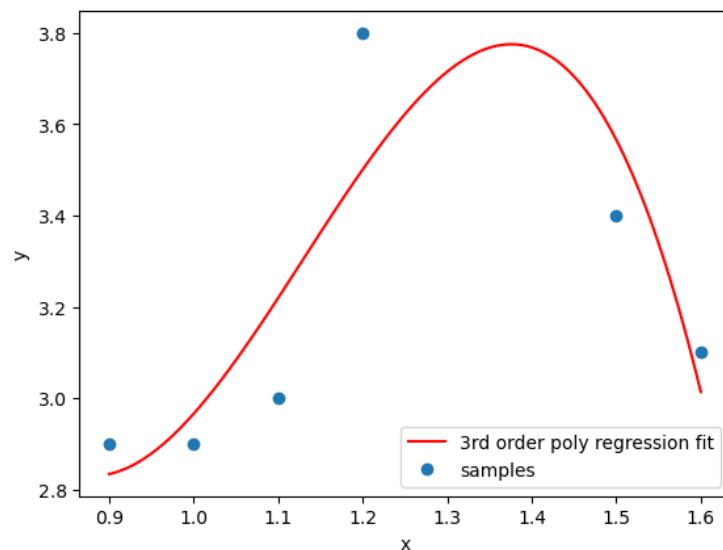
- (i) Perform a linear regression, sketch the result of line fitting, calculate the mean squared error for the training set, and predict a test point {x=0.5}. **0.0890**
- (ii) Perform a 3<sup>rd</sup>-order polynomial regression, sketch the result of line fitting, calculate the mean squared error for the training set, and predict the same test point {x=0.5}. **0.03024**
- (iii) Does it have a unique 6<sup>th</sup>-order polynomial regression solution? If so, proceed to solve it. Otherwise, apply ridge regression with different regularization parameter  $\lambda$ : [0.0001, 0.01, 1]. **X**

**Answer:** See TUT6\_Q6.ipynb

(i)  $w = \begin{bmatrix} 2.5828 \\ 0.4936 \end{bmatrix}, MSE = 0.08896, y_t = 2.8296$



$$(ii) w = \begin{bmatrix} 22.6841 \\ -57.2428 \\ 53.2438 \\ -15.7193 \end{bmatrix}, MSE = 0.03024, y_t = 5.4088$$



(iii) No, because  $P^T P$  is not invertible.

**When  $\lambda = 0.0001$ :**  $MSE = 0.0330, y_t = 2.8114$

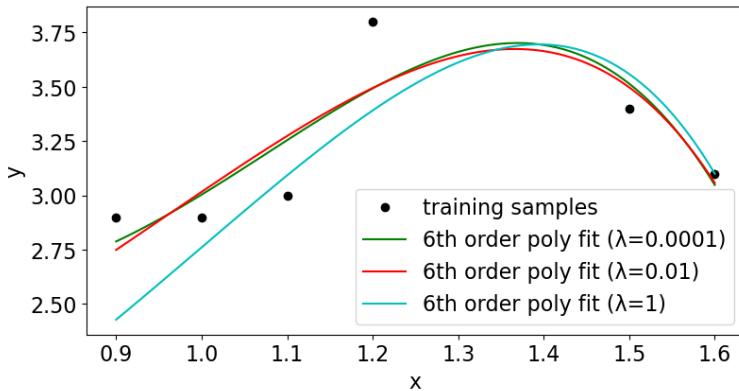
$$w = [4.6439, -4.0828, -1.3906, 3.7044, 3.2392, -3.8597, 0.7490]^T,$$

