

EE2211

Lec1: What is machine learning?	5
AI, Machine Learning, and Deep Learning.....	7
When do we need machine learning?.....	8
Application of Machine Learning.....	9
Types of Machine Learning.....	11
Supervised Learning.....	11
Regression.....	12
Classification.....	14
Unsupervised Learning.....	16
Clustering.....	16
Reinforcement Learning.....	18
Example.....	20
Walking Through A Toy Example: Token Classification.....	20
Step 1: Feature Extraction.....	21
Step 2: Nearest Neighbour Classifier.....	22
Inductive vs. Deductive Reasoning.....	23
Summary and Practise Questions.....	24
Lec2: Data Engineering	25
Types of Data.....	25
Levels/Scales of Measurement.....	25
Nominal Data.....	27
Ordinal Data.....	27
Interval Data.....	28
Ratio Data.....	29
Example.....	29
Based on Numerical/Categorical.....	31
Other Aspects.....	31
Data Wrangling and Cleaning (Data Preprocessing).....	32
Formatting Data.....	33
Binary Coding.....	33
Normalisation.....	34
Data Cleaning.....	35
Data Integrity and Visualisation.....	37
Data Integrity.....	37
Data Visualization.....	39
Why Visualization is Necessary.....	42
Summary and Practice Question.....	42
Lec3: Introduction to Linear Algebra, Probability and Statistics	43
Introduction to Linear Algebra.....	43
Notations, Vectors, Matrices.....	43
Linear dependence and independence.....	47
Systems of Linear Equations.....	48
Causality.....	49

Causality is a Statistical Relationship.....	52
Correlation vs Causality.....	52
Correlation does not imply causation!.....	53
Simpson's paradox.....	54
Example.....	54
Random Variable.....	55
Axioms of Probability.....	56
Discrete Random Variable.....	57
Continuous Random Variable.....	59
Example 1.....	61
Example 2.....	62
Example 3.....	63
Bayes' Rule.....	63
Example.....	64
Summary and Practice Question.....	64
Lec4: Systems of linear equations.....	65
Vector and Matrix Operation.....	65
Dot Product = Inner Product.....	65
Inverse of Matrix.....	66
Determinant and Cofactor to find Inverse.....	67
Systems of Linear Equations.....	68
Types of System.....	69
Square or even-determined system.....	71
Over-determined system.....	72
Under-determined system.....	75
Set.....	80
Function.....	81
The inner product function.....	82
Linearity and Superposition.....	83
Affine Function.....	84
Class Quiz.....	84
Summary.....	85
Lec5: Least Squares and Linear Regression.....	86
Functions: Maximum and Minimum.....	86
Derivative and Gradient.....	87
Differentiation of a scalar function w.r.t. a vector.....	89
$\nabla = (\partial \partial x, \partial \partial y, \partial \partial z)$	89
$\nabla f = (\partial f \partial x, \partial f \partial y, \partial f \partial z)$	89
Differentiation of a vector function w.r.t. a vector.....	90
Properties.....	90
Linear Regression.....	91
Learning (Training).....	95
Least Squares Regression.....	95
Example 1: one dimension input, 1 feature.....	97
Example 2: 3 dimension input, 3 features.....	99
Learning of Vectored Function (Multiple Outputs).....	100

Example 3.....	102
Example 4.....	105
Class Quiz.....	107
Summary.....	108
Lec6: Ridge regression & Polynomial regression.....	109
Linear Regression for Classification.....	109
Binary Classification.....	109
Example.....	110
Multi-Categorical/Class Classification.....	112
Ridge Regression.....	116
Derivation.....	117
Primal and Dual Form.....	117
Polynomial Regression.....	122
Example.....	123
Polynomial Expansion.....	123
Example.....	125
Polynomial Regression for Classification.....	129
In Class Quiz.....	130
Lec7: Over-fitting, bias/variance trade-off.....	131
Overfitting, Underfitting & Model Complexity.....	131
Overfitting.....	131
Underfitting.....	132
Good Fit.....	133
Overfitting / Underfitting Schematic.....	133
Feature Selection.....	134
Selecting Features With Pearson's r.....	135
Regularisation.....	136
Regularisation Example.....	137
Bias-Variance Trade-off.....	138
Bias-Variance Decomposition Theorem.....	139
Example.....	140
Test 1.....	141
Test 2.....	141
Test 3.....	142
Test 4.....	142
Test 5.....	143
All together.....	143
Average.....	144
Conclusion.....	144
In Class Quiz.....	144
Python and Conda Stuff.....	144
Tutorial 1.....	146
Tutorial 2.....	146
Q1.....	147
Q2.....	148
Q3.....	152

Q5.....	155
Tutorial 3.....	155
Q5.....	155
Q10-Q11.....	156
Tutorial 4.....	156
Q1.....	156
Q2.....	157
Q6.....	158
Tutorial 5.....	160
Q1.....	160
Q2.....	161
Q3.....	162
Q4.....	163
Without removing duplicates VS removing the duplicates.....	163
Q5.....	165
Q6.....	166
Q8.....	168
Q9.....	168
Tutorial 6.....	169
Q7.....	169
Q8.....	169
Functions.....	170
Admin.....	172

Lec1: What is machine learning?

Learning is any process by which a system improves performance from experience.

- Herbert Simon

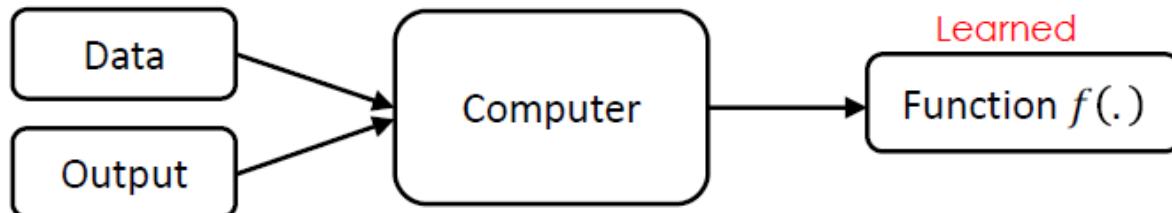
A computer program is said to learn

- from **experience E**
- with respect to some class of **tasks T**
- and **performance measure P**,

if its performance at tasks in T, as measured by P, improves with experience E.

- Tom Mitchell

Machine Learning (Supervised Learning)



Data Output



Cat

⋮



Dog

When applied

$f(\cdot)$



) → Cat !

New image

Machine Learning: field of study that gives computers the ability to learn without being explicitly programmed

- Arthur Samuel

- Computer with machine learning algorithm
- Output also called label
- Data output pair
- the function is said to be learned



Cat

⋮



Dog

$f(\text{cat}) = \text{'cat'}$

$f(\text{dog}) = \text{'dog'}$

- - $f(\text{data}) = \text{output}$

- Can predict after learnt

ARTIFICIAL INTELLIGENCE

Any technique which enables computers to mimic human behavior



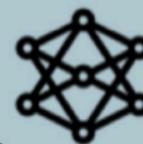
MACHINE LEARNING

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



DEEP LEARNING

A subset of ML which make the computation of multi-layer neural networks feasible



Timeline: 1950's, 1960's, 1970's, 1980's, 1990's, 2000's, 2010s

Example of AI but not ML: Deductive Reasoning

NUS is in Singapore, Singapore is in Asia \rightarrow NUS is in Asia

- ML can be used to achieve AI
- Deep learning is a subset of ML with specific architecture called neural network
- Deductive reasoning
 - If given A is in B and B is in C, then can reason that A is in C

When do we need machine learning?

Lack of human expertise
(Navigating on Mars)



Involves huge amount of data
(Genomics)



Learning is not always useful:

No need to “learn” to calculate payroll!

My Salary = Days_of_work * Daily Salary + Bonus

Application of Machine Learning

Task T, Performance P, Experience E

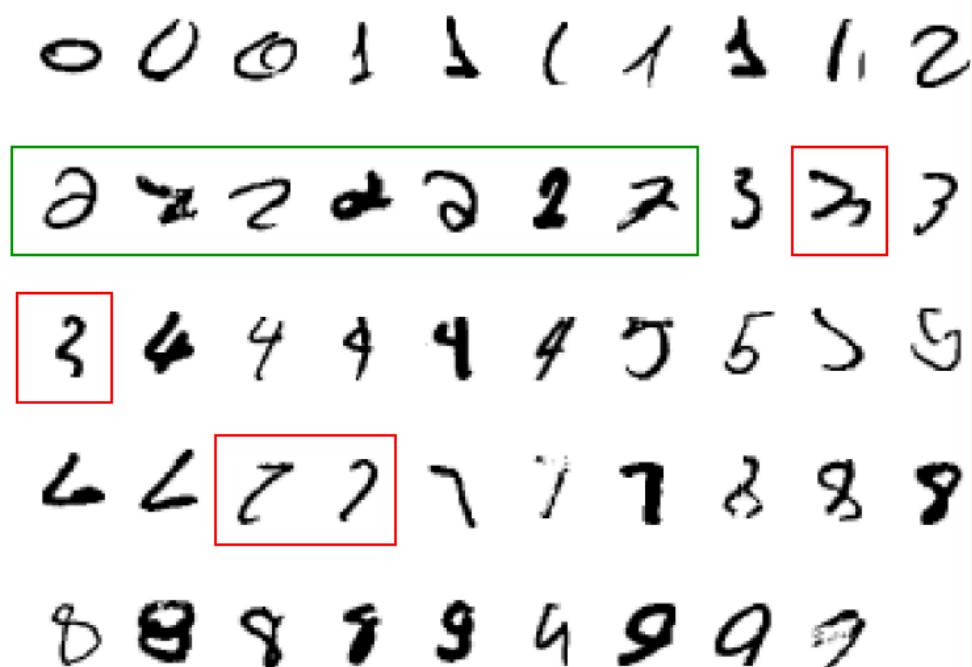
T: Digit Recognition

P: Classification Accuracy

E: Labelled Images

4 “four”

3 “three”



Labels -> Supervision!

- used by USPS, United States Postal Service
- By giving the labelled images the supervision is given to the model

Task T, Performance P, Experience E

T: Email Categorization

P: Classification Accuracy

E: Email Data, Some Labelled



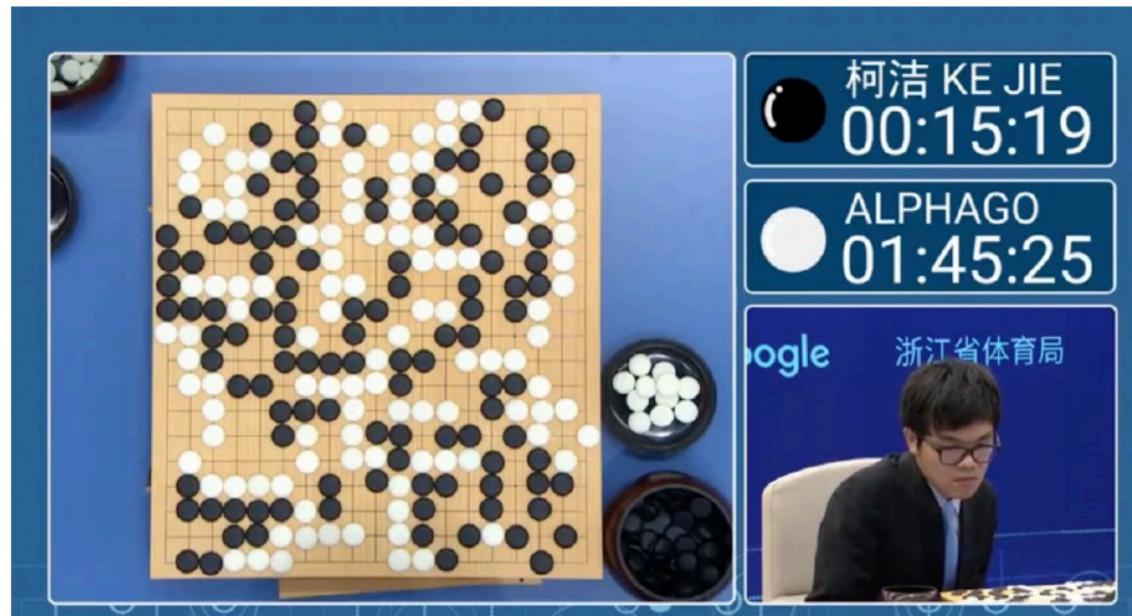
- Only some email data is labelled, user does not label all mails only some

Task T, Performance P, Experience E

T: Playing Go Game

P: Chances of Winning

E: Records of Past Games



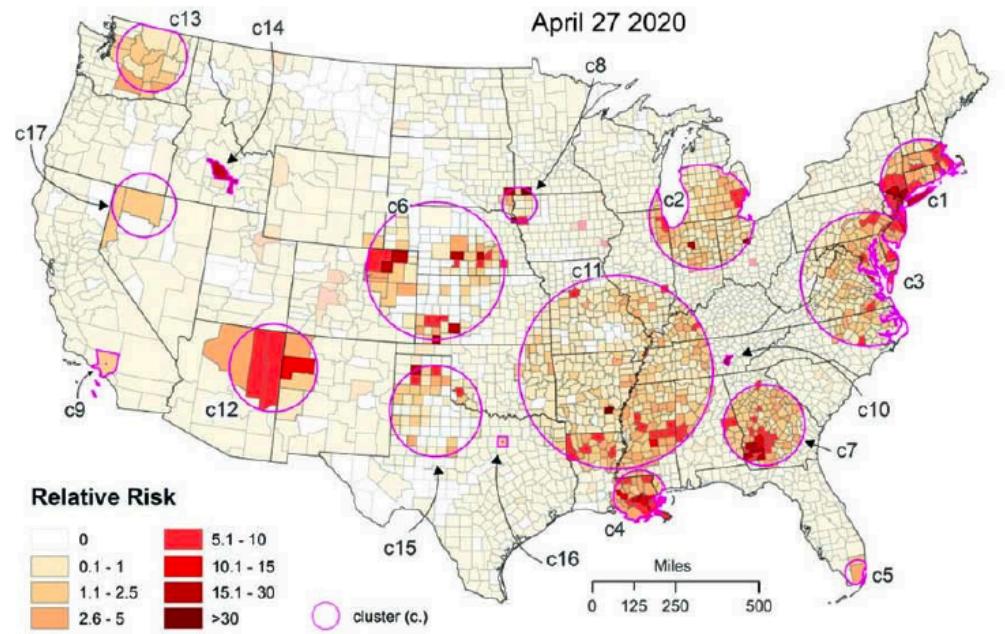
Task T, Performance P, Experience E

T: Identifying Covid-19 Clusters

P: Small Internal Distances

Larger External Distances

E: Records of Patients



- Since model is finding out clusters

- Between cluster should be big, picking two counties far apart from each other should belong in different cluster
 - small internal distance
- But between neighbouring county should be in one cluster, so picking two counties close to each other they should belong in the same cluster
 - larger external distance
- Cannot be supervised as only have records of patients not cluster

Types of Machine Learning

Supervised Learning

Input:

1) Training Samples,
2) Desired Output
(Teacher/Supervision)

Output:

A rule that maps input to output

Unsupervised Learning

Input:

Samples

Output:

Underlying patterns in data

Reinforcement Learning

Input:

Sequence of States,
Actions, and
Delayed Rewards

Output:

Action Strategy: a rule
that maps the
environment to action

Supervised Learning

Supervised Learning

Input:

1) Training Samples,
2) Desired Output
(Teacher/Supervision)

Output:

A rule that maps input to output

Supervised Learning

Data



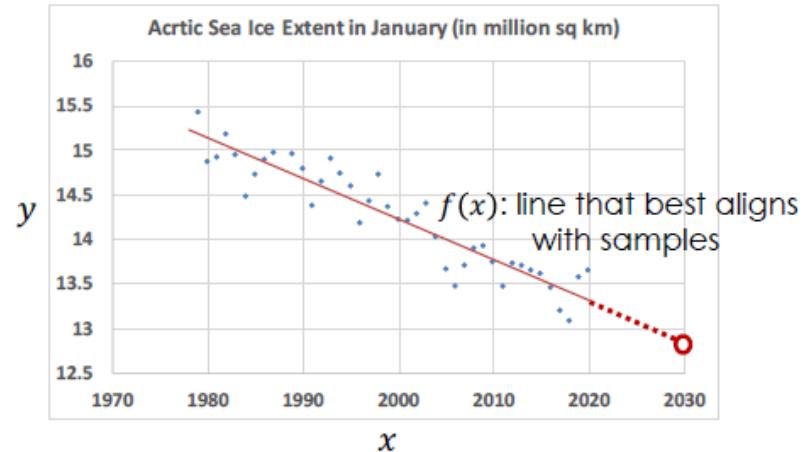
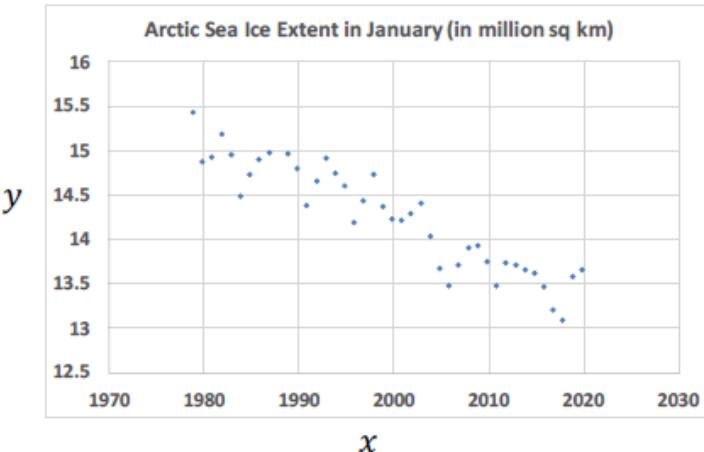
- If y is continuous, the type of learning is called regression

- if y is categorical, the type of learning is called classification

Regression

Regression

- Given $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- Learn a function $f(\mathbf{x})$ to predict real-valued y given \mathbf{x}



- if it's a line, then linear regression. If not line, regression but not linear
- to predict to the future
- Regression also belongs to the class of supervised learning methods.

Regression is just like classification except that y_i are no longer restricted to belong to a finite set. Rather it can take on uncountably many values. In the case of regression, we typically do not say that y_i is the label.

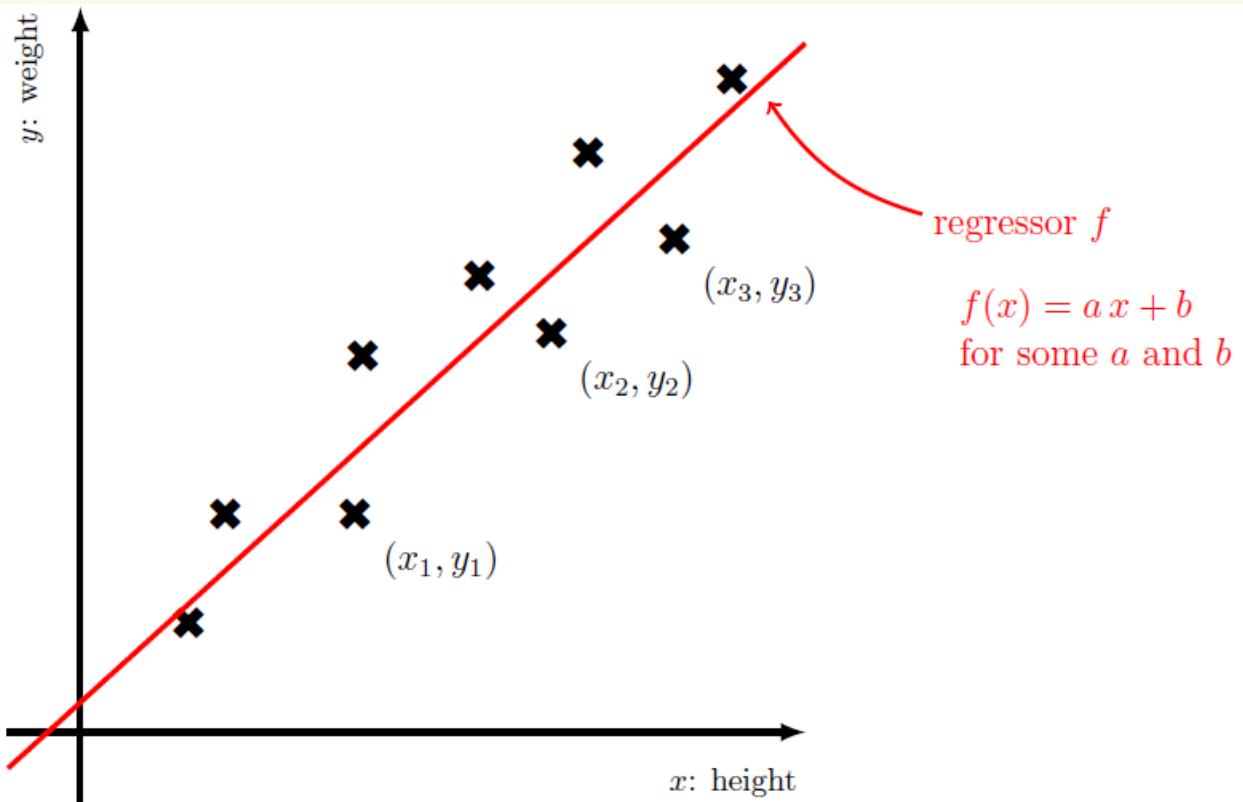


Figure 1.2: Height and weight regression example

Rather we use the term *target variable* or *outcome variable* or *dependent variable* to refer to the y_i 's. For instance y_i could take any value in the interval $[0, 1]$ (which is uncountable) or it could take values in \mathbb{R} . In either case, the data examples in the dataset \mathcal{D} , denoted as (\mathbf{x}_i, y_i) belong to $\mathbb{R}^d \times \mathbb{R}$, which means that each feature vector or training sample $\mathbf{x}_i \in \mathbb{R}^d$ (d -dimensional Euclidean space) and each target variable $y_i \in \mathbb{R}$ is a real number. We wish to learn a regressor $f : \mathbb{R}^d \rightarrow \mathbb{R}$ which is a function that brings us from the feature space to the set of real numbers. The regressor f should be designed such that given a new feature vector \mathbf{x}' , its corresponding target value y' is “well predicted”, in some precise mathematical sense, by $f(\mathbf{x}')$.

- Say each feature vector $x_i \in \mathbb{R}$ is scalar (so $d = 1$) and represents the height of a student. The target variable y_i represents the weight of the same student. Given the data of $m = 600$ students in EE2211 captured in the dataset $\mathcal{D} = \{(x_i, y_i) : 1 \leq i \leq 600\}$, I am tasked to learn a regressor $f : \mathbb{R} \rightarrow \mathbb{R}$ so that I can use the height of Alice x' in a new class EE9999 to predict her weight $y' = f(x')$. See Fig. 1.2.
- Say each feature vector $\mathbf{x}_i \in \mathbb{R}^3$ is a 3-dimensional vector and $x_{i,1}$ represents the air pressure at location i , $x_{i,2}$ represents the amount of rainfall at location i and $x_{i,3}$ represents the “amount of greenery” at location i . Each y_i corresponds to the temperature at location i . The temperature can take on uncountably many values—it is a real number. Given the dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i) : 1 \leq i \leq m\}$, we would like to learn a regressor $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ such that if I give you the feature vector of new location \mathbf{x}' , you can tell me the temperature $y' = f(\mathbf{x}')$.

Classification

The most canonical task is classification. This problem takes the following form. We have a *dataset* \mathcal{D} which consists of a certain number m of *data examples* $(\mathbf{x}_i, y_i), i = 1, \dots, m$. Each data example (\mathbf{x}_i, y_i) consists of a *feature vector* (also called *training sample* or *training example*)

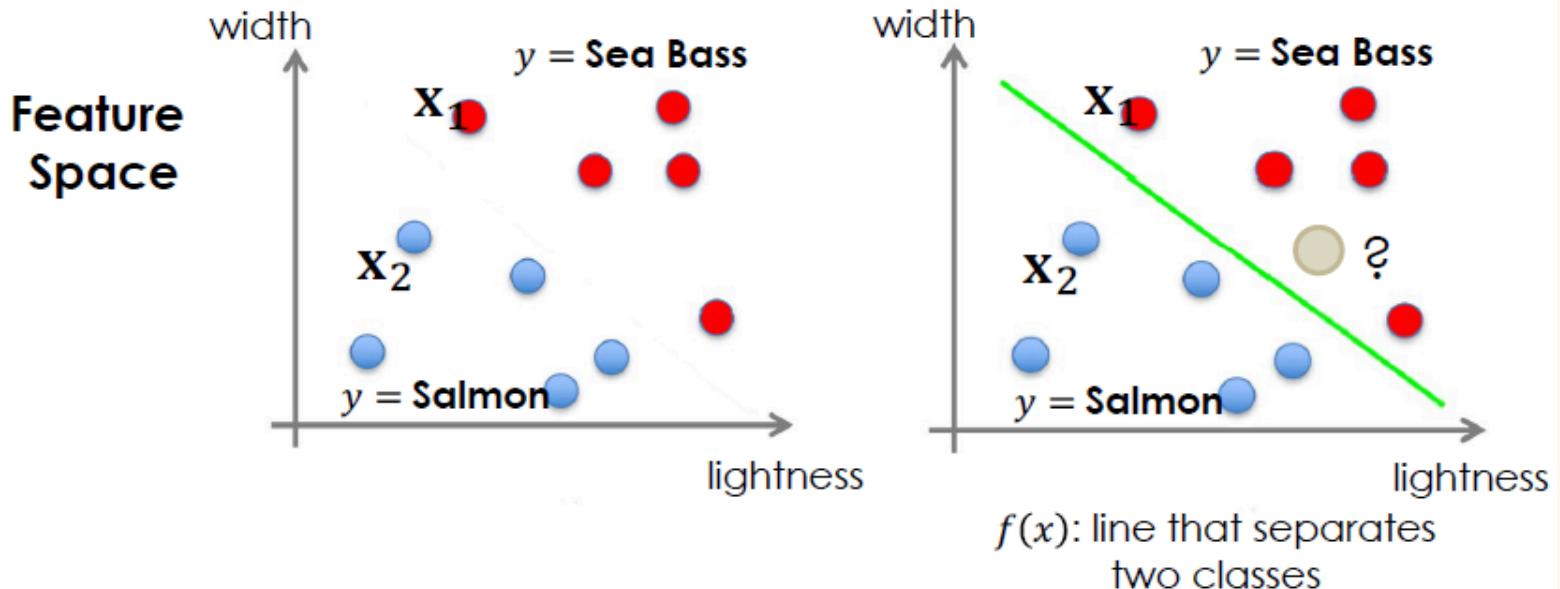
$$\mathbf{x}_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,d} \end{bmatrix} \quad \text{or} \quad \mathbf{x}_i = [x_{i,1} \ x_{i,2} \ \dots \ x_{i,d}]^\top, \quad (1.1)$$

(usually an element of d -dimensional Euclidean space \mathbb{R}^d) and a *label* y_i . The crux of classification is that the label y_i can only take on *finitely many values* so $y_i \in \{1, 2, \dots, c\}$ for some integer $c \geq 2$, where c is the number of classes. We do not allow y_i to take on infinitely many values.

The classification problem consists in finding a *classifier* which is a function $f : \mathbb{R}^d \rightarrow \{1, 2, \dots, c\}$ such that it accurately predicts labels given new samples, called *test samples*. This function f is constructed or learned based on the dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i) : 1 \leq i \leq m\}$.

Classification

- Given $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
- Learn a function $f(\mathbf{x})$ to predict categorical y given \mathbf{x}



- The two axis is the feature of the category, the plane is called feature space
- $f(x)$ learnt is to separate the two category
- Classification belongs to the class of supervised learning methods.
- Example

- We have records of $m = 100$ emails. Let $d = 2$. Thus, there are two components in each \mathbf{x}_i . Furthermore the first component of \mathbf{x}_i , denoted as $x_{i,1}$, counts the number of times the word “love” appears. The second component of \mathbf{x}_i , denoted as $x_{i,2}$, counts the number of times the word “money” appears. Each label $y_i \in \{0, 1\}$ where $y_i = 1$ means that the i^{th} email is spam while $y_i = 0$ means that the i^{th} email is non-spam. Based on the dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i) : 1 \leq i \leq 100\}$ of $m = 100$ emails, we are tasked to learn a classifier $f : \mathbb{R}^2 \rightarrow \{0, 1\}$ such that we can accurately predict whether the next email you receive (which is of course not in the training dataset) is spam or not. See Fig. 1.1. In this case in which there are only two classes, we refer to this as *binary classification*.

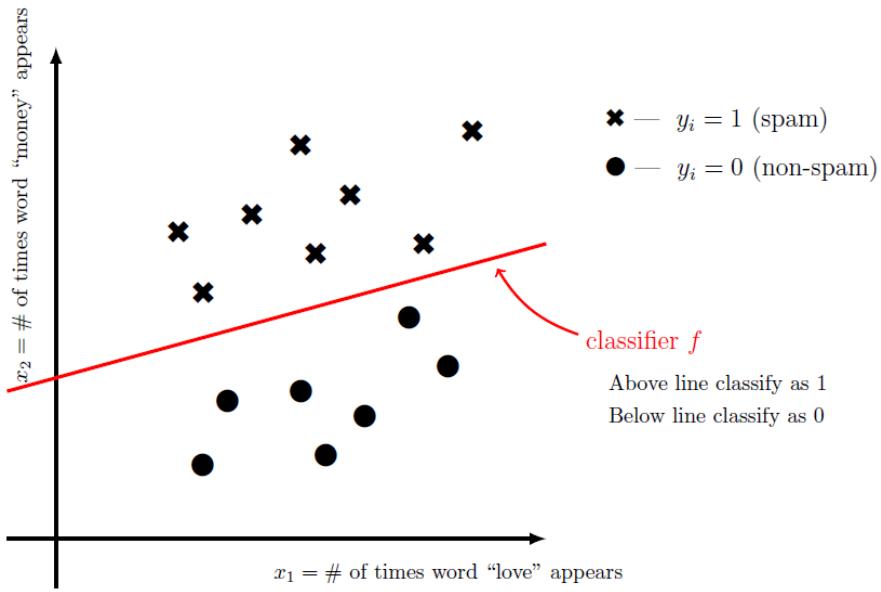


Figure 1.1: Spam classification

-

- Consider an image classification problem. We are given a dataset of size $m = 10^6$, namely, $\mathcal{D} = \{(\mathbf{x}_i, y_i) : 1 \leq i \leq 10^6\}$ where each \mathbf{x}_i represents an image that is vectorized into a vector. For example, each image contains $5 \times 5 = 25$ pixels and for the sake of simplicity, each pixel value is 0 or 1. Thus $\mathbf{x}_i \in \{0, 1\}^{5 \times 5} \cong \{0, 1\}^{25}$ and $d = 25$. Of course, we can have more complicated images that have multiple quantization levels (not restricted to binary) and colors as well, but the treatment will

○

be the same. To each image, there is a label y_i which can take on 10 values so $c = 10$. These 10 values signify what object is in the image, e.g., $y_i \in \{\text{dog, cat, \dots, snake}\}$. We would now like to design a image classifier $f : \{0, 1\}^{25} \rightarrow \{\text{dog, cat, \dots, snake}\}$ that “does well” (in some sense) on new images (or test images) that contain one of the 10 animals. Since there are more than two classes in this scenario, we refer to this as *multi-class* classification.

-

- Is the following a classification problem? There are a total of $M \geq 2$ football teams. Each football team $i \in \{1, 2, \dots, M\}$ plays L games against every other football team j . Say there are no draws. The number of times i beats j is denoted as $b_{ij} \in \{0, 1, \dots, L\}$ so necessarily $b_{ij} + b_{ji} = L$ (convince yourself that this is the case). We set $b_{ii} = 0$ for all i . Clearly, if one team wins all its games against other team, it is the best. When this is not case, which is common, it is not so clear what is the best strategy to *rank* the team. More precisely, given the matrix $\mathbf{B} = [b_{ij}]_{1 \leq i \leq M, 1 \leq j \leq M}$ we would like to rank the teams from top (best) to bottom (worst). This is a non-standard machine learning problem. However, notice that the total number of rankings is $M!$, which is finite, so in some sense this is like a classification problem. In what way is it not a classification problem? See [XFTF19] for some discussion of how we can use machine learning to learn the skills of tennis players.

Unsupervised Learning

Unsupervised Learning

Input:

Samples

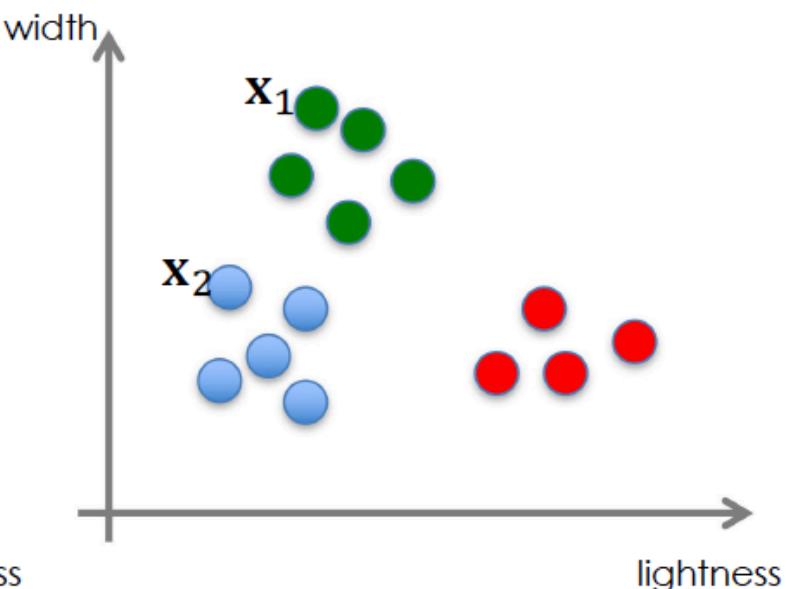
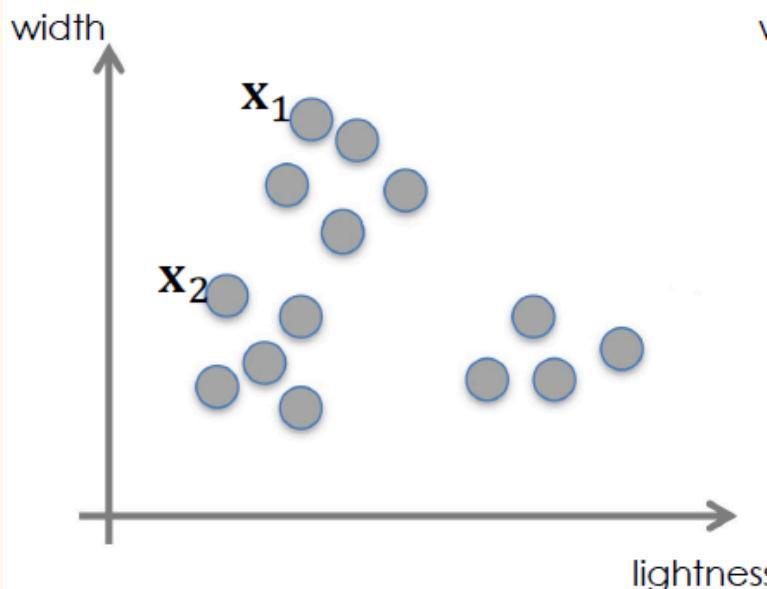
Output:

Underlying patterns in data

Clustering

Clustering

- Given $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, without labels
- Output Hidden Structure Behind



No Label/Supervision is given!

- Only clustering in syllabus for unsupervised
 - only categorical in syllabus, continuous exist
- Can only give x , cannot give y , there is no label/desire output
- Output hidden structure behind, cluster the groups
- small internal distance and larger external distance

In clustering, the labels or target variables y_i are no longer available and we *only* have access to the feature vectors $\mathcal{D} = \{\mathbf{x}_i : 1 \leq i \leq m\}$. There is often reason to believe that these feature vectors can be grouped or clustered into different clusters.

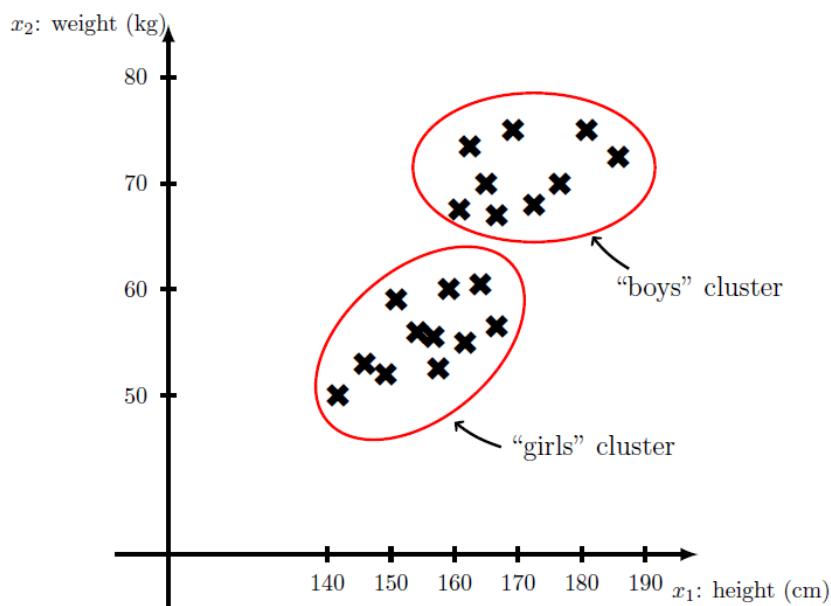


Figure 1.3: Height and weight clustering example

- Say the vector $\mathbf{x}_i = [x_{i,1}, x_{i,2}]^\top$ represents the height and weight of the i^{th} student in this semester's EE2211 cohort. There is reason to believe that girls are shorter and lighter than boys. So given $\mathcal{D} = \{\mathbf{x}_i : 1 \leq i \leq m\}$, can we assign each partition this set into 2 groups automatically? See Fig. 1.3. Note that while the clustering algorithm knows that there are 2 clusters, it does not and cannot assign semantic meanings (such as "boys" and "girls") to the clusters it discovers. It can only tell us that \mathcal{D} is partitioned into two non-empty disjoint subsets \mathcal{D}_1 and \mathcal{D}_2 (i.e., $\mathcal{D}_1 \cap \mathcal{D}_2 = \emptyset$ and $\mathcal{D}_1 \cup \mathcal{D}_2 = \mathcal{D}$).
- Say each $\mathbf{x}_i = [x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4}]^\top$ and $x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4}$ respectively represent the number of times the words "winner", "victory", "stock", "prices" appeared in the i^{th} article in a bag of m news articles. There is reason to believe that sports articles would contain more of the first 2 words but finance articles would contain more of the last 2 words. Hence, it seems possible to design a clustering algorithm to group the m feature vectors into two clusters, one representing sports articles while the other representing finance articles.