**When  $\lambda = 0.01$ :**  $MSE = 0.0360, y_t = 1.8481$

$$w = [1.2471, 0.8001, 0.5515, 0.4051, 0.2553, 0.0220, -0.2649]^T,$$

**When  $\lambda = 1$ :**  $MSE = 0.0739, y_t = 1.3768,$

$$w = [0.7708, 0.7170, 0.6438, 0.5342, 0.3618, 0.0864, -0.3535]^T$$



### Question 7 (Lecture 8)

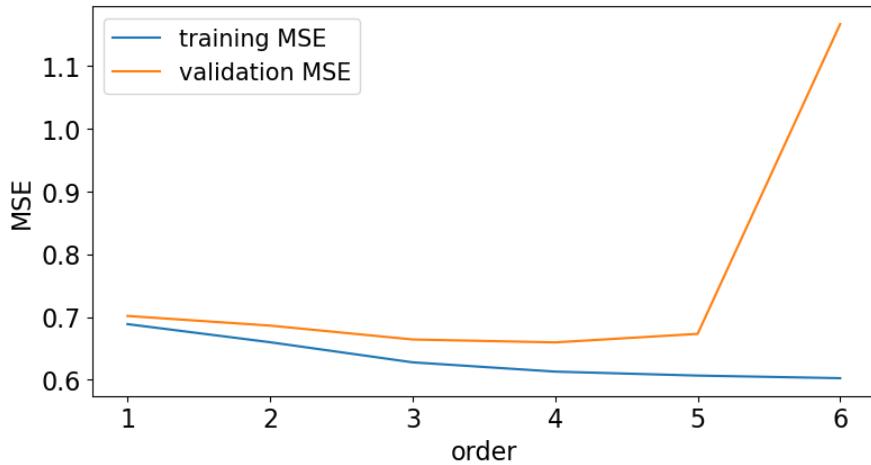
The California Housing Dataset is used for predicting the median house value for a given district in California, based on features from the 1990 U.S. Census. Get the dataset by “from sklearn.datasets import fetch\_california\_housing” and load the dataset as DataFrame by “fetch\_california\_housing(as\_frame=True)”. Use Python to perform the following tasks.

- (i) Split the database into three sets: 60% of samples for training, 20% of samples for validation, and 20% of samples for testing. Hint: you might want to utilize “from sklearn.model\_selection import train\_test\_split” for the splitting.
- (ii) Use Pearson’s correlation as a feature selection metric to select the top 2 features that have the highest absolute correlations with the target output. Use the selected top 2 features to construct the training, validation and test datasets. Hint: you might want to utilize “.corr()” for calculating Pearson’s correlation coefficients.
- (iii) Apply the polynomial model from orders 1 to 6 to the dataset generated from part (ii) without L2 regularization unless  $\mathbf{P}^T \mathbf{P}$  is not invertible (then, use  $\lambda = 0.00001$ ). Plot Mean Squared Errors (MSE) over polynomial orders for both training and the validation sets. Which polynomial order provides the best MSE in the training and validation sets? Which polynomial order should be selected? Use the selected polynomial order to evaluate the MSE for the test set. Hint: you might want to utilize “from sklearn.metrics import mean\_squared\_error” for calculating MSE.
- (iv) Use regularization parameter  $\lambda = 1$  for all polynomial orders and repeat the same analyses. Compare the plots of (iii) and (iv). What do you see?

**Answer:** See TUT6\_Q7.ipynb

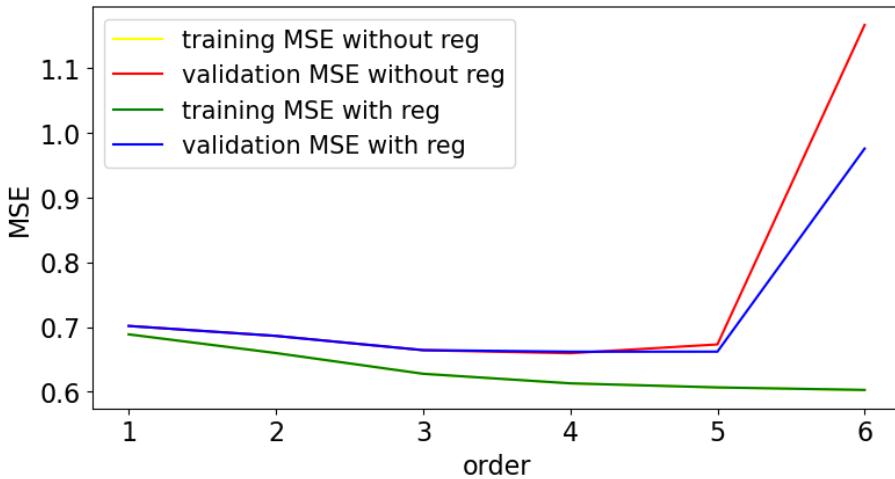
Note: using different “random\_state” in “train\_test\_split” may have different results. The following answers are obtained using “random\_state=42”

- (ii) Top 2 features are “MedInc”, “AveRooms”
- (iii) MSE for training set: [0.6888, 0.6599, 0.6279, 0.6132, 0.6068, **0.6027**], order 6  
MSE for validation set: [0.7017, 0.6864, 0.6643, **0.6597**, 0.6732, 1.1663], order 4  
Order 4 should be selected according to the validation set.  
MSE for test set using order 4 is 4.4618.



- (iv) MSE for training set: [0.6888, 0.6599, 0.6279, 0.6132, 0.6069, **0.6030**], order 6
- MSE for validation set: [0.7017, 0.6864, 0.6645, **0.6620**, 0.6621, 0.9755], order 4
- Order 4 should be selected according to the validation set.

MSE for test set using order 4 is 4.4618.



Regularization helps generalization (prevent overfitting) for higher polynomial order.

### Question 8 (Lecture 9)

The Breast Cancer Wisconsin Diagnostic Dataset is used for classifying tumors as either malignant (labelled 0) or benign (labelled 1), based on digitized measurements of breast mass cell nuclei. Get the dataset by “from sklearn.datasets import load\_breast\_cancer” and load the dataset as DataFrame by “load\_breast\_cancer (as\_frame=True)”. Use Python to perform the following tasks.

- (i) Split the database into three sets: 80% of samples for training, 10% of samples for validation, and 10% of samples for testing. Hint: you might want to utilize “from sklearn.model\_selection import train\_test\_split” for the splitting.
- (ii) Use Pearson’s correlation as a feature selection metrics. First, select features that have

absolute correlations greater than 0.5. Second, among the selected features, remove features that have absolute correlations with each other greater than 0.9. Use the selected features to construct the training, validation and test datasets. Hint: you might want to utilize “.corr()” for calculating Pearson’s correlation coefficients.

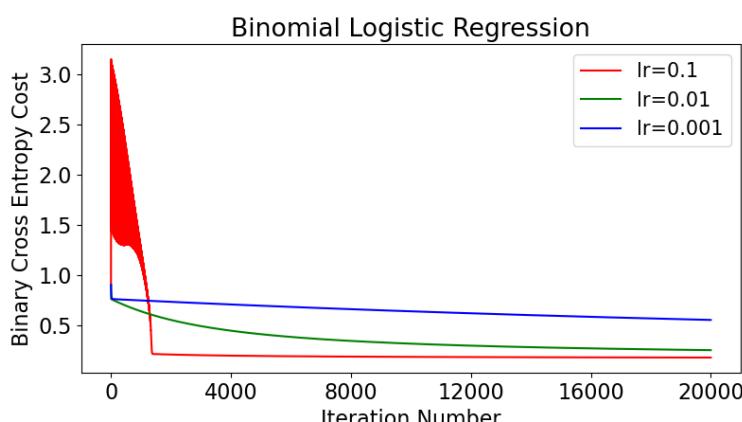
- (iii) Apply linear regression for classification to the dataset generated in part (ii) without regularization unless  $\mathbf{X}^T \mathbf{X}$  is not invertible (then, use  $\lambda = 0.00001$ ). The threshold is set to 0.5 for class prediction. Compute the training, validation and test accuracies. Demonstrate the confusion matrix for the test set. Hint: you might want to utilize “from sklearn.metrics import accuracy\_score, confusion\_matrix” for evaluation metrics.
- (iv) Apply logistic regression to the dataset generated from part (ii) without regularization. The threshold is set to 0.5 for class prediction. The number of iterations for training is set to 20000. Try different learning rates: 0.1, 0.01 and 0.001. Plot the cost over training iterations for each learning rate. Which learning rate provides the best accuracy in the training and the validation sets? Which learning rate should be selected? Use the selected learning rate to compute the accuracy and confusion matrix for the test set.

**Answer:** See TUT6\_Q8.ipynb

Note: using different “random\_state” in “train\_test\_split” may have different results. The following answers are obtained using “random\_state=42”.

- (ii) Final selected features: 'mean radius', 'mean compactness', 'mean concavity', 'radius error', 'worst compactness', 'worst concavity', 'worst concave points'
- (iii) Training accuracy is 0.9341; validation accuracy is 0.9123; test accuracy is 0.9474.
- (iv) Confusion matrix for test set is  $\begin{bmatrix} 14 & 2 \\ 1 & 40 \end{bmatrix}$ .
- Training accuracy: [0.9187, 0.9143, 0.8308]; learning rate 0.1.
- Validation accuracy: [0.9298, 0.9123, 0.7895]; learning rate 0.1.
- Learning rate 0.1 should be selected according to the validation accuracy.
- Test accuracy using learning rate 0.1: 0.9825.