Reinforcement Learning

Reinforcement Learning

Input:

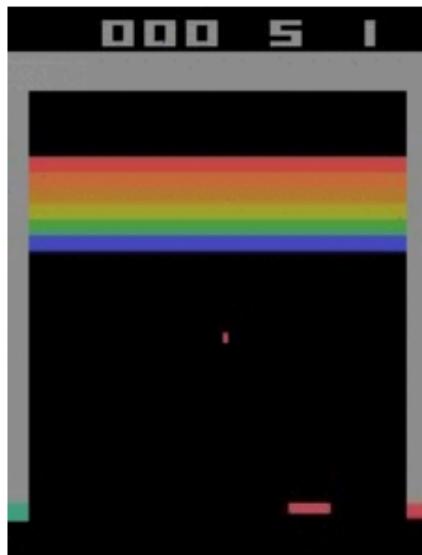
Sequence of States,
Actions, and
Delayed Rewards

Output:

Action Strategy: a rule
that maps the
environment to action

- When task requires a sequence of state or actions

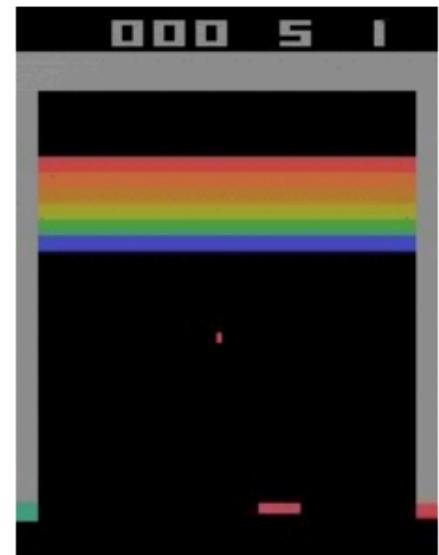
Breakout Game



Initial Performance



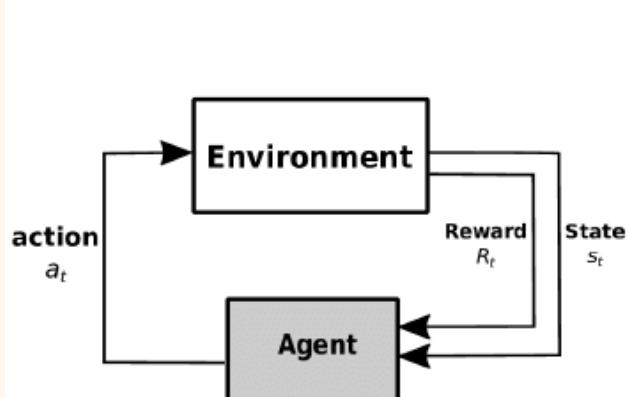
Training 15 minutes



Training 30 minutes

- Its a gif of the game playing better after training

- Given sequence of states S and actions A with (delayed) rewards R
- Output a policy $\pi(a, s)$, to guide us what action a to take in state s



S : Ball Location, Paddle Location, Bricks

A : left, right

R :

positive reward
Knocking a brick, clearing all bricks

negative reward
Missing the ball

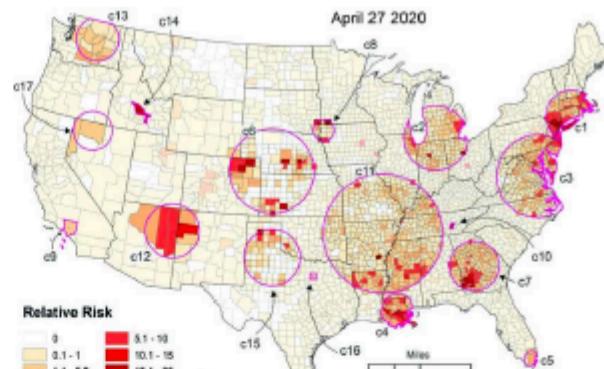
zero reward
Cases in between

- Intuitively, reinforcement learning with its actions and rewards is like supervision learning. But it's given a different label because the nature is so different

Example

0 0 0 1 1 1 1 1 1 2
 2 2 2 2 2 2 3 3 3
 3 4 4 4 4 4 5 5 5
 6 6 7 7 7 7 8 8 8
 8 8 9 9 9 9 9 9

Supervised



Unsupervised

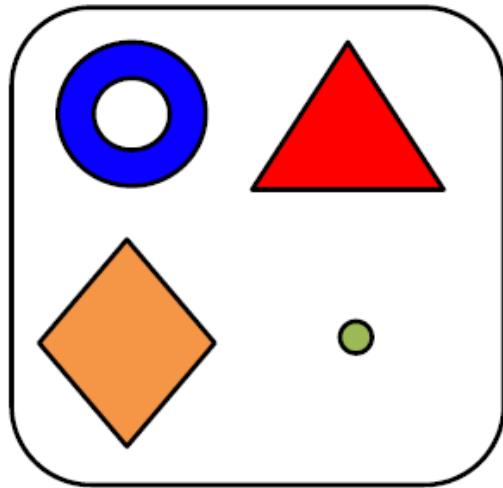


Supervised

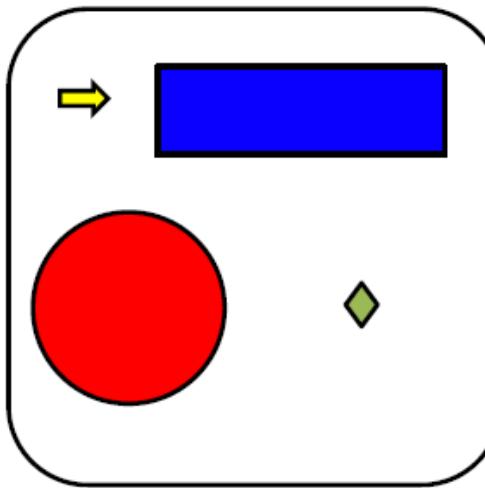


Reinforcement

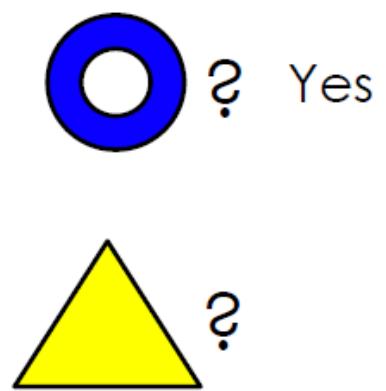
Walking Through A Toy Example: Token Classification



Yes



No



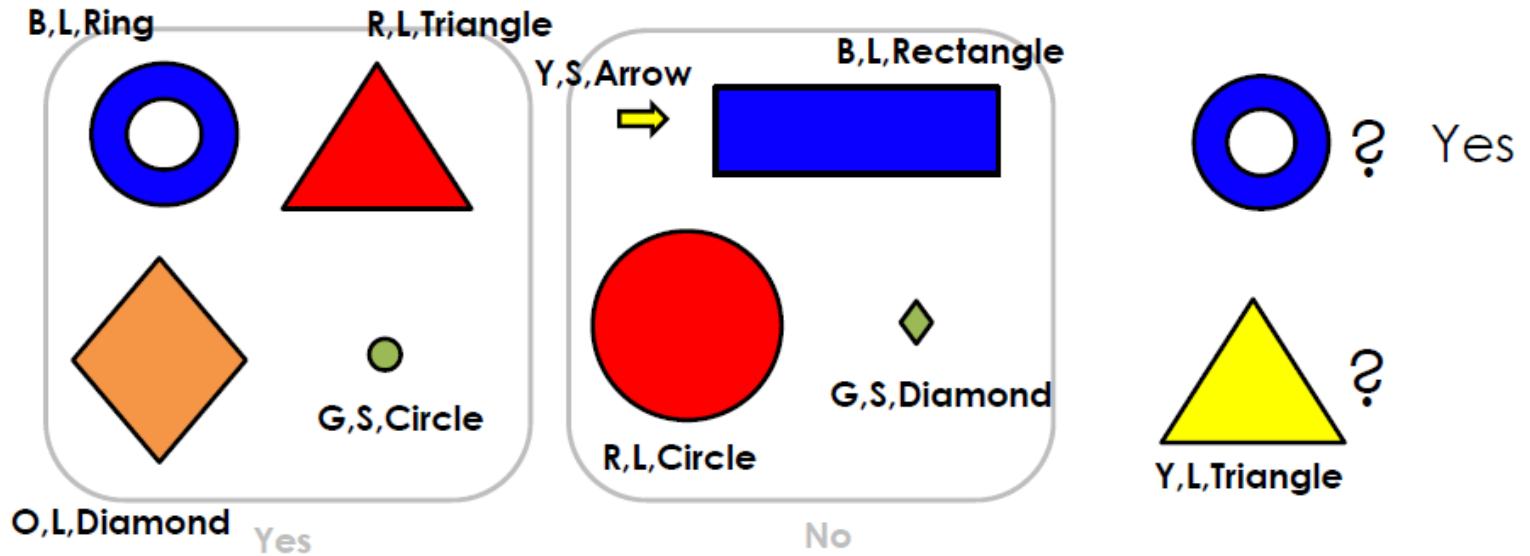
- Issue comes out when the testing sample (yellow triangle) does not overlap with the training sample (the set of all shapes with labels)
- Need to train a classification system

Step 1: Feature Extraction
Extract Attributes of Samples

Step 2: Sample Classification
Decide Label for a Sample

- requires a way to describe each sample, by extracting out their attributes
- based on the feature extracted, conduct the classification
- Sample classification can have many ways to do
 - Nearest Neighbour Classifier is just one way

Step 1: Feature Extraction



- Extract the feature of colour, size, shape (can have more)

Feature Extraction

	Color	Size	Shape	Label
○	Blue	Large	Ring	Yes
▲	Red	Large	Triangle	Yes
◆	Orange	Large	Diamond	Yes
●	Green	Small	Circle	Yes
→	Yellow	Small	Arrow	No
■	Blue	Large	Rectangle	No
●	Red	Large	Circle	No
◆	Green	Small	Diamond	No
▼	Yellow	Large	Triangle	?

- To check the testing sample, check the testing sample against all the other training sample, and build up a Similarity table
- feature extraction = feature selection (for EE2211)

Step 2: Nearest Neighbour Classifier

Similarity

	Color	Size	Shape	Total
	0	1	0	1
	0	1	1	2
	0	1	0	1
	0	0	0	0
	0	0	0	0
	1	0	0	1
	0	1	0	1
	0	1	0	1
	0	0	0	0

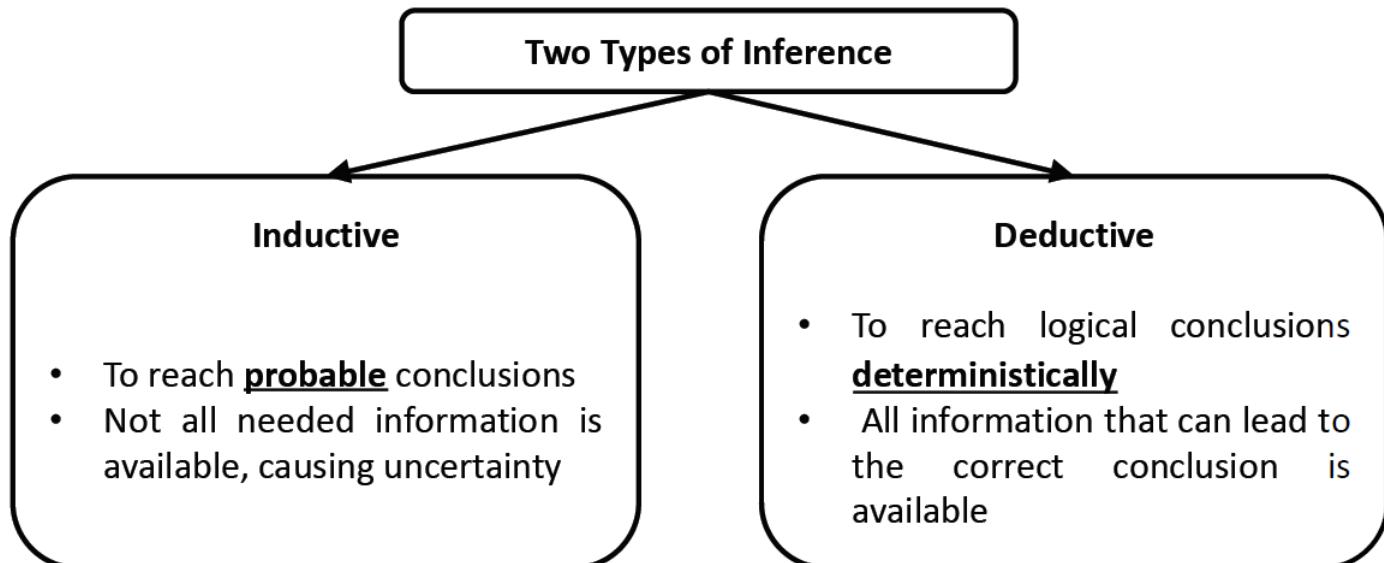
Nearest Neighbor Classifier:

- 1) Find the “nearest neighbor” of a sample in the feature space
- 2) Assign the label of the nearest neighbor to the sample

- Yellow big triangle is the most similar to the red big triangle, hence, given the label of red big triangle. Most similar in the feature space
- Nearest neighbour = more similar
- If attributes is changed the result may differ
 - need to extract as many feature as possible to avoid draw in similarity
 - if has draw choose any
- In EE2211, the newly labelled will not be learn, only the experience from the input/output pair given
 - done in advanced technique
- Nearest Neighbour Classifier does not requires training, it's just comparing against lookup table

Inductive vs. Deductive Reasoning

- Main Task of Machine Learning: to make inference



Probability and Statistics



Rule-based reasoning

NUS is in Singapore, Singapore is in Asia ->
NUS is in Asia

- Inductive
 - Do not have all the information needed, can only give probable conclusion
 - Result have probability/confidence
 - Has uncertainty
 - If face is far away hard to tell
 - EE2211 ML: is inductive
 - Guessing population from sample, small picture to big picture
- Deductive
 - Have all the information needed
 - Conclusion is deduced from information and hence 100% confidence and certainty
 - Draw conclusion deterministically
 - With the big picture get the small picture

Summary and Practise Questions

Three Components in ML Definition

Task T, Performance P, Experience E

Two Types of Supervised Learning

Classification, Regression

Three Types of in ML

Supervised Learning

Unsupervised Learning

Reinforcement Learning

One Type of Unsupervised Learning

Clustering

Inductive and Deductive

Inductive: Probable

Deductive: Rule-based

Example of a Classifier Model

Nearest Neighbor Classifier

Which of the following statement is true?

- A. Nearest Neighbor Classifier is an example of unsupervised learning
- B. Nearest Neighbor Classifier is an example of deductive learning
- C. Nearest Neighbor Classifier is an example of feature selection
- D. None of the above is correct.

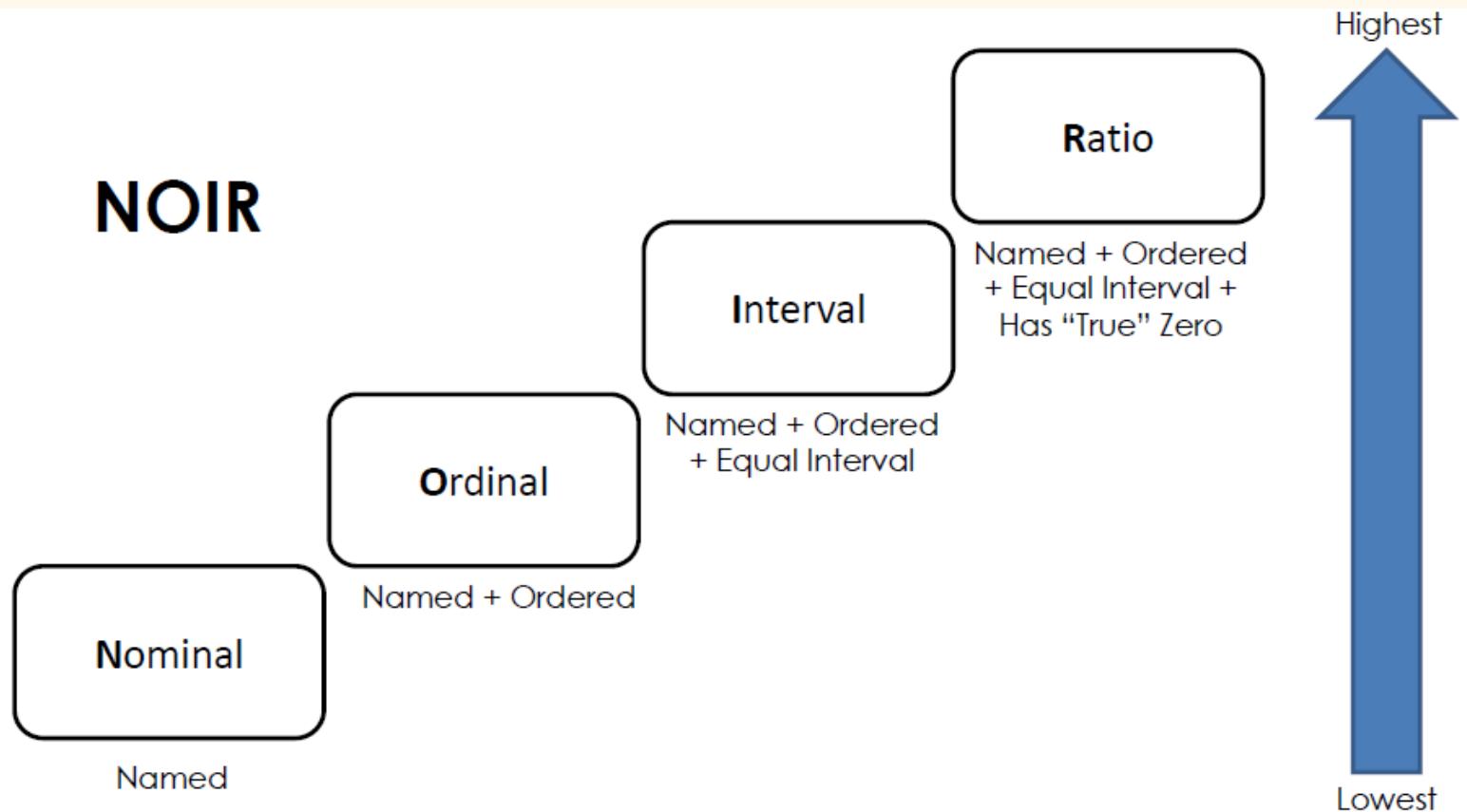
- A is wrong, nearest neighbour classifier is supervised
- B is wrong, nearest neighbour classifier is not with full certainty
- C is wrong, nearest neighbour classifier is the second step, the first step is feature extraction/selection
 - other classification can be used
 - feature extraction = feature selection (for EE2211)

Lec2: Data Engineering

Types of Data

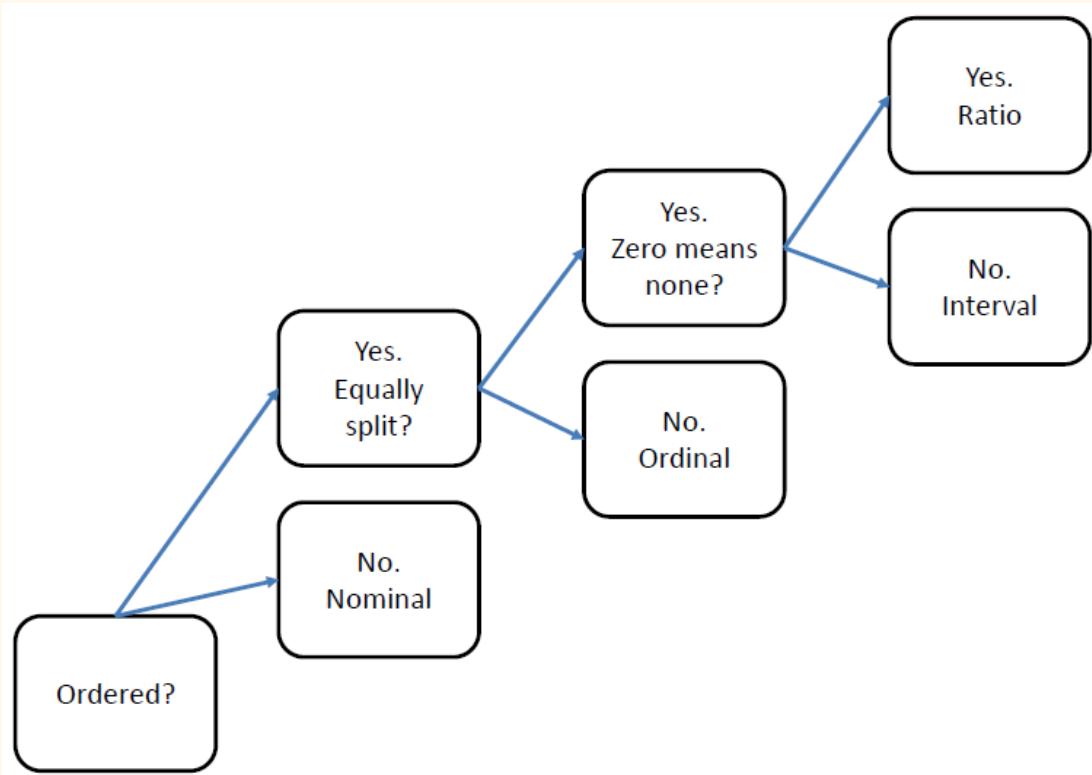
- Based on [Levels/Scales of Measurement](#)
 - Nominal Data
 - Ordinal Data
 - Interval Data
 - Ratio Data
- Based on [Numerical/Categorical](#)
 - Numerical, also known as Quantitative
 - Categorical, also known as Qualitative
- Other aspects
 - Available or Missing Data

Levels/Scales of Measurement



- Going up means, Ordinal is Named + Ordered

We can estimate	Nominal	Ordinal	Interval	Ratio
Frequency Distribution	Yes	Yes	Yes	Yes
Median	No	Yes	Yes	Yes
Add or subtract	No	No	Yes	Yes
Mean, standard deviation	No	No	Yes	Yes
Ratios	No	No	No	Yes



Nominal Data

- Lowest Level of Measurement
- Discrete Categories
- **NO** natural order
- Estimating a **mean**, **median**, or **standard deviation**, would be meaningless.
- Possible Measure: **mode**, **frequency distribution**
- Example:

Gender	Occupation
	
man	Doctor
	
woman	Police
	
	Teacher

Ordinal Data

- **Ordered** Categories
- Relative Ranking
- Unknown “distance” between categories: orders matter but not the difference between values
- Possible Measure: **mode**, **frequency distribution** + **median**
- Example:
 - Evaluate the difficulty level of an exam
 - 1: Very Easy, 2: Easy, 3: About Right, 4: Difficult, 5: Very Difficult
 - can sort means can have median

Interval Data

- Ordered Categories
- Well-defined “unit” measurement:
 - Distances between points on the scale are measurable and well-defined
 - Can measure differences!
- **Equal Interval** (between two consecutive unit points)
- **Zero is arbitrary** (not absolute), in many cases human-defined
 - If the variable equals zero, it does not mean there is none of that variable
- **Ratio is meaningless**
- Possible Measure: mode, frequency distribution + median + mean, standard deviation, addition/subtraction
- Example:
 - Temperature measured in Celsius
 - For instance: 10 degrees C, 28 degrees C
 - Year of someone's birth
 - For instance: 1990, 2005, 2010, 2022
- because zero is arbitrary defined, the ratio of two value is meaningless
- 20 degC is not twice as energetic as 10degC



Ratio Data

- Most precise and **highest** level of measurement
- Ordered
- Equal Intervals
- **Natural Zeros:**
 - If the variable equals zero, it means there is none of that variable
 - Not arbitrary
- Possible Measure: mode, frequency distribution + median + mean, standard deviation, addition/subtraction + multiplication and division (ratio)
- Example:
 - Weights
 - 10 KG, 20 KG, 30 KG
 - Time
 - 10 Seconds, 1 Hour, 1 Day
- well defined ratio, 20 kg twice as heavy as 10 kg



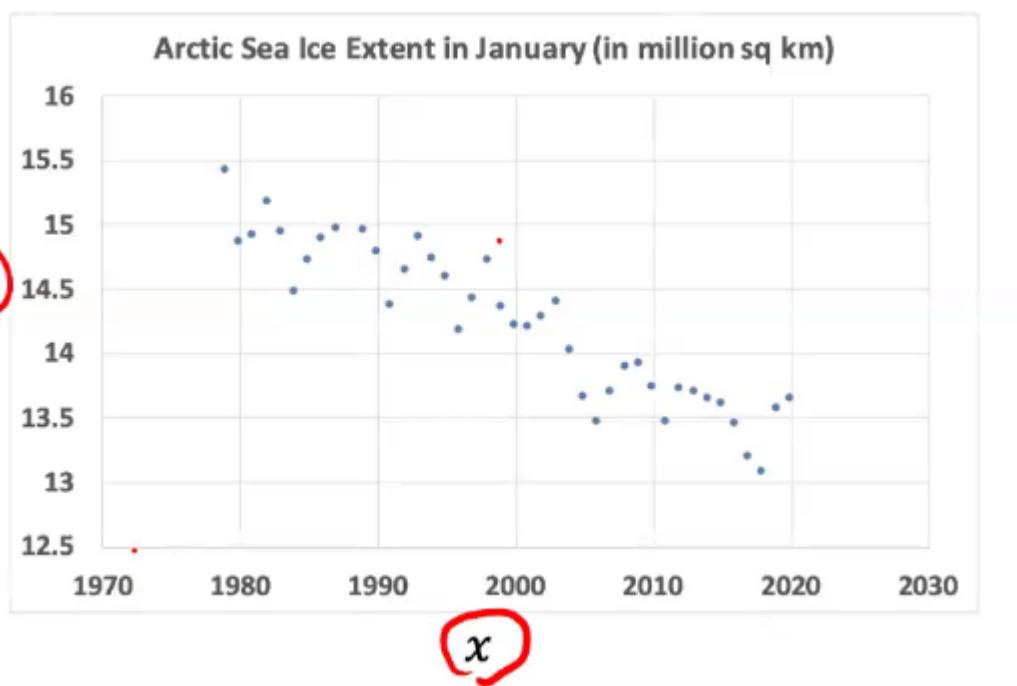
Example

1. **Favorite Restaurant**
 - Mcdonald's, Burger King, Subway, KFC, ...
2. **Weight of luggage measured in KG**
3. **SAT Scores: note that, SAT ranges is [400, 1600]**
4. **Size of Packed Eggs in supermarkets**
 - Small, Medium, Large, Extra Large, ...
5. **Military rank**
 - General, Major, Captain, ...
6. **Number of people in a household**
 - 1, 2, 3, 4, 5, ...
7. **Credit Score in United States: the range is [300, 850]**



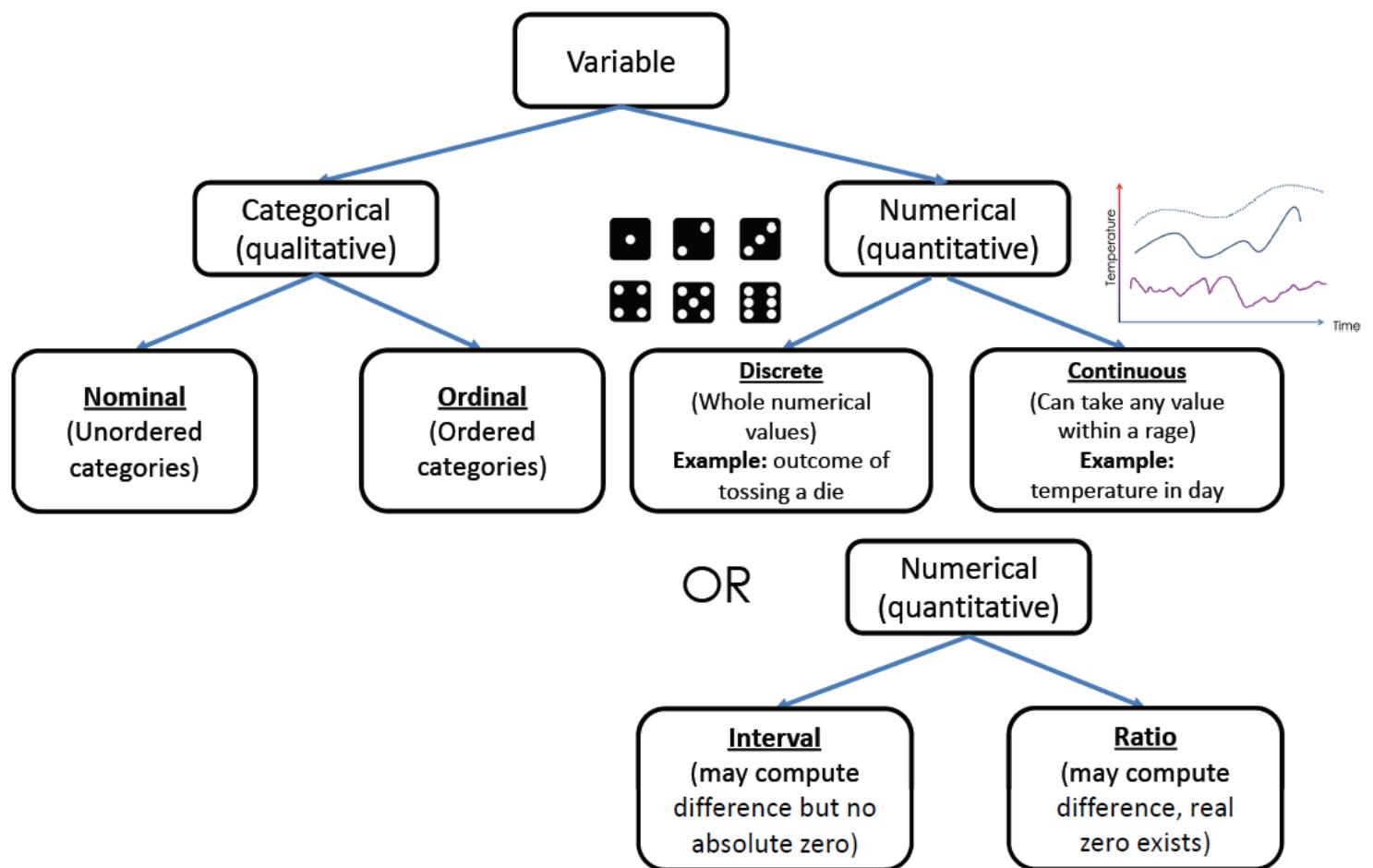
1. Nominal

2. Ratio
3. Interval
4. Ordinal
5. Ordinal
6. Ratio
7. Interval



- x is interval, y is ratio

Based on Numerical/Categorical



Other Aspects

- Missing data: data that is missing and you do not know the mechanism.
 - You should use a single common code for all missing values (for example, “NA”), rather than leaving any entries blank.

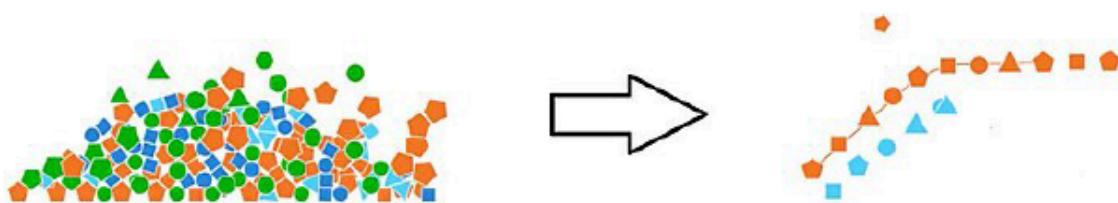
NUS student	Age	Country of birth
Olivia Tan	20	Singapore
Hendra Setiawan	19	Indonesia
John Smith	19	NA

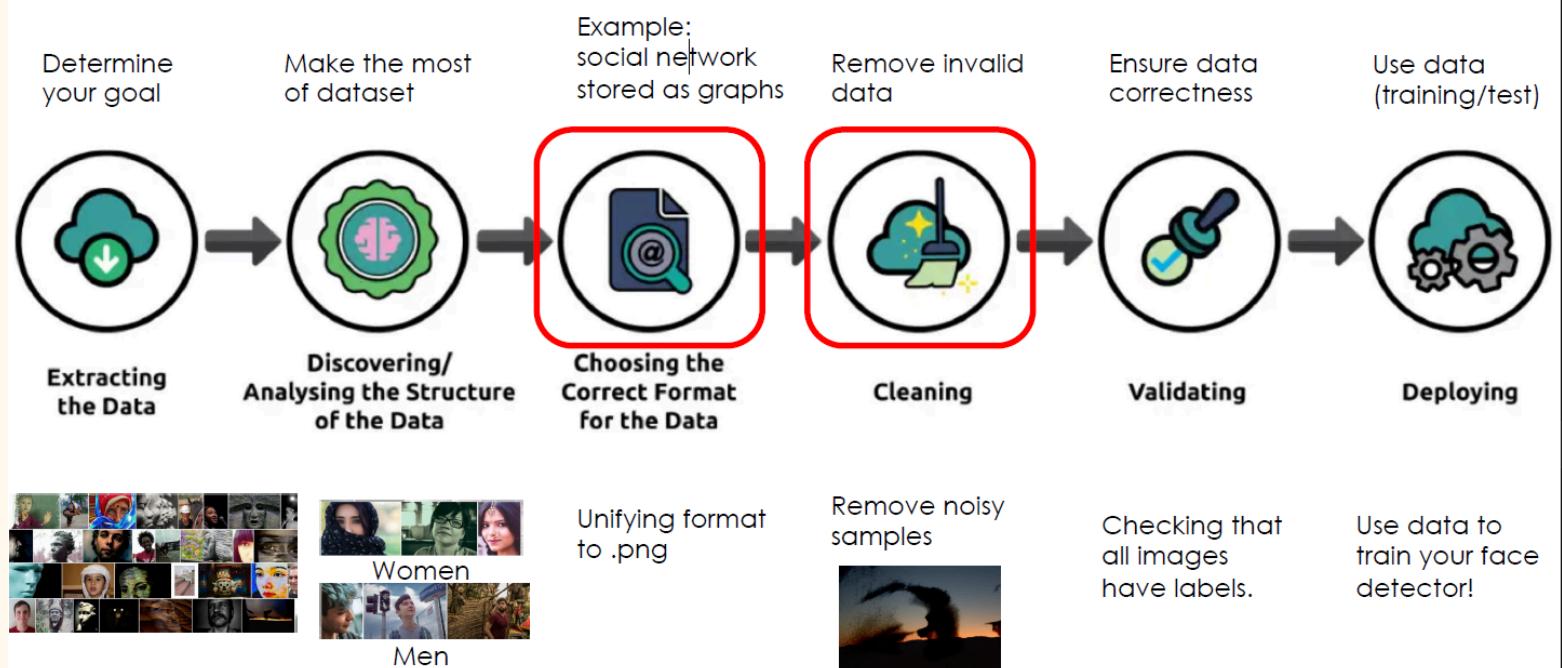
Data Wrangling and Cleaning (Data Preprocessing)



- Prepare data for ML/feature extraction
- Cleaning is subset of wrangling

- The process of transforming and mapping data from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.
- In short, transforms data to gain insight
- It is a general process!





Collect Human Face Images for Face Detector

Credit: <https://understandingdata.com/what-is-data-wrangling/>

- second step is to group the data

Formatting Data

Binary Coding

- Binary Coding** to convert categories into binary form
 - One-hot encoding: unify several entities within one vector
 - Example: the color of a pixel can be red, yellow, or green
 - Very common in classification tasks!

$$\text{red} = [1, 0, 0]$$

$$\text{yellow} = [0, 1, 0]$$

$$\text{green} = [0, 0, 1]$$

- Format Labels/output
- One-hot encoding is one type of binary coding
 - Put stuff in one vector

! $\text{red} = [1, 0, 0]$

$\text{yellow} = [0, 1, 0]$

$\text{green} = [0, 0, 1]$

- o One-hot memes, no inbetween, only one attribute can be present

■ $\boxed{[1, 0, 0]}$

Normalisation

Often we have feature vectors in which features are on different scales. We let m be the number of training samples and d be the number of features as usual. Say the feature vectors are

$$\mathbf{x}_1 = \begin{bmatrix} x_{1,1} \\ x_{1,2} \\ \vdots \\ x_{1,d} \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} x_{2,1} \\ x_{2,2} \\ \vdots \\ x_{2,d} \end{bmatrix}, \quad \dots \quad \text{and} \quad \mathbf{x}_m = \begin{bmatrix} x_{m,1} \\ x_{m,2} \\ \vdots \\ x_{m,d} \end{bmatrix}. \quad (2.2)$$

The first feature $x_{i,1}$ for $1 \leq i \leq m$ could be the height of students measured in centimeters so the dynamic range is $[x_{\min,1}, x_{\max,1}] = [140, 190]$. The second feature $x_{i,2}$ for $1 \leq i \leq m$ could measure the (American) shoe size of the students so it ranges from $[x_{\min,2}, x_{\max,2}] = [6, 13]$. So even if both features are deemed equally “important”, unfortunately, any machine learning method would place more importance on the first feature because of its larger values, which is not ideal. Thus, we have to scale or normalize the features so that their dynamic ranges are roughly the same. We can use (at least) two methods to do so:

• Normalization

- Linear Scaling:

scale each variable to $[0, 1]$

$$x_i = \frac{x_i^{\text{raw}} - x^{\min}}{x^{\max} - x^{\min}}, \quad i = 1, 2, \dots, M$$

- Z-score standardization:

each independent dimension of data is normally distributed

$$x_i = \frac{x_i^{\text{raw}} - E[X]}{\sigma(X)}, \quad i = 1, 2, \dots, M.$$

- Format data/input

- Linear scaling

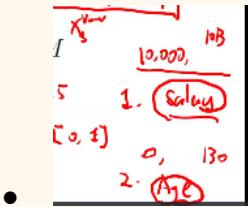
input raw $1, 2, 3, 5, 8, 9$

$$x_i = \frac{x_i^{\text{raw}} - x^{\min}}{x^{\max} - x^{\min}}, \quad i = 1, 2, \dots, M$$

$$x_3 = \frac{3 - 1}{9 - 1} = \frac{2}{8} = 0.25$$

- o convert into 0-1

- any range to 0-1
- uses for converting two data of very different range into comparable range



- - example is age and salary

$$\bar{x}_{i,1} := \frac{x_{i,1} - x_{\min,1}}{x_{\max,1} - x_{\min,1}}$$

-

- Z-score

- must be normally distributed

$$\hat{x}_1 = \frac{1}{m} \sum_{i=1}^m x_{i,1}, \quad \text{and} \quad \hat{\sigma}_1 = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (x_{i,1} - \hat{x}_1)^2}$$

-

$$\bar{x}_{i,1} := \frac{x_{i,1} - \hat{x}_1}{\hat{\sigma}_1}$$

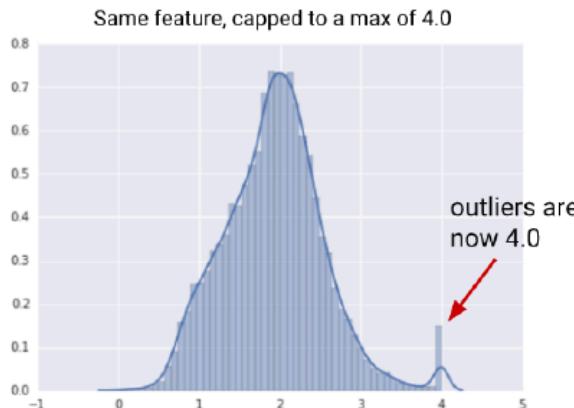
-

Data Cleaning

- The process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database.

- Example:

- Clipping outliers



- Handling missing features

Students	Year of Birth	Gender	Height	GPA
Tan Ah Kow	1995	M	1.72	4.2
Ahmad Abdul	X NA	M	1.65	4.1
John Smith	1995	M	1.75	X NA
Chen Lulu	1995	F	X NA	4.0
Raj Kumar	1995	M	1.73	4.5
Li Xiuxiu	1994	F	1.70	3.8

- must make sure outlier is outlier before taking them out, otherwise it's still valuable data
- handling missing features has several ways

Data Cleaning: Handling missing features

1. Removing the examples with missing features from the dataset

- Can be done if the dataset is big enough so we can sacrifice some training examples

2. Using a learning algorithm that can deal with missing feature values

- Example: random forest

3. Using a data imputation technique

- Only can remove if dataset is big enough
- imputation means fill in the data

Data Cleaning: Handling missing features: Imputation

- Method 1. Replace the missing value of a feature by an average value of this feature in the dataset:

$$\hat{x}^{(j)} \leftarrow \frac{1}{N} \sum_{i=1}^N x_i^{(j)}$$

- Method 2. Highlight the missing value

- Replace the missing value with a value outside the normal range of values.
- For example, if the normal range is [0, 1], then you can set the missing value to -1.
- Enforce the learning algorithm to learn what is best to do when the feature has a value significantly different from regular values.