Confusion matrix for test set is  $\begin{bmatrix} 15 & 1 \\ 0 & 41 \end{bmatrix}$



## Question 9 (Lecture 9)

The Iris dataset is used to classify 3 species of iris flower: Setosa, Versicolor, and Virginica. Get the dataset by “from skelarn.datasets import load\_iris”. Use Python to perform the following tasks.

- (i) Split the database into three sets: 80% of samples for training, 10% of samples for validation, and 10% of samples for testing. Hint: you might want to utilize “from sklearn.model\_selection import train\_test\_split” for the splitting.
- (ii) Construct the target output using one-hot encoding. Hint: you might want to utilize “from sklearn.preprocessing import OneHotEncoder”.
- (iii) Apply linear regression for classification to the dataset generated in part (ii) without regularization unless  $\mathbf{X}^T \mathbf{X}$  is not invertible (then, use  $\lambda = 0.00001$ ). Compute the training, validation and test accuracies. Demonstrate the confusion matrix for the test set. Hint: you might want to utilize “from sklearn.metrics import accuracy\_score, confusion\_matrix” for evaluation metrics.
- (iv) Perform logistic regression for classification with different learning rates: 0.1, 0.01 and 0.001. The number of iterations for training is set to 20000. Plot the cost over training iterations for each learning rate. Which learning rate provides the best accuracy in the training and the validation sets? Which learning rate should be selected? Use the selected learning rate to compute the accuracy and confusion matrix for the test set.

**Answer:** See TUT6\_Q9.ipynb

Note: using different “random\_state” in “train\_test\_split” may have different results. The following answers are obtained using “random\_state=42”.

(iii) Training accuracy: 0.8417; Validation accuracy: 0.8; Test accuracy: 0.9333

Confusion matrix for test set is  $\begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 1 \\ 0 & 0 & 4 \end{bmatrix}$ .

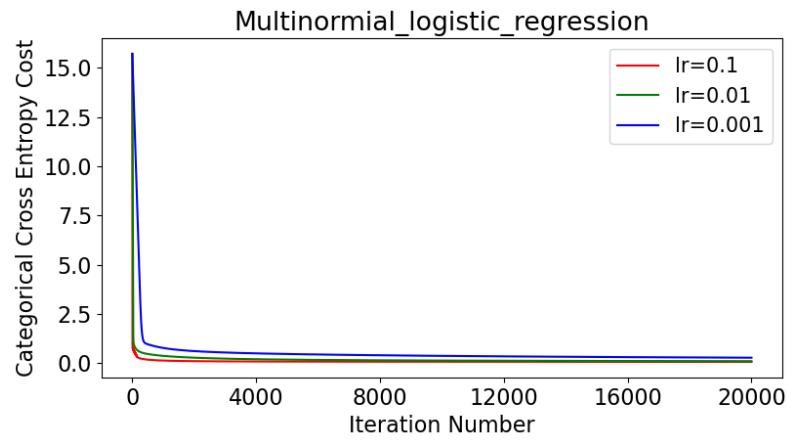
(iv) Training accuracy: [0.975, 0.975, 0.975]

Validation accuracy: [1, 1, 1]

Learning rate 0.1 should be selected as it converges the fastest although both validation and training accuracy are the same for all the learning rates.

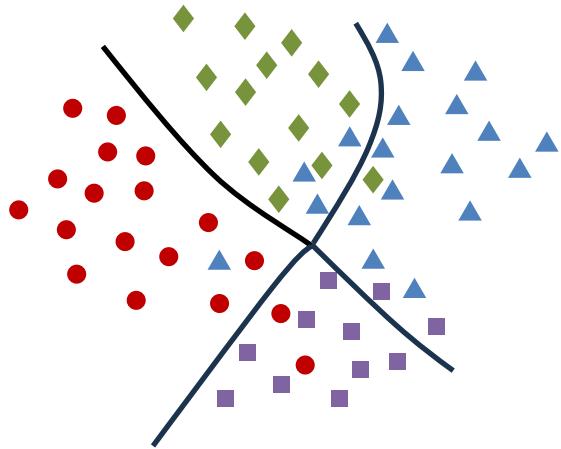
Test accuracy for learning rate 0.1: 0.9333

Confusion matrix for test set is:  $\begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 1 \\ 0 & 0 & 4 \end{bmatrix}$



### Question 10 (Lecture 9)

Tabulate the confusion matrix for the following classification problem. The class-1, class-2, class-3, and class-4 data points are indicated by circles, diamonds, triangles, and squares.



Answer:

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$ (circles)	16	0	0	2
$P_2$ (diamonds)	0	13	1	0
$P_3$ (triangles)	1	3	15	0
$P_4$ (squares)	0	0	2	9

**NATIONAL UNIVERSITY OF SINGAPORE  
DEPARTMENT OF ELECTRICAL & COMPUTER ENGINEERING**

**ACADEMIC YEAR 2025-2026  
SEMESTER 1**

**EE2213: Introduction to AI**

**TUTORIAL 7: L10-L11**

---

**Question 1** (Lecture 10)

What type of cluster does the K-means clustering method form?

- (a) Well-separated clusters
- (b) Center-based clusters
- (c) Contiguous clusters
- (d) Density-based clusters

**Answer:** (b)

**Question 2** (Lecture 10)

The K-means clustering method can only find local optima.

- (a) True
- (b) False

**Answer:** (a). Because the cost function is non-convex.

**Question 3** (Lecture 10)

In order to determine an optimal  $K$  from a possible set of  $K$  values for K-means clustering, Elbow method which selects the  $K$  with the minimum within-cluster variance can be applied.

- (a) True
- (b) False

**Answer:** (b). Should choose  $K$  corresponding to “bend” or “knee”, instead of the minimum within-cluster variance.

**Question 4** (Lecture 10)

Prove that the membership degree sums to 1 for each data point in Fuzzy C-means clustering.

**Answer:**

Need to prove  $\sum_{k=1}^C w_{ik} = 1$

$$\begin{aligned}
\sum_{k=1}^C w_{ik} &= \sum_{k=1}^C \frac{1}{\sum_{j=1}^C \left( \frac{\|x_i - c_k\|}{\|x_i - c_j\|} \right)^{\frac{2}{r-1}}} = \sum_{k=1}^C \frac{1}{\sum_{j=1}^C \frac{(\|x_i - c_k\|)^{\frac{2}{r-1}}}{(\|x_i - c_j\|)^{\frac{2}{r-1}}}} \\
&= \sum_{k=1}^C \frac{1}{(\|x_i - c_k\|)^{\frac{2}{r-1}} \sum_{j=1}^C \frac{1}{(\|x_i - c_j\|)^{\frac{2}{r-1}}}} \\
&= \sum_{k=1}^C \frac{1}{(\|x_i - c_k\|)^{\frac{2}{r-1}} \sum_{j=1}^C \frac{1}{(\|x_i - c_j\|)^{\frac{2}{r-1}}}} \\
&= \sum_{k=1}^C \frac{1}{(\|x_i - c_k\|)^{\frac{2}{r-1}} \sum_{k=1}^C \frac{1}{(\|x_i - c_k\|)^{\frac{2}{r-1}}}} = 1
\end{aligned}$$

### Question 5 (Lecture 10)

Consider the following data points:  $\mathbf{x}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ ,  $\mathbf{x}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ ,  $\mathbf{x}_3 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ ,  $\mathbf{x}_4 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ . Both the K-means and Fuzzy C-means clustering are initialized with centers at  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .

- (i) What are the two cluster centers found by K-means clustering upon convergence, or after a maximum of 100 iterations? (0.45)
- (ii) What are the two cluster centers found by Fuzzy C-means clustering upon convergence, or after a maximum of 100 iterations? Assume fuzzier  $r=2$ . (1 p.5)

**Answer:** See TUT\_Q5.ipynb

(i) Initialization:  $\mathbf{c}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ ,  $\mathbf{c}_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