- Method 1 is less commonly used, as it's not the best practice
- Method 2, is used to make the algorithm able to learn what to do when encounter -1 in the future
 - -1 is only for the algo, when finding mean and whatnot, ignore the -1

Data Integrity and Visualisation

Data Integrity

- Data integrity is the maintenance and the assurance of data accuracy and consistency;
 - A critical aspect to the design, implementation, and usage of any system that stores, processes, or retrieves data.
 - Very broad concept!
- Example:
 - In a dataset, numeric columns/cells should not accept alphabetic data.
 - A binary entry should only allow binary inputs

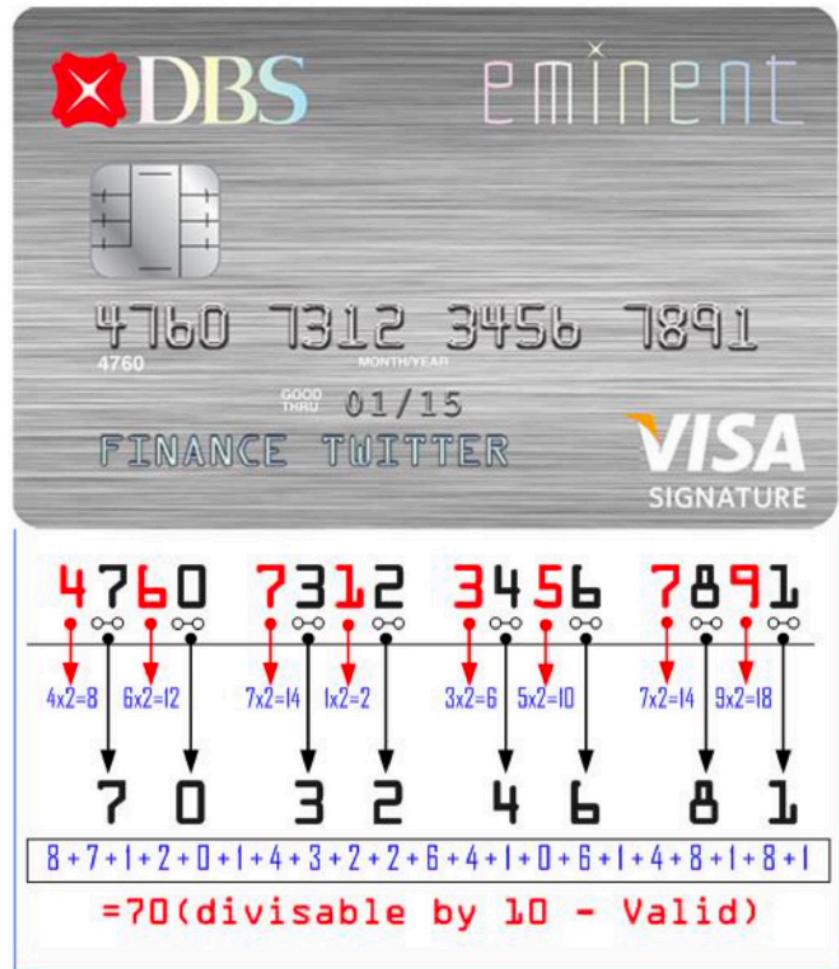
We can only select one of these

Organization	User Type	Is_Emergency ↑	External Profile Entered	Subject Areas		Bid	Relevance	Candidate Suggestion Rank	Tpms Rank	Quota	Number Of Assignments
				Primary	Secondary						
National University of Singapore	Student, >3 times as reviewer for CVPR, ICCV, or ECCV			Machine learning	3D from single images; Adversarial attack and defense; Computer vision theory; Explainable computer vision; Self- & semi-& meta- & unsupervised learning; Transfer/ low-shot/ long-tail learning; Vision + graphics	Not Entered	0.08	1	1434		4
Zhejiang University	Faculty/Researcher, 3-10 times as reviewer for CVPR, ICCV, or ECCV	No		Transfer/ low-shot/ long-tail learning	Efficient learning and inferences; Explainable computer vision; Image and video synthesis and generation; Recognition; detection, categorization,	Not Entered	0.16	7			2



4760 7312 3456 7891

Issuer Number Bank Number Account Number Check Digits



$$8 + 7 + 1 + 2 + 0 + 1 + 4 + 3 + 2 + 6 + 4 + 1 + 0 + 6 + 1 + 4 + 8 + 1 + 8 + 1$$

=70 (divisible by 10 - Valid)

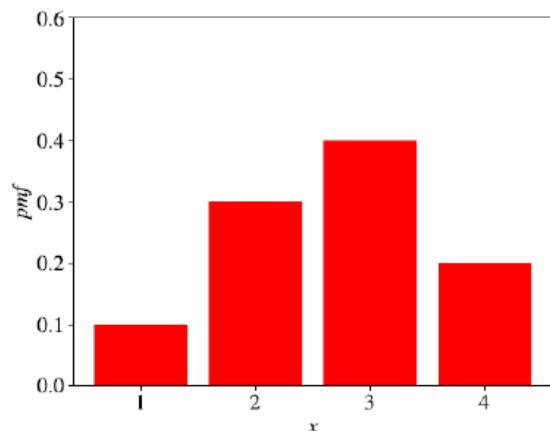


Chart Types

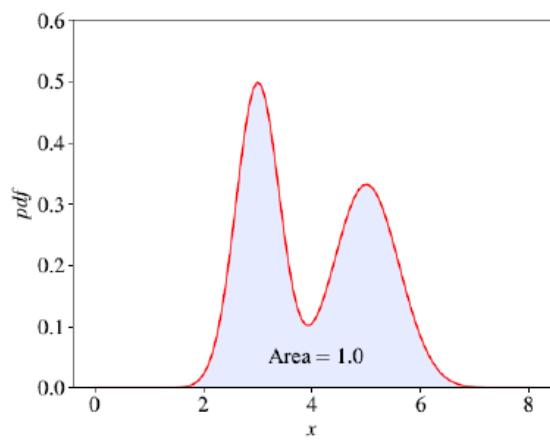


www.adatis.co.uk

Visualization: Distribution



Probability Mass Function



Probability Density Function

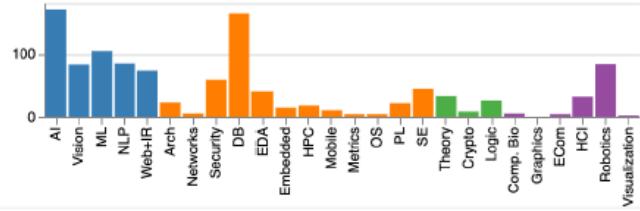
Visualization: Bars

CSRankings: Computer Science Rankings

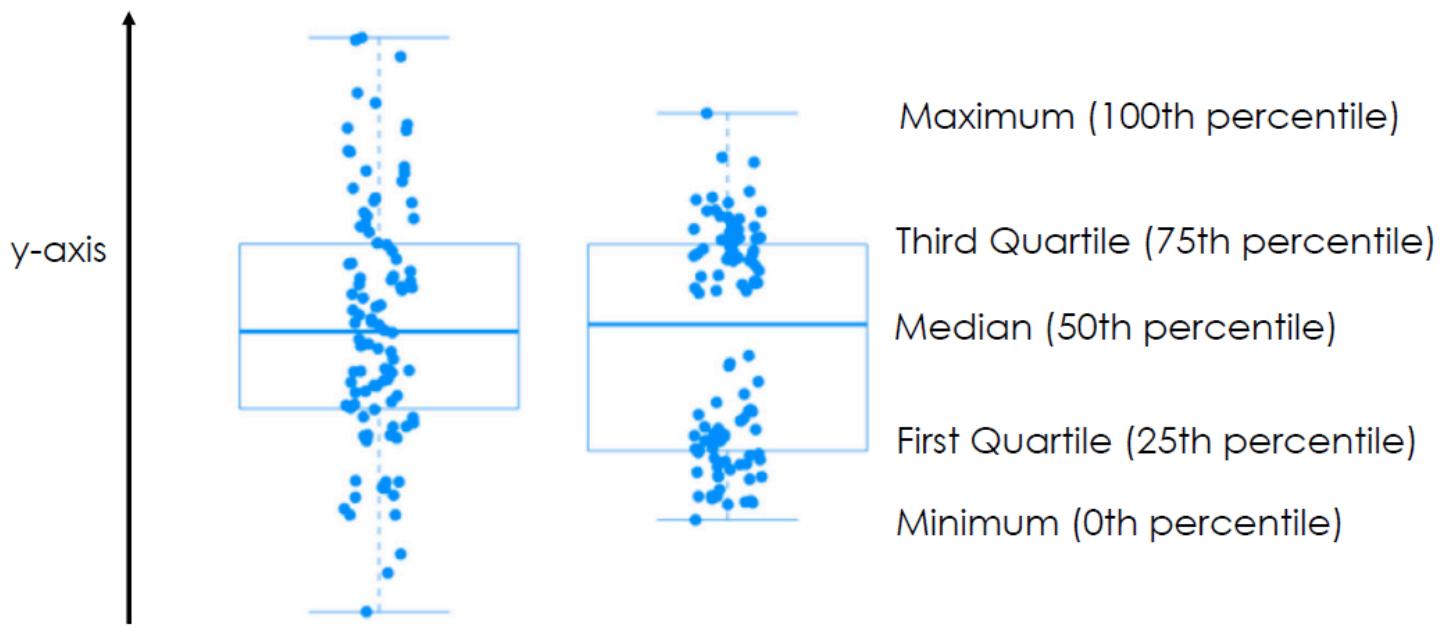
CSRankings is a metrics-based ranking of top computer science institutions around the world. Click on a triangle (\blacktriangleright) to expand areas or institutions. Click on a name to go to a faculty member's home page. Click on a chart icon (the after a name or institution) to see the distribution of their publication areas as a bar chart. Click on a Google Scholar icon (to see publications, and click on the DBLP logo (to go to a DBLP entry. Applying to grad school? Read this first.

Rank institutions in by publications from to

12	Georgia Institute of Technology	9.1	94
13	University of Maryland - College Park	8.2	83
14	University of Wisconsin - Madison	7.6	65
15	Columbia University	7.4	55
15	National University of Singapore	7.4	66

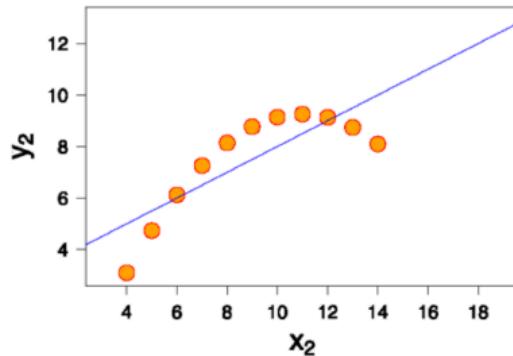
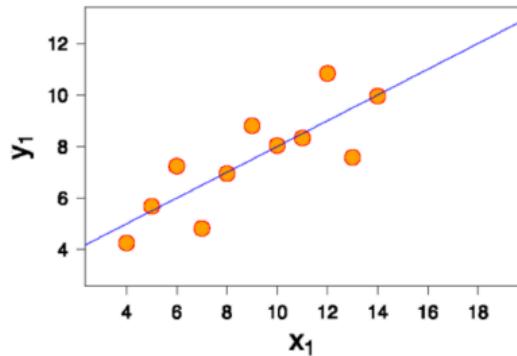


Visualization: Boxplots

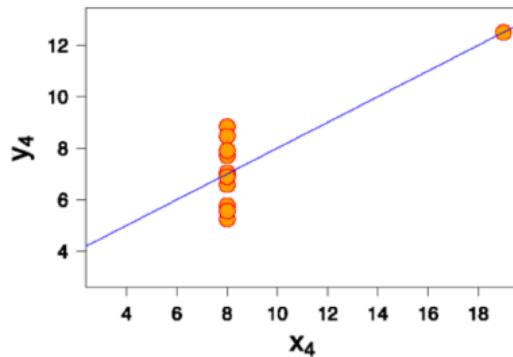
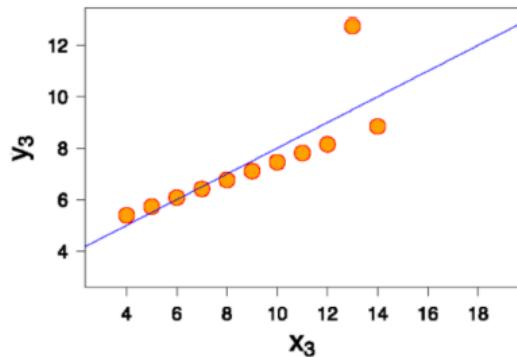


- The first quartile (Q_1) is defined as the middle number between the smallest number (i.e., Minimum) and the median of the data set.
- The third quartile (Q_3) is the middle number between the median and the highest value (i.e., Maximum) of the data set.

Why Visualization is Necessary



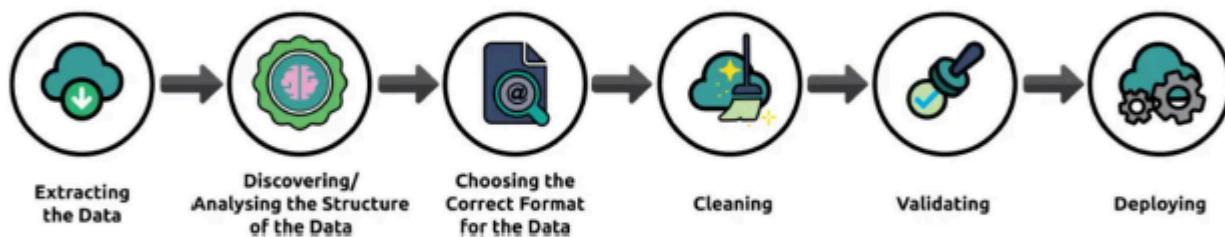
Four datasets with **identical** means, variances and regression lines!



Hence, we need visualization to show their difference!

Summary and Practice Question

- Types of data
 - NOIR
- Data wrangling and cleaning



- Data integrity and visualization
 - Integrity: Design
 - Visualization: Graphical Representation

Color	Size	Shape
Blue	Large	Ring
Red	Large	Triangle
Orange	Large	Diamond
Green	Small	Circle
Yellow	Small	Arrow
Blue	Large	Rectangle
Red	Large	Circle
Green	Small	Diamond

What are the NOIR data types of *color*, *size*, and *shape* in the table?

- Colour - Nominal
- Size - Ordinal
- Shape - Nominal

Lec3: Introduction to Linear Algebra, Probability and Statistics

Introduction to Linear Algebra

Notations, Vectors, Matrices

- A **scalar** is a simple numerical value, like 15 or -3.25
 - Focus on **real** numbers
- **Variables** or **constants** that take scalar values are denoted by an *italic* letter, like x or a

- A **vector** is an ordered list of scalar values
 - Denoted by a **bold character**, e.g. \mathbf{x} or \mathbf{a}

- In many books, vectors are written column-wise:

$$\mathbf{a} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -2 \\ 5 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

- The three vectors above are two-dimensional, or have two elements

- We denote an **entry** or **attribute** of a vector as an italic value with an index, e.g. $a^{(j)}$ or $x^{(j)}$.

- The index j denotes a specific dimension of the vector, the position of an attribute in the list

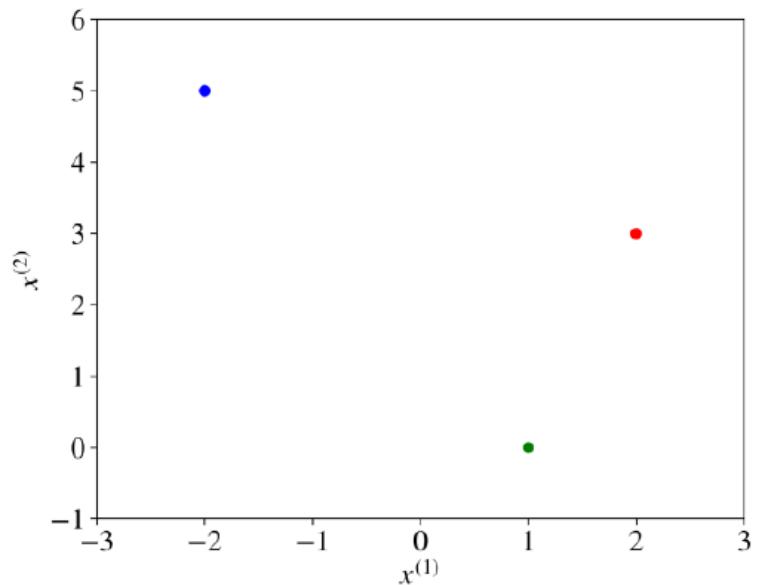
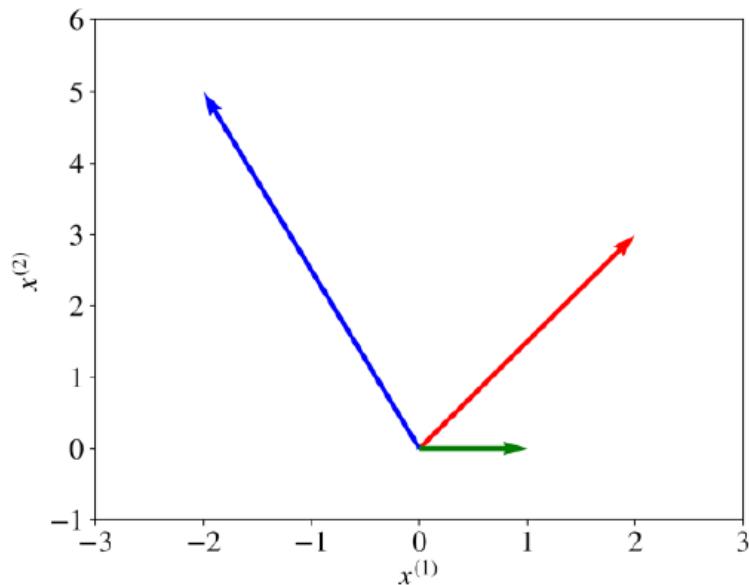
$$\mathbf{a} = \begin{bmatrix} a^{(1)} \\ a^{(2)} \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad \text{or more commonly} \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

- Note:

- $x^{(j)}$ is not to be confused with the power operation, e.g., x^2 (squared)
- Square of an indexed attribute of a vector is denoted as $(x^{(j)})^2$.

- **Vectors** can be visualized as, in a multi-dimensional space,
 - arrows that point to some directions, or
 - points

Illustrations of three two-dimensional vectors, $\mathbf{a} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} -2 \\ 5 \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$



- A **matrix** is a rectangular array of numbers arranged in rows and columns
 - Denoted with bold capital letters, such as **X** or **W**
 - An example of a matrix with two rows and three columns:
$$\mathbf{X} = \begin{bmatrix} 2 & 4 & -3 \\ 21 & -6 & -1 \end{bmatrix}$$
- A **set** is an unordered collection of unique elements
 - When an element x belongs to a set S , we write $x \in S$.
 - A special set denoted **R** includes all real numbers from minus infinity to plus infinity
- Note:
 - For elements in matrix **X**, we shall use the indexing $x_{1,1}$ where the first and second indices indicate the row and the column position.
 - Usually, for input data, rows represent samples and columns represent features
- set is unordered with unique elements

$S = \{ 3, 2 \}$

or **W** Feature 1 Feature 2 Feature 3

$$\mathbf{X} = \begin{bmatrix} 2 & 4 & -3 \\ 21 & -6 & -1 \end{bmatrix} \begin{matrix} \text{Sample 1} \\ \text{Sample 2} \end{matrix}$$

- **Capital Sigma:** the summation over a collection $\{x_1, x_2, x_3, x_4, \dots, x_m\}$ is denoted by:

$$\sum_{i=1}^m x_i = x_1 + x_2 + \dots + x_{m-1} + x_m$$

- **Capital Pi:** the product over a collection $\{x_1, x_2, x_3, x_4, \dots, x_m\}$ is denoted by:

$$\prod_{i=1}^m x_i = x_1 \cdot x_2 \cdot \dots \cdot x_{m-1} \cdot x_m$$

Linear dependence and independence

- A collection of d -vectors $\mathbf{x}_1, \dots, \mathbf{x}_m$ (with $m \geq 1$) is called **linearly dependent** if

$$\beta_1 \mathbf{x}_1 + \dots + \beta_m \mathbf{x}_m = 0$$

holds for some β_1, \dots, β_m that are **not all zero**.

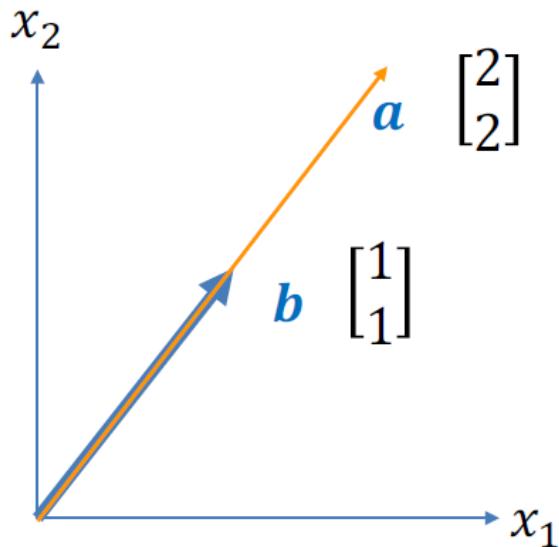
- A collection of d -vectors $\mathbf{x}_1, \dots, \mathbf{x}_m$ (with $m \geq 1$) is called **linearly independent** if it is not linearly dependent, which means that

$$\beta_1 \mathbf{x}_1 + \dots + \beta_m \mathbf{x}_m = 0$$

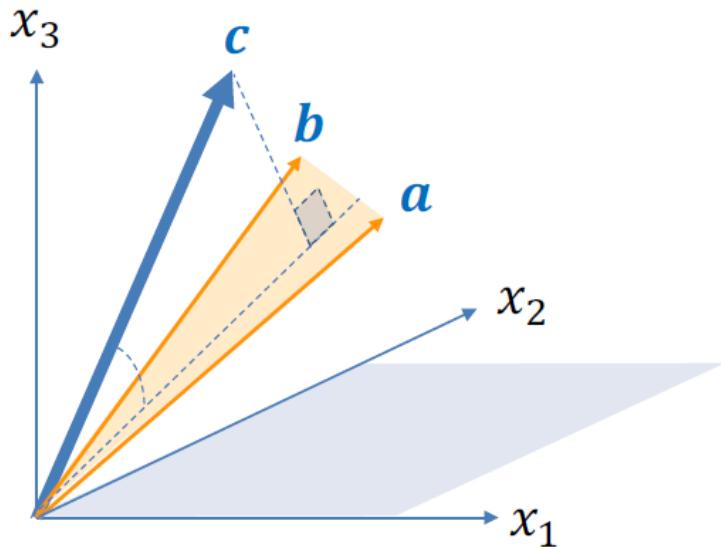
only holds for $\beta_1 = \dots = \beta_m = 0$.

Note: If all rows or columns of a square matrix \mathbf{X} are **linearly independent**, then \mathbf{X} is **invertible**.

Geometry of dependency and independency



$$\beta_1 \mathbf{a} + \beta_2 \mathbf{b} = 0$$



$$\beta_1 \mathbf{a} + \beta_2 \mathbf{b} \neq \beta_3 \mathbf{c}$$

Systems of Linear Equations

These equations can be written compactly in matrix-vector notation:

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

Where

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,d} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

Note:

- The data matrix $\mathbf{X} \in \mathcal{R}^{m \times d}$ and the target vector $\mathbf{y} \in \mathcal{R}^m$ are given
- The unknown vector of parameters $\mathbf{w} \in \mathcal{R}^d$ is to be learnt
- The rank(\mathbf{X}) corresponds to the maximal number of linearly independent columns/rows of \mathbf{X} .

- Learn the solution to $Xw = y$, to map X to y , hence leaning w
- The principled way for computing rank is to do Echelon Form
 - <https://stattrek.com/matrix-algebra/echelon-transform.aspx#MatrixA>
- For small-size matrices, however, the rank is in many cases easy to estimate
- What is the rank of

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

$$\begin{bmatrix} 1 & 1 \\ 100 & 100 \end{bmatrix}$$

$$\begin{bmatrix} 1 & -2 & 3 \\ 0 & -3 & 3 \\ 1 & 1 & 0 \end{bmatrix}$$

Causality

≡ GEA1000 Notes

- Causality, or causation is:
 - The influence by which one event or process (i.e., **cause**) contributes to another (i.e. **effect**),
 - The **cause** is partly responsible for the **effect**, and the **effect** is partly dependent on the **cause**
- Causality relates to an extremely very wide domain of subjects: philosophy, science, management, humanity.
- Causality research is extremely complex
 - Researcher can never be completely certain that there are **no other factors** influencing the causal relationship,
 - In most cases, we can only say "probably" causal.
- Partially or partly, cannot be completely sure that there is anything else that something is affecting the relationship

- (**Probable**) causal relations or non-causal?

- New web design implemented ? Web page traffic increased Yes
- Your height and weight ? Gets A in EE2211 No
- Uploaded new app store images ? Downloads increased by 2X Yes
- One works hard and attends lectures/tutorials ? Gets A in EE2211 Yes
- Your favorite color ? Your GPA in NUS No

- One popular way to causal data analysis is **Randomized Controlled Trial (RCT)**
 - A study design that randomly assigns participants into an experimental group or a control group.
 - As the study is conducted, the only expected difference between two groups is the outcome variable being studied.
- Example:
 - To decide whether smoking and lung cancer has a causal relation, we put participants into experimental group (people who smoke) and control group (people who don't smoke), and check whether they develop lung cancer eventually.
- RCT is sometimes infeasible to conduct, and also has moral issues.

Causality is very difficult to study, but it's possible

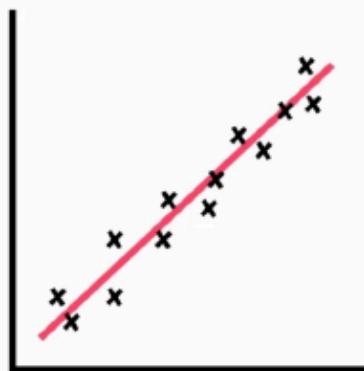
	Experiments	Observational Studies
Assignment	By researchers	Participants assign themselves
Randomisation	Preferable	Not possible
Ethical issues	Possible (if intervention may be harmful)	Unlikely
Confounders	Unlikely (if randomisation is done on enough participants)	Usually exist many
Possible to show causation	Yes (in the ideal case)	Very difficult
Able to show association	Yes	

Causality is a Statistical Relationship

- Decades of data show a clear causal relationship between smoking and cancer.
- If one smokes, it is a sure thing that his/her risk of cancer will increase.
- But it is not a sure thing that one will get cancer.
- The relationship is not deterministic.

Correlation vs Causality

- In statistics, **correlation** is any **statistical relationship**, whether causal or not, between two random variables.
- Correlations are useful because they can indicate a **predictive relationship** that can be exploited in practice.



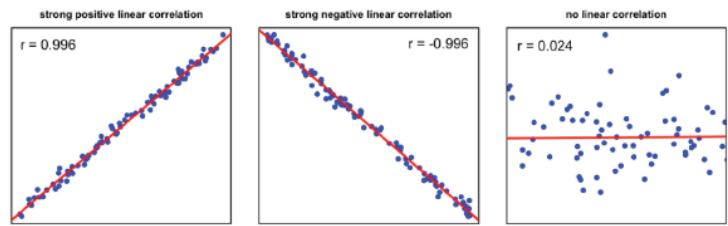
Positive
Correlation



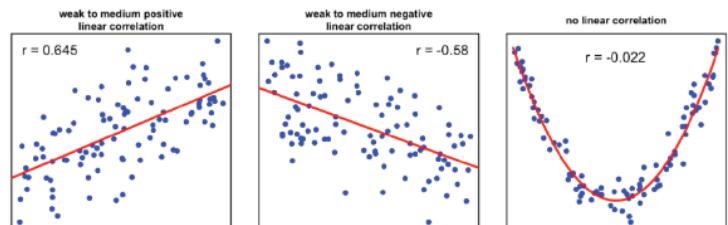
Negative
Correlation

- Linear correlation coefficient, r , which is also known as the Pearson Coefficient.

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{s_{xy}}{s_x s_y},$$



Strong linear relationship	$r > 0.9$
Medium linear relationship	$0.7 < r \leq 0.9$
Weak linear relationship	$0.5 < r \leq 0.7$
No or doubtful linear relationship	$0 < r \leq 0.5$

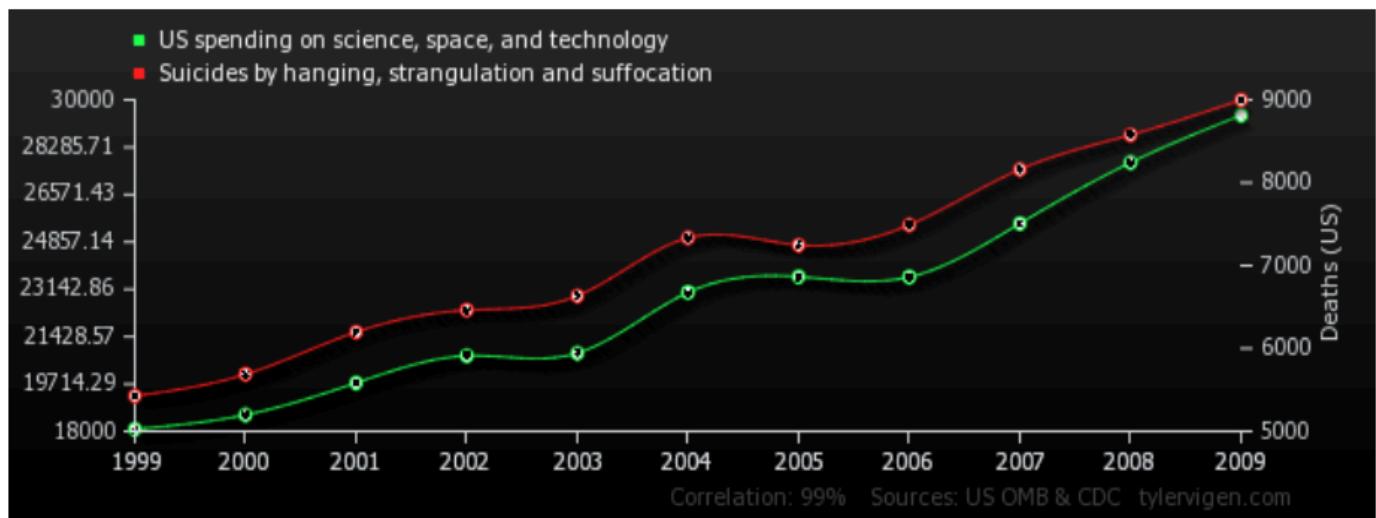


The same holds for negative values.

<https://www.geo.fu-berlin.de/en/v/soga/>

Correlation does not imply causation!

- Some great examples of correlations that can be calculated but are clearly not causally related appear at <http://tylervigen.com/>

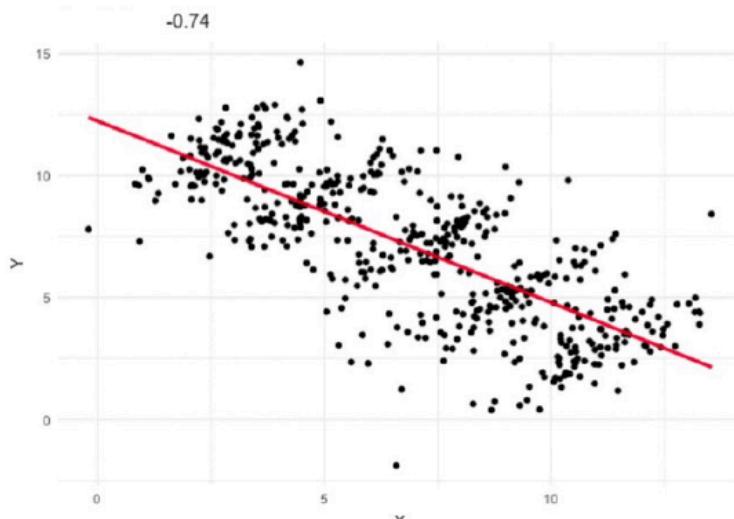
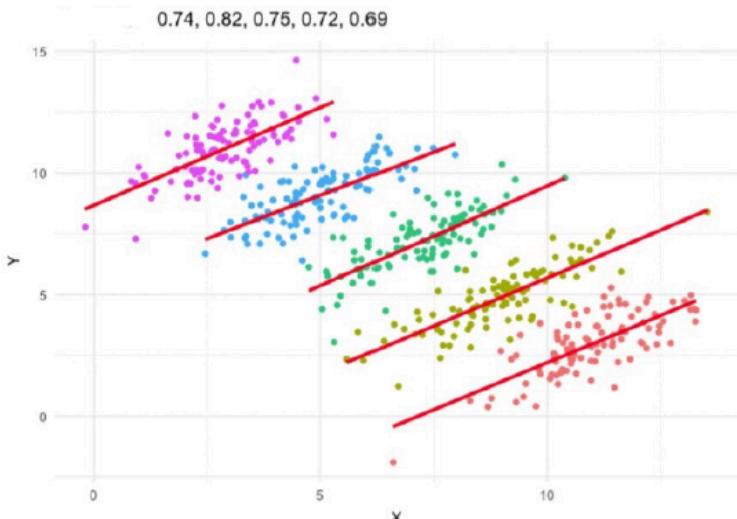


<https://tylervigen.com>

Simpson's paradox

Simpson's paradox: GEA1000 Notes

- **Simpson's paradox** is a phenomenon in probability and statistics, in which a trend appears in several different groups of data but disappears or reverses when these groups are combined.



The same set of samples!

Example

- Batting Average in professional baseball game
- Two well-known players, Derek Jeter and David Justice

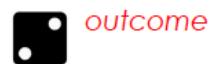
Batter \ Year	1995		1996		Combined	
Batter	# of wins	# of games	# of wins	# of games	# of wins	# of games
Derek Jeter	12/48	.250	183/582	.314	195/630	.310
David Justice	104/411	.253	45/140	.321	149/551	.270

Random Variable

GEA1000 Notes

EE2012 Notes

- We describe a *random experiment* by describing its procedure and observations of its *outcomes*.
- *Outcomes* are mutual exclusive in the sense that only one outcome occurs in a specific trial of the random experiment.
 - This also means an outcome is not decomposable.
 - All unique outcomes form a *sample space*.
- A subset of sample space S , denoted as A , is an *event* in a random experiment $A \subset S$, that is meaningful to an application.
 - Example of an event: faces with numbers no greater than 3
- Some books used $P(\cdot)$ and $p(\cdot)$ to distinguish between the probability of discrete random variable and the probability of continuous random variables respectively.
- We shall use $Pr(\cdot)$ for both the above cases



- A **random variable**, usually written as an *italic* capital letter, like X , is a variable whose possible values are numerical outcomes of a random event.
- There are two types of random variables: **discrete** and **continuous**.

Axioms of Probability

Assuming events $A \subseteq S$ and $B \subseteq S$, the probabilities of events related with and must satisfy,

1. $Pr(A) \geq 0$
2. $Pr(S) = 1$
3. If $A \cap B = \emptyset$, then $Pr(A \cup B) = Pr(A) + Pr(B)$
*otherwise, $Pr(A \cup B) = Pr(A) + Pr(B) - Pr(A \cap B)$

Two Basic Rules

- Sum Rule

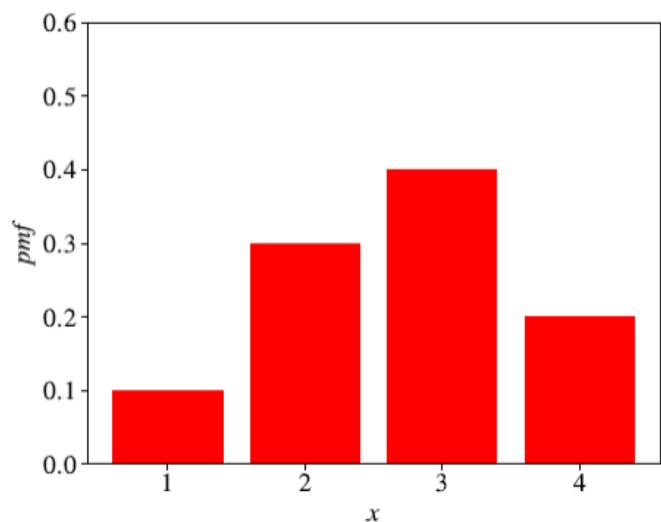
$$\Pr(X = x) = \sum_Y \Pr(X = x, Y = y_i)$$

- Product Rule

$$\Pr(X = x, Y = y) = \Pr(Y = y|X = x) P(X = x)$$

Discrete Random Variable

- A **discrete random variable (DRV)** takes on only a countable number of distinct values such as **red**, **orange**, **blue** or 1, 2, 3.
- The **probability distribution** of a discrete random variable is described by a list of probabilities associated with each of its possible values.
- This list of probabilities is called a **probability mass function (pmf)**.
 - Like a histogram, except that here the probabilities sum to 1



A probability mass function

- Let a **discrete** random variable X have k possible values $\{x_i\}_{i=1}^k$.
- The **expectation** of X denoted as $E(x)$ is given by,

$$E(x) \stackrel{\text{def}}{=} \sum_{i=1}^k [x_i \cdot \Pr(X = x_i)] \\ = x_1 \cdot \Pr(X = x_1) + x_2 \cdot \Pr(X = x_2) + \cdots + x_k \cdot \Pr(X = x_k)$$

where $\Pr(X = x_i)$ is the probability that X has the value x_i according to the **pmf**.

- The expectation of a random variable is also called the **mean, average** or **expected value** and is frequently denoted with the letter μ .
- Another important statistic is the **standard deviation**, defined as,

$$\sigma \stackrel{\text{def}}{=} \sqrt{E[(X - \mu)^2]} .$$

- **Variance**, denoted as σ^2 or $\text{var}(X)$, is defined as,

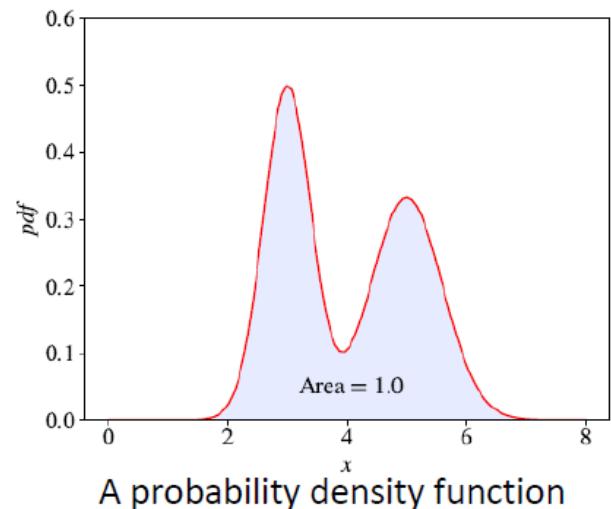
$$\sigma^2 = E[(X - \mu)^2]$$
- For a **discrete random variable**, the standard deviation is given by

$$\sigma = \sqrt{\Pr(X = x_1)(x_1 - \mu)^2 + \cdots + \Pr(X = x_k)(x_k - \mu)^2}$$

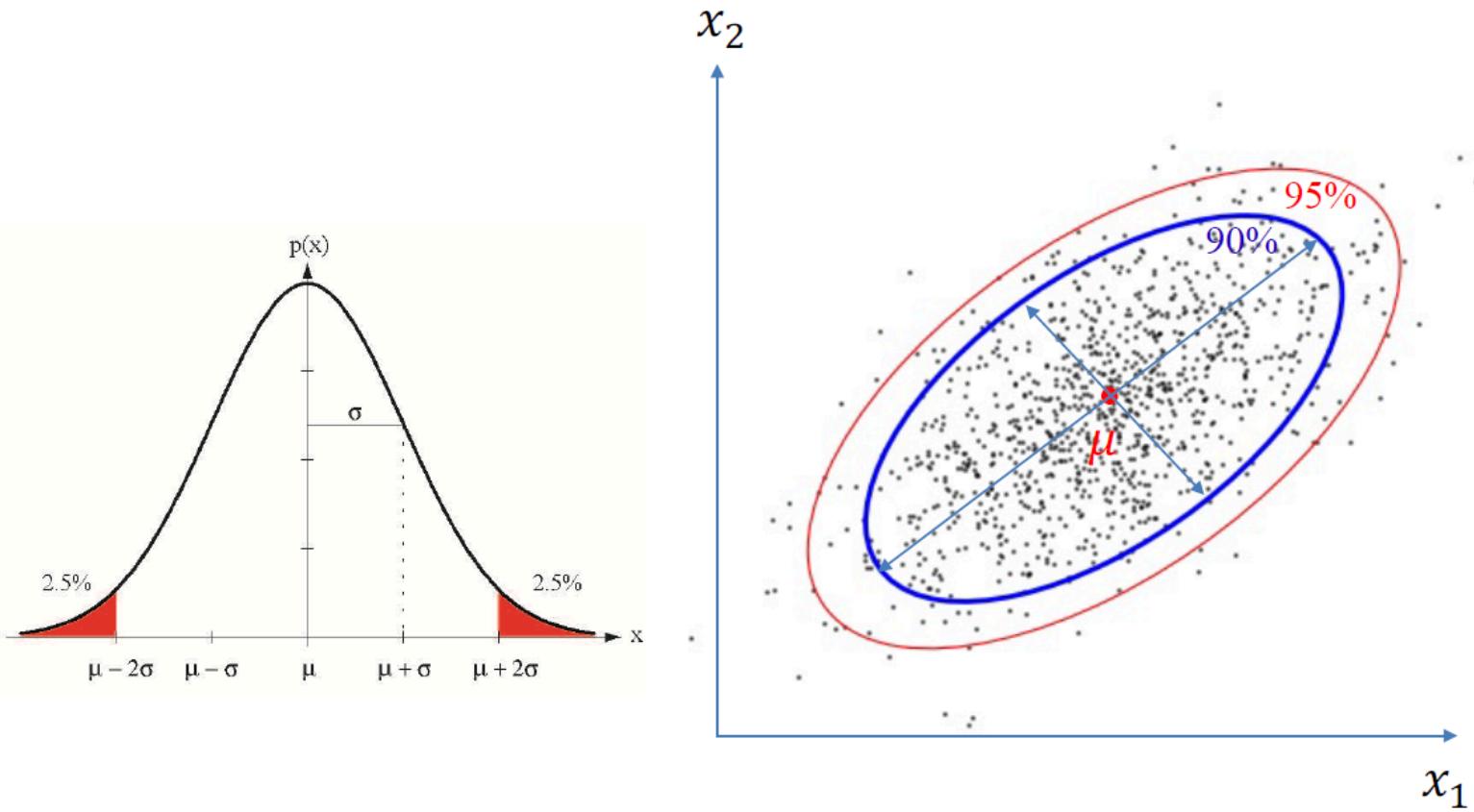
 where $\mu = E(X)$.

Continuous Random Variable

- A **continuous random variable (CRV)** takes an infinite number of possible values in some interval.
 - Examples include height, weight, and time.
 - The number of values of a continuous random variable X is infinite, the probability $\Pr(X = c)$ for any c is 0
 - Therefore, instead of the list of probabilities, the probability distribution of a CRV (a continuous probability distribution) is described by a **probability density function (pdf)**.
 - The pdf is a function whose range is **nonnegative** and the **area under the curve is equal to 1**.



- The **expectation** of a continuous random variable X is given by $E[x] \stackrel{\text{def}}{=} \int_R x f_X(x) dx$ where f_X is the **pdf** of the variable X and \int_R is the integral of function $x f_X$.
- The **variance** of a continuous random variable X is given by $\sigma^2 \stackrel{\text{def}}{=} \int_R (X - \mu)^2 f_X(x) dx$
- **Integral** is an equivalent of the **summation** over all values of the function when the function has a continuous domain.
- It equals the **area under the curve** of the function.
- The property of the pdf that the **area under its curve is 1** mathematically means that $\int_R f_X(x) dx = 1$



Example 1

- **Independent random variables**
- Consider tossing a fair coin twice, what is the probability of having (H,H)? Assuming a coin has two sides, H=head and T=Tail
 - $\Pr(x=H, y=H) = \Pr(x=H)\Pr(y=H) = (1/2)(1/2) = 1/4$

Example 2

- **Dependent random variables**
- Given 2 balls with different colors (**Red** and Black), what is the probability of first drawing B and then **R**? Assuming we are drawing the balls **without replacement**.
- The space of outcomes of taking two balls sequentially without replacement:
 $B-R$, $R-B$
 - Thus having **B-R** is $1/2$.
- Mathematically:
 - $\Pr(x=B, y=R) = \Pr(y=R | x=B) \Pr(x=B) = 1 \times (1/2) = 1/2$

Conditional Probability

Example 3

- **Dependent random variables**
- Given 3 balls with different colors (**R,G,B**), and we draw 2 balls. What is the probability of first having **B** and then **G**, if we draw **without replacement**?
- The space of outcomes of taking two balls sequentially without replacement:

R-G | G-B | B-R

R-B | G-R | B-G Thus, $\Pr(y=G, x=B) = 1/6$

- Mathematically:

$$\begin{aligned}\Pr(y=G, x=B) &= \Pr(y=G \mid x=B) \Pr(x=B) \\ &= (1/2) \times (1/3) \\ &= 1/6\end{aligned}$$

Bayes' Rule

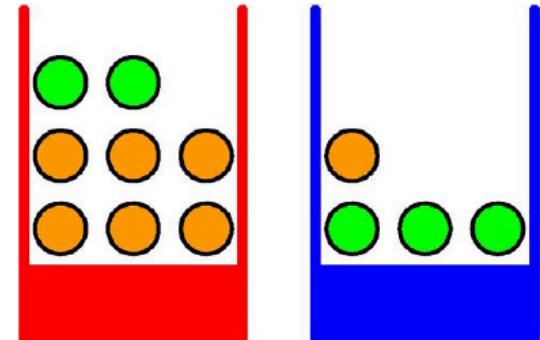
- The conditional probability $\Pr(Y = y \mid X = x)$ is the probability of the random variable **Y** to have a specific value **y**, given that another random variable **X** has a specific value of **x**.
- The **Bayes' Rule** (also known as the **Bayes' Theorem**):

$$\Pr(Y = y \mid X = x) = \frac{\text{likelihood} \quad \text{prior}}{\text{posterior} \quad \text{evidence}} \quad \Pr(X = x \mid Y = y) \Pr(Y = y)$$

Example

- Drawing a sample of fruit from a box
 - First pick a box, and then draw a sample of fruit from it
 - B: variable for Box, can be r (red) or b (blue)
 - F: variable for Fruit, can be o (orange) or a (apple)

- $\Pr(B=r)=0.4$ prior
- $\Pr(F=o | B=r)= 0.75$ likelihood
- $\Pr(F=o)= 0.45$ evidence



$$\begin{aligned} \Pr(B=r | F=o) &= \Pr(F=o | B=r) * \Pr(B=r) / \Pr(F=o) \\ &= 0.75 * 0.4 / 0.45 = 2/3 \quad \text{posterior} \end{aligned}$$

Summary and Practice Question

Suppose the random variable X has the following probability mass function (pmf) listed in the table below. k is unknown.

X	1	2	3	4	5
Pr[X]	0.1	0.05	0.05	0.6	k

What is the probability that X takes a value of odd numbers?

$$\begin{aligned} k &= 0.2 \\ \Pr(X = \text{odd}) &= 0.35 \end{aligned}$$

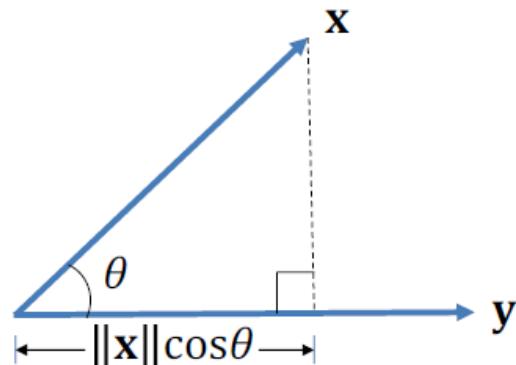
Lec4: Systems of linear equations

Vector and Matrix Operation

Dot Product = Inner Product

Dot Product or Inner Product of Vectors:

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



Geometric definition:

$$\mathbf{x} \cdot \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos\theta$$

where θ is the angle between \mathbf{x} and \mathbf{y} ,
and $\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}}$ is the Euclidean length of vector \mathbf{x}

E.g. $\mathbf{a} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$, $\mathbf{c} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \rightarrow \mathbf{a} \cdot \mathbf{c} = 2*1 + 3*0 = 2$

Orthogonal Projection of a Vector to a Subspace	
$S = \{u_1, u_2, \dots, u_k\}$	is a orthonormal basis for a subspace $V \subseteq \mathbb{R}^n$, if $u_i \in \mathbb{R}^n$, $u_i \perp u_j$ for $i \neq j$ and $\ u_i\ = 1$
requires the \Rightarrow	$w_p = (w \cdot u_1)u_1 + (w \cdot u_2)u_2 + \dots + (w \cdot u_k)u_k$ = projection of w to V
orthonormal basis	for any different orthonormal basis for V , projection of w to V is the same unique.
of V to find the projection	Example: $\{(1, 0)\}$ & $\{\frac{1}{\sqrt{2}}(1, 0), \frac{1}{\sqrt{2}}(0, 1)\}$ are orthonormal basis for the $x-y$ plane in \mathbb{R}^3
1) get by writing Gram-Schmidt	Let $w = (1, 0)$ then $w_p = ((1, 0) \cdot (1, 0))(1, 0) + ((1, 0) \cdot (0, 1))(0, 1)$ $w_p = (1, 0) \cdot \frac{1}{\sqrt{2}}(1, 0) + (1, 0) \cdot \frac{1}{\sqrt{2}}(0, 1) = (1, 0)$
2) least square approximation	$\hat{z} = (1, 0)$

- The dot product is the to find the projection of y to x or x to y

Inverse of Matrix

if A is $n \times n$, $(A I) \rightarrow (R I)$ $R \neq I$ $\Rightarrow R$ has zero row non pivot col. per steed pivots in non trivial sol	$AB = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}$
Inverse of Matrix	invertible } only for square invertible } singular } matrices
if A is non-singular \rightarrow	A square matrix A of order n is invertible if there exists a square matrix B of order n such that
if A is invertible	$AB = I_n = BA$. A square matrix is singular if it is not invertible.
$A\vec{x} = \vec{b}$	Uniqueness of inverse - if B & C are both inverses of a square matrix A , then $B=C$ $\hookrightarrow BA = I = AC \Rightarrow B = BI = BIA(C) = (BA)C = IC = C \Rightarrow B = C \#$
$\vec{x} = A^{-1}\vec{b}$	\hookrightarrow since inverse is unique, it is given a notation A^{-1} for inverse of A . $AA^{-1} = I = A^{-1}A$ Not $\frac{1}{A}$
\therefore a unique solution exists & system is consistent	for non-square matrix, if left inverse exists then there is no right inverse and vice versa. Right & Left inverse need not be unique. \hookrightarrow if both right & left inverse exist then A must be square & right = left inverse and are unique.
Proof	Cancellation Law for matrices - if A is invertible and $\begin{matrix} AB = AC \\ BA = CA \end{matrix}$ then $B = C$ $\because A^{-1}AB = A^{-1}AC \Rightarrow B = C \therefore$ this is only true if A^{-1} exists \therefore inverse of symmetric matrix is symmetric
i) $(A^{-1})^{-1} = A$	i) $(A^{-1})^{-1} = A$
ii) $I = I^T = (AA^{-1})^T$ $I = (A^{-1})^T(A^{-1})$	ii) for any non-zero $a \in \mathbb{R}$, (aA) is invertible with inverse $(aA)^{-1} = \frac{1}{a}A^{-1}$ $A = A^T$ $(A^{-1})^T = (A^T)^{-1}$
$(A^{-1})^T = (A^T)^{-1}$	iii) A^T is invertible with inverse $(A^T)^{-1} = (A^{-1})^T$, inverse of the transpose is the transpose of the inverse. $A = (A^{-1})^{-1}$ $= (A^{-1})^T$
$(AB)^{-1} = B^{-1}A^{-1}$	iv) $(AB)^{-1} = B^{-1}A^{-1}$, for $(AB)^{-1}$ to exist A, B must be invertible. $\Rightarrow (A_1A_2 \dots A_k)^{-1} = A_k^{-1} \dots A_2^{-1}A_1^{-1}$ $A^{-1} = (A^{-1})^T$
$(AB)C = I = C(AB)$	<u>Finding Inverse For 2×2</u> $\hookrightarrow C = AB \text{ is invertible}$ $\hookrightarrow AB \text{ is invertible}$ $\hookrightarrow A, B, C \text{ all invertible}$ $\hookrightarrow AB \text{ must be invertible}$ \hookrightarrow inverse X of A satisfies $AX = I$ $\hookrightarrow (A I) \xrightarrow{\text{REF}} (R B)$ $\therefore A^{-1} \text{ is symmetric}$
$ABC = I$ $A^{-1}ABC = A^{-1}I$ $B^{-1}A^{-1} = A^{-1}A^{-1}$ $C^{-1} = B^{-1}A^{-1}$	$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ \hookrightarrow can explain only elementary matrix
$A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$	$\det A = ad - bc$ A^{-1} exists if $ad - bc \neq 0$ \hookrightarrow if $R = I$ then $R = A^{-1}$ $\therefore R$ has a zero row & determinant of $A = 0$, A^{-1} doesn't exist.

Determinant computation $\det(A) = \sum_{j=1}^k (-1)^{i+j} a_{ij} M_{ij}$

Example: 3x3 matrix, use the first row ($i = 1$)

$$\begin{aligned}
 |A| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} \square & \square & \square \\ \square & e & f \\ \square & h & i \end{vmatrix} - b \begin{vmatrix} \square & \square & \square \\ d & \square & f \\ g & \square & i \end{vmatrix} + c \begin{vmatrix} \square & \square & \square \\ d & e & \square \\ g & h & \square \end{vmatrix} \\
 &= a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\
 &= a(ei - fh) - b(di - fg) + c(dh - eg)
 \end{aligned}$$

Consider a 3x3 matrix

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}. \quad \text{The minor of } a_{22} = \begin{vmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{vmatrix}$$

Its cofactor matrix is

$$\mathbf{C} = \begin{pmatrix} + \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} & - \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} & + \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix} \\ - \begin{vmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} \end{vmatrix} & + \begin{vmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{vmatrix} & - \begin{vmatrix} a_{11} & a_{12} \\ a_{31} & a_{32} \end{vmatrix} \\ + \begin{vmatrix} a_{12} & a_{13} \\ a_{22} & a_{23} \end{vmatrix} & - \begin{vmatrix} a_{11} & a_{13} \\ a_{21} & a_{23} \end{vmatrix} & + \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \end{pmatrix}.$$

$\text{adj}(\mathbf{A}) = \mathbf{C}^T$
 $\det(\mathbf{A}) = \sum_{j=1}^k a_{ij} C_{ij} = (-1)^{i+j} a_{ij} M_{ij}$
 $\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A})$

Systems of Linear Equations

Consider a system of m linear equations with d variables or unknowns w_1, \dots, w_d :

$$\begin{aligned} x_{1,1}w_1 + x_{1,2}w_2 + \cdots + x_{1,d}w_d &= y_1 \\ x_{2,1}w_1 + x_{2,2}w_2 + \cdots + x_{2,d}w_d &= y_2 \\ &\vdots \\ x_{m,1}w_1 + x_{m,2}w_2 + \cdots + x_{m,d}w_d &= y_m. \end{aligned}$$

- m is number of sample
- x is number of feature

These equations can be written compactly in matrix-vector notation:

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

Where

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \cdots & x_{m,d} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

Note:

- The data matrix $\mathbf{X} \in \mathcal{R}^{m \times d}$ and the target vector $\mathbf{y} \in \mathcal{R}^m$ are given
- The unknown vector of parameters $\mathbf{w} \in \mathcal{R}^d$ is to be learnt
- For number of feature d , is X , and number of sample m
- w is the classifier to map $R^d \rightarrow R$, which is y

Types of System

A set of linear equations can have no solution, one solution, or multiple solutions:

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

Where

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,d} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}.$$

\mathbf{X} is Square	Even-determined	$m = d$	Equal number of equations and unknowns
\mathbf{X} is Tall	Over-determined	$m > d$	More number of equations than unknowns
\mathbf{X} is Wide	Under-determined	$m < d$	Fewer number of equations than unknowns

\mathbf{X} is Square	Even-determined	$m = d$	One unique solution in general	$\hat{\mathbf{w}} = \mathbf{X}^{-1}\mathbf{y}$
\mathbf{X} is Tall	Over-determined	$m > d$	No exact solution in general; An approximated solution	$\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ Left-inverse
\mathbf{X} is Wide	Under-determined	$m < d$	Infinite number of solutions in general; Unique constrained solution	$\hat{\mathbf{w}} = \mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{y}$ Right-inverse

```

28 X = np.array([[1, 2], [0, 6], [1, 0]])
29 leftinverseX = helper.leftInverse(X)
30 print("leftinverseX =\n", leftinverseX)
31 rightinverseX = helper.rightInverse(X)
32 print("rightinverseX =\n", rightinverseX)
33
34 Xpadded = helper.paddingOfOnes(X)
35 leftinverseX = helper.leftInverse(Xpadded)
36 print("leftinverseX =\n", leftinverseX)
37 rightinverseX = helper.rightInverse(Xpadded)
38 print("rightinverseX =\n", rightinverseX)
39
40 print("inv(Xpadded) =\n", inv(Xpadded))
•

```

```

leftinverseX =
[[ 0.47368421 -0.15789474  0.52631579]
 [ 0.02631579  0.15789474 -0.02631579]]
rightinverseX =
[[ 0.          0.          1.        ]
 [-0.25       0.1875     0.25      ]]
leftinverseX =
[[[-3.   1.   3.  ]
 [ 3.  -1.  -2.  ]
 [ 0.5  0.   -0.5 ]]]
rightinverseX =
[[[-3.000000e+00  1.000000e+00  3.000000e+00]
 [ 3.000000e+00 -1.000000e+00 -2.000000e+00]
 [ 5.000000e-01  4.4408921e-16 -5.000000e-01]]]
inv(Xpadded) =
[[[-3.   1.   3.  ]
 [ 3.  -1.  -2.  ]
 [ 0.5  0.   -0.5 ]]]

```

- If square and invertible, left and right inverse give the same thing
- Left inverse for over determined is actually the least square approximation
 - There is no solution in column space of \mathbf{X} , hence find the closest one in $\text{Col}(\mathbf{X})$

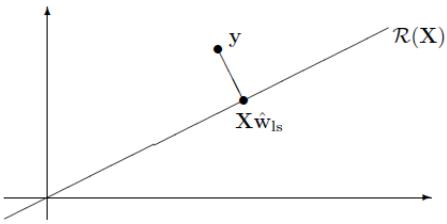


Figure 4.1: The least squares problem consists in finding $\hat{\mathbf{w}}_{ls}$, the point such that $\mathbf{X}\hat{\mathbf{w}}_{ls} \in \mathcal{R}(\mathbf{X})$ is closest to a given point \mathbf{y} .

- This means $\mathbf{X}^T\mathbf{X}$ has to have non zero determinant OR \mathbf{X} is full rank
- Right inverse for under determined is actually the least norm solution
 - There is infinitely many solution, hence find the one closest to origin, having the least norm

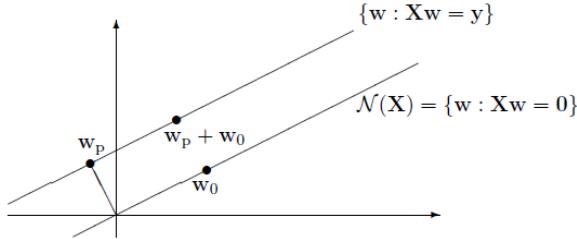


Figure 4.2: The least norm problem consists of finding the particular solution $\hat{\mathbf{w}}_p$ that minimizes the norm.

- All other solutions $\hat{\mathbf{w}}_p + \mathbf{w}_0$ where $\mathbf{w}_0 \in \mathcal{N}(\mathbf{X})$ have norms that are at least as large as that of $\hat{\mathbf{w}}_p$.

- This means $\mathbf{X}\mathbf{X}^T$ has to have non zero determinant OR \mathbf{X} is full rank

The matrix $\mathbf{X} \in \mathbb{R}^{m \times d}$ and vector $\mathbf{y} \in \mathbb{R}^m$ are given and $\mathbf{w} \in \mathbb{R}^d$ is to be found. As mentioned, if the matrix \mathbf{X} is square and full rank, \mathbf{X}^{-1} exists and so we can solve for \mathbf{w} by simple matrix inversion and multiplication $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$. However, most of the time in engineering, $m \neq d$ and more care is needed to discuss the existence and uniqueness of solutions to the linear system in (4.6). For this, we appeal to the Rouché-Capelli Theorem. We need the notion of the *augmented matrix*

$$\tilde{\mathbf{X}} = [\mathbf{X} \quad \mathbf{y}] \in \mathbb{R}^{m \times (d+1)}. \quad (4.8)$$

Note that the augmented matrix $\tilde{\mathbf{X}}$ has rank at least as large as that of \mathbf{X} , i.e., $\text{rank}(\mathbf{X}) \leq \text{rank}(\tilde{\mathbf{X}})$. This is because $\tilde{\mathbf{X}}$ has more columns than \mathbf{X} so the dimension of its column space must be as large as that of \mathbf{X} .

Theorem 4.1 (Rouché-Capelli Theorem). *For the linear system in (4.6), the following hold:*

(i) *The system in (4.6) admits a unique solution if and only if*

$$\text{rank}(\mathbf{X}) = \text{rank}(\tilde{\mathbf{X}}) = d; \quad (4.9)$$

(ii) *The system in (4.6) has no solution if and only if*

$$\text{rank}(\mathbf{X}) < \text{rank}(\tilde{\mathbf{X}}); \quad (4.10)$$

(iii) *The system in (4.6) has infinitely many solutions if and only if*

$$\text{rank}(\mathbf{X}) = \text{rank}(\tilde{\mathbf{X}}) < d. \quad (4.11)$$

Proof sketch (Only the \Leftarrow directions). For part (i), we note that the condition that $\text{rank}(\mathbf{X}) = \text{rank}(\tilde{\mathbf{X}})$ means that \mathbf{y} is in the column space of \mathbf{X} . This means that there exists $\{w_i\}_{i=1}^d \subset \mathbb{R}$ such that $\sum_{i=1}^d w_i \underline{x}_i = \mathbf{y}$ where the \underline{x}_i 's are the d columns of \mathbf{X} . Since $\text{rank}(\mathbf{X}) = d$, $\{\underline{x}_i\}_{i=1}^d$ span \mathbb{R}^d and the representation $\sum_{i=1}^d w_i \underline{x}_i = \mathbf{y}$ is unique (see discussion after Definition 4.3), which means there is a unique solution.

For part (ii), the condition that $\text{rank}(\mathbf{X}) < \text{rank}(\tilde{\mathbf{X}})$ means that \mathbf{y} is not in the column space of \mathbf{X} so there is no solution.

For part (iii), since $\text{rank}(\mathbf{X}) < d$, $\{\underline{x}_i\}_{i=1}^d$ do not span \mathbb{R}^d and the dimension of the nullspace of \mathbf{X} is non-zero. This means that if $\hat{\mathbf{w}}_p$ is a particular solution so is $\hat{\mathbf{w}}_p + \mathbf{w}_0$ where $\mathbf{w}_0 \in \mathcal{N}(\mathbf{X})$. Hence, are infinitely many solutions. \square

Square or even-determined system

1. Square or even-determined system: $m = d$

- Equal number of equations and unknowns, i.e., $\mathbf{X} \in \mathcal{R}^{d \times d}$
- **One unique solution** if \mathbf{X} is invertible or all rows/columns of \mathbf{X} are linearly independent
- If all rows or columns of \mathbf{X} are linearly independent, then \mathbf{X} is invertible.

Solution:

If \mathbf{X} is invertible (or $\mathbf{X}^{-1}\mathbf{X} = \mathbf{I}$), then pre-multiply both sides by \mathbf{X}^{-1}

$$\begin{aligned} \mathbf{X}^{-1}\mathbf{X}\mathbf{w} &= \mathbf{X}^{-1}\mathbf{y} \\ \Rightarrow \quad \hat{\mathbf{w}} &= \mathbf{X}^{-1}\mathbf{y} \end{aligned}$$

(Note: we use a *hat* on top of \mathbf{w} to indicate that it is a specific point in the space of \mathbf{w})

- Number of sample, $m =$ number of features, d

Example 1

$$\begin{aligned} w_1 + w_2 &= 4 \\ w_1 - 2w_2 &= 1 \end{aligned}$$

(1)

(2)

Two unknowns

Two equations

$$\mathbf{X} \quad \mathbf{w} \quad \mathbf{y}$$

$$\begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

$$\hat{\mathbf{w}} = \mathbf{X}^{-1} \mathbf{y}$$

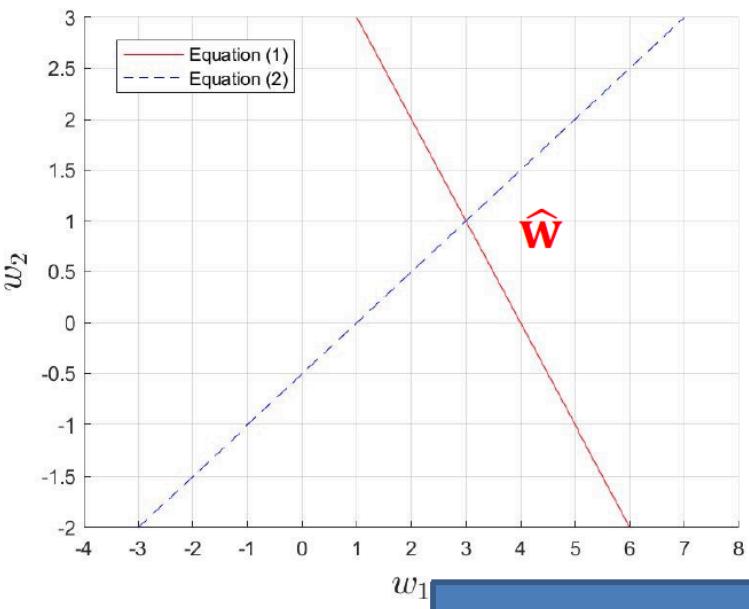
$$= \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix}^{-1} \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

$$= \frac{-1}{3} \begin{bmatrix} -2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

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

$$\text{adj}(\mathbf{A}) = \mathbf{C}^T = \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$\det(\mathbf{A}) = ad - bc$$



Python demo 3

2. Over-determined system: $m > d$

- More equations than unknowns
- X is non-square (tall) and hence not invertible
- Has no exact solution in general *
- An approximated solution is available using the left inverse

If the **left-inverse** of X exists such that $X^\dagger X = I$, then pre-multiply both sides by X^\dagger results in

$$\begin{aligned} X^\dagger X w &= X^\dagger y \\ \Rightarrow \hat{w} &= X^\dagger y \end{aligned}$$

Definition:

A matrix B that satisfies $B_{d \times m} A_{m \times d} = I$ is called a **left-inverse** of A .

The **left-inverse** of X : $X^\dagger = (X^T X)^{-1} X^T$ given $X^T X$ is invertible.

Note: * exception: when $\text{rank}(X) = \text{rank}([X, y])$, there is a solution.

- Number of sample, $m >$ number of features, d
- No exact solution in general, only approximation (least square solution)

Let $V \subseteq \mathbb{R}^n$ and $S = \{u_1, u_2, \dots, u_k\}$ be a basis for V . Then the orthogonal projection of any vector $w \in \mathbb{R}^n$ onto $V = A(A^T A)^{-1} A^T w$, where $A = (u_1, u_2, \dots, u_k) \quad n \times k \quad n \geq k$

Since S is basis, columns of A are linearly independent ($\text{rank}(A) = k$, A is full ranked) so $A^T A$ is invertible

$\therefore \hat{u} = (A^T A)^{-1} A^T w$ is the unique solution to $A^T A \hat{u} = A^T w$

\circ $\therefore \text{proj}_V w = A \hat{u} = A (A^T A)^{-1} A^T w$

\circ Except for $\text{rank}(X) = \text{rank}([X, y])$ which just means y is in the column space of X , hence there is an exact solution

- Lots of noisy measurement/sample to find something

<u>Case 1</u>	<u>Full Rank Matrices</u>	no of $R =$
Rank = no of col	let A be $m \times n$, $m \geq n$. Then A has full rank, $\text{rank}(A) = \text{column} = n$, iff RREF(A) has the form:	n
	$R = \begin{pmatrix} I_n \\ 0_{(m-n) \times n} \end{pmatrix}$	$\leftarrow n$ first columns, n non-zero rows, $m-n$ zero rows
1) $\text{rank}(A) = n$.		$\rightarrow R$ has n non-zero rows with n entries \Rightarrow spans \mathbb{R}^n
2) The rows of A spans \mathbb{R}^n , $\text{Row}(A) = \mathbb{R}^n$.		All first columns
3) The columns of A are linearly independent.		$\rightarrow A\vec{x} = \vec{0}$
4) Homogeneous system $A\vec{x} = \vec{0}$ has only the trivial solution, $\text{Null}(A) = \{\vec{0}\}$.		$\rightarrow ((A^T A)^{-1} A^T) A\vec{x} = ((A^T A)^{-1} A^T) \vec{0}$
5) $A^T A$ is an invertible matrix of order n .		$\rightarrow I_n = \vec{0}$
6) A has a left inverse. $= (A^T A)^{-1} A^T$		$\vec{x} = \vec{0}$
⑤ \rightarrow ⑥: $I = (A^T A)^{-1} A^T A = ((A^T A)^{-1} A^T) A \Rightarrow$ left inverse of $A = (A^T A)^{-1} A^T \neq$		
⑥ \rightarrow ⑤: Suppose C is a left inverse of A , then if $\vec{u} \in \mathbb{R}^n$ such that $A\vec{u} = \vec{0}$,		
$\vec{u} = I\vec{u} = C A\vec{u} = C\vec{0} = \vec{0} \Rightarrow$ nullspace of A is trivial		
since $\text{null}(A^T A) = \text{null}(A)$		\Rightarrow nullspace of $A^T A$ is trivial $\Rightarrow A^T A$ is invertible

- to many samples for little features, cannot find inverse
 - w cannot fit the constraint needed
 - so need to be an approximation

Example 2

$$w_1 + w_2 = 1 \quad (1)$$

$$w_1 - w_2 = 0 \quad (2)$$

$$w_1 = 2 \quad (3)$$

Two unknowns
Three equations

$$\begin{matrix} \mathbf{X} & \mathbf{w} & \mathbf{y} \end{matrix} \quad \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$$

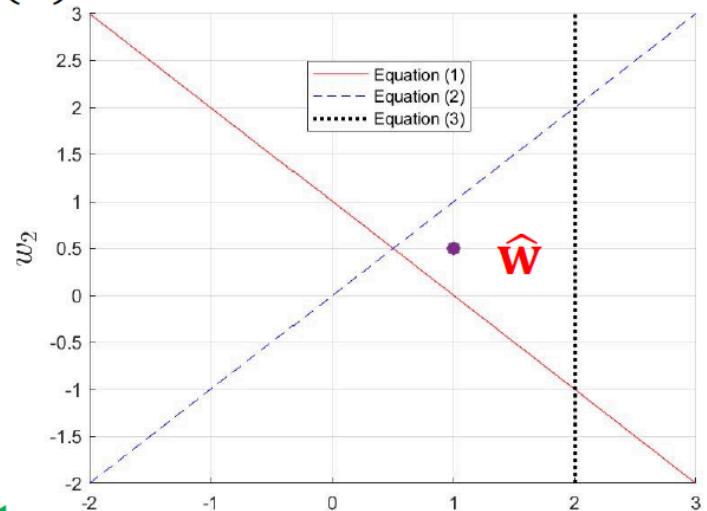
No exact solution

Approximated solution

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

$$= \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$

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



Python demo 4

- Consider the following over-determined system in which $m = 3$ and $d = 2$:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 \\ 4 & 3 \\ 5 & 6 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}. \quad (4.12)$$

The augmented matrix is

$$\tilde{\mathbf{X}} = [\mathbf{X} \quad \mathbf{y}] = \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 2 \\ 5 & 6 & 3 \end{bmatrix}. \quad (4.13)$$

In this case $\text{rank}(\mathbf{X}) = 2$ and $\text{rank}(\tilde{\mathbf{X}}) = 3$. This is case (ii) of the Rouché-Capelli Theorem and there is no solution. This is the usual case for over-determined systems. Note that in Python, you can find the rank of a matrix (2D array) \mathbf{A} using `np.linalg.matrix_rank(A)`.

Consider the following over-determined system in which $m = 3$ and $d = 2$:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 \\ 4 & 3 \\ 5 & 6 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 4 \\ 10 \\ 17 \end{bmatrix}. \quad (4.14)$$

In this case $\text{rank}(\mathbf{X}) = 2$ and $\text{rank}(\tilde{\mathbf{X}}) = 2$. This is case (i) of the Rouché-Capelli Theorem and there is a unique solution even though the system is over-determined. Note that \mathbf{y} is one times the first column of \mathbf{X} plus two times the second column of \mathbf{X} , so it is in the linear span of the columns of \mathbf{X} .

Consider the following over-determined system in which $m = 3$ and $d = 2$:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 \\ 4 & 2 \\ 6 & 3 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 8 \\ 16 \\ 24 \end{bmatrix}. \quad (4.15)$$

In this case $\text{rank}(\mathbf{X}) = 1$ and $\text{rank}(\tilde{\mathbf{X}}) = 1$ and both these ranks are less than $d = 2$. This is case (iii) of the Rouché-Capelli Theorem and there are infinitely many solutions even though the system is over-determined. Note that the three columns of $\tilde{\mathbf{X}}$ are collinear.

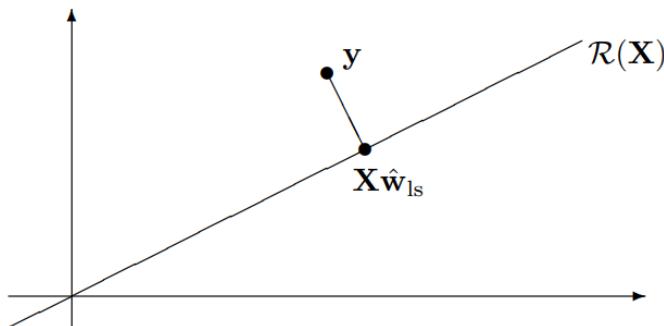


Figure 4.1: The least squares problem consists in finding $\hat{\mathbf{w}}_{ls}$, the point such that $\mathbf{X}\hat{\mathbf{w}}_{ls} \in \mathcal{R}(\mathbf{X})$ is closest to a given point \mathbf{y} .

3. Under-determined system: $m < d$

- More unknowns than equations
- Infinite number of solutions in general *

If the **right-inverse** of X exists such that $XX^\dagger = I$, then the d -vector $w = X^\dagger y$ (one of the infinite cases) satisfies the equation $Xw = y$, i.e.,

$$\begin{aligned} Xw &= y \quad \Rightarrow \quad XX^\dagger y = y \\ &\quad \Rightarrow \quad Iy = y \end{aligned}$$

Definition:

A matrix B that satisfies $A_{m \times d} B_{d \times m} = I$ is called a **right-inverse** of A . The **right-inverse** of X : $X^\dagger = X^T (XX^T)^{-1}$ given XX^T is invertible.

If X is right-invertible, we can find a unique constrained solution.

Note: * exception: no solution if the system is inconsistent $\text{rank}(X) < \text{rank}([X, y])$

- Number of sample, $m <$ number of features, d
 - Few patients but uses high resolution MRI
- Infinite number of solution in general
 - But can find unique constraint solution

A unique solution is yet possible by constraining the search using
 $w = X^T a$

If XX^T is invertible, let $w = X^T a$, then

$$\begin{aligned} X X^T a &= y \\ \Rightarrow \hat{a} &= (X X^T)^{-1} y \\ \Rightarrow \hat{w} &= X^T \hat{a} = X^T \underbrace{(X X^T)^{-1} y}_{X^\dagger} \end{aligned}$$

right-inverse

- By constraint w to be the form of $X^T a$, then from infinite solution to unique solution
 - This is finding the least norm solution
- Or exception when $\text{rank}(X) < \text{rank}([X, y])$, meaning y is linearly independent of the column space of X , so confirm no solution

Case 2

Rank = no of rows Let A be $m \times n$, $n \geq m$. Then A has a full rank, $\text{rank}(A) = \text{no of rows} = m$, iff $R = \text{RREF}$ has the form:

$$R = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & \cdots & 0 \end{pmatrix} \leftarrow m \text{ non-zero rows, } m \text{ pivot column, } n-m \text{ non-pivot column} \quad \text{full Rank} = 3$$

* columns before the 1st pivot column can be zero column also $\begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$

1) $\text{rank}(A) = m$.

$\Rightarrow R$ has m pivot column with m entries \Rightarrow spans \mathbb{R}^m

2) The columns of A spans \mathbb{R}^m , $\text{Col}(A) = \mathbb{R}^m$

All non-zero

3) The rows of A are linearly independent.

4) The linear system $Ax = b$ is consistent for every $b \in \mathbb{R}^m$.

5) AA^T is an invertible matrix of order m .

6) A has a right inverse. $= A^T(AA^T)^{-1}$

③ \rightarrow ④: $AA^T(AA^T)^{-1} = I \Rightarrow A^T(AA^T)^{-1}$ is the right inverse of A . For any b , $A(A^T(AA^T)^{-1})^T b = b$, $1_A = A^T(AA^T)^{-1} b$ \therefore consistent &

④ \rightarrow ③: Suppose $Ax = b$ is consistent for every $b \in \mathbb{R}^m$. Let b_i be a soln to $Ax = e_i$, e_i = i -th vector in the standard basis, for $i=1, \dots, m$.

Let $B = (b_1 \ b_2 \ \dots \ b_m)$, then $AB = A(b_1 \ b_2 \ \dots \ b_m) = e_1 \ e_2 \ \dots \ e_m = I_m \Rightarrow B$ is a right inverse of A #

To satisfy both case 1 & 2, A must be a square matrix with $R = I$. \rightarrow RREF (full rank square matrix) = I

Case 1 & 2 are mutually exclusive for non-square matrix \Rightarrow can have at most a left or a right inverse not both.

↳ if $r > c$, then can't have right inverse; if A has right inverse, then $\text{rank}(A) = r$ but $r > c \therefore \text{rank}(A) \neq r$

- Too many features, too little samples
 - o too many possible solution that is possible and most are not good

Example 3 $w_1 + 2w_2 + 3w_3 = 2 \quad (1)$

$$w_1 - 2w_2 + 3w_3 = 1 \quad (2)$$

Three unknowns

Two equations

X **w** **y**

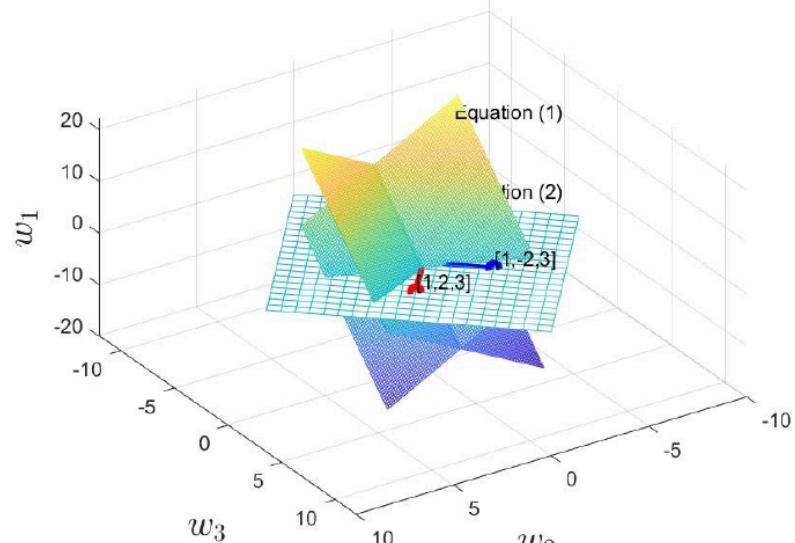
$$\begin{bmatrix} 1 & 2 & 3 \\ 1 & -2 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Infinitely many solutions along the intersection line

Here $\mathbf{X}\mathbf{X}^T$ is invertible

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

$$= \begin{bmatrix} 1 & 1 \\ 2 & -2 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 14 & 6 \\ 6 & 14 \end{bmatrix}^{-1} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.15 \\ 0.25 \\ 0.45 \end{bmatrix}$$



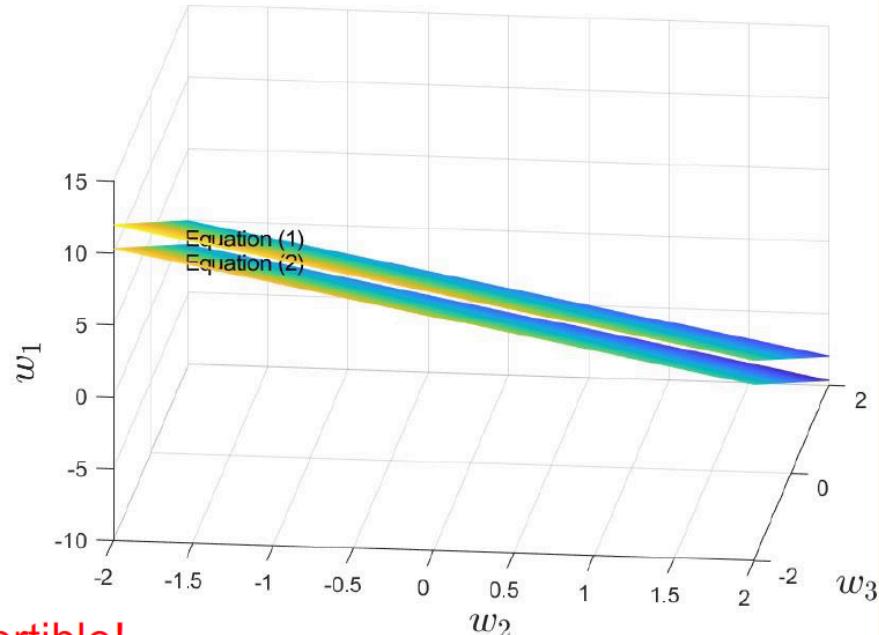
Constrained solution

- The two original planes, intersect in a line, infinite solution
- Putting in the constraint, introduce one more plane and hence only one point of intersection

Example 4 $w_1 + 2w_2 + 3w_3 = 2 \quad (1)$ Three unknowns
 $3w_1 + 6w_2 + 9w_3 = 1 \quad (2)$ Two equations

$$\mathbf{X} \quad \mathbf{w} \quad \mathbf{y}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



Both \mathbf{XX}^T and $\mathbf{X}^T\mathbf{X}$ are not invertible!

There is no solution for the system

Let A be $m \times n$, $n \geq m$. Then A has a full rank, $\text{rank}(A) = \text{no of rows} = m$, iff $A = \text{REF}$ has the form:

$$A = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & \dots & 0 & \dots & 0 & \dots \\ 0 & 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 & \dots \\ 0 & 0 & 1 & \dots & 0 & \dots & 0 & \dots & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & \dots & 0 & \dots & 0 & \dots \\ 0 & 0 & 0 & \dots & 0 & \dots & 1 & \dots & 0 & \dots \\ 0 & 0 & 0 & \dots & 0 & \dots & 0 & \dots & 1 & \dots \end{pmatrix} \leftarrow m \text{ non-zero rows, } m \text{ pivot column, } n-m \text{ non-pivot column}$$

* columns before the 1st pivot column can be zero column also

- \mathbf{X} is not full rank so cannot find right inverse

Consider the following under-determined system in which $m = 2$ and $d = 3$:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 & 3 \\ 4 & 2 & 5 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 10 \\ 7 \end{bmatrix}. \quad (4.16)$$

In this case, $\text{rank}(\mathbf{X}) = 2$ and $\text{rank}(\tilde{\mathbf{X}}) = 2$ but $d = 3$. This is case (iii) of the Rouché-Capelli Theorem and there are infinitely many solutions. This is the usual case for under-determined systems.

Consider the following under-determined system in which $m = 2$ and $d = 3$:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 & 3 \\ 4 & 2 & 6 \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}. \quad (4.17)$$

In this case, $\text{rank}(\mathbf{X}) = 1$ and $\text{rank}(\tilde{\mathbf{X}}) = 2$ because $\mathbf{y} \notin \mathcal{R}(\mathbf{X})$. This is case (ii) of the Rouché-Capelli Theorem and there is no solution. Note that \mathbf{y} boosts the rank of \mathbf{X} by 1 in the augmented matrix $\tilde{\mathbf{X}}$, i.e., \mathbf{y} is not in the column space of \mathbf{X} , which is the ray $\{[t, 2t]^\top : t \in \mathbb{R}\}$.

Least Norm Solution

4.4 Least Norm Solution for $m < d$

Now we consider the case in which \mathbf{X} is wide ($m < d$) and full rank. This means that $\text{rank}(\mathbf{X}) = m$; equivalently, all rows are linearly independent (full row rank). This *under-determined* situation also occurs a lot in engineering.

Example 4.3. For example, in control engineering (a field of study within mechanical and electrical engineering), one often considers the following discrete-time state-space system (e.g., describing the dynamics of a robot operating over a quantized time interval):

$$v_{i+1} = av_i + bw_i, \quad i = 0, 1, \dots, d-1, \quad (4.21)$$

where v_i is the state of the system at time i and w_i is our control. We assume the system starts at the origin $v_0 = 0$. We desire to design the w_i 's such that the terminal state $v_d = y$ (for some given y) while minimizing the cost of the control $\sum_{i=0}^{d-1} w_i^2$. After some algebra, this can be rewritten as

$$\text{minimize } \|\mathbf{w}\|^2 \quad \text{subject to} \quad y = [b \ ab \ a^2b \ \dots \ a^{d-1}b] \begin{bmatrix} w_{d-1} \\ w_{d-2} \\ \vdots \\ w_0 \end{bmatrix}. \quad (4.22)$$

This is exactly an under-determined problem if we make the identifications $\mathbf{X} = [b \ ab \ a^2b \ \dots \ a^{d-1}b]$ (for obvious reasons, this is called the d -step reachability matrix) and $\mathbf{w} = [w_{d-1} \ w_{d-2} \ \dots \ w_0]^\top$ and y is scalar (i.e., $m = 1$).