1<sup>st</sup> iteration: Assignment:

$$\begin{aligned}
\text{dist}(\mathbf{x}_1, \mathbf{c}_1) &= 0, \text{dist}(\mathbf{x}_1, \mathbf{c}_2) = 1, \mathbf{x}_1 \in \text{Cluster 1} \\
\text{dist}(\mathbf{x}_2, \mathbf{c}_1) &= 1, \text{dist}(\mathbf{x}_2, \mathbf{c}_2) = 0, \mathbf{x}_2 \in \text{Cluster 2} \\
\text{dist}(\mathbf{x}_3, \mathbf{c}_1) &= 1, \text{dist}(\mathbf{x}_3, \mathbf{c}_2) = 1.41, \mathbf{x}_3 \in \text{Cluster 1} \\
\text{dist}(\mathbf{x}_4, \mathbf{c}_1) &= 1.41, \text{dist}(\mathbf{x}_4, \mathbf{c}_2) = 1, \mathbf{x}_4 \in \text{Cluster 2}
\end{aligned}$$

Centroid Update:  $\mathbf{c}_1 = \frac{1}{2}(\mathbf{x}_1 + \mathbf{x}_3) = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}$ ,  $\mathbf{c}_2 = \frac{1}{2}(\mathbf{x}_2 + \mathbf{x}_4) = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$

2<sup>nd</sup> iteration: Assignment:

$$\begin{aligned}
\text{dist}(\mathbf{x}_1, \mathbf{c}_1) &= 0.5, \text{dist}(\mathbf{x}_1, \mathbf{c}_2) = 1.12, \mathbf{x}_1 \in \text{Cluster 1} \\
\text{dist}(\mathbf{x}_2, \mathbf{c}_1) &= 1.12, \text{dist}(\mathbf{x}_2, \mathbf{c}_2) = 0.5, \mathbf{x}_2 \in \text{Cluster 2} \\
\text{dist}(\mathbf{x}_3, \mathbf{c}_1) &= 0.5, \text{dist}(\mathbf{x}_3, \mathbf{c}_2) = 1.12, \mathbf{x}_3 \in \text{Cluster 1} \\
\text{dist}(\mathbf{x}_4, \mathbf{c}_1) &= 1.12, \text{dist}(\mathbf{x}_4, \mathbf{c}_2) = 0.5, \mathbf{x}_4 \in \text{Cluster 2} \\
\text{Cluster assignment is unchanged. Converged.}
\end{aligned}$$

Two cluster centers found:  $\mathbf{c}_1 = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}$ ,  $\mathbf{c}_2 = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$

(ii) Two cluster centers found:  $\mathbf{c}_1 = \begin{bmatrix} 0.0449 \\ 0.4999 \end{bmatrix}$ ,  $\mathbf{c}_2 = \begin{bmatrix} 0.9551 \\ 0.4999 \end{bmatrix}$

$$\text{Membership matrix } \mathbf{W} = \begin{bmatrix} 0.8219 & 0.1781 \\ 0.1781 & 0.8219 \\ 0.8217 & 0.1783 \\ 0.1783 & 0.8217 \end{bmatrix}$$

### Question 6 (Lecture 10)

Load the wine dataset by “from sklearn.datasets import load\_wine”. Assume that the class labels are not given. Apply K-means clustering algorithm to group all the data with different K from 2 to 10. Which K is the best according to the Elbow method?

**Answer:** See TUT7\_Q6.ipynb

Choose the K with “bend” or “knee” (Abrupt change in the slope): either 3 or 4.

### Question 7 (Lecture 11)

In neural networks, what is the purpose of introducing nonlinear activation functions like sigmoid and ReLU?

- (a) Ensure the output values to be between 0 and 1.
- ~~(b) Enable the model to learn complex patterns.~~
- (c) Speed up the gradient calculation in backpropagation.

**Answer:** (b)

### Question 8 (Lecture 11)

A multi-layer perceptron (MLP) of 2 layers has been constructed as:

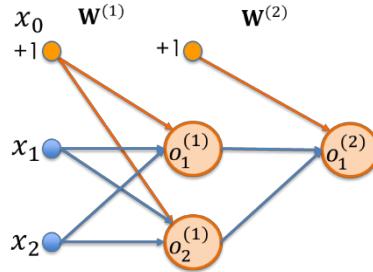
$$f_{\mathbf{W}}(\mathbf{x}) = [1, \sigma(\mathbf{x}^T \mathbf{W}^{(1)})] \mathbf{W}^{(2)}$$

where  $\sigma$  is the ReLU function,  $\mathbf{W}^{(1)} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix}$ ,  $\mathbf{W}^{(2)} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$ .