We return to the equation $\mathbf{X}\mathbf{w} = \mathbf{y}$ in which $m < d$ and the matrix \mathbf{X} has full row rank. It is clear that $(\mathbf{X}\mathbf{X}^\top)^{-1}$ exists and

$$\hat{\mathbf{w}}_p = \mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y} \quad (4.23)$$

is a solution to the equation (check!). The subscript p indicates that this is a *particular* solution to the linear system. From the usual case of the Rouché-Capelli Theorem (case (iii)), we know that there are infinitely many solutions. Where are these infinitely many solutions? Let $\mathbf{w}_0 \in \mathcal{N}(\mathbf{X})$ be any vector in the nullspace of \mathbf{X} . Note that the nullspace has positive dimension because $\text{rank}(\mathbf{X}) = m < d$ so \mathbf{w}_0 can be chosen to be a non-zero vector. Then $\hat{\mathbf{w}}_p + \mathbf{w}_0$ is also a solution to (4.6) (check!). In the following, we argue that among all the solutions, $\hat{\mathbf{w}}_p$ has a special place in our hearts because it is the *least norm solution* to (4.6).

Suppose that \mathbf{w} is any solution to $\mathbf{X}\mathbf{w} = \mathbf{y}$. Then since $\hat{\mathbf{w}}_p$ is also a solution, we have $\mathbf{X}(\mathbf{w} - \hat{\mathbf{w}}_p) = \mathbf{0}$ and

$$(\mathbf{w} - \hat{\mathbf{w}}_p)^\top \hat{\mathbf{w}}_p \stackrel{(4.23)}{=} (\mathbf{w} - \hat{\mathbf{w}}_p)^\top \mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y} = (\mathbf{X}(\mathbf{w} - \hat{\mathbf{w}}_p))^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y} = 0. \quad (4.24)$$

This means that $\mathbf{w} - \hat{\mathbf{w}}_p$ is orthogonal to $\hat{\mathbf{w}}_p$. By the Pythagorean theorem,

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}_p) + \hat{\mathbf{w}}_p\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}_p\|^2 + \|\hat{\mathbf{w}}_p\|^2 \geq \|\hat{\mathbf{w}}_p\|^2, \quad (4.25)$$

which shows that $\hat{\mathbf{w}}_p$ is the least norm solution to $\mathbf{X}\mathbf{w} = \mathbf{y}$, so we also denote this as

$$\hat{\mathbf{w}}_{ln} = \mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y} \quad (4.26)$$

See the geometry of the least norm problem in Fig. 4.2.

Finally, we say a few words about the matrix $\mathbf{X}^\dagger = \mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top)^{-1}$. This matrix is called the *pseudo-inverse* of the full-rank wide matrix \mathbf{X} . It is also known as the *right-inverse* of \mathbf{X} (why?). The right-inverse of a matrix with full row rank can be implemented in Python as `np.linalg.pinv(X)`.

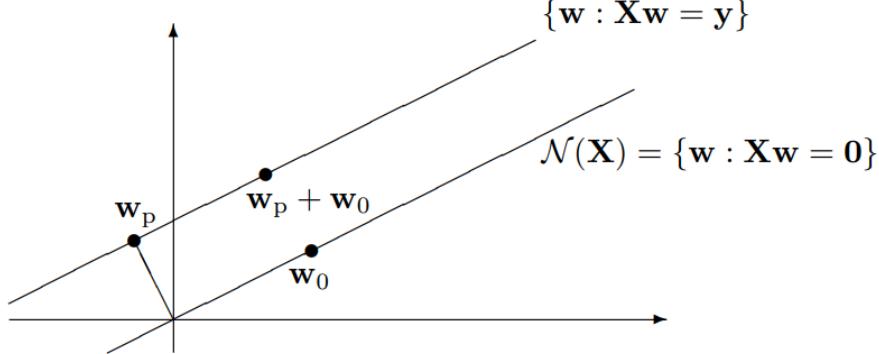


Figure 4.2: The least norm problem consists of finding the particular solution $\hat{\mathbf{w}}_p$ that minimizes the norm. All other solutions $\hat{\mathbf{w}}_p + \mathbf{w}_0$ where $\mathbf{w}_0 \in \mathcal{N}(\mathbf{X})$ have norms that are at least as large as that of $\hat{\mathbf{w}}_p$.

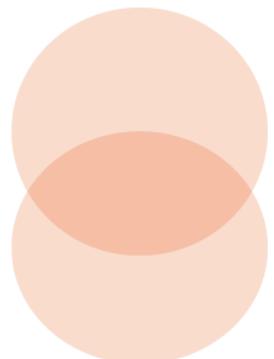
Set

- A **set** is an **unordered** collection of unique elements
 - Denoted as a calligraphic capital character e.g., $\mathcal{S}, \mathcal{R}, \mathcal{N}$ etc
 - When an element x belongs to a set S , we write $x \in S$
- A set of numbers can be **finite** - include a fixed amount of values
 - Denoted using accolades, e.g. $\{1, 3, 18, 23, 235\}$ or $\{x_1, x_2, x_3, x_4, \dots, x_d\}$
- A set can be **infinite** and include all values in some interval
 - If a set of real numbers includes all values between a and b , **including a and b** , it is denoted using square brackets as $[a, b]$
 - If the set **does not include the values a and b** , it is denoted using parentheses as (a, b)
- Examples:
 - The special set denoted by \mathcal{R} includes all real numbers from minus infinity to plus infinity
 - The set $[0, 1]$ includes values like 0, 0.0001, 0.25, 0.9995, and 1.0

- **Intersection** of two sets:

$$\mathcal{S}_3 \leftarrow \mathcal{S}_1 \cap \mathcal{S}_2$$

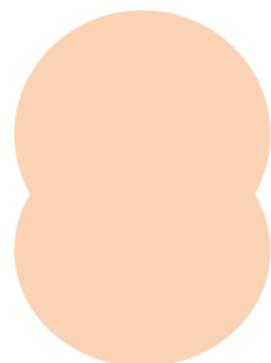
Example: $\{1, 3, 5, 8\} \cap \{1, 8, 4\} = \{1, 8\}$



- **Union** of two sets:

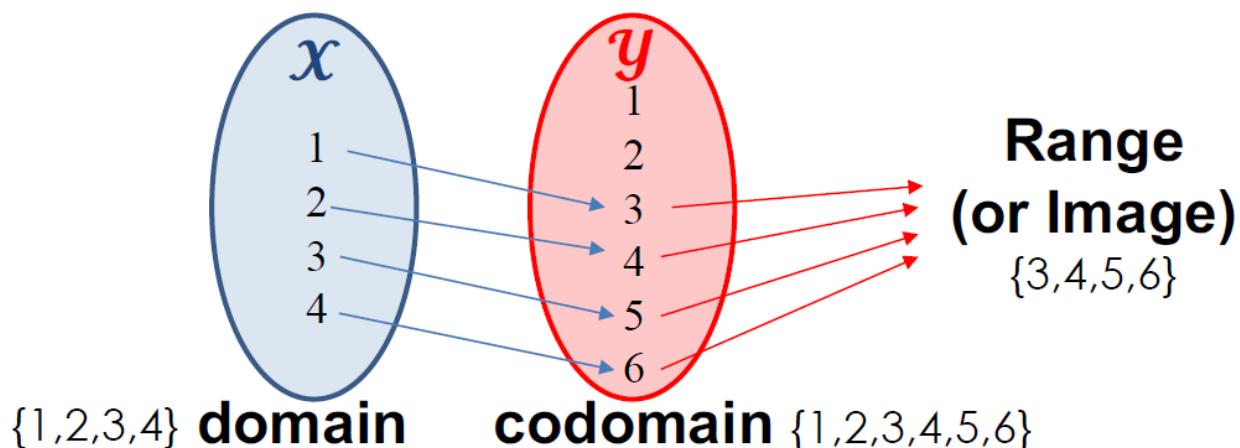
$$\mathcal{S}_3 \leftarrow \mathcal{S}_1 \cup \mathcal{S}_2$$

Example: $\{1, 3, 5, 8\} \cup \{1, 8, 4\} = \{1, 3, 4, 5, 8\}$



Function

- 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
- If the function is called f , this relation is denoted $y = f(x)$
 - The element x is the **argument** or **input** of the function
 - y is the value of the function or the **output**
- The symbol used for representing the input is the **variable** of the function
 - $f(x)$ f is a function of the variable x ; $f(x, w)$ f is a function of the variable x and w



- A **scalar function** can have vector argument
 - E.g. $y = f(\mathbf{x}) = x_1 + x_2 + 2x_3$
- A **vector function**, denoted as $\mathbf{y} = \mathbf{f}(\mathbf{x})$ is a function that returns a vector \mathbf{y}
 - Input argument can be a **vector** $\mathbf{y} = \mathbf{f}(\mathbf{x})$ or a **scalar** $y = f(x)$
 - E.g. $\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} -x_1 \\ x_2 \end{bmatrix}$
 - E.g. $\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} -2x_1 \\ 3x_1 \end{bmatrix}$

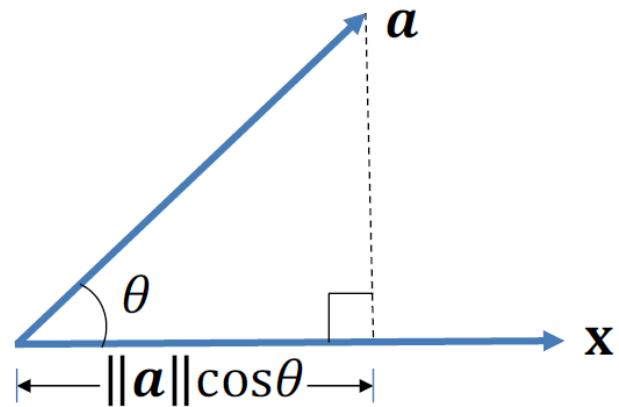
- The notation $f: \mathcal{R}^d \rightarrow \mathcal{R}$ means that f is a function that maps **real** d -vectors to **real** numbers
 - i.e., f is a scalar-valued function of d -vectors
- If \mathbf{x} is a d -vector argument, then $f(\mathbf{x})$ denotes the value of the function f at \mathbf{x}
 - i.e., $f(\mathbf{x}) = f(x_1, x_2, \dots, x_d)$, $\mathbf{x} \in \mathcal{R}^d$, $f(\mathbf{x}) \in \mathcal{R}$
- Example: we can define a function $f: \mathcal{R}^4 \rightarrow \mathcal{R}$ by

$$f(\mathbf{x}) = x_1 + x_2 - x_4^2$$

The inner product function

- Suppose \mathbf{a} is a d -vector. We can define a scalar valued function f of d -vectors, given by

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \dots + a_d x_d \quad (1)$$
 for any d -vector \mathbf{x}
- The inner product of its d -vector argument \mathbf{x} with some (fixed) d -vector \mathbf{a}
- We can also think of f as forming a **weighted sum** of the elements of \mathbf{x} ; the elements of \mathbf{a} give the weights



- Weighted sum, \mathbf{a} can be called weight or coefficient
- \mathbf{a} is to be found to map \mathbf{x} to a scalar

Linearity and Superposition

A function $f: \mathcal{R}^d \rightarrow \mathcal{R}$ is **linear** if it satisfies the following two properties:

- **Homogeneity**
 - For any d -vector \mathbf{x} and any scalar α , $f(\alpha\mathbf{x}) = \alpha f(\mathbf{x})$
 - **Scaling** the (vector) argument is the same as scaling the function value
- **Additivity**
 - For any d -vectors \mathbf{x} and \mathbf{y} , $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$
 - **Adding** (vector) arguments is the same as adding the function values

Superposition and linearity

- The inner product function $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ defined in equation (1) (slide 42) satisfies the property

$$\begin{aligned}f(\alpha\mathbf{x} + \beta\mathbf{y}) &= \mathbf{a}^T(\alpha\mathbf{x} + \beta\mathbf{y}) \\&= \mathbf{a}^T(\alpha\mathbf{x}) + \mathbf{a}^T(\beta\mathbf{y}) \\&= \alpha(\mathbf{a}^T\mathbf{x}) + \beta(\mathbf{a}^T\mathbf{y}) \\&= \alpha f(\mathbf{x}) + \beta f(\mathbf{y})\end{aligned}$$

for all d -vectors \mathbf{x} , \mathbf{y} , and all scalars α , β .

- This property is called **superposition**, which consists of **homogeneity** and **additivity**
- A function that satisfies the superposition property is called **linear**
- If a function f is **linear**, superposition extends to linear combinations of any number of vectors:

$$f(\alpha_1\mathbf{x}_1 + \cdots + \alpha_k\mathbf{x}_k) = \alpha_1 f(\mathbf{x}_1) + \cdots + \alpha_k f(\mathbf{x}_k)$$

for any d vectors $\mathbf{x}_1 + \cdots + \mathbf{x}_k$, and any scalars $\alpha_1 + \cdots + \alpha_k$.

Affine Function

A linear function plus a constant is called an affine function

A linear function $f: \mathcal{R}^d \rightarrow \mathcal{R}$ is **affine** if and only if it can be expressed as $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b$ for some d -vector \mathbf{a} and scalar b , which is called the **offset (or bias)**

Example:

$$f(\mathbf{x}) = 2.3 - 2x_1 + 1.3x_2 - x_3$$

This function is affine, with $b = 2.3$, $\mathbf{a}^T = [-2, 1.3, -1]$.

- Linear function plus DC offset constant get non linear function, get affine function

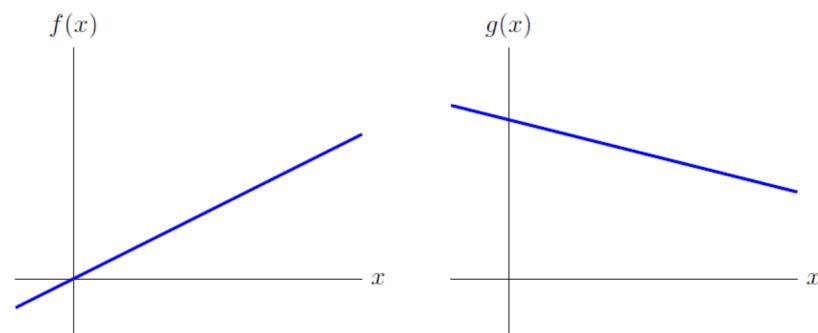


Figure 2.1 Left. The function f is linear. Right. The function g is affine, but not linear.

Class Quiz

100 elderly participated in a research study on brain imaging and cognitive decline in ageing. Each subject has 10 brain measures. If we would like to use these measures to predict cognition, this is an over-determined system.

Yes, overdetermined as no of sample > no of feature

If each of the 100 older adults has 100 brain and body measures, this is an even-determined system. We can usually find a unique solution.

Yes, even determine as no of sample = no of feature

If only a subset of 40 older adults have completed 100 brain and body measures, this is an under-determined system, i.e. more equations than unknowns. To predict cognition, we need to use right inverse.

No, under determined as no of sample < no of feature, find right inverse, but it's less equations than unknown

Summary

Linear Functions

A function $f: \mathcal{R}^d \rightarrow \mathcal{R}$ is **linear** if it satisfies the following **two properties**:

- **Homogeneity** $f(\alpha \mathbf{x}) = \alpha f(\mathbf{x})$ **Scaling**
- **Additivity** $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$ **Adding**

Inner product function

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = a_1 x_1 + a_2 x_2 + \cdots + a_d x_d$$

Affine function

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b \quad \text{scalar } b \text{ is called the offset (or bias)}$$

- Operations on Vectors and Matrices
 - Dot-product, matrix inverse
- Systems of Linear Equations $\mathbf{Xw} = \mathbf{y}$
 - Matrix-vector notation, linear dependency, invertible
 - Even-, over-, under-determined linear systems
- Set and Functions

Assignment 1 (week 7 Wed)
Tutorial 4

X is Square	Even-determined	$m = d$	One unique solution in general	$\hat{\mathbf{w}} = \mathbf{X}^{-1} \mathbf{y}$
X is Tall	Over-determined	$m > d$	No exact solution in general; An approximated solution	$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ Left-inverse
X is Wide	Under-determined	$m < d$	Infinite number of solutions in general; Unique constrained solution	$\hat{\mathbf{w}} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ Right-inverse

- Scalar and vector functions
- Inner product function
- Linear and affine functions

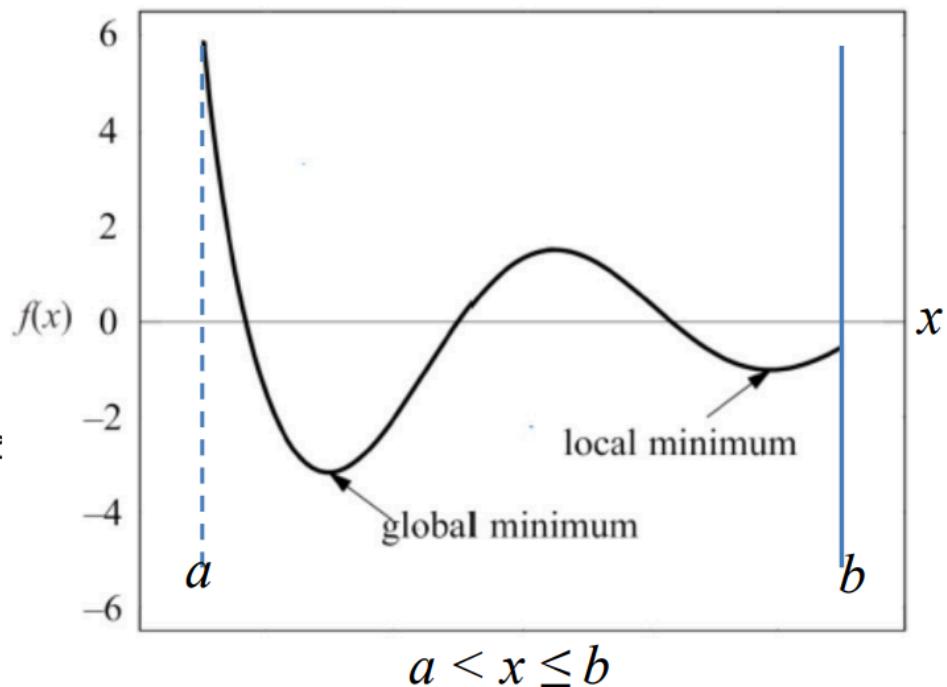
python package *numpy*
Inverse: *numpy.linalg.inv(X)*
Transpose: *X.T*

Lec5: Least Squares and Linear Regression

Functions: Maximum and Minimum

- $f(x)$ has a **local minimum** at $x = c$ if $f(x) \geq f(c)$ for every x in some open interval around $x = c$
- $f(x)$ has a **global minimum** at $x = c$ if $f(x) \geq f(c)$ for all x in the domain of f

A local and a global minima of a function



Note: An **interval** is a set of real numbers with the property that any number that lies between two numbers in the set is also included in the set.

An **open interval** does not include its endpoints and is denoted using parentheses. E.g. $(0, 1)$ means “all numbers greater than 0 and less than 1”.

- optimisation problem

Max and Arg Max

- Given a set of values $\mathcal{A} = \{a_1, a_2, \dots, a_m\}$,
- The operator $\max_{a \in \mathcal{A}} f(a)$ returns the highest value $f(a)$ for all elements in the set \mathcal{A}
- The operator $\arg \max_{a \in \mathcal{A}} f(a)$ returns the element of the set \mathcal{A} that maximizes $f(a)$
- When the set is **implicit** or **infinite**, we can write

$$\max_a f(a) \quad \text{or} \quad \arg \max_a f(a)$$

E.g. $f(a) = 3a$, $a \in [0,1]$ $\rightarrow \max_a f(a) = 3$ and $\arg \max_a f(a) = 1$

Min and Arg Min operate in a similar manner

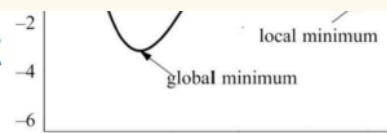
Note: **arg max** returns a value from the **domain** of the function and **max** returns from the **range (codomain)** of the function.

- $f(a)$ is what we want to optimise

Derivative and Gradient

Optimise in terms of cost functions, so need to find the maximum or minimum of the cost functions

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



- If the derivative is a constant value, e.g. 5 or -3
 - The function f grows (or decreases) constantly at any point x of its domain
- When the derivative f' is a function
 - If f' is positive at some x , then the function f grows at this point
 - If f' is negative at some x , then the function f decreases at this point
 - The derivative of zero at x means that the function's slope at x is horizontal (e.g. maximum or minimum points)
- The process of finding a derivative is called **differentiation**.
- Gradient** is the generalization of derivative for functions that take several inputs (or one input in the form of a vector or some other complex structure).

Differentiation of a scalar function w.r.t. a vector

The gradient of a function is a vector of **partial derivatives**

Differentiation of a **scalar** function w.r.t. a **vector**

If $f(\mathbf{x})$ is a **scalar** function of d variables, \mathbf{x} is a $d \times 1$ vector.

Then differentiation of $f(\mathbf{x})$ w.r.t. \mathbf{x} results in a $d \times 1$ vector

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

T : $\mathbb{R}^3 \rightarrow \mathbb{R}^1$
 (x, y, z)
 3x1

This is referred to as the **gradient** of $f(\mathbf{x})$ and often written as $\nabla_{\mathbf{x}}f$.

E.g. $f(\mathbf{x}) = ax_1 + bx_2 \quad \nabla_{\mathbf{x}}f = \begin{bmatrix} a \\ b \end{bmatrix}$

Ref: Duda, Hart, and Stork, "Pattern Classification", 2001 (Appendix)

$$\nabla = (\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z})$$

$$\nabla f = (\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z})$$

- f takes in a vector and returns a scalar
- differentiating f with respect to the vector it takes in is nothing but the gradient of $f = \nabla f$
- Imagine $V(x, y, z)$, then $\nabla V = -\mathbf{E}$
- Partial derivative of f with respect to all the variable in \mathbf{x}

Differentiation of a vector function w.r.t. a vector

If $\mathbf{f}(\mathbf{x})$ is a vector function of size $h \times 1$ and \mathbf{x} is a $d \times 1$ vector.
Then differentiation of $\mathbf{f}(\mathbf{x})$ results in a $h \times d$ matrix

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

The matrix is referred to as the **Jacobian** of $\mathbf{f}(\mathbf{x})$

- \mathbf{f} is vector valued function of vector
- equivalent to the matrix multiplication of \mathbf{f} as a column vector ($h \times 1$) and ∇ as a row vector ($1 \times d$)
 - first row is the partial derivative of the first scalar function in \mathbf{f}

Properties

Derivative and Gradient

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$$

Some Vector-Matrix Differentiation Formulae

$$\frac{d(\mathbf{Ax})}{d\mathbf{x}} = \mathbf{A} \quad \begin{matrix} 2 \times 1 \\ 3 \times 1 \end{matrix} \rightarrow \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 \end{bmatrix}$$

$$\frac{d(\mathbf{b}^T \mathbf{x})}{d\mathbf{x}} = \underline{\mathbf{b}} \quad \frac{d(\mathbf{y}^T \mathbf{Ax})}{d\mathbf{x}} = \underline{\mathbf{A}^T \mathbf{y}}$$

$$\frac{d(\mathbf{x}^T \mathbf{Ax})}{d\mathbf{x}} = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}$$

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{x} = \underline{a_1 x_1 + a_2 x_2 + \dots + a_d x_d}$$

<https://math.stackexchange.com/questions/312077/differentiate-fx-xtax>

- last one is hard lmao

Linear Regression

$$\sum_{i=1}^m (f_{\mathbf{w}, b}(\mathbf{x}_i) - y_i)^2$$

- All model-based learning algorithms have a **loss function**
- What we do to find the best model is **to minimize the objective** known as the **cost function**
- **Cost function** is a sum of **loss functions** over training set plus possibly some model complexity penalty (regularization)
- In linear regression, the cost function is given by the *average loss*, also called the **empirical risk** because we do not have all the data (e.g. testing data)
 - The average of all penalties is obtained by applying the model to the training data
- Loss function is the penalty of being inaccurate
- Empirical risk, because, sample may not be representative the real world distribution, the best possible thing doable is to minimise the empirical risk

Linear regression is a popular regression learning algorithm that learns a model which is a linear combination of features of the input example.

$$\mathbf{X}\mathbf{w} = \mathbf{y}, \quad \mathbf{X} \in \mathcal{R}^{m \times d}, \mathbf{w} \in \mathcal{R}^{d \times 1}, \mathbf{y} \in \mathcal{R}^{m \times 1}$$

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,d} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Problem Statement: To predict the unknown y for a given \mathbf{x} (**testing**)

- We have a collection of labeled examples (**training**) $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$
 - m is the size of the collection
 - \mathbf{x}_i is the d -dimensional feature vector of example $i = 1, \dots, m$ (input)
 - y_i is a real-valued target (1-D)
 - Note:
 - when y_i is **continuous** valued, it is a **regression problem**
 - when y_i is **discrete** valued, it is a **classification problem**
- We want to build a model $f_{\mathbf{w}, b}(\mathbf{x})$ as a linear combination of features of example \mathbf{x} : $f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$
where \mathbf{w} is a d -dimensional vector of parameters and b is a real number.
- The notation $f_{\mathbf{w}, b}$ means that the model f is parametrized by two values: \mathbf{w} and b

Ref: [Book4] Stephen Boyd and Lieven Vandenberghe, "Introduction to Applied Linear Algebra", Cambridge University Press, 2018 (chp.14)

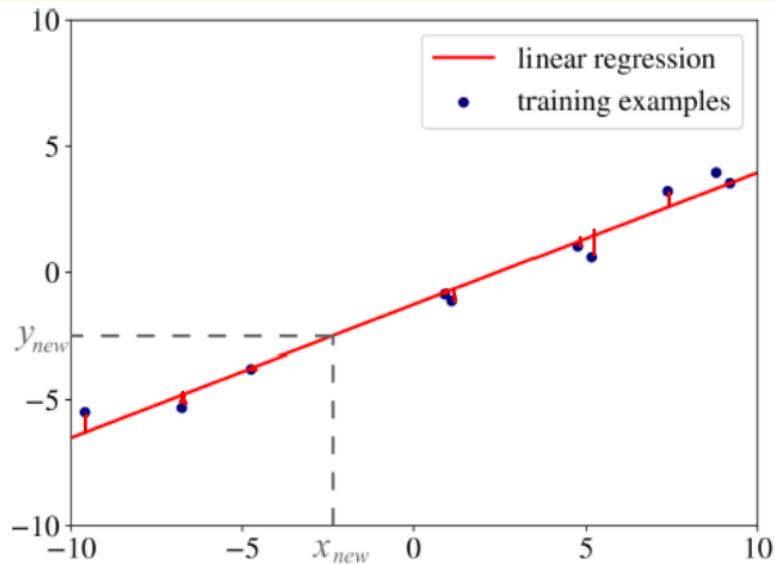
- during training cannot use test data, only training data, the use of test data is called data leakage
- a model is a function that turns \mathbf{x} into y
 - parametrized by \mathbf{w} and b , so these are the values to be found

Learning objective function

- To find the optimal values for w^* and b^* which **minimizes** the following expression:

$$\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x_i) - y_i)^2$$

- In mathematics, the expression we minimize or maximize is called an **objective function**, or, simply, an **objective**



$(f_w(x_i) - y_i)^2$ is called the **loss function**: a measure of the difference between $f_w(x_i)$ and y_i or a penalty for misclassification of example i .

- Objective function = cost function, minimise or maximise this to find the best regression, ie, the best value of w and b
 - sum of squared error or mean squared error
 - sum of the loss functions is the objective function
- Compare the value predicted by the model and the ground truth value
- Bias term is needed to shift the line in

Learning objective function (using simplified notation hereon)

- To find the optimal values for \mathbf{w}^* which **minimizes** the following expression:

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

with $f_{\mathbf{w}}(\mathbf{x}_i) = \mathbf{x}^T \mathbf{w}$,

where we define $\mathbf{w} = [b, w_1, \dots, w_d]^T = [w_0, w_1, \dots, w_d]^T$,

and $\mathbf{x}_i = [1, x_{i,1}, \dots, x_{i,d}]^T = [x_{i,0}, x_{i,1}, \dots, x_{i,d}]^T, i = 1, \dots, m$

- This particular choice of the loss function is called **squared error loss**

Note: The normalization factor $\frac{1}{m}$ can be omitted as it does not affect the optimization.

- putting b inside w to simplify the expression, adjust the x accordingly put a bunch of 1

Learning (Training)

- Consider the set of feature vector \mathbf{x}_i and target output y_i indexed by $i = 1, \dots, m$, a linear model $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ can be stacked as

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{X}) &= \mathbf{X}\mathbf{w} & \longleftrightarrow & \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} \\ \vdots \\ \mathbf{x}_m^T \mathbf{w} \end{bmatrix} & & \text{Learning target vector} \end{aligned}$$

where $\mathbf{x}_i^T \mathbf{w} = [1, x_{i,1}, \dots, x_{i,d}] \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$

Note: The **bias/offset term** is responsible for translating the line/plane/hyperplane away from the origin.

- shove all the sample together and form matrix X

Least Squares Regression

In vector-matrix notation, the minimization of the objective function can be written compactly using $\mathbf{e} = \mathbf{X}\mathbf{w} - \mathbf{y}$:

$$\begin{aligned} J(\mathbf{w}) &= \mathbf{e}^T \mathbf{e} \\ &= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= (\mathbf{w}^T \mathbf{X}^T - \mathbf{y}^T)(\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y} \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y}. \end{aligned}$$

Note: when $f_{\mathbf{w}}(\mathbf{X}) = \mathbf{X}\mathbf{w}$, then

$$\sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}).$$

- minimise sum of square of error \Rightarrow least squares regression
- $J(w)$ is a scalar

Differentiating $J(w)$ with respect to w and setting the result to 0 :

$$\frac{\partial}{\partial w} J(w) = 0$$

$$\begin{aligned}\frac{\partial}{\partial w} (w^T X^T X w - 2y^T X w + y^T y) &= 0 \\ \Rightarrow 2X^T X w - 2X^T y &= 0 \\ \Rightarrow X^T X w &= X^T y\end{aligned}$$

\Rightarrow Any minimizer \hat{w} of $J(w)$ must satisfy $X^T(Xw - y) = 0$

If $X^T X$ is invertible, then

Learning/training:

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

Prediction/testing:

$$\hat{f}_w(X_{new}) = X_{new} \hat{w}$$

- differentiate with respect to w to minimise the cost, to find the best w

Let $V \subseteq \mathbb{R}^n$ and $S = \{u_1, u_2, \dots, u_k\}$ be a basis for V . Then the orthogonal projection of any vector $w \in \mathbb{R}^n$ onto $V = A(A^T A)^{-1} A^T w$, where $A = (u_1 \ u_2 \ \dots \ u_k)$ $n \times k$

Since S is basis, columns of A are linearly independent ($\text{Rank}(A) = k$, A is full ranked) so $A^T A$ is invertible

$\therefore u = (A^T A)^{-1} A^T w$ is the unique solution to $A^T A u = A^T w$

\therefore project of w onto $\text{Col}(A) = V = Au = A(A^T A)^{-1} A^T w$

- the least square is the left inverse

- $X^T X$ must be invertible

- note that w include the b

Example 1: one dimension input, 1 feature

Example 1 **Training set** $\{(x_i, y_i)\}_{i=1}^m$

X	w	y	
$\begin{bmatrix} 1 & -9 \\ 1 & -7 \\ 1 & -5 \\ 1 & 1 \\ 1 & 5 \\ 1 & 9 \end{bmatrix}$	$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$	$\begin{bmatrix} -6 \\ -6 \\ -4 \\ -1 \\ 1 \\ 4 \end{bmatrix}$	$\{x = -9\} \rightarrow \{y = -6\}$ $\{x = -7\} \rightarrow \{y = -6\}$ $\{x = -5\} \rightarrow \{y = -4\}$ $\{x = 1\} \rightarrow \{y = -1\}$ $\{x = 5\} \rightarrow \{y = 1\}$ $\{x = 9\} \rightarrow \{y = 4\}$

$$\begin{bmatrix} 1 & -9 \\ 1 & -7 \\ 1 & -5 \\ 1 & 1 \\ 1 & 5 \\ 1 & 9 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} -6 \\ -6 \\ -4 \\ -1 \\ 1 \\ 4 \end{bmatrix}$$

This set of linear equations has no exact solution

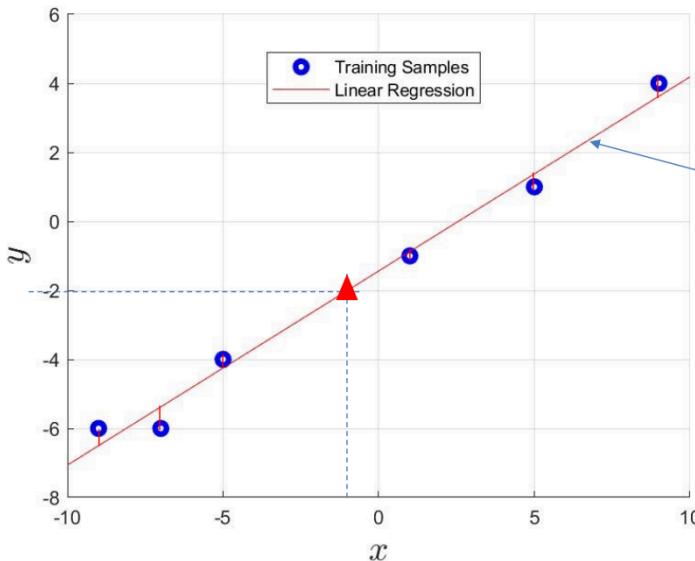
However, $X^T X$ is invertible

Least square approximation

$$\hat{w} = X^\dagger y = (X^T X)^{-1} X^T y$$

$$= \begin{bmatrix} 6 & -6 \\ -6 & 262 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ -9 & -7 & -5 & 1 & 5 & 9 \end{bmatrix} \begin{bmatrix} -6 \\ -6 \\ -4 \\ -1 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} -1.4375 \\ 0.5625 \end{bmatrix}$$

- remember to add 1 to get the bias in (augmentation)
- then shove the x into X along with the padded 1
- Overdetermined system, likely to have left inverse
 - no exact solution, has least square solution



$$\hat{y} = \mathbf{X}\hat{\mathbf{w}}$$

$$= \mathbf{X} \begin{bmatrix} -1.4375 \\ 0.5625 \end{bmatrix}$$

$$y = -1.4375 + 0.5625x$$

Prediction:

Test set

$$\{x = -1\} \rightarrow \{y = ?\}$$

$$\hat{y} = [1 \quad -1] \begin{bmatrix} -1.4375 \\ 0.5625 \end{bmatrix}$$

$$= -2$$

Linear Regression on one-dimensional samples

Python demo 1

```

1  # EE2211 Lecture 5 Demo 1 linear regression
2  import pandas as pd
3  #import matplotlib.pyplot as plt
4  import numpy as np
5  from numpy.linalg import inv
6  from sklearn.metrics import mean_squared_error
7  ##
8  X = np.array([[1, -9], [1, -7], [1, -5], [1, 1], [1, 5], [1, 9]])
9  Y = np.array([[ -6], [-6], [-4], [-1], [1], [4]])
10 w = inv(X.T @ X) @ X.T @ Y
11 print(w)
12
13 Xnew = np.array([1, -1])
14 Ynew = Xnew@w
15 print(Ynew)
16
17 ## difference
18 print("Mean squared error between Y and Xw")
19 Ytest=X@w
20 MSE = np.square(np.subtract(Ytest,Y)).mean()
21 print(MSE)
22 MSE = mean_squared_error(Ytest,Y)
23 print(MSE)
24

```

- manually do MSE or use built in function
- remember to pad 1 for X

```

[[[-1.4375]
 [ 0.5625]]
 [-2.]
Mean squared error between Y and XW
0.1666666666666666
0.1666666666666666

```

Example 2: 3 dimension input, 3 features

Example 2	$\{(\mathbf{x}_i, y_i)\}_{i=1}^m$	$\{x_1 = 1, x_2 = 1, x_3 = 1\} \rightarrow \{y = 1\}$
Training set		$\{x_1 = 1, x_2 = -1, x_3 = 1\} \rightarrow \{y = 0\}$
		$\{x_1 = 1, x_2 = 1, x_3 = 3\} \rightarrow \{y = 2\}$
		$\{x_1 = 1, x_2 = 1, x_3 = 0\} \rightarrow \{y = -1\}$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 3 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \\ -1 \end{bmatrix}$$

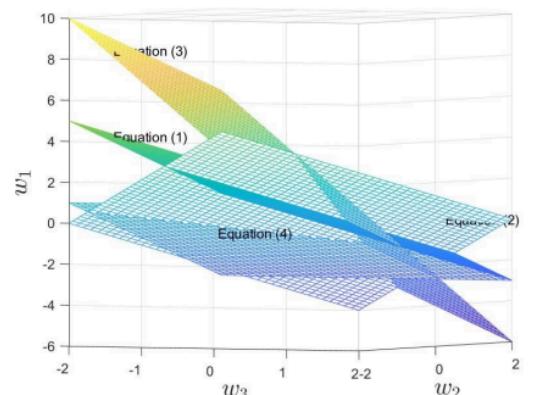
This set of linear equations has no exact solution

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

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

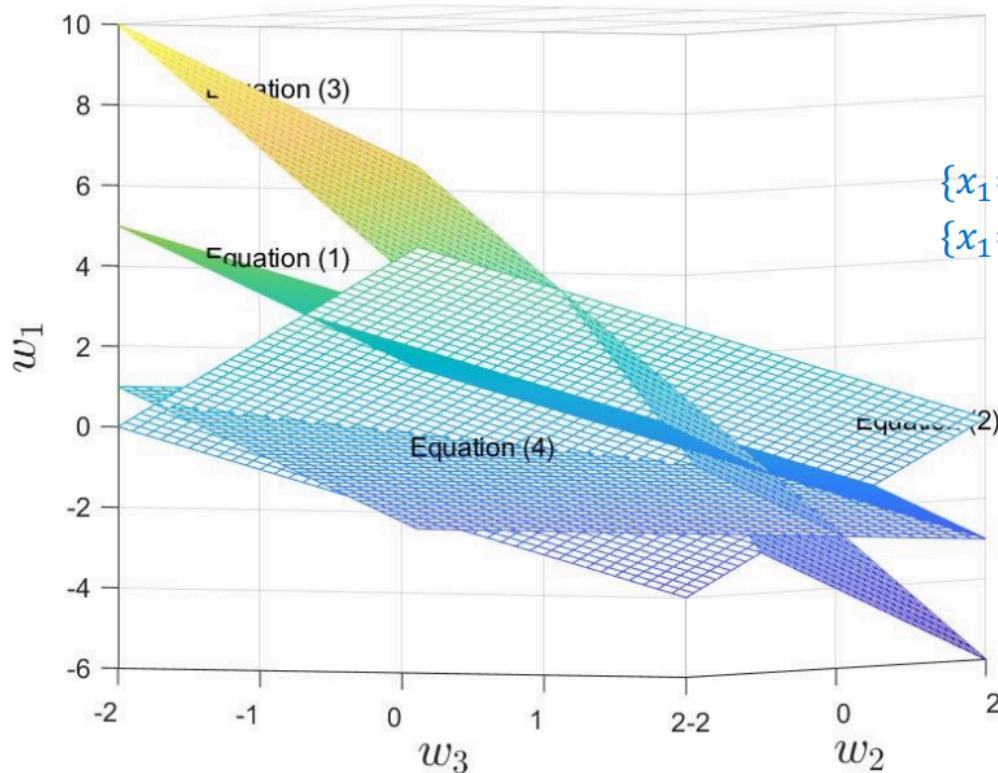
$$= \begin{bmatrix} 4 & 2 & 5 \\ 2 & 4 & 3 \\ 5 & 3 & 11 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & 1 & 3 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 2 \\ -1 \end{bmatrix} = \begin{bmatrix} -0.7500 \\ 0.1786 \\ 0.9286 \end{bmatrix}$$

Least square approximation



- in this case the first feature is all one (or a constant), so no need to pad ones
 - so w_1 used as the offset

The four linear equations



Prediction:

Test set

$$\{x_1 = 1, x_2 = 6, x_3 = 8\} \rightarrow \{y = ?\}$$

$$\{x_1 = 1, x_2 = 0, x_3 = -1\} \rightarrow \{y = ?\}$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new}\hat{\mathbf{w}}$$

$$\begin{aligned}\hat{\mathbf{y}} &= \begin{bmatrix} 1 & 6 & 8 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} -0.7500 \\ 0.1786 \\ 0.9286 \end{bmatrix} \\ &= \begin{bmatrix} 7.7500 \\ -1.6786 \end{bmatrix}\end{aligned}$$

Learning of Vectored Function (Multiple Outputs)

For one sample: a linear model $\mathbf{f}_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{W}$ Vector function

For m samples: $\mathbf{F}_{\mathbf{w}}(\mathbf{X}) = \mathbf{X}\mathbf{W} = \mathbf{Y}$

$$\begin{array}{l} \text{Sample 1} \longrightarrow \mathbf{x}_1^T \\ \vdots \\ \text{Sample } m \longrightarrow \mathbf{x}_m^T \end{array} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix} \mathbf{W} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m,1} & \dots & x_{m,d} \end{bmatrix} \underbrace{\begin{bmatrix} w_{0,1} & \dots & w_{0,h} \\ w_{1,1} & \dots & w_{1,h} \\ \vdots & \ddots & \vdots \\ w_{d,1} & \dots & w_{d,h} \end{bmatrix}}_h$$

$$\begin{array}{l} \text{Sample 1's output} \longrightarrow \begin{bmatrix} y_{1,1} & \dots & y_{1,h} \end{bmatrix} \\ \vdots \\ \text{Sample } m \text{'s output} \longrightarrow \underbrace{\begin{bmatrix} y_{m,1} & \dots & y_{m,h} \end{bmatrix}}_h \end{array} \underbrace{m}_h$$

$$\mathbf{X} \in \mathcal{R}^{m \times (d+1)}, \mathbf{W} \in \mathcal{R}^{(d+1) \times h}, \mathbf{Y} \in \mathcal{R}^{m \times h}$$

- w is no longer a vector, it's a matrix W
- it's solving the same system of equation but with different y, the features are the same, to predict the different output, the weights must necessarily be different
- But now must be least squares for all the output predictions

Objective: $\sum_{i=1}^m (\mathbf{f}_w(\mathbf{x}_i) - \mathbf{y}_i)^2 = \mathbf{E}^T \mathbf{E}$

Least Squares Regression of Multiple Outputs

In matrix notation, the sum of squared errors cost function can be written compactly using $\mathbf{E} = \mathbf{XW} - \mathbf{Y}$:

$$\begin{aligned} J(\mathbf{W}) &= \text{trace}(\mathbf{E}^T \mathbf{E}) \\ &= \text{trace}[(\mathbf{XW} - \mathbf{Y})^T (\mathbf{XW} - \mathbf{Y})] \end{aligned}$$

- With error matrix $\mathbf{E} = \mathbf{XW} - \mathbf{Y}$, \mathbf{E} is $m \times h$
- Minimising J is hard, trace of J is easier
- $\mathbf{E}^T \mathbf{E}$ is $h \times h$ with the trace of $\mathbf{E}^T \mathbf{E}$ being the sum of error squared

Least Squares Regression of Multiple Outputs

$$\begin{aligned} J(\mathbf{W}) &= \text{trace}(\mathbf{E}^T \mathbf{E}) \\ &= \text{trace}\left(\begin{bmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_h^T \end{bmatrix} [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_h]\right) \\ &= \text{trace}\left(\begin{bmatrix} \mathbf{e}_1^T \mathbf{e}_1 & \mathbf{e}_1^T \mathbf{e}_2 & \dots & \mathbf{e}_1^T \mathbf{e}_h \\ \mathbf{e}_2^T \mathbf{e}_1 & \mathbf{e}_2^T \mathbf{e}_2 & \dots & \mathbf{e}_2^T \mathbf{e}_h \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{e}_h^T \mathbf{e}_1 & \mathbf{e}_h^T \mathbf{e}_2 & \dots & \mathbf{e}_h^T \mathbf{e}_h \end{bmatrix}\right) = \sum_{k=1}^h \mathbf{e}_k^T \mathbf{e}_k \end{aligned}$$

○

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

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

Prediction/testing: $\widehat{\mathbf{F}}_w(\mathbf{X}_{new}) = \mathbf{X}_{new} \widehat{\mathbf{W}}$

- - least square estimate for multiple outputs is the similar as the for single output

- The regression is the same as the outputs are independent, and can just do the same regression for the output separately. But doing together can save computation resource

Example 3:

Training set $\{x_1 = 1, x_2 = 1, x_3 = 1\} \rightarrow \{y_1 = 1, y_2 = 0\}$

$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^m \quad \{x_1 = 1, x_2 = -1, x_3 = 1\} \rightarrow \{y_1 = 0, y_2 = 1\}$

$\{x_1 = 1, x_2 = 1, x_3 = 3\} \rightarrow \{y_1 = 2, y_2 = -1\}$

$\{x_1 = 1, x_2 = 1, x_3 = 0\} \rightarrow \{y_1 = -1, y_2 = 3\}$

X **W** **Y**

$$\text{Bias} \rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 3 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \\ w_{3,1} & w_{3,2} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 2 & -1 \\ -1 & 3 \end{bmatrix}$$

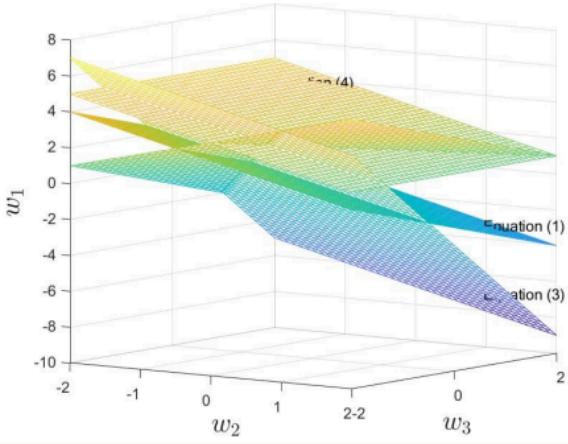
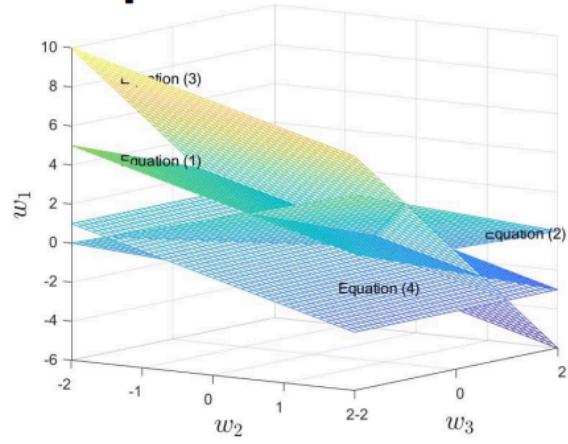
This set of linear equations has NO exact solution

$$\hat{\mathbf{W}} = \mathbf{X}^\dagger \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \mathbf{X}^T \mathbf{X} \text{ is invertible}$$

Least square approximation

$$= \begin{bmatrix} 4 & 2 & 5 \\ 2 & 4 & 3 \\ 5 & 3 & 11 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & 1 & 3 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 2 & -1 \\ -1 & 3 \end{bmatrix} = \begin{bmatrix} -0.75 & 2.25 \\ 0.1786 & 0.0357 \\ 0.9286 & -1.2143 \end{bmatrix}$$

Example 3



Prediction:

Test set: two new samples

$$\{x_1 = 1, x_2 = 6, x_3 = 8\} \rightarrow \{y_1 = ?, y_2 = ?\}$$

$$\{x_1 = 1, x_2 = 0, x_3 = -1\} \rightarrow \{y_1 = ?, y_2 = ?\}$$

$$\hat{\mathbf{Y}} = \mathbf{X}_{new} \hat{\mathbf{W}}$$

$$\begin{aligned} \text{Bias} &= \begin{bmatrix} 1 & 6 & 8 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} -0.75 & 2.25 \\ 0.1786 & 0.0357 \\ 0.9286 & -1.2143 \end{bmatrix} \\ &= \begin{bmatrix} 7.75 & -7.25 \\ -1.6786 & 3.4643 \end{bmatrix} \end{aligned}$$

Python demo 2

```

27 # EE2211 Lecture 5 Demo 2
28 import pandas as pd
29 import matplotlib.pyplot as plt
30 import numpy as np
31 from numpy.linalg import inv
32 from sklearn.metrics import mean_squared_error
33 ##
34 X = np.array([[1, 1, 1], [1, -1, 1], [1, 1, 3], [1, 1, 0]])
35 Y = np.array([[1, 0], [0, 1], [2, -1], [-1, 3]])
36 w = inv(X.T @ X) @ X.T @ Y
37 print("the estimated w")
38 print(w)
39
40 Xnew = np.array([[1, 6, 8], [1, 0, -1]])
41 Ynew = Xnew@w
42 print("testing Ynew")
43 print(Ynew)
44
45 ## difference
46 print("Mean squared error between Y and Xw")
47 Ytest=X@w
48 MSE = np.square(np.subtract(Ytest,Y)).mean()
49 print(MSE)
50 MSE = mean_squared_error(Ytest,Y)
51 print(MSE)
52

```

```

the estimated w
[[ -0.75      2.25      ]
 [  0.17857143  0.03571429]
 [  0.92857143 -1.21428571]]
testing Ynew
[[ 7.75      -7.25      ]
 [ -1.67857143  3.46428571]]
Mean squared error between Y and Xw
0.3035714285714286
0.30357142857142855

```

Example 4

The values of feature x and their corresponding values of multiple outputs target y are shown in the table below.

Based on the least square regression, what are the values of w ?

Based on the current mapping, when $x = 2$, what is the value of y ?

x	[3]	[4]	[10]	[6]	[7]
y	[0, 5]	[1.5, 4]	[-3, 8]	[-4, 10]	[1, 6]

$$\hat{W} = X^{\dagger}Y = (X^T X)^{-1} X^T Y = \begin{bmatrix} 1.9 & 3.6 \\ -0.4667 & 0.5 \end{bmatrix}$$

Python demo 3

$$\widehat{Y}_{new} = X_{new} \hat{W} = [1 \quad 2] \hat{W} = [0.9667 \quad 4.6]$$

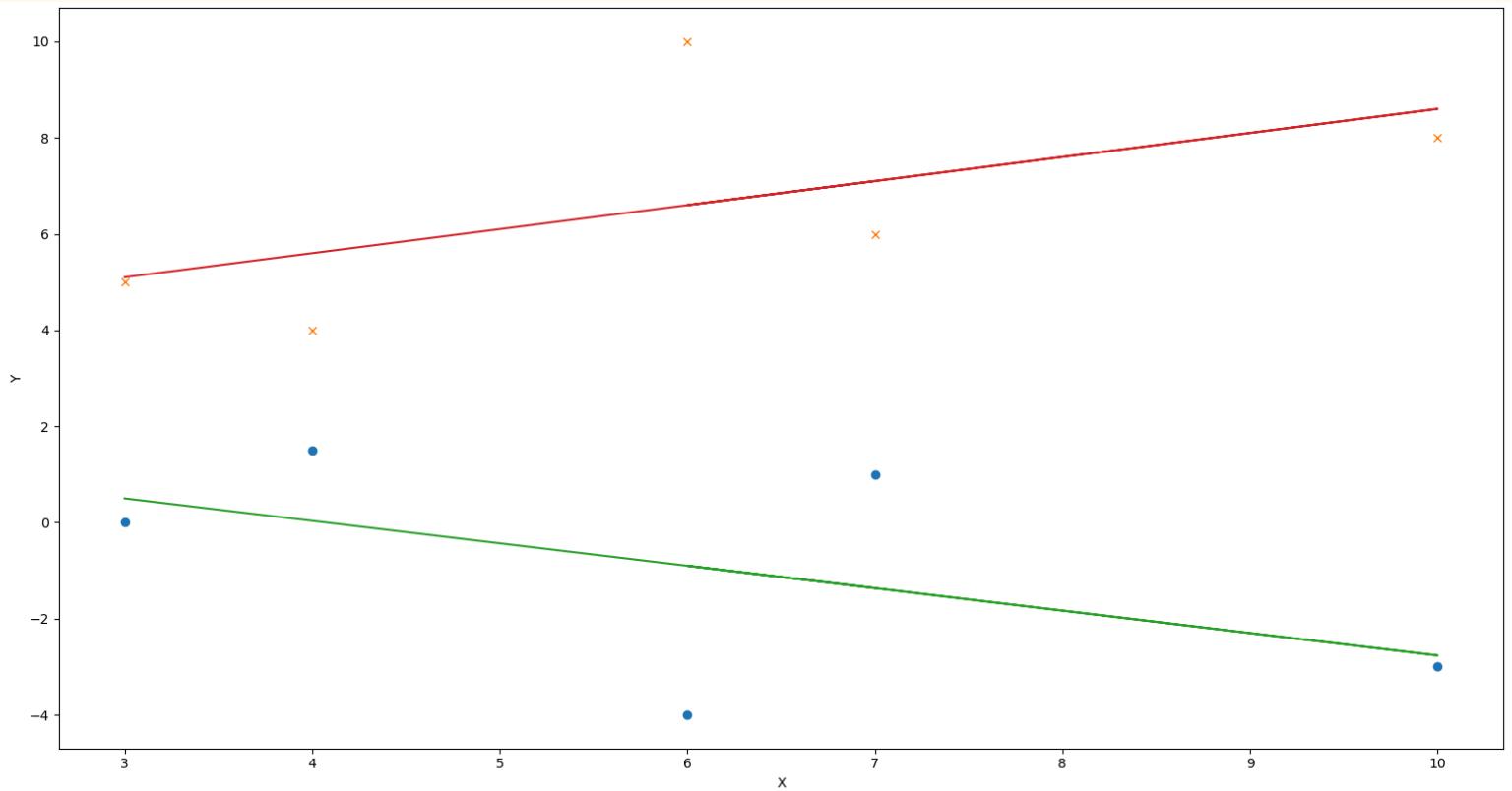
Prediction

```

53 # EE2211 Lecture 5 Demo 3
54 import pandas as pd
55 import matplotlib.pyplot as plt
56 import numpy as np
57 from numpy.linalg import inv
58 from sklearn.metrics import mean_squared_error
59 ##
60 X = np.array([[1, 3], [1, 4], [1, 10], [1, 6], [1, 7]])
61 Y = np.array([[0, 5], [1.5, 4], [-3, 8], [-4, 10], [1, 6]])
62 w = inv(X.T @ X) @ X.T @ Y
63 print('W')
64 print(w)
65
66 Xnew = np.array([[1, 2]])
67
68 Ynew = Xnew@w
69
70 print('Ynew')
71 print(Xnew)
72 print(Xnew.shape)
73 print(Ynew)
74
75 Ytest= X@w
76 print('Ytest')
77 print(Ytest)
78 plt.plot(X[:,1], Y[:,0], 'o', label = 'Y1')
79 plt.plot(X[:,1], Y[:,1], 'x', label = 'Y2')
80 plt.plot(X[:,1], Ytest[:,0])
81 plt.plot(X[:,1], Ytest[:,1])
82 plt.xlabel('X')
83 plt.ylabel('Y')
84 plt.show()

```

W	$\begin{bmatrix} 1.9 & 3.6 \\ -0.46666667 & 0.5 \end{bmatrix}$
Ynew	$\begin{bmatrix} [1 2] \\ (1, 2) \end{bmatrix}$
Ytest	$\begin{bmatrix} [0.96666667 4.6] \\ [[0.5 5.1] [0.03333333 5.6] [-2.76666667 8.6] [-0.9 6.6] [-1.36666667 7.1]] \end{bmatrix}$



- green is Y_1 and red is Y_2

Class Quiz

Suppose $y = 3x + 5$, this is a linear function.

Jin Yueming 42

True

69%

False

31%

- its affine

Suppose $g(\mathbf{x})$ is a scalar function of d variables where \mathbf{x} is a $d \times 1$ vector, the outcome of differentiation of $g(\mathbf{x})$ w.r.t. \mathbf{x} is a $d \times 1$ vector.

Jin Yueming 43

True

91%

False

9%

Summary

- Notations, Vectors, Matrices
- Operations on Vectors and Matrices
 - Dot-product, matrix inverse
- Systems of Linear Equations $f_w(\mathbf{X}) = \mathbf{Xw} = \mathbf{y}$
 - Matrix-vector notation, linear dependency, invertible
 - Even-, over-, under-determined linear systems
- Functions, Derivative and Gradient
 - Inner product, linear/affine functions
 - Maximum and minimum, partial derivatives, gradient
- Least Squares, Linear Regression
 - Objective function, loss function
 - Least square solution, training/learning and testing/prediction
 - Linear regression with multiple outputs

Midterm (L1 to L5)
Trial quiz

Learning/training $\hat{\mathbf{w}} = (\mathbf{X}_{train}^T \mathbf{X}_{train})^{-1} \mathbf{X}_{train}^T \mathbf{y}_{train}$

Prediction/testing $\mathbf{y}_{test} = \mathbf{X}_{test} \hat{\mathbf{w}}$

- Classification
- Ridge Regression
- Polynomial Regression

Python packages: numpy, pandas, matplotlib.pyplot, numpy.linalg, and sklearn.metrics (for mean_squared_error), numpy.linalg.pinv

Lec6: Ridge regression & Polynomial regression

Linear Regression for Classification

- We have a collection of labeled examples
 - m is the size of the collection
 - \mathbf{x}_i is the d -dimensional feature vector of example $i = 1, \dots, m$
 - y_i is discrete target label (e.g., $y_i \in \{-1, +1\}$ or $\{0, 1\}$ for binary classification problems)
 - Note:
 - when y_i is **continuous** valued \rightarrow a **regression problem**
 - when y_i is **discrete** valued \rightarrow a **classification problem**
- Linear model: $f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$ or in compact form $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ (having the offset term absorbed into the inner product)

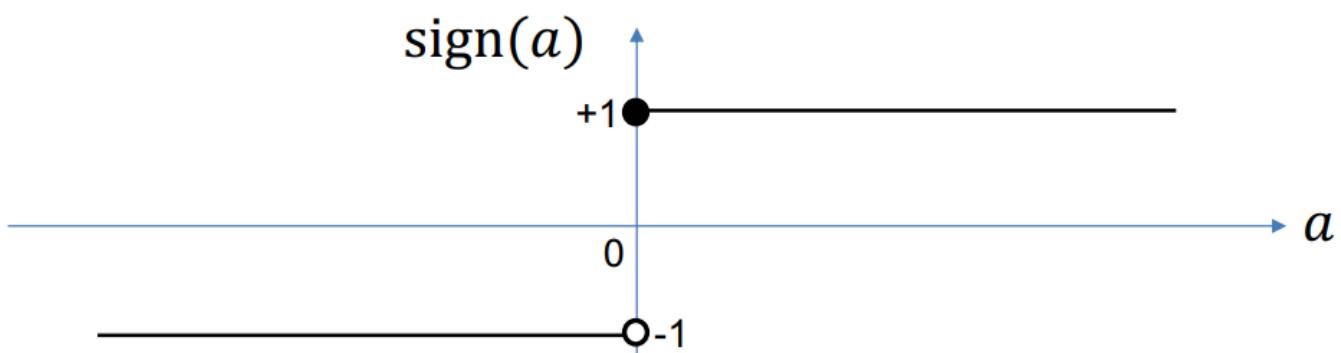
Binary Classification

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

Learning: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, $y_i \in \{-1, +1\}, i = 1, \dots, m$

Prediction: $\hat{f}_{\mathbf{w}}^c(\mathbf{x}_{new}) = \text{sign}(\mathbf{x}_{new}^T \hat{\mathbf{w}})$ for each row \mathbf{x}_{new}^T of \mathbf{X}_{new}

$$\text{sign}(a) = +1 \text{ for } a \geq 0 \text{ and } -1 \text{ for } a < 0$$



- If classified by 0, 1, take threshold as 0.5
 - or use -1 +1
 - don't use larger number incase large number overshadows smaller number
- Use sign function for binary

Example

Example 1 Training set $\{x_i, y_i\}_{i=1}^m$

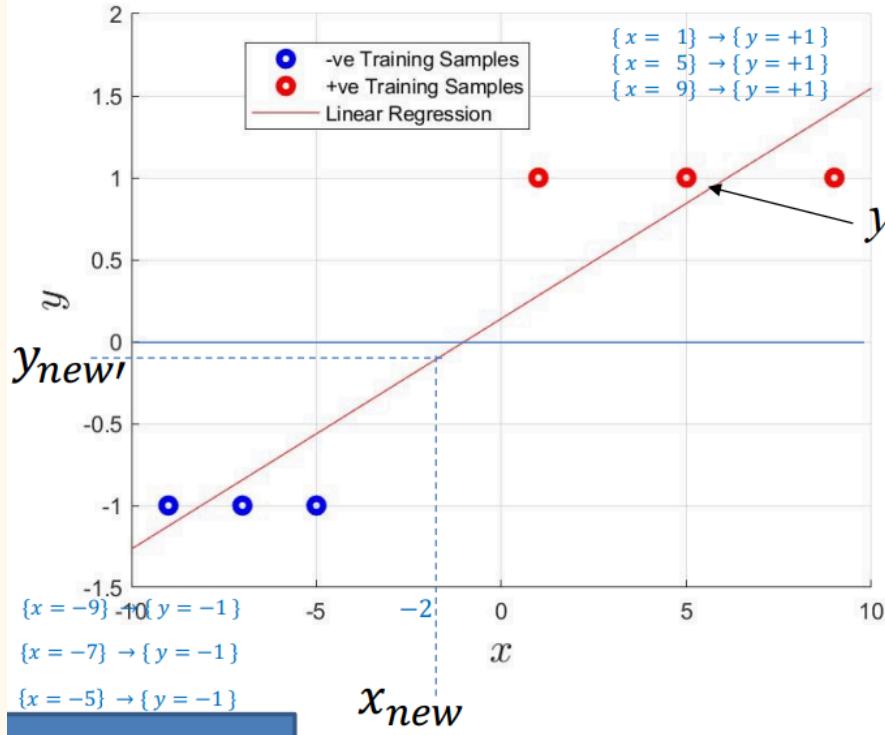
X	w	y	
Bias			
$\begin{bmatrix} 1 & -9 \\ 1 & -7 \\ 1 & -5 \\ 1 & 1 \\ 1 & 5 \\ 1 & 9 \end{bmatrix}$	$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$	$\begin{bmatrix} -1 \\ -1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$	
			$\{x = -9\} \rightarrow \{y = -1\}$ $\{x = -7\} \rightarrow \{y = -1\}$ $\{x = -5\} \rightarrow \{y = -1\}$ $\{x = 1\} \rightarrow \{y = +1\}$ $\{x = 5\} \rightarrow \{y = +1\}$ $\{x = 9\} \rightarrow \{y = +1\}$

This set of linear equations has NO exact solution

$$\hat{\mathbf{w}} = \mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \mathbf{X}^T \mathbf{X} \text{ is invertible}$$

$$= \begin{bmatrix} 6 & -6 \\ -6 & 262 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ -9 & -7 & -5 & 1 & 5 & 9 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.1406 \\ 0.1406 \end{bmatrix} \quad \text{Least square approximation}$$

Example 1



```
# EE2211 Lecture 6 Demo 1 Binary classification
import numpy as np
from numpy.linalg import inv

X = np.array([[1,-9], [1,-7], [1,-5], [1,1], [1,5], [1, 9]])
y = np.array([-1, -1, -1, 1, 1, 1])
## Linear regression for classification
w = inv(X.T @ X) @ X.T @ y
print("Estimated w")
print(w)
xt = np.array([[1,-2]])
y_predict = xt @ w
print("Predicted y")
print(y_predict)
y_class_predict = np.sign(y_predict)
print("Predicted y class")
print(y_class_predict)
```

$$\hat{y} = \text{sign}(\mathbf{X}\hat{\mathbf{w}}) \\ = \text{sign}(\mathbf{X} \begin{bmatrix} 0.1406 \\ 0.1406 \end{bmatrix})$$

Prediction:

Test set $\{x = -2\} \rightarrow \{y = ?\}$

$$y_{new} = \hat{f}_w^c(\mathbf{x}_{new}) = \text{sign}(\mathbf{x}_{new}\hat{\mathbf{w}})$$

Bias
 $= \text{sign}([1 \quad -2] \begin{bmatrix} 0.1406 \\ 0.1406 \end{bmatrix})$

$= \text{sign}(-0.1406) = -1$

```
Estimated w
[[0.140625]
 [0.140625]]
Predicted y
[[-0.140625]]
Predicted y class
[[-1.]]
```

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

$$\text{Learning: } \widehat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}, \quad \mathbf{Y} \in \mathbf{R}^{m \times C}$$

Prediction: $\hat{f}_w^c(x_{new}) = \arg \max_{k=1,\dots,c} (\mathbf{x}_{new}^T \widehat{\mathbf{W}}(:, k))$ for each \mathbf{x}_{new}^T of \mathbf{X}_{new}

Each row (of $i = 1, \dots, m$) in \mathbf{Y} has an **one-hot** encoding/assignment:

e.g., target for class-1 is labelled as $y_i^T = [1, 0, 0, \dots, 0]$ for the i th sample,

target for class-2 is labelled as $\mathbf{y}_i^T = [0, \textcolor{red}{1}, 0, \dots, 0]$ for the j th sample,

target for class-C is labelled as $\mathbf{y}_m^T = [0, 0, \dots, 0, 1]$ for the m th sample.

- Use one hot encoding and multiple output linear regression
 - One column for each class
 - Then for prediction, the column with the biggest value to determine which class it is

Training set	$\{x_i, y_i\}_{i=1}^m$	$\{x_1 = 1, x_2 = 1\} \rightarrow \{y_1 = 1, y_2 = 0, y_3 = 0\}$	Class 1
		$\{x_1 = -1, x_2 = 1\} \rightarrow \{y_1 = 0, y_2 = 1, y_3 = 0\}$	Class 2
		$\{x_1 = 1, x_2 = 3\} \rightarrow \{y_1 = 1, y_2 = 0, y_3 = 0\}$	Class 1
		$\{x_1 = 1, x_2 = 0\} \rightarrow \{y_1 = 0, y_2 = 0, y_3 = 1\}$	Class 3

$$\text{Bias} \rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 3 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \\ w_{2,1} & w_{2,2} & w_{2,3} \\ w_{3,1} & w_{3,2} & w_{3,3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

This set of linear equations has NO exact solution.

$$\hat{\mathbf{W}} = \mathbf{X}^\dagger \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \text{if } \mathbf{X}^T \mathbf{X} \text{ is invertible}$$

Least square approximation

$$= \begin{bmatrix} 4 & 2 & 5 \\ 2 & 4 & 3 \\ 5 & 3 & 11 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & 1 & 3 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0.2857 & -0.5 & 0.2143 \\ 0.2857 & 0 & -0.2857 \end{bmatrix}$$

Test set \mathbf{X}_{new} $\{x_1 = 6, x_2 = 8\} \rightarrow \{\text{class 1, 2, or 3?}\}$
 $\{x_1 = 0, x_2 = -1\} \rightarrow \{\text{class 1, 2, or 3?}\}$

$$\widehat{\mathbf{Y}} = \mathbf{X}_{new} \widehat{\mathbf{W}} = \begin{bmatrix} 1 & 6 & 8 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0.2857 & -0.5 & 0.2143 \\ 0.2857 & 0 & -0.2857 \end{bmatrix}$$

Category prediction:

$$\begin{aligned} \widehat{f}_w^c(\mathbf{X}_{new}) &= \arg \max_{k=1, \dots, C} (\widehat{\mathbf{Y}}(:, k)) \\ &= \arg \max_{k=1, \dots, C} \left(\begin{bmatrix} 4 & -2.50 & -0.50 \\ -0.2587 & 0.50 & 0.7857 \end{bmatrix} \right) \\ &= \begin{bmatrix} 1 \\ 3 \end{bmatrix} \quad \begin{array}{l} \xrightarrow{\text{Class 1}} \\ \xrightarrow{\text{Class 3}} \end{array} \end{aligned}$$

For each row of \mathbf{Y} , the **column position** of the **largest** number (across all columns for that row) determines the **class label**.

E.g. in the first row, the maximum number is 4 which is in column 1. Therefore, the resulting predicted class is 1.

Python demo 2

```

# EE2211 Lecture 6 Demo 2 Multi-class classification
import numpy as np
from numpy.linalg import inv
from sklearn.preprocessing import OneHotEncoder
X = np.array([[1, 1, 1], [1, -1, 1], [1, 1, 3], [1, 1, 0]])
y_class = np.array([1, 2, 1, 3])
y_onehot = np.array([[1, 0, 0], [0, 1, 0], [1, 0, 0], [0, 0, 1]])
print("One-hot encoding manual")
print(y_class)
print(y_onehot)

print("One-hot encoding function")
onehot_encoder=OneHotEncoder(sparse=False)
print(onehot_encoder)
Ytr_onehot = onehot_encoder.fit_transform(y_class)
print(Ytr_onehot)

#reshaped = y_class.reshape(len(y_class), 1)
#print(reshaped)
#Ytr_onehot = onehot_encoder.fit_transform(reshaped)

```

```

## Linear Classification
print("Estimated W")
W = inv(X.T @ X) @ X.T @ Ytr_onehot
print(W)
X_test = np.array([[1, 6, 8], [1, 0, -1]])
yt_est = X_test@W;
print("Test")
print(yt_est)
yt_class = [[1 if y == max(x) else 0 for y in x] for x in yt_est ]
print("class label test")
print(yt_class)

```

```
One-hot encoding manual
```

```
[[1]
 [2]
 [1]
 [3]]
[[1 0 0]
 [0 1 0]
 [1 0 0]
 [0 0 1]]
```

```
One-hot encoding function
```

```
OneHotEncoder(sparse=False)
```

```
[[1. 0. 0.]
 [0. 1. 0.]
 [1. 0. 0.]
 [0. 0. 1.]]
```

```
Estimated W
```

```
[[ 1.11022302e-16  5.00000000e-01  5.00000000e-01]
 [ 2.85714286e-01 -5.00000000e-01  2.14285714e-01]
 [ 2.85714286e-01  2.77555756e-17 -2.85714286e-01]]
```

```
Test
```

```
[[ 4.          -2.5          -0.5         ]
 [-0.28571429  0.5          0.78571429]]
```

```
class label test
```

```
[[1, 0, 0], [0, 0, 1]]
```

Ridge Regression

Recall Linear regression

Objective: $\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$

The learning computation: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

We cannot guarantee that the matrix $\mathbf{X}^T \mathbf{X}$ is invertible

Ridge regression: shrinks the regression coefficients w by imposing a penalty on their size

Objective: $\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^m (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \sum_{j=1}^d w_j^2$
 $= \operatorname{argmin} (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$

Here $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage: the larger the value of λ , the greater the amount of shrinkage.

Note: m samples & d parameters

- Least square solution for linear regression requires $\mathbf{X}^T \mathbf{X}$ must be invertible, but this is not guarantee for any given sample
- Ridge regression introduce additional objective term to shrink the w
 - first part is minimise square error, second part is to minimise size of w
 - reduce size of weights
- Ridge regression to avoid overfitting issue
- Multiple output ridge regression is out of scope

Using a linear model:

$$\min_{\mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

Solution:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} ((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}) &= \mathbf{0} \\ \Rightarrow 2\mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} &= \mathbf{0} \\ \Rightarrow \mathbf{X}^T \mathbf{X}\mathbf{w} + \lambda \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\mathbf{w} &= \mathbf{X}^T \mathbf{y} \end{aligned}$$

where \mathbf{I} is the $d \times d$ identity matrix

Here on, we shall focus on single column of output \mathbf{y} in derivations in the sequel

Learning:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- Adding a small value along the diagonal means even if $\mathbf{X}^T \mathbf{X}$ is not invertible, $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is.

Primal and Dual Form

- dual and prime give the same answer but one is computationally less expensive
- Always use the one that requires to find the inverse of the smaller square matrix

```
## dual solution m < d (without ridge)
w_dual = P.T @ inv(P @ P.T) @ y
print("Under-determined system")
print("Unique constrained solution, no ridge")
print(w_dual)

print("*****")
print("Approximation with dual ridge regression")
print(P.shape)
reg_L2 = 0.0001*np.identity(P.shape[0]) #number of rows of P = Dual I
print(reg_L2)
w_dual_ridge = P.T @ (inv(P @ P.T + reg_L2)) @ y
print(w_dual_ridge)

print("*****")
## primal ridge
print("Approximation with primal ridge regression")
print(P.shape)
reg_L = 0.0001*np.identity(P.shape[1]) #number of columns of P = Primal I
print(reg_L)
w_primal_ridge = inv(P.T @ P + reg_L) @ P.T @ y
print(w_primal_ridge)
```

<pre>Under-determined system Unique constrained solution, no ridge [[-1.] [1.] [1.] [1.] [-4.] [1.]]</pre>	<pre>***** Approximation with dual ridge regression (4, 6) [[0.0001 0. 0. 0. 0.] [0. 0.0001 0. 0. 0.] [0. 0. 0.0001 0. 0.] [0. 0. 0. 0.0001 0.] [0. 0. 0. 0. 0.0001]]</pre>
<pre>***** Approximation with primal ridge regression (4, 6) [[0.0001 0. 0. 0. 0.] [0. 0.0001 0. 0. 0.] [0. 0. 0.0001 0. 0.] [0. 0. 0. 0.0001 0.] [0. 0. 0. 0. 0.0001]]</pre>	<pre>[[-0.9993004] [0.99940033] [0.99940033] [0.99940033] [-3.9790115] [0.99940033]]</pre>

Ridge Regression in Primal Form (when $m > d$)

$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$ is invertible for $\lambda > 0$,

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

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$

- Primal form used for $m > d$, over determined
 - always find inverse in smaller in dimension
- \mathbf{X} is $m \times d$
- $\mathbf{X}^T \mathbf{X}$ is $d \times d$
- \mathbf{I} is $d \times d$

Ridge Regression in Dual Form (when $m < d$)

$(\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})$ is invertible for $\lambda > 0$,

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

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{X}_{new}) = \mathbf{X}_{new} \hat{\mathbf{w}}$

- Dual form used for $m < d$, under determined
 - always find inverse in smaller in dimension
- \mathbf{X} is $m \times d$
- $\mathbf{X}^T \mathbf{X}$ is $m \times m$
- \mathbf{I} is $m \times m$

Derivation as homework (see tutorial 6).

Hint: start off with $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$ and make use of $\mathbf{w} = \mathbf{X}^T \mathbf{a}$ and $\mathbf{a} = \lambda^{-1}(\mathbf{y} - \mathbf{X} \mathbf{w})$, $\lambda > 0$

Given $(X^T X + \lambda I) w = X^T y$

derive $\hat{w} = X^T (X X^T + \lambda I)^{-1} y$

Hint: $w = X^T a$ where

$$a = \lambda^{-1} (y - Xw), \lambda > 0$$

$$\begin{aligned} : X^T X w + \lambda \cancel{I} w &= X^T y \\ \lambda w &= X^T y - X^T X w \\ &= X^T (y - Xw) \end{aligned}$$

$$\begin{aligned} w &= \cancel{\lambda^{-1}} X^T \underline{(y - Xw)} \\ &= X^T a \quad \text{--- (1)} \end{aligned}$$

$$: a = \lambda^{-1} (y - Xw)$$

$$\begin{aligned} \lambda a &= y - Xw \\ &= y - X X^T a \end{aligned}$$

$$X X^T a + \lambda a = y$$

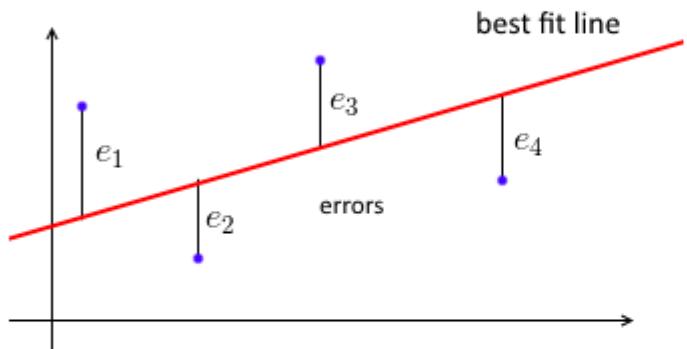
$$(X X^T + \lambda) a = y$$

$$a = \underline{(X X^T + \lambda I)}^{-1} y \quad \text{--- (2)}$$

Short answer: no difference between Primal and Dual - it's only about the way of arriving to the solution. Kernel ridge regression is essentially the same as usual ridge regression, but uses the kernel trick to go non-linear.

Linear Regression

First of all, a usual Least Squares Linear Regression tries to fit a straight line to the set of data points in such a way that the sum of squared errors is minimal.



We parametrize the best fit line with \mathbf{w} and for each data point (\mathbf{x}_i, y_i) we want $\mathbf{w}^T \mathbf{x}_i \approx y_i$. Let $e_i = y_i - \mathbf{w}^T \mathbf{x}_i$ be the error - the distance between predicted and true values. So our goal is to

minimize the sum of squared errors $\sum e_i^2 = \|\mathbf{e}\|^2 = \|X\mathbf{w} - \mathbf{y}\|^2$ where $X = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$ - a data matrix with each \mathbf{x}_i being a row, and $\mathbf{y} = (y_1, \dots, y_n)$ a vector with all y_i 's.

Thus, the objective is $\min_{\mathbf{w}} \|X\mathbf{w} - \mathbf{y}\|^2$, and the solution is $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$ (known as "Normal Equation").

For a new unseen data point \mathbf{x} we predict its target value \hat{y} as $\hat{y} = \mathbf{w}^T \mathbf{x}$.

Ridge Regression

When there are many correlated variables in linear regression models, the coefficients \mathbf{w} can become poorly determined and have lots of variance. One of the solutions to this problem is to restrict weights \mathbf{w} so they don't exceed some budget C . This is equivalent to using L_2 -regularization, also known as "weight decay": it will decrease the variance at the cost of sometimes missing the correct results (i.e. by introducing some bias).

The objective now becomes $\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$, with λ being the regularization parameter.

By going through the math, we obtain the following solution:

$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$. It's very similar to the usual linear regression, but here we add λ to each diagonal element of $\mathbf{X}^T \mathbf{X}$.

Note that we can re-write \mathbf{w} as $\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$ (see [here](#) for details). For a new unseen data point \mathbf{x} we predict its target value \hat{y} as $\hat{y} = \mathbf{x}^T \mathbf{w} = \mathbf{x}^T \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$. Let $\boldsymbol{\alpha} = (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$. Then $\hat{y} = \mathbf{x}^T \mathbf{X}^T \boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \cdot \mathbf{x}^T \mathbf{x}_i$.

Ridge Regression Dual Form

We can have a different look at our objective - and define the following quadratic program problem:

$$\min_{\mathbf{e}, \mathbf{w}} \sum_{i=1}^n e_i^2 \text{ s.t. } e_i = y_i - \mathbf{w}^T \mathbf{x}_i \text{ for } i = 1 \dots n \text{ and } \|\mathbf{w}\|^2 \leq C.$$

It's the same objective, but expressed somewhat differently, and here the constraint on the size of \mathbf{w} is explicit. To solve it, we define the Lagrangian $\mathcal{L}_p(\mathbf{w}, \mathbf{e}; C)$ - this is the primal form that contains primal variables \mathbf{w} and \mathbf{e} . Then we optimize it w.r.t. \mathbf{e} and \mathbf{w} . To get the dual formulation, we put \mathbf{e} and \mathbf{w} back to $\mathcal{L}_p(\mathbf{w}, \mathbf{e}; C)$.

So, $\mathcal{L}_p(\mathbf{w}, \mathbf{e}; C) = \|\mathbf{e}\|^2 + \beta^T (\mathbf{y} - X\mathbf{w} - \mathbf{e}) - \lambda (\|\mathbf{w}\|^2 - C)$. By taking derivatives w.r.t. \mathbf{w} and \mathbf{e} , we obtain $\mathbf{e} = \frac{1}{2} \beta$ and $\mathbf{w} = \frac{1}{2\lambda} X^T \beta$. By letting $\alpha = \frac{1}{2\lambda} \beta$, and putting \mathbf{e} and \mathbf{w} back to $\mathcal{L}_p(\mathbf{w}, \mathbf{e}; C)$, we get dual Lagrangian

$\mathcal{L}_d(\alpha, \lambda; C) = -\lambda^2 \|\alpha\|^2 + 2\lambda \alpha^T \mathbf{y} - \lambda \|X^T \alpha\| - \lambda C$. If we take a derivative w.r.t. α , we get $\alpha = (X X^T - \lambda I)^{-1} \mathbf{y}$ - the same answer as for usual Kernel Ridge regression. There's no need to take a derivative w.r.t λ - it depends on C , which is a regularization parameter - and it makes λ regularization parameter as well.

Next, put α to the primal form solution for \mathbf{w} , and get $\mathbf{w} = \frac{1}{2\lambda} X^T \beta = X^T \alpha$. Thus, the dual form gives the same solution as usual Ridge Regression, and it's just a different way to come to the same solution.