Suppose that this MLP is used to train a dataset  $\mathbf{X} = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 5 & 1 \end{bmatrix}$  with the target  $\mathbf{y} = \begin{bmatrix} 0.1 \\ 0.7 \end{bmatrix}$ . The cost function is  $J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N (f_{\mathbf{W}}(\mathbf{x}_i) - y_i)^2$  for optimization. The learning rate  $\eta$  is 0.1. What will  $\mathbf{W}^{(1)}$  and  $\mathbf{W}^{(2)}$  be after the 1<sup>st</sup> iteration of backpropagation?

**Answer:** See TUT\_Q8.ipynb

1. Architecture of MLP:



## 2. Learning:

1<sup>st</sup> iteration: Step 1: Forward Pass:

$$\begin{aligned}
 \mathbf{f}_{\mathbf{W}(0)}(\mathbf{x}) &= \left[ \mathbf{1}, \sigma^{(1)}(\mathbf{x} \mathbf{W}^{(1)}(0)) \right] \mathbf{W}^{(2)}(0) \\
 &= \left[ \mathbf{1}, \text{ReLU} \left( \begin{bmatrix} 1 & 2 & 1 \\ 1 & 5 & 1 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} \right) \right] \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 2 \\ 5 \end{bmatrix}
 \end{aligned}$$

Step 2: Backward Pass:

$$\mathcal{L} = (f_{\mathbf{W}}(\mathbf{x}) - y)^2$$

(1) Output layer:

$$\text{formula from the lecture notes: } \frac{\partial \mathcal{L}}{\partial w_{i,j}^{(L)}} = \epsilon_j^{(L)} a_i^{(L)},$$

$$\text{where } \epsilon_j^{(L)} = \sum_{m=1}^{D_L} \mathcal{L}'(o_m^{(L)}) \sigma^{(L)\prime}(z_j^{(L)})_m$$

$$\epsilon_1^{(2)} = \mathcal{L}'(o_1^{(2)}) = \mathcal{L}'(f_{\mathbf{W}}(\mathbf{x})) = 2(f_{\mathbf{W}}(\mathbf{x}) - y)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(2)}} = 2(f_{\mathbf{W}}(\mathbf{x}) - y) \mathbf{a}^{(2)}$$

$$\frac{dJ}{d\mathbf{W}^{(2)}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}_i}{\partial \mathbf{W}^{(2)}} = \frac{1}{N} \sum_{i=1}^2 2(f_{\mathbf{W}}(\mathbf{x}_i) - y_i) \mathbf{a}^{(2)} = \frac{1}{N} (\mathbf{A}^{(2)})^T 2(f_{\mathbf{W}}(\mathbf{x}) - \mathbf{y})$$

$$\frac{dJ(\mathbf{W}^{(2)}(0))}{d\mathbf{W}^{(2)}} = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 4 \end{bmatrix}^T \left( \begin{bmatrix} 2 \\ 5 \end{bmatrix} - \begin{bmatrix} 0.1 \\ 0.7 \end{bmatrix} \right) = \begin{bmatrix} 6 \\ 0 \\ 19.1 \end{bmatrix}$$

$$\mathbf{W}^{(2)}(1) = \mathbf{W}^{(2)}(0) - \eta \frac{dJ(\mathbf{W}^{(2)}(0))}{d\mathbf{W}^{(2)}} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} - 0.1 \begin{bmatrix} 6 \\ 0 \\ 19.1 \end{bmatrix} = \begin{bmatrix} 0.4 \\ -1 \\ -0.91 \end{bmatrix}$$

(2) Hidden layer:

Since it is ReLU activation,

$$\frac{dJ(\mathbf{W}^{(1)}(0))}{d\mathbf{W}^{(1)}} = \frac{1}{N} (\mathbf{A}^{(1)})^T \mathbf{E}^{(1)},$$

$$\begin{aligned} \text{where } \mathbf{E}^{(1)} &= \delta_{\mathbf{Z}^{(1)}} * \left( \mathbf{E}^{(2)} \left( \mathbf{W}_{1,:}^{(2)}(0) \right)^T \right) \\ &= \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} * \left( 2 \left( \begin{bmatrix} 2 \\ 5 \end{bmatrix} - \begin{bmatrix} 0.1 \\ 0.7 \end{bmatrix} \right) [-1, 1] \right) = \begin{bmatrix} 0 & 3.8 \\ 0 & 8.6 \end{bmatrix}, \end{aligned}$$