<https://stats.stackexchange.com/questions/92672/difference-between-primal-dual-and-kernel-ridge-regression>

Polynomial Regression

Motivation: nonlinear decision surface

- Based on the sum of products of the variables
- E.g. when the input dimension is $d=2$,

a polynomial function of degree = 2 is:

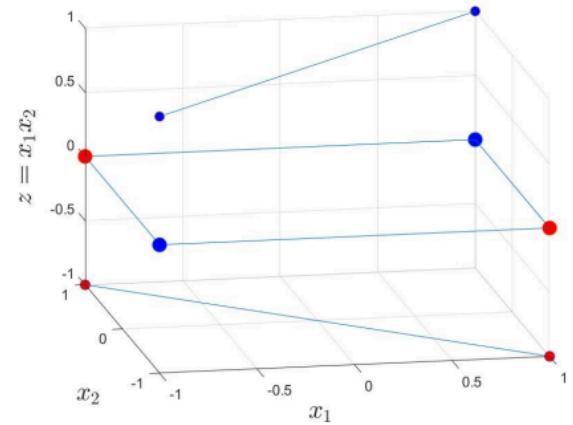
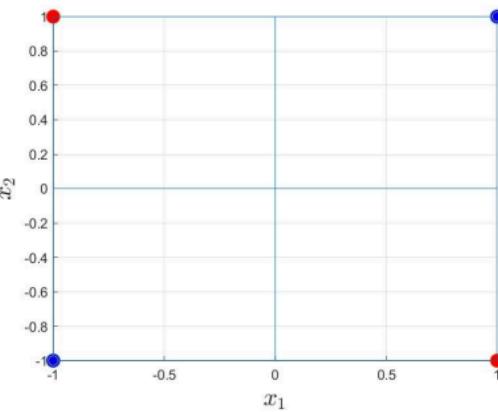
$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2.$$

Example

XOR problem

$$\begin{aligned} \mathbf{x}_1 &= [+1 \quad +1]^T & y_1 &= +1 \\ \mathbf{x}_2 &= [-1 \quad +1]^T & y_2 &= -1 \\ \mathbf{x}_3 &= [+1 \quad -1]^T & y_3 &= -1 \\ \mathbf{x}_4 &= [-1 \quad -1]^T & y_4 &= +1 \end{aligned}$$

$$f_{\mathbf{w}}(\mathbf{x}) = x_1 x_2$$



- Must need non linear regression to fit the situation, not possible to split into non linear, not linearly separable
- Pic on left shows how the axis of $x_1 x_2$ separate the blue and red dots

Polynomial Expansion

- The linear model $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ can be written as

$$\begin{aligned} f_{\mathbf{w}}(\mathbf{x}) &= \mathbf{x}^T \mathbf{w} \\ &= \sum_{i=0}^d x_i w_i, \quad x_0 = 1 \\ &= w_0 + \sum_{i=1}^d x_i w_i. \end{aligned}$$

- By including additional terms involving the products of pairs of components of \mathbf{x} , we obtain a quadratic model:

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

2nd order: $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2$

3rd order: $f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2 + \sum_{i=1}^d \sum_{j=1}^d \sum_{k=1}^d w_{ijk} x_i x_j x_k, \quad d = 2$

$$\text{C}_{n+r}^r = \frac{(n+r-1)!}{r! (n-1)!}$$

Ref: Duda, Hart, and Stork, "Pattern Classification", 2001 (Chp.5)

- d is number of feature
- 2nd order has 6 parameter/weight
 - $x_1 x_2$ is the same as $x_2 x_1$ need to combine
- 3rd order has 10 parameter/weight
 - same as 2nd order, + 4
- nCr, ????

Polynomial Regression

Generalized Linear Discriminant Function

- In general:

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

*1
C(3,1)
C(3,2)
C(3,3)*

Weierstrass Approximation Theorem: Every continuous function defined on a closed interval $[a, b]$ can be uniformly approximated as closely as desired by a polynomial function.

- Suppose f is a continuous real-valued function defined on the real interval $[a, b]$.
 - For every $\varepsilon > 0$, there exists a polynomial p such that for all x in $[a, b]$, we have $|f(x) - p(x)| < \varepsilon$.
- (Ref: https://en.wikipedia.org/wiki/Stone%20%26%20Weierstrass_theorem)

Notes:

- For high dimensional input features (large *d value*) and high *polynomial order*, the number of polynomial terms becomes explosive! (i.e., grows exponentially)
- For high dimensional problems, polynomials of order larger than 3 is seldom used.

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

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{P}\mathbf{w} \quad (\text{Note: } \mathbf{P} \triangleq \mathbf{P}(\mathbf{X}) \text{ for symbol simplicity})$$

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

$$\text{where } \mathbf{p}_l^T \mathbf{w} = [1, x_{l,1}, \dots, x_{l,d}, \dots, x_{l,i} x_{l,j}, \dots, x_{l,i} x_{l,j} x_{l,k}, \dots]$$

$l = 1, \dots, m$; d denotes the dimension of input features; m denotes the number of samples

$$\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \\ \vdots \\ w_{ij} \\ \vdots \\ w_{ijk} \end{bmatrix}$$

Example 3

Training set $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^m$

2nd order polynomial model

$$\{x_1 = 0, x_2 = 0\} \rightarrow \{y = -1\}$$

$$\{x_1 = 1, x_2 = 1\} \rightarrow \{y = -1\}$$

$$\{x_1 = 1, x_2 = 0\} \rightarrow \{y = +1\}$$

$$\{x_1 = 0, x_2 = 1\} \rightarrow \{y = +1\}$$

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_{12} x_1 x_2 + w_{11} x_1^2 + w_{22} x_2^2$$

$$= [1 \ x_1 \ x_2 \ x_1 x_2 \ x_1^2 \ x_2^2] \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_{12} \\ w_{11} \\ w_{22} \end{bmatrix}$$

Python
demo 3

Stack the 4 training samples as a matrix

$$\mathbf{P} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & x_{1,1}x_{1,2} & x_{1,1}^2 & x_{1,2}^2 \\ 1 & x_{2,1} & x_{2,2} & x_{2,1}x_{2,2} & x_{2,1}^2 & x_{2,2}^2 \\ 1 & x_{3,1} & x_{3,2} & x_{3,1}x_{3,2} & x_{3,1}^2 & x_{3,2}^2 \\ 1 & x_{4,1} & x_{4,2} & x_{4,1}x_{4,2} & x_{4,1}^2 & x_{4,2}^2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

- Number of rows in P = number of samples
- Number of columns in P = number of terms in the polynomial expansion of with d feature up till order O
- Sub in the value of x1 and x2, and solve the linear system $\mathbf{y} = \mathbf{P}\mathbf{w}$

Example 3 (cont'd)

Training set

2nd order polynomial model

$$\mathbf{P} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & x_{1,1}x_{1,2} & x_{1,1}^2 & x_{1,2}^2 \\ 1 & x_{2,1} & x_{2,2} & x_{2,1}x_{2,2} & x_{2,1}^2 & x_{2,2}^2 \\ 1 & x_{3,1} & x_{3,2} & x_{3,1}x_{3,2} & x_{3,1}^2 & x_{3,2}^2 \\ 1 & x_{4,1} & x_{4,2} & x_{4,1}x_{4,2} & x_{4,1}^2 & x_{4,2}^2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$$\hat{\mathbf{w}} = \mathbf{P}^T (\mathbf{P} \mathbf{P}^T)^{-1} \mathbf{y}$$

$$= \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 6 & 3 & 3 \\ 1 & 3 & 3 & 1 \\ 1 & 3 & 1 & 3 \end{bmatrix}^{-1} \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 1 \\ -4 \\ 1 \\ 1 \end{bmatrix}$$

Python
demo 3

Example 3 (cont'd) Prediction

Test set

- Test point 1: $\{x_1 = 0.1, x_2 = 0.1\} \rightarrow \{y = \text{class } -1 \text{ or } +1?\}$
- Test point 2: $\{x_1 = 0.9, x_2 = 0.9\} \rightarrow \{y = \text{class } -1 \text{ or } +1?\}$
- Test point 3: $\{x_1 = 0.1, x_2 = 0.9\} \rightarrow \{y = \text{class } -1 \text{ or } +1?\}$
- Test point 4: $\{x_1 = 0.9, x_2 = 0.1\} \rightarrow \{y = \text{class } -1 \text{ or } +1?\}$

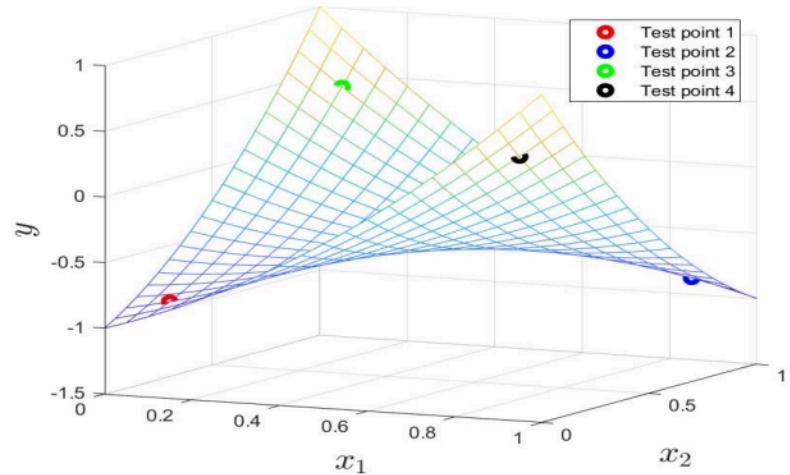
$$\hat{\mathbf{y}} = \mathbf{P}_{new} \hat{\mathbf{w}}$$

$$= \begin{bmatrix} 1 & 0.1 & 0.1 & 0.01 & 0.01 & 0.01 \\ 1 & 0.9 & 0.9 & 0.81 & 0.81 & 0.81 \\ 1 & 0.1 & 0.9 & 0.09 & 0.01 & 0.81 \\ 1 & 0.9 & 0.1 & 0.09 & 0.81 & 0.01 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ 1 \\ -4 \\ 1 \\ 1 \end{bmatrix}$$

$$[1 \ x_1 \ x_2 \ x_1x_2 \ x_1^2 \ x_2^2]$$

$$\hat{f}_w^c(\mathbf{P}(\mathbf{X}_{new})) = \text{sign}(\hat{\mathbf{y}}) = \text{sign}\left(\begin{bmatrix} -0.82 \\ -0.82 \\ 0.46 \\ 0.46 \end{bmatrix}\right)$$

$$= \begin{bmatrix} -1 \\ -1 \\ +1 \\ +1 \end{bmatrix} \begin{array}{l} \text{-----} \rightarrow \text{Class } -1 \\ \text{-----} \rightarrow \text{Class } -1 \\ \text{-----} \rightarrow \text{Class } +1 \\ \text{-----} \rightarrow \text{Class } +1 \end{array}$$



- non linear decision surface, each value of x_1 and x_2 mapped onto the surface, then compare to threshold value 0 to get which class it is

```
#EE2211 Lecture 6 Demo 3 Polynomial regression
import numpy as np
from numpy.linalg import inv
from numpy.linalg import matrix_rank
from sklearn.preprocessing import PolynomialFeatures
X = np.array([[0, 0], [1, 1], [1, 0], [0, 1]])
y = np.array([[-1], [-1], [1], [1]])
## Generate polynomial features
order = 2
poly = PolynomialFeatures(order)
print(poly)
P = poly.fit_transform(X)
print("matrix P")
print(P)

print("*****")
#print(matrix_rank(P))
#PY = np.vstack((P.T, y.T))
#print(matrix_rank(PY.T))
```

```
PolynomialFeatures()
matrix P
[[1. 0. 0. 0. 0. 0.]
 [1. 1. 1. 1. 1. 1.]
 [1. 1. 0. 1. 0. 0.]
 [1. 0. 1. 0. 0. 1.]]
*****
```

- input to polynomial fit transform doesn't need to pad 1, it automatically does it

```

## dual solution m < d (without ridge)
w_dual = P.T @ inv(P @ P.T) @ y
print("Under-determined system")
print("Unique constrained solution, no ridge")
print(w_dual)

print("*****")
print("Approximation with dual ridge regression")
print(P.shape)
reg_L2 = 0.0001*np.identity(P.shape[0]) #number of rows of P = Dual I
print(reg_L2)
w_dual_ridge = P.T @ (inv(P @ P.T + reg_L2)) @ y
print(w_dual_ridge)

print("*****")
## primal ridge
print("Approximation with primal ridge regression")
print(P.shape)
reg_L = 0.0001*np.identity(P.shape[1]) #number of columns of P = Primal I
print(reg_L)
w_primal_ridge = inv(P.T @ P + reg_L) @ P.T @ y
print(w_primal_ridge)

```

```

Under-determined system
Unique constrained solution, no ridge
[[-1.]
 [ 1.]
 [ 1.]
 [ 1.]
 [-4.]
 [ 1.]]
*****
Approximation with dual ridge regression
(4, 6)
[[0.0001 0.    0.    0.    0.    ]
 [0.    0.0001 0.    0.    0.    ]
 [0.    0.    0.0001 0.    0.    ]
 [0.    0.    0.    0.0001 0.    ]
 [0.    0.    0.    0.    0.0001]]
[[-0.9993004 ]
 [ 0.99940033]
 [ 0.99940033]
 [ 0.99940033]
 [-3.99790115]
 [ 0.99940033]]
*****
Approximation with primal ridge regression
(4, 6)
[[0.0001 0.    0.    0.    0.    0.    ]
 [0.    0.0001 0.    0.    0.    0.    ]
 [0.    0.    0.0001 0.    0.    0.    ]
 [0.    0.    0.    0.0001 0.    0.    ]
 [0.    0.    0.    0.    0.0001 0.    ]
 [0.    0.    0.    0.    0.    0.0001]]
[[-0.9993004 ]
 [ 0.99940033]
 [ 0.99940033]
 [ 0.99940033]
 [-3.99790115]
 [ 0.99940033]]

```

```

#EE2211 Lecture 6 Demo 3 Testing/prediction
import numpy as np
from numpy.linalg import inv
from numpy.linalg import matrix_rank
from sklearn.preprocessing import PolynomialFeatures
X = np.array([[0, 0], [1, 1], [1, 0], [0, 1]])
y = np.array([-1, -1, 1, 1])
## Generate polynomial features
order = 2
poly = PolynomialFeatures(order)
print(poly)
P = poly.fit_transform(X)
print("matrix P")
print(P)
print("Under-determined system")
#print(matrix_rank(P))
#PY = np.vstack((P.T, y.T))
#print(matrix_rank(PY.T))

## dual solution m < d (without ridge)
w_dual = P.T @ inv(P @ P.T) @ y
print("Unique constrained solution, no ridge")
print(w_dual)

#testing
print("Prediction")
Xnew= np.array([ [0.1, 0.1], [0.9, 0.9], [0.1, 0.9], [0.9, 0.1]])
Pnew = poly.fit_transform(Xnew)
Ynew=Pnew@w_dual
print(Ynew)
print(np.sign(Ynew))

```

```

PolynomialFeatures()
matrix P
[[1.  0.  0.  0.  0.  0.]
 [1.  1.  1.  1.  1.  1.]
 [1.  1.  0.  1.  0.  0.]
 [1.  0.  1.  0.  0.  1.]]
Under-determined system
Unique constrained solution, no ridge
[[-1.]
 [ 1.]
 [ 1.]
 [ 1.]
 [-4.]
 [ 1.]]
Prediction
[[ -0.82]
 [ -0.82]
 [  0.46]
 [  0.46]]
[[-1.]
 [-1.]
 [ 1.]
 [ 1.]]

```

Ridge Regression in Primal Form ($m > d$)

For $\lambda > 0$,

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

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{P}(\mathbf{X}_{new})) = \mathbf{P}_{new} \hat{\mathbf{w}}$

Ridge Regression in Dual Form ($m < d$)

For $\lambda > 0$,

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

Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{P}(\mathbf{X}_{new})) = \mathbf{P}_{new} \hat{\mathbf{w}}$

Note: Change \mathbf{X} to \mathbf{P} with reference to slides 15/16; m & d refers to the size of \mathbf{P} (not \mathbf{X})

For Regression Applications

- Learn **continuous** valued y using either primal form or dual form
- Prediction: $\hat{f}_{\mathbf{w}}(\mathbf{P}(\mathbf{X}_{new})) = \mathbf{P}_{new} \hat{\mathbf{w}}$

For Classification Applications

- Learn **discrete** valued y ($y \in \{-1, +1\}$) or \mathbf{Y} (one-hot) using either primal form or dual form
- Binary Prediction: $\hat{f}_{\mathbf{w}}^c(\mathbf{P}(\mathbf{X}_{new})) = \text{sign}(\mathbf{P}_{new} \hat{\mathbf{w}})$
- Multi-Category Prediction: $\hat{f}_{\mathbf{w}}^c(\mathbf{P}(\mathbf{X}_{new})) = \arg \max_{k=1, \dots, C} (\mathbf{P}_{new} \hat{\mathbf{W}}(:, k))$

In Class Quiz

Suppose we would like to build up one full third-order polynomial model for only one input feature and single output, the polynomial model has three parameters to learn.

Jin Yueming
15

True



False



- One variable (feature) 3rd order means $y = c + w_1x + w_2x^2 + w_3x^3$
- Four parameters/weights

If there are four data samples with two input features each, the full second-order polynomial model is an over-determined system.

Jin Yueming
14

True



False



- $f_w(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_{12}x_1x_2 + w_{11}x_1^2 + w_{22}x_2^2$
- 6 columns, 4 rows, under determined

The polynomial model can be used to solve problems with nonlinear decision boundary.

Jin Yueming

True



False



Lec7: Over-fitting, bias/variance trade-off

- Linear is special case of polynomial => use “**P**” instead of “**X**” from now on
- Training/Learning (primal) on training set

$$\hat{\mathbf{w}} = (\mathbf{P}_{train}^T \mathbf{P}_{train})^{-1} \mathbf{P}_{train}^T \mathbf{y}_{train}$$

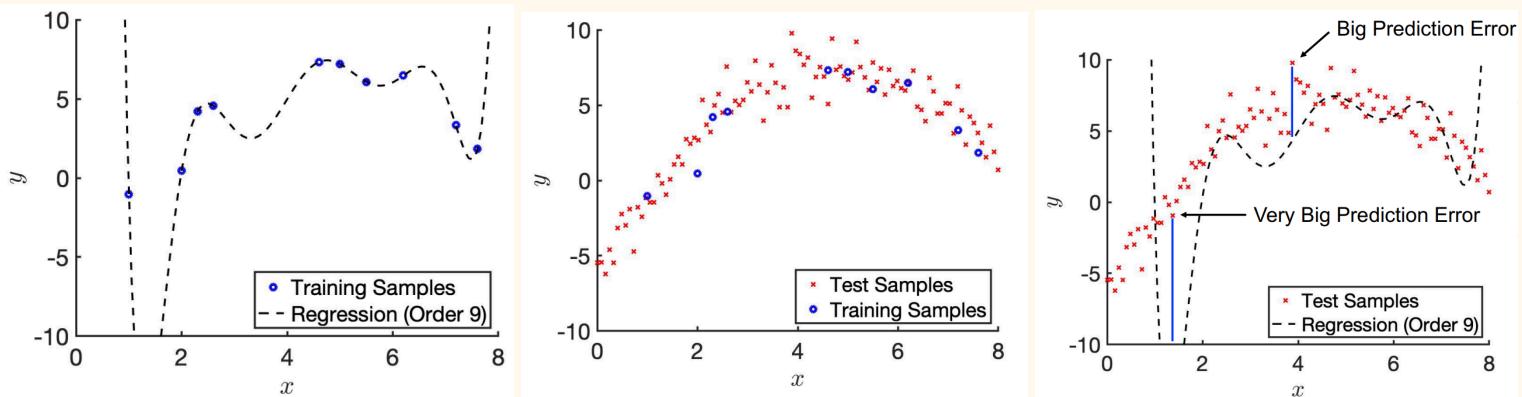
- Prediction/Testing/Evaluation on test set

$$\hat{\mathbf{y}}_{test} = \mathbf{P}_{test} \hat{\mathbf{w}}$$

- There should be **zero** overlap between training & test sets
- Important goal of regression: prediction on **new unseen data**, i.e., test set
- Why is test set important for evaluation?
 - Ideally training set >> test set

Overfitting, Underfitting & Model Complexity

Overfitting



- Regression fits the training data well but not the test data
- has very big prediction error when using test data (MSE is big)
- Because training data set is small while trying to fit a very complicated model of order 9, a lot parameter
- Training good, test bad \Rightarrow over fitting

- **Overfitting** occurs when model predicts the training data well, but predicts new data (e.g., from test set) poorly

- **Reason 1**

- Model is too complex for the data
- Previous example: Fit order 9 polynomial to 10 data points

- **Reason 2**

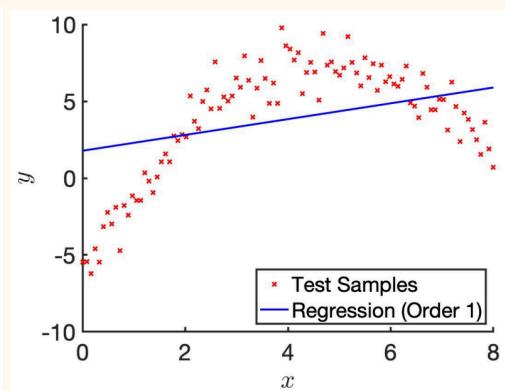
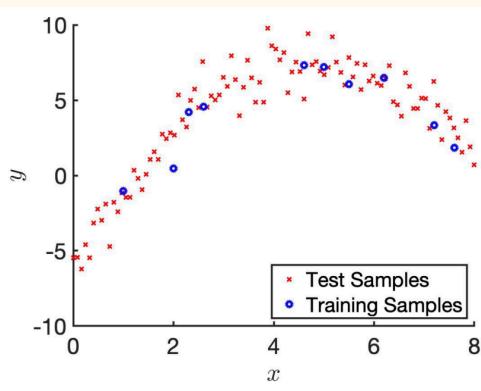
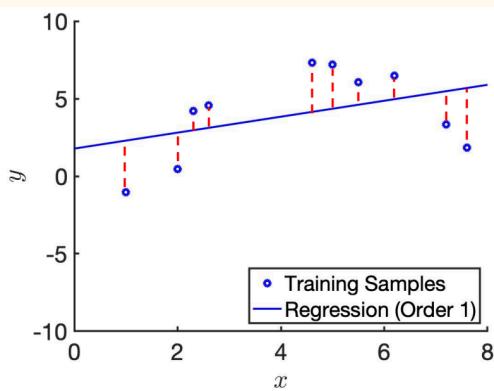
- Too many features but number of training samples too small
- Even linear model can overfit, e.g., linear model with 9 input features (i.e., w is 10-D) and 10 data points in training set => data might not be enough to estimate 10 unknowns well

- **Solutions**

- Use **simpler models** (e.g., lower order polynomial)
- Use **regularization** (see next part of lecture)

- Can use less features also

Underfitting



- Regression is bad for the training data and the test data => under fitting
- Means model is too simple

- **Underfitting** is the inability of trained model to predict the targets in the training set

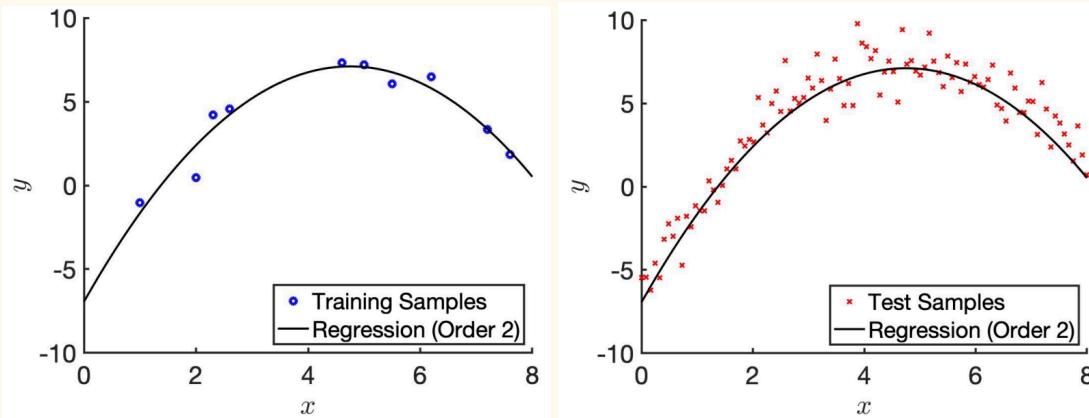
- **Reason 1**

- Model is too simple for the data
- Previous example: Fit order 1 polynomial to 10 data points that came from an order 2 polynomial
- **Solution:** Try more complex model

- **Reason 2**

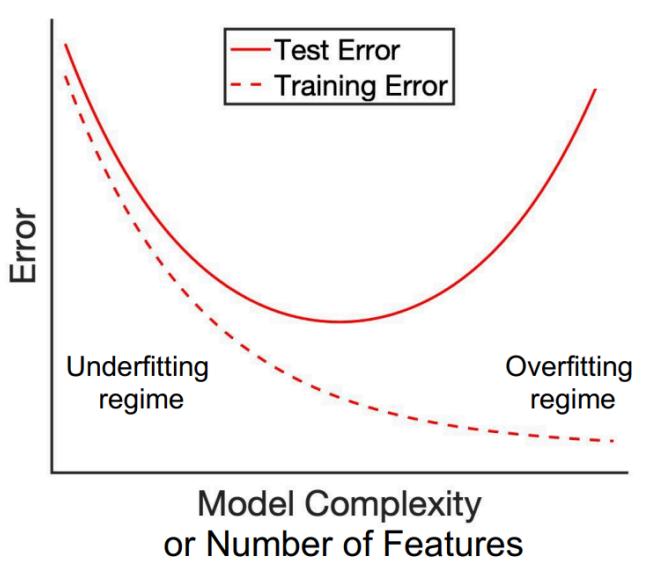
- Features are not informative enough
- **Solution:** Try to develop more informative features
 - Need more features, or have features that doesn't tell much
 - Be better at feature selection

Good Fit



- Regression does not pass through all the training data, the MSE is actually larger than the 9 order overfitting case
- But this is due to quadratic only having 3 parameters, so the regression is also still good
- The MSE to the test data is also very good

Overfitting / Underfitting Schematic



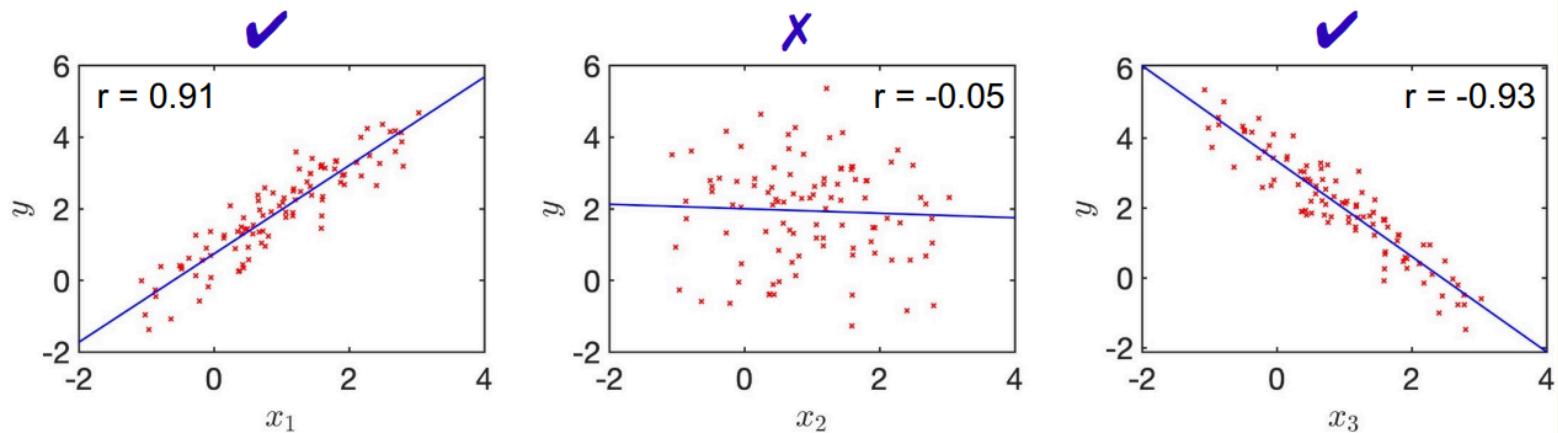
- When both training and test error are high, is when model complexity is too low or number of features is not enough \Rightarrow underfitting
- When training error is low but test error is high, is when model complexity is too high or number of features is too much \Rightarrow over fitting
 - if the model is very complicated and has a lot of parameter, the parameter can be chosen to fit the training data very well
- note the test error has a u-shape indicating there is a sweet spot between under and over fitting

Feature Selection

- Less features might reduce overfitting
 - Want to discard useless features & keep good features, so perform feature selection
 - Feature selection procedure
 - Step 1: feature selection in **training** set
 - Step 2: fit model using selected features in **training** set
 - Step 3: evaluate trained model using **test** set
 - Very common mistake
 - Feature selection with test set (or full dataset) leads to inflated performance
 - Do not perform feature selection with test data
- Do not train with test set
 - Do not feature select from test set
 - Both will lead to information leakage
 - If need to compare two models, give both the same training set and test with the same test set

Selecting Features With Pearson's r

- Given features x , we want to predict target y
- Assume x & y both continuous
- Compute Pearson's correlation coefficient between each feature & target y in the training set
 - Pearson's correlation r measures linear relationship between two variables



- Two options
 - Option 1: Pick K features with largest absolute correlations
 - Option 2: Pick all features with absolute correlations $> C$
 - K & C are “magic” numbers set by the ML practitioner
- Other metrics besides Pearson's correlation are possible

Regularisation

- Regularization is an umbrella term that includes methods forcing learning algorithm to build less complex models.
- **Motivation 1:** Solve an ill-posed problem
 - For example, estimate 10th order polynomial with just 5 datapoints
- **Motivation 2:** Reduce overfitting
- For example, in previous lecture, we added $\lambda \mathbf{w}^T \mathbf{w}$:

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

- Minimizing with respect to \mathbf{w} , primal solution is

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

- For $\lambda > 0$, matrix becomes invertible (**Motivation 1**)
- $\hat{\mathbf{w}}$ might also perform better in test set, i.e., reduces overfitting (**Motivation 2**) – will show example later
 - ill posed problem means trying to solve under determine system
 - with λ , \mathbf{w} reduces fitting by making sure features that does not contribute much have low weights, so can fix over fitting
- Consider minimization from previous slide

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

Cost function quantifying data
fitting error in training set

Regularization

- Consider minimization from previous slide

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

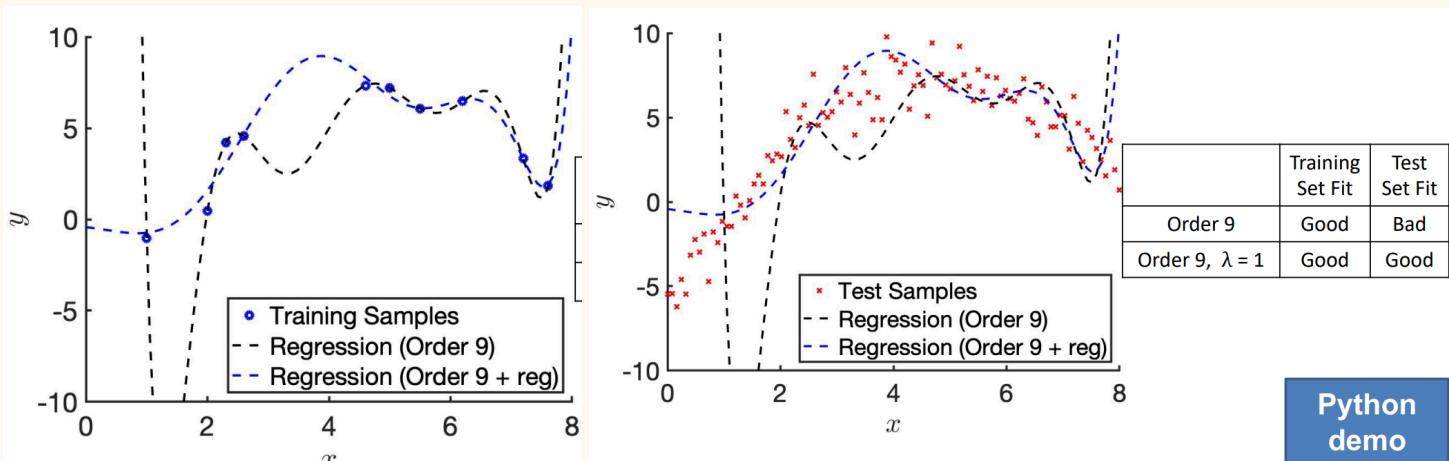
- $\mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_d^2$ L2 - Regularization

- Encourage w_0, \dots, w_d to be small (also called shrinkage or weight-decay) => constrain model complexity
- More generally, most machine learning algorithms can be formulated as the following optimization problem

$$\underset{\mathbf{w}}{\operatorname{argmin}} \text{Data-Loss}(\mathbf{w}) + \lambda \text{Regularization}(\mathbf{w})$$

- **Data-Loss(w)** quantifies fitting error to training set given parameters \mathbf{w} : smaller error => better fit to training data
- **Regularization(w)** penalizes more complex models
 - Because its squared, its called L2 - Regularization
 - Bigger λ means regularisation dominates, forces model to be less complex at the cost of increasing MSE to training data

Regularisation Example

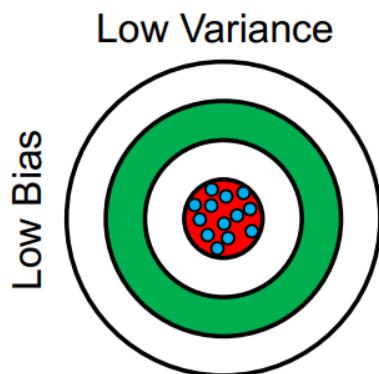


- Can see that with Ridge regression, the blue curve does not fit the training data as well anymore. But it fit the overall data much better

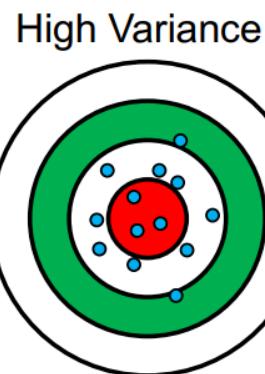
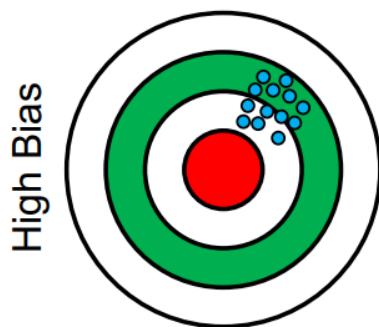
Bias-Variance Trade-off

- Suppose we are trying to predict red target below:

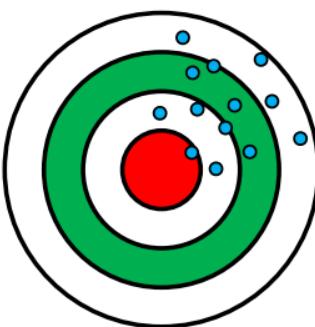
Low Bias: blue predictions on average close to red target
Low Variance: low variability among predictions



High Bias: blue predictions on average not close to red target
Low Variance: Low variability among predictions



Low Bias: blue predictions on average close to red target
High Variance: large variability among predictions

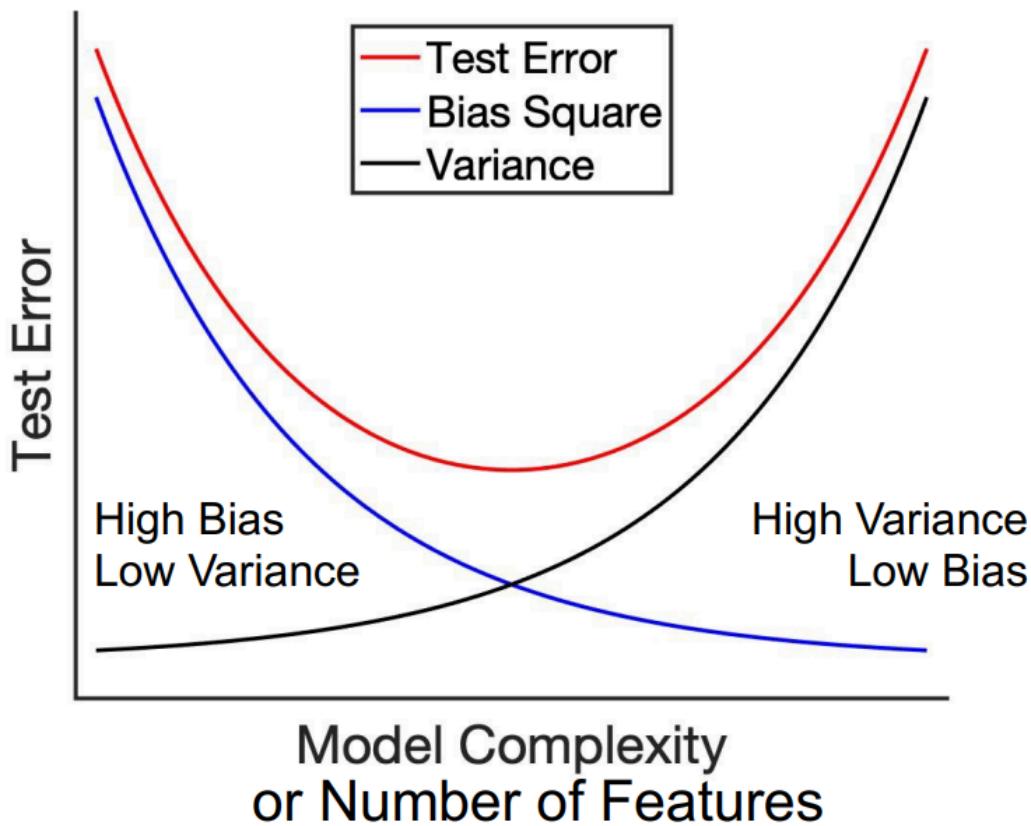


- Both bias and variance contribute to prediction error

Bias-Variance Decomposition Theorem

- **Test error** = Bias Squared + Variance + Irreducible Noise
 - Mathematical details in optional uploaded material (won't be tested)
- “**Variance**” refers to variability of prediction models across different training sets
 - In previous example, every time the training set of 10 samples changes, the trained model changes
 - “Variance” quantifies variability across trained models
- “**Bias**” refers to how well an average prediction model will perform
 - In previous example, every time the training set of 10 samples changes, the trained model changes
 - If we average the trained models, how well will this average trained model perform?
- “**Irreducible Noise**” reflects the fact that even if we are perfect modelers, it might not be possible to predict target y with 100% accuracy from feature(s) x

- Test error = Bias Squared + Variance + Irreducible Noise

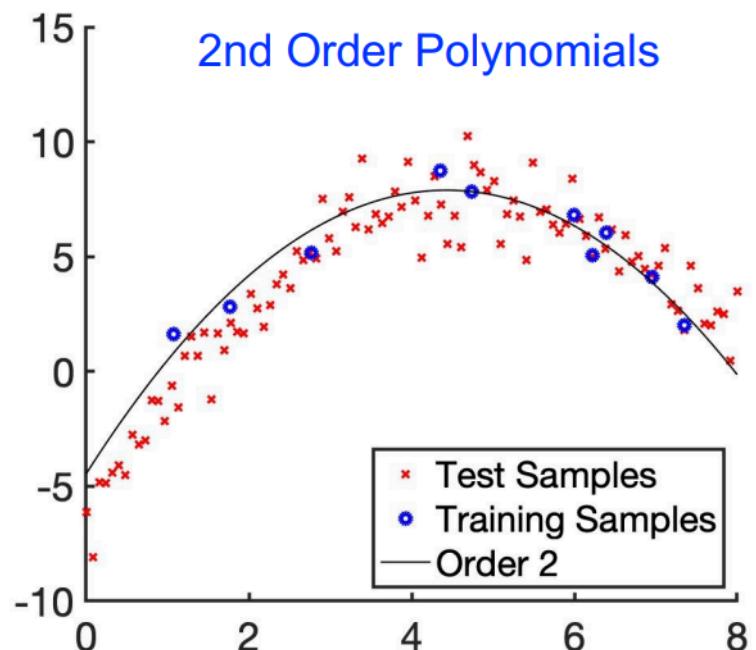
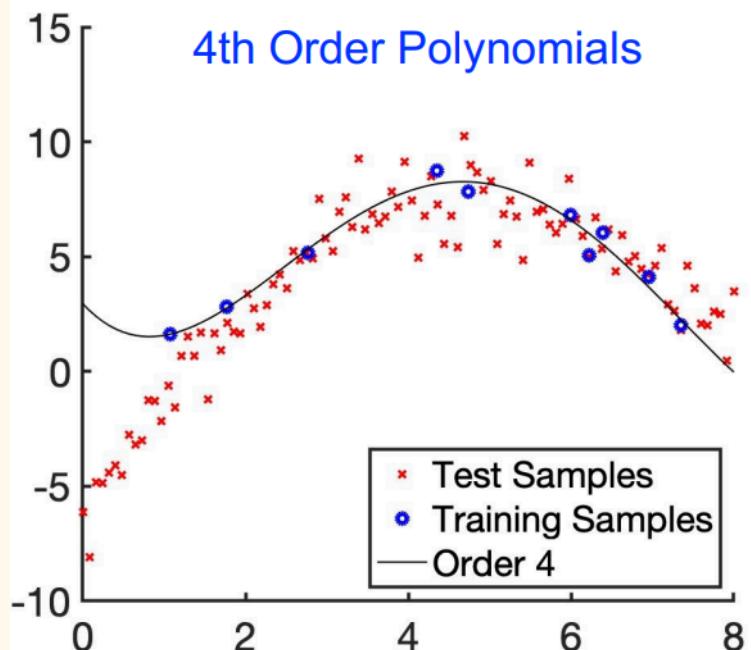


- test error with u shape can be explained here
- With simple models, predictions has high bias but low variance
- With complex models, predictions has low bias but high variance
- This is the bias variance trade off