$$\mathbf{A}^{(1)} = \mathbf{X} = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 5 & 1 \end{bmatrix}$$

$$\text{therefore, } \frac{dJ(\mathbf{W}^{(1)}(0))}{d\mathbf{W}^{(1)}} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 2 & 5 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 3.8 \\ 0 & 8.6 \end{bmatrix} = \begin{bmatrix} 0 & 6.2 \\ 0 & 25.3 \\ 0 & 6.2 \end{bmatrix}$$

$$\begin{aligned} \text{then, } \mathbf{W}^{(1)}(1) &= \mathbf{W}^{(1)}(0) - \eta \frac{dJ(\mathbf{W}^{(1)}(0))}{d\mathbf{W}^{(1)}} \\ &= \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 & 6.2 \\ 0 & 25.3 \\ 0 & 6.2 \end{bmatrix} = \begin{bmatrix} -1 & -0.62 \\ 0 & -1.53 \\ 1 & -1.62 \end{bmatrix} \end{aligned}$$

### Question 9 (Lecture 11)

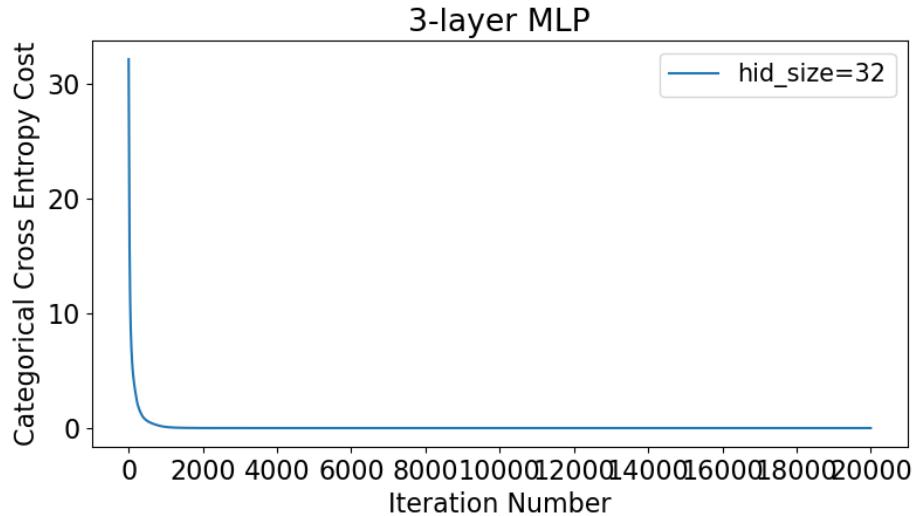
The Digits dataset is a classic, lightweight dataset ideal for testing simple machine learning models like MLPs, especially on CPU. It has 10 classes (digits 0 to 9). Each sample is a grayscale image with a dimension of  $8 \times 8$ . Get the dataset by “from sklearn.datasets import load\_digits”. Use Python to do the following tasks.

- (i) Split the database into three sets: 70% of samples for training, 15% of samples for validation, and 15% of samples for testing. Normalize the features using z-score standardization. Hint: you may use “from sklearn.preprocessing import StandardScaler”.
- (ii) Construct the target output using one-hot encoding.
- (iii) Use a 3-layer MLP for the digit classification. Suppose that ReLU is the activation function for the hidden layers and Softmax is the activation function for the output layer. Assume that the number of the neurons for each hidden layer is 32. The learning rate is 0.01. The number of iterations for training is 20000. The objective function is  $J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N L_{CCE}(\mathbf{f}_W(\mathbf{x}_i), \mathbf{y}_i)$ , where  $L_{CCE}(\mathbf{f}_W(\mathbf{x}_i), \mathbf{y}_i)$  is the categorical cross-entropy loss. Plot the cost over number of iterations. Compute the training accuracy and validation accuracy.
- (iv) Change the hidden layer size to 36 and do the same thing as in (iii). Which one is better? Use the model with the better hidden layer size to compute the accuracy of the test set.

**Answer:** See TUT7\_Q9

Note: Different initialization of the weight parameters and data splitting may have different results.

(iii) Training accuracy = 1, Validation accuracy = 0.8741



(iv) Training accuracy = 1, Validation accuracy = 0.8259,

Hidden layer size 32 is better according to the validation accuracy.

Test accuracy for model with hidden layer size of 32 = 0.8444.