Example

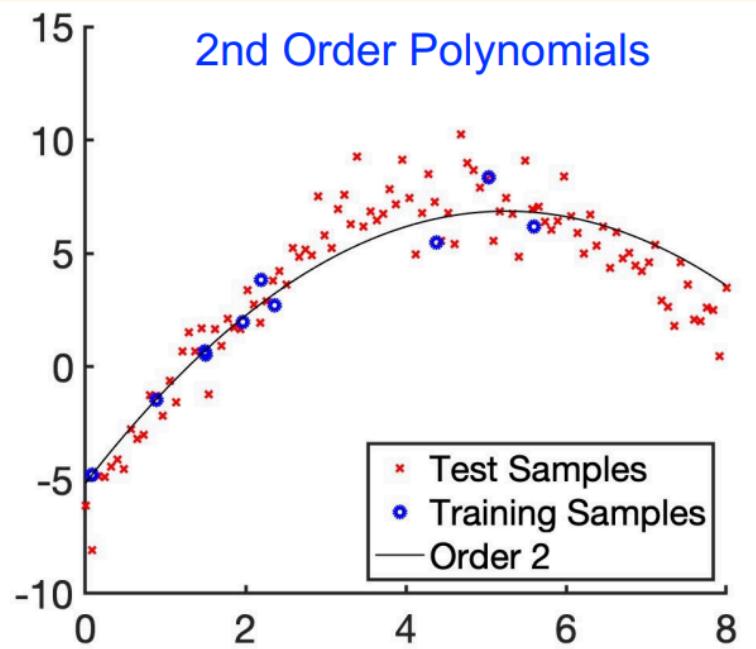
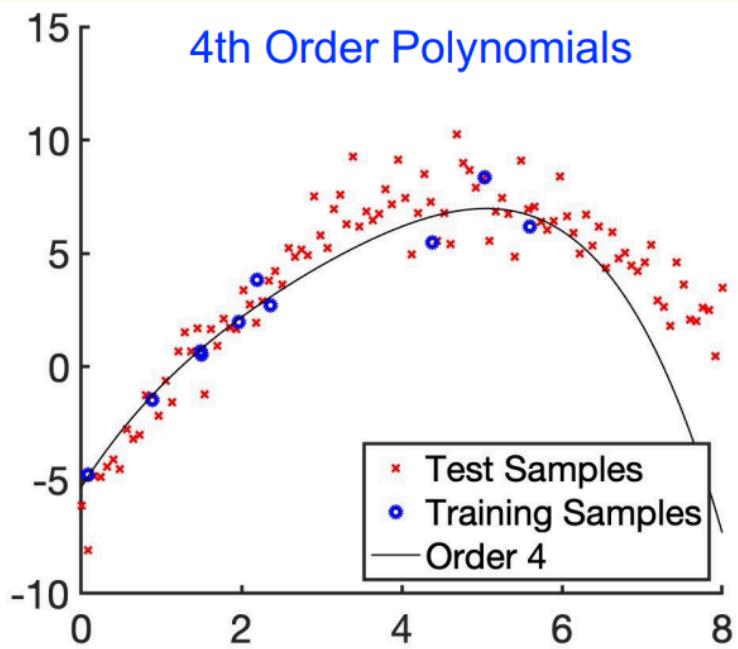
- Simulate data from order 2 polynomial (+ noise)
- Randomly sample 10 training samples each time
- Fit with order 2 polynomial
- Fit with order 4 polynomial
- The data points are the same for all
- The same two models is trained on the same data set each time the experiment is repeated
 - each time the sample is randomly chosen
- Note the variance the 4th order one has

Test 1



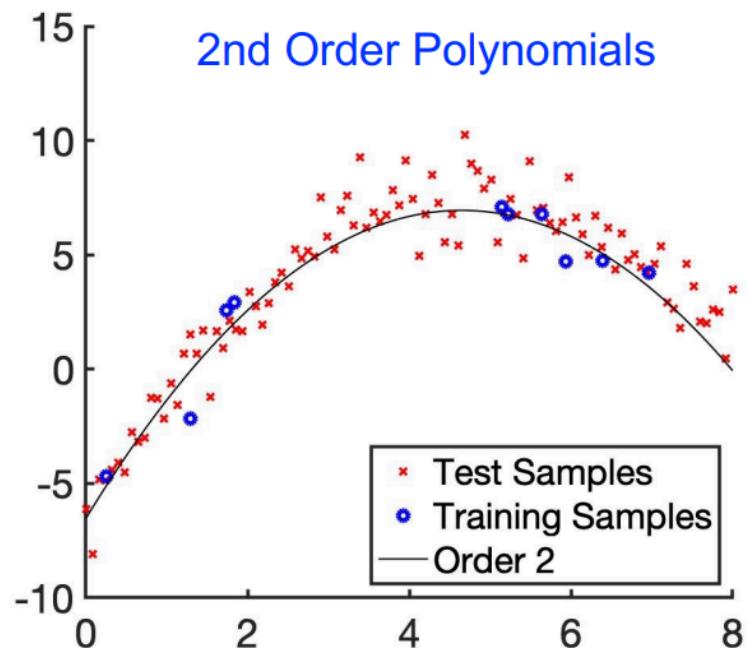
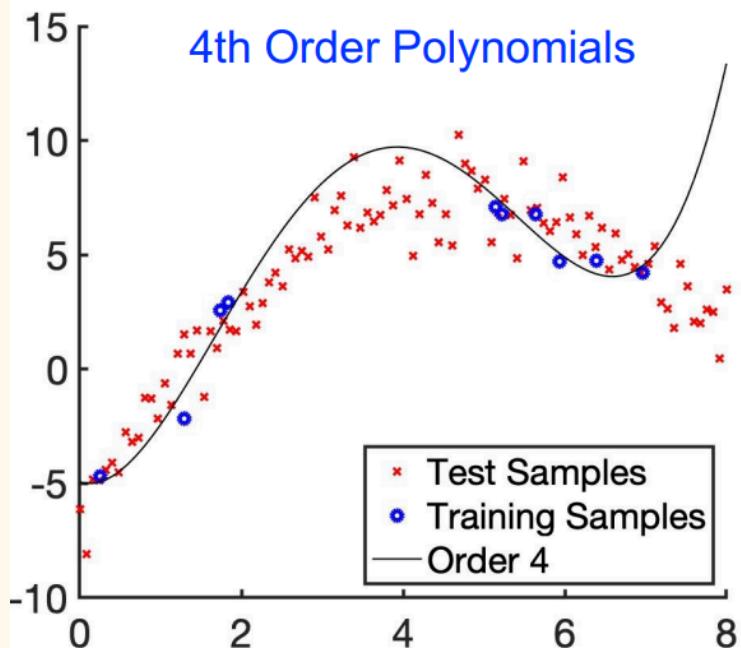
- 4th order has big error on LHS
- 2nd order fit training and test well
- Both fit training quite well

Test 2



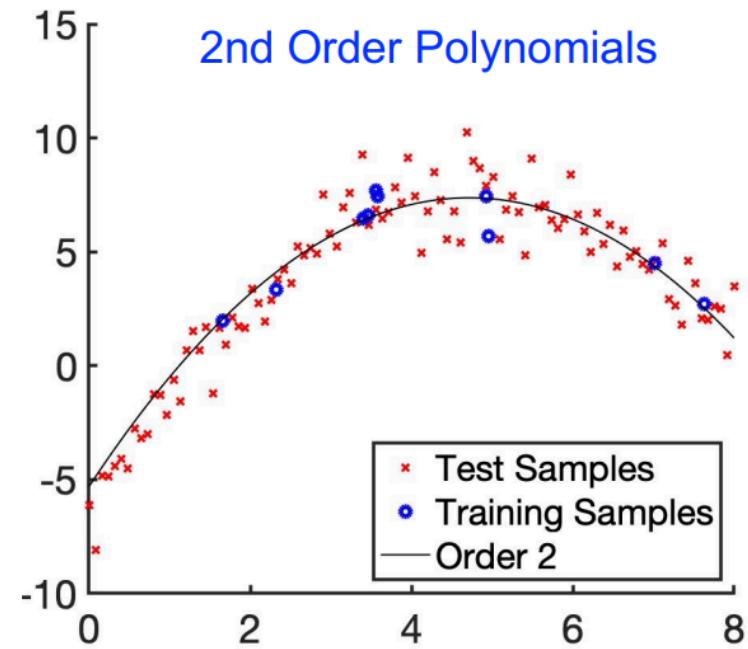
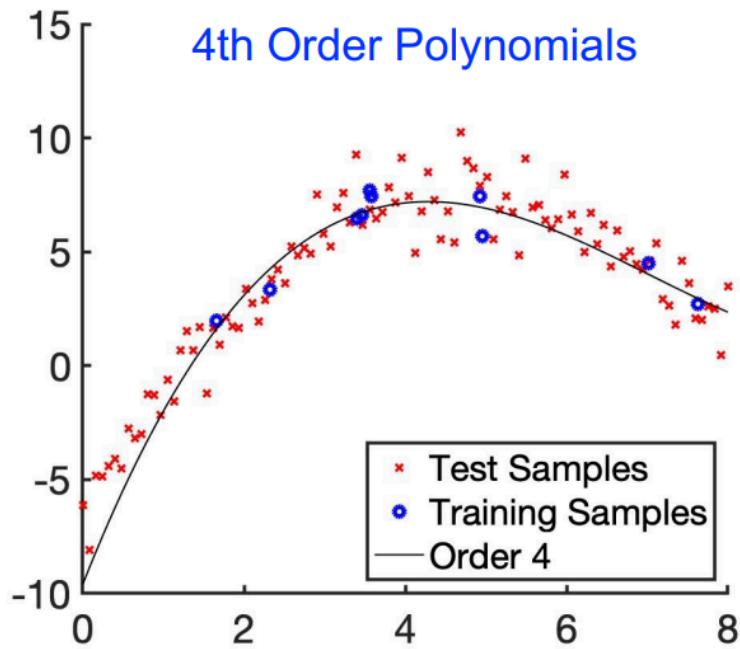
- 4th order has big error on RHS
- 2nd order fit training and test well
- Both fit training quite well

Test 3



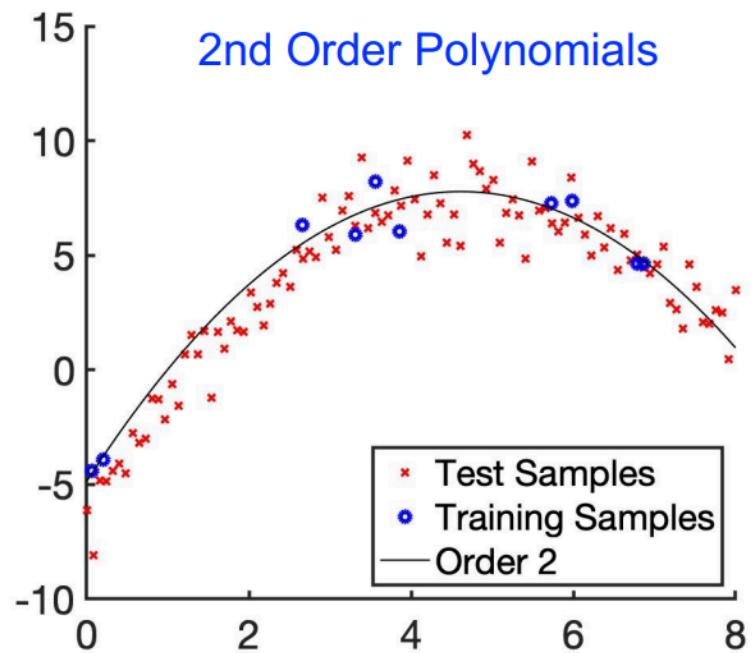
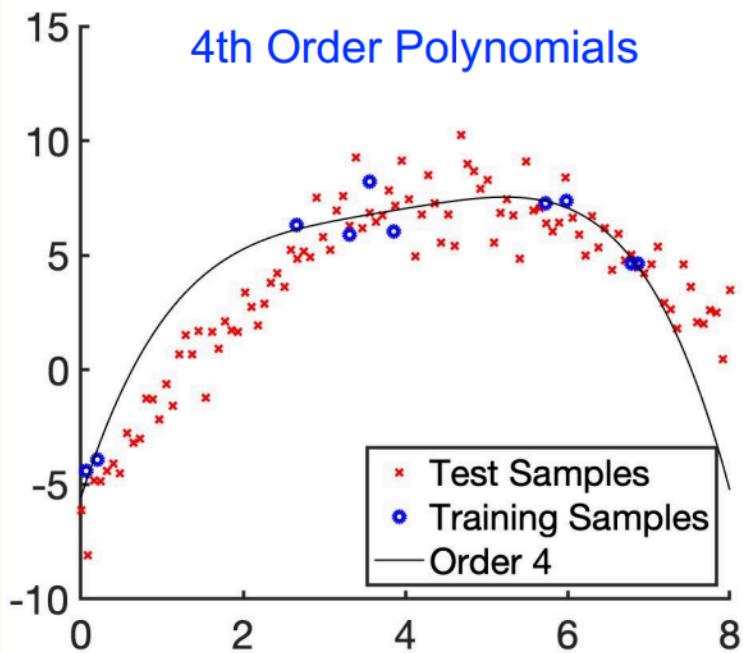
- 4th order has big error on RHS
- 2nd order fit training and test well
- Both fit training quite well

Test 4



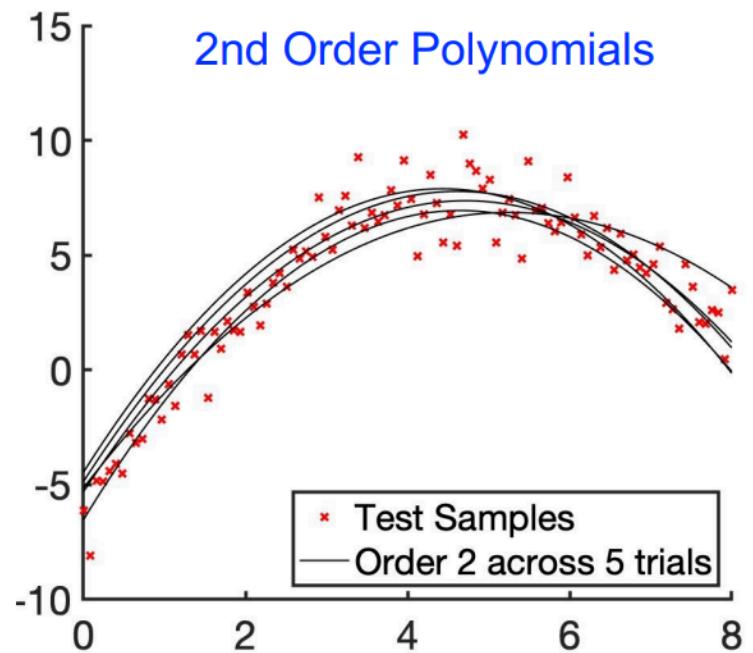
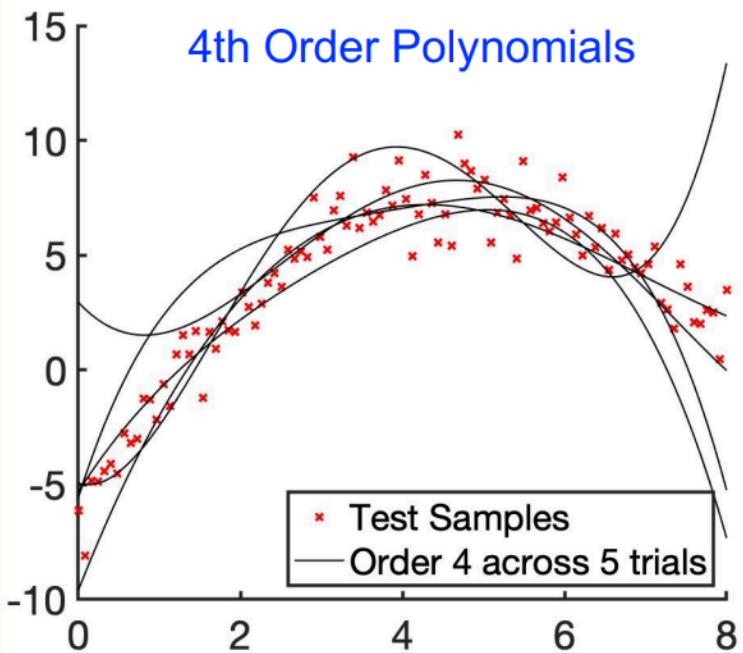
- Both fit training and test quite well

Test 5



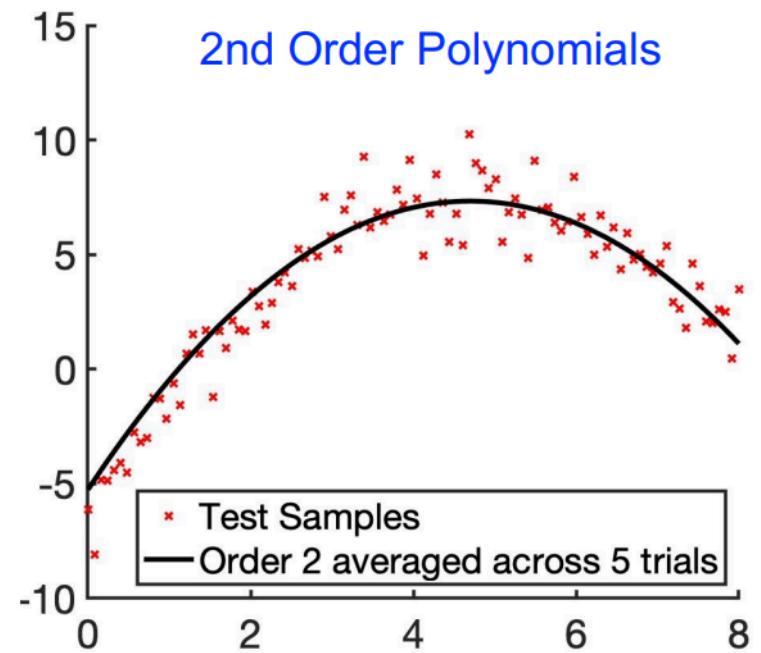
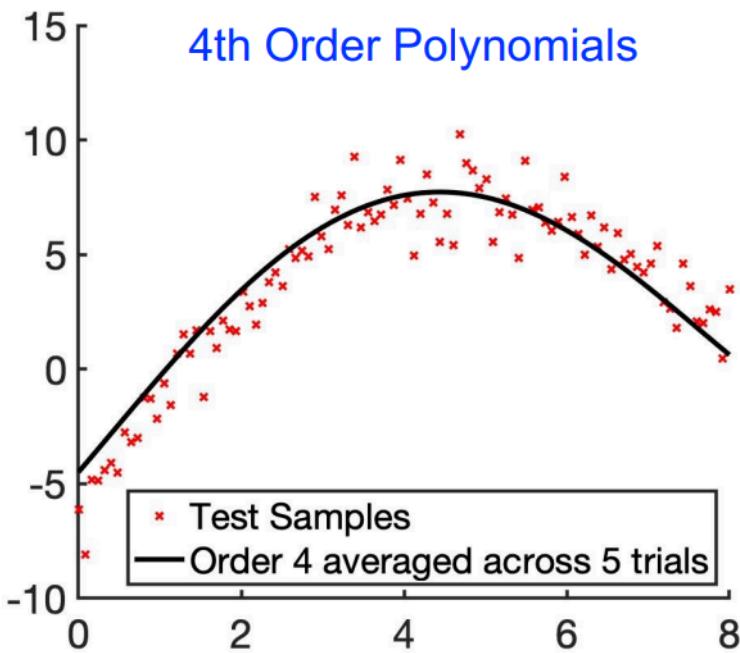
- 4th order has big error on LHS
- 2nd order fit training and test well
- Both fit training quite well

All together



- 4th varies a lot, high variance
- 2nd all quite similar, low variance

Average



- Both averages follows the test data well
- Both has low bias

Conclusion

- Fit with order 2 polynomial: low variance, low bias
- Fit with order 4 polynomial: high variance, low bias
- 2nd order performs better and has lower test error

Order 2 Achieves Lower Test Error

In Class Quiz

An algorithm achieves 10% error in the training set and 50% error in the test set. This situation is described as

26

Overfitting

81%

Underfitting

19%

Python and Conda Stuff

need to run the conda commands in the conda prompt

2.1 Installation

After installing Anaconda^[3], you will probably find `python=3.8.8` in your base(root) environment as shown in the Anaconda Navigator under the environment tab. In this module we will use

```
python=3.8 numpy=1.24
```

So let's create a new environment for EE2211 and we will call it `ee2211`. To do this we enter the following command in the terminal (for Windows user search Anaconda Powershell Prompt, for MacOS user just search for terminal)

```
conda create --name ee2211 python=3.8 numpy=1.24 matplotlib scikit-learn pandas
```

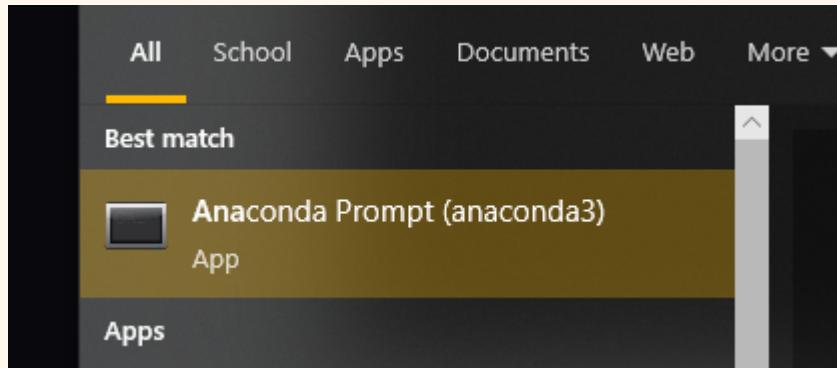
Now in the conda navigator you should find a new environment we just created `ee2211`, with the version we want :-)

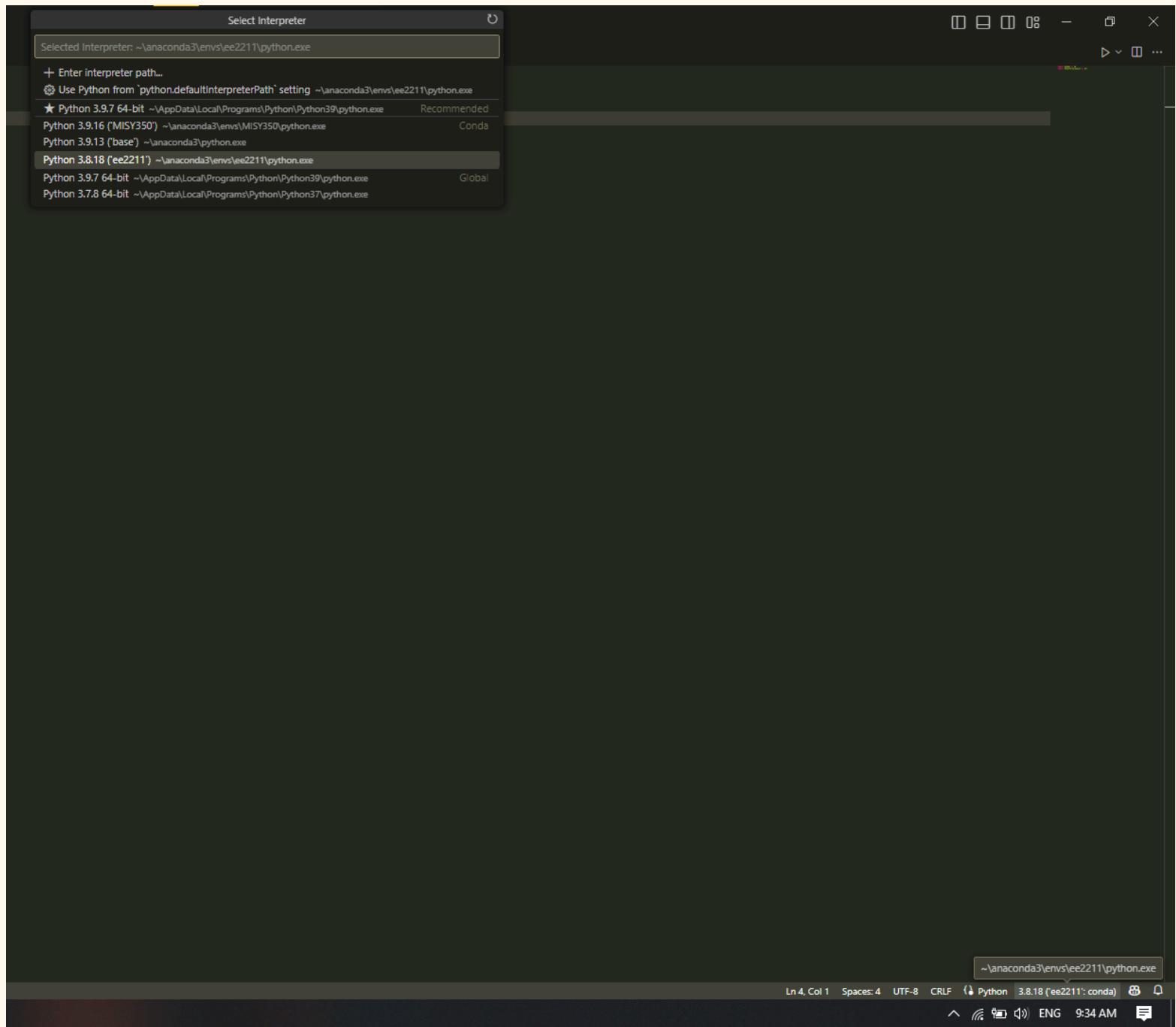
Go ahead and activate the `ee2211` environment, by the following command

```
conda activate ee2211
```

In order to view the Jupyter notebook, let's install it in `ee2211`

```
conda install -c conda-forge jupyterlab
```





Tutorial 1

Tutorial 2

pandas mean python data analysis

matplotlib to plot

sklearn for machine learning

csv easy to read and create, used for data processing and manipulation

pandas reading csv will add index to the data

```
csv_file = pd.read_csv('GovernmentExpenditureonEducation.csv')
print(csv_file)
```

```

● year total_expenditure_on_education
0 1981 942517
1 1982 1358430
2 1983 1611647
3 1984 1769728
4 1985 1812376
5 1986 1641893
6 1987 1654115
7 1988 1604473
8 1989 1765250
9 1990 2056374
10 1991 2816371
11 1992 2597894
12 1993 2902886
13 1994 3318956
14 1995 3443857
15 1996 3771955
16 1997 4440254

```

https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf

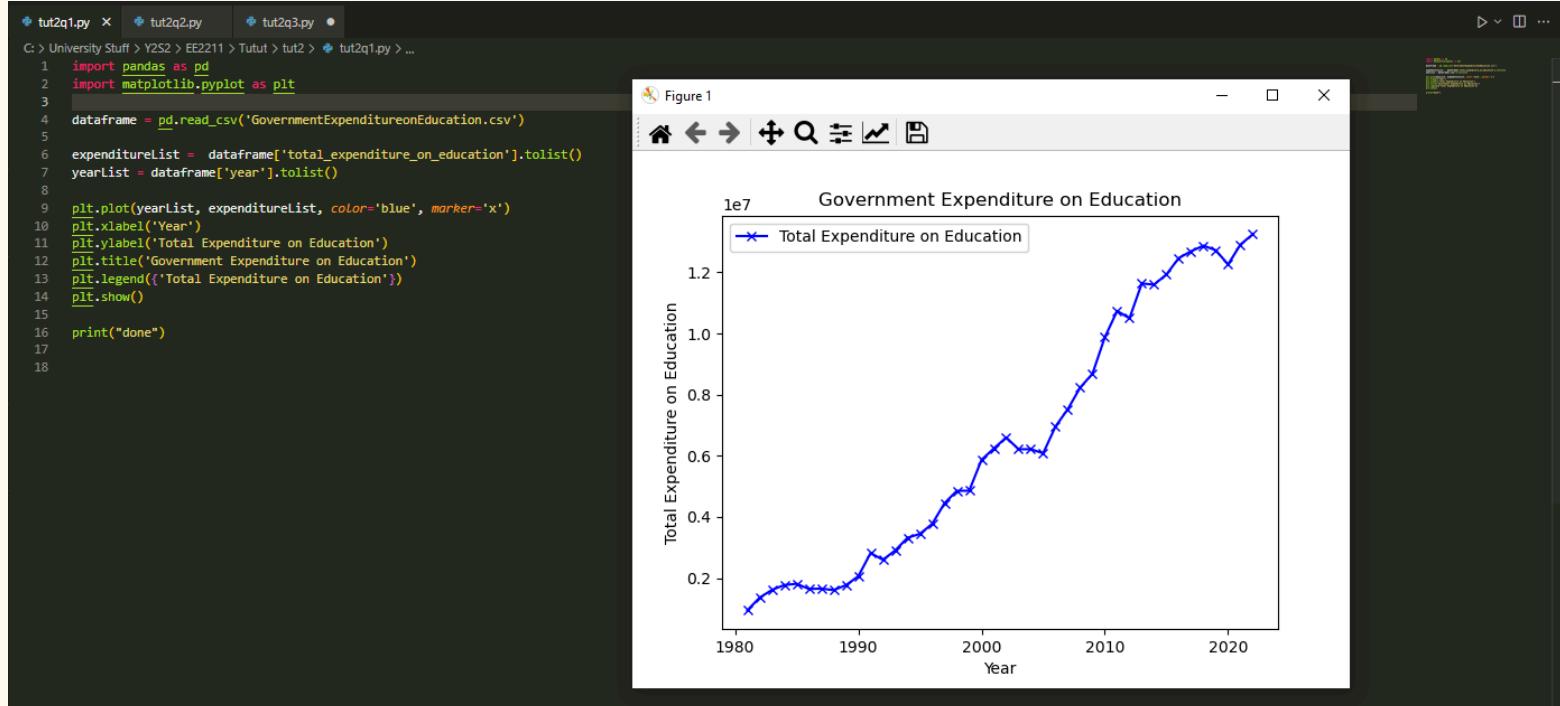
https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.plot.html

Q1

(Data Reading and Visualization, simple data structure)

Question 1:

A Comma Separated Values (CSV) file is a plain text file that contains a list of data. These files are often used for exchanging data between different applications. Download the file “[government-expenditure-on-education.csv](#)” from <https://data.gov.sg/dataset/government-expenditure-on-education>. Plot the educational expenditure over the years. (Hint: you might need “import pandas as pd” and “import matplotlib.pyplot as plt”).



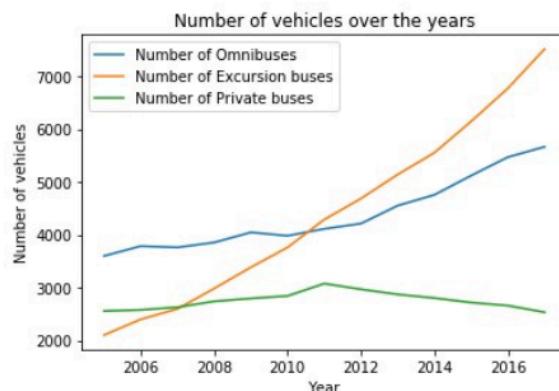
- if absolute file path not used, need to run in the same directory
- use `.tolist()` to make column into list in python

Q2

(Data Reading and Visualization, slightly more complicated data structure)

Question 2:

Download the CSV file from <https://data.gov.sg/dataset/annual-motor-vehicle-population-by-vehicle-type>. Extract and plot the number of Omnibuses, Excursion buses and Private buses over the years as shown below. (Hint: you might need “import pandas as pd” and “import matplotlib.pyplot as plt”.)



The screenshot shows a Jupyter Notebook interface with two code cells and their corresponding plots.

Code Cell 1 (Left):

```
# Display the plot
plt.show()

Replace 'column1' and 'column2' with the names of the columns you want to plot. If you're not sure what the column names are, you can print them out with print(csv_file.columns).

Here's how your code would look:

import pandas as pd
import matplotlib.pyplot as plt

csv_file = pd.read_csv('GovernmentExpenditureonEducation.csv')

# Replace 'column1' and 'column2' with your column names
plt.plot(csv_file['column1'], csv_file['column2'])

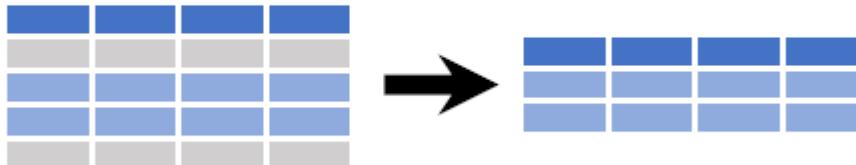
plt.show()

print("done")
```

Code Cell 2 (Right):

```
C:\>University Stuff\Y2S2\EE2211>Tutut>Nt2> tut2q1.py ...  
C:\>University Stuff\Y2S2\EE2211>Tutut>Nt2> tut2q2.py ...  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
17
```

Subset Observations - rows



`df[df.Length > 7]`

Extract rows that meet logical criteria.

- extract the index out instead of getting the whole list

```
import pandas as pd
import matplotlib.pyplot as plt

dataframe = pd.read_csv('AnnualMotorVehiclePopulationbyVehicleType.csv')

privateBusIndexList = dataframe[dataframe['type'] == "Private buses"].index.tolist()

print(privateBusIndexList)

privateBusIndexList = dataframe.loc[dataframe['type'] == "Private buses"].index

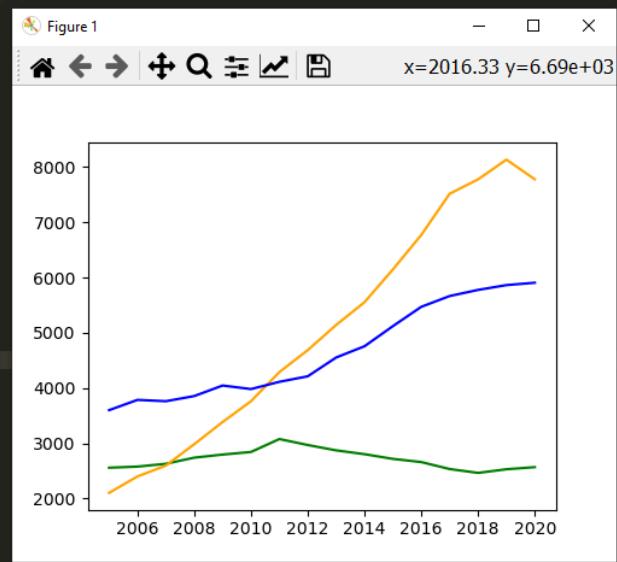
print(privateBusIndexList)
```

- PS C:\University Stuff\Y2S2\EE2211\Tutut\tut2> [132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 216, 237, 258, 279, 300, 321]
Index([132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 216, 237, 258, 279,
 300, 321],
 dtype='int64')
- PS C:\University Stuff\Y2S2\EE2211\Tutut\tut2> []
- with .index, the data type is in a diff form

```

tut2q1.py tut2q2.py * tut2q3.py
C:\> University Stuff > Y2S2 > EE2211 > Tutut > tut2 > tut2q2.py > ...
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4
5 dataframe = pd.read_csv('AnnualMotorVehiclePopulationbyVehicleType.csv')
6
7 # privateBusIndexList = dataframe[dataframe['type'] == "Private buses"].index.tolist()
8 # print(privateBusIndexList)
9
10 # privateBusIndexList = dataframe.loc[dataframe['type'] == "Private buses"].index
11 # excursionBusIndexList = dataframe.loc[dataframe['type'] == "Excursion buses"].index
12 # omniBusIndexList = dataframe.loc[dataframe['type'] == "Omnibuses"].index
13
14 # privateBusDataframe = dataframe.loc[privateBusIndexList]
15 # excursionBusDataframe = dataframe.loc[excursionBusIndexList]
16 # omniBusDataframe = dataframe.loc[omniBusIndexList]
17
18 # or without extracting the index first
19 privateBusDataframe = dataframe.loc[dataframe['type'] == "Private buses"]
20 excursionBusDataframe = dataframe.loc[dataframe['type'] == "Excursion buses"]
21 omniBusDataframe = dataframe.loc[dataframe['type'] == "Omnibuses"]
22
23 plt.plot(privateBusDataframe['year'], privateBusDataframe['number'], color='green')
24 plt.plot(excursionBusDataframe['year'], excursionBusDataframe['number'], color='orange')
25 plt.plot(omniBusDataframe['year'], omniBusDataframe['number'], color='blue')
26
27 plt.show()
28
29
30
31 # privateBusNumber = []
32 # privateBusYear = []
33
34 # excursionBusNumber = []
35 # excursionBusYear = []
36
37 # omniBusNumber = []

```

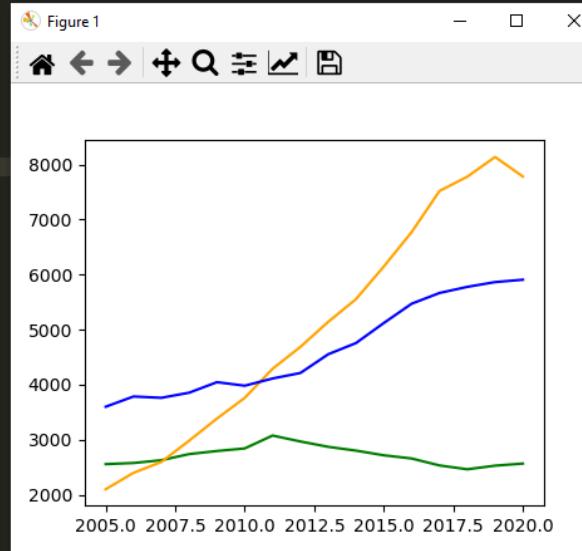


- without extracting index first

```

tut2q1.py tut2q2.py * tut2q3.py
C:\> University Stuff > Y2S2 > EE2211 > Tutut > tut2 > tut2q2.py > ...
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4
5 dataframe = pd.read_csv('AnnualMotorVehiclePopulationbyVehicleType.csv')
6
7 # privateBusIndexList = dataframe[dataframe['type'] == "Private buses"].index.tolist()
8 # print(privateBusIndexList)
9
10 privateBusIndexList = dataframe.loc[dataframe['type'] == "Private buses"].index
11 excursionBusIndexList = dataframe.loc[dataframe['type'] == "Excursion buses"].index
12 omniBusIndexList = dataframe.loc[dataframe['type'] == "Omnibuses"].index
13
14 privateBusDataframe = dataframe.loc[privateBusIndexList]
15 excursionBusDataframe = dataframe.loc[excursionBusIndexList]
16 omniBusDataframe = dataframe.loc[omniBusIndexList]
17
18 # or without extracting the index first
19 # privateBusDataframe = dataframe.loc[dataframe['type'] == "Private buses"]
20 # excursionBusDataframe = dataframe.loc[dataframe['type'] == "Excursion buses"]
21 # omniBusDataframe = dataframe.loc[dataframe['type'] == "Omnibuses"]
22
23 plt.plot(privateBusDataframe['year'], privateBusDataframe['number'], color='green')
24 plt.plot(excursionBusDataframe['year'], excursionBusDataframe['number'], color='orange')
25 plt.plot(omniBusDataframe['year'], omniBusDataframe['number'], color='blue')
26
27 plt.show()
28
29
30
31 # privateBusNumber = []
32 # privateBusYear = []
33
34 # excursionBusNumber = []
35 # excursionBusYear = []
36
37 # omniBusNumber = []

```

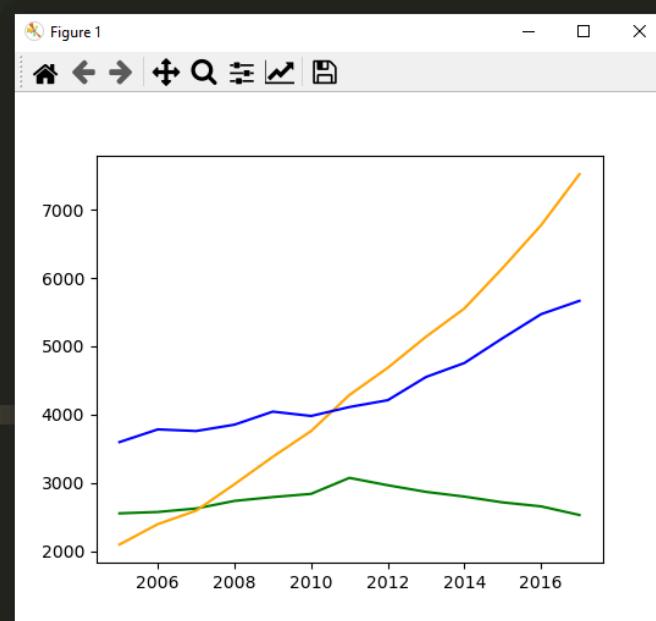


- extracting index first

Cut away the year after 2017

tut2q1.py tut2q2.py tut2q3.py

```
C: > University Stuff > Y2S2 > EE2211 > Tutut > tut2 > tut2q2.py > ...
10 privateBusIndexList = dataframe.loc[dataframe['type'] == "Private buses"].index
11 excursionBusIndexList = dataframe.loc[dataframe['type'] == "Excursion buses"].index
12 omniBusIndexList = dataframe.loc[dataframe['type'] == "Omnibuses"].index
13
14 privateBusDataframe = dataframe.loc[privateBusIndexList]
15 excursionBusDataframe = dataframe.loc[excursionBusIndexList]
16 omniBusDataframe = dataframe.loc[omniBusIndexList]
17
18 # or without extracting the index first
19 # privateBusDataframe = dataframe.loc[dataframe['type'] == "Private buses"]
20 # excursionBusDataframe = dataframe.loc[dataframe['type'] == "Excursion buses"]
21 # omniBusDataframe = dataframe.loc[dataframe['type'] == "Omnibuses"]
22
23 plt.plot(privateBusDataframe['year'], privateBusDataframe['number'], color='green')
24 plt.plot(excursionBusDataframe['year'], excursionBusDataframe['number'], color='orange')
25 plt.plot(omniBusDataframe['year'], omniBusDataframe['number'], color='blue')
26
27 plt.show()
28
29 yearCut = 2017
30
31 print("cutting away years after {}".format(yearCut))
32 privateBusDataframe = privateBusDataframe.loc[privateBusDataframe['year'] <= yearCut]
33 excursionBusDataframe = excursionBusDataframe.loc[excursionBusDataframe['year'] <= yearCut]
34 omniBusDataframe = omniBusDataframe.loc[omniBusDataframe['year'] <= yearCut]
35
36 plt.plot(privateBusDataframe['year'], privateBusDataframe['number'], color='green')
37 plt.plot(excursionBusDataframe['year'], excursionBusDataframe['number'], color='orange')
38 plt.plot(omniBusDataframe['year'], omniBusDataframe['number'], color='blue')
39
40 plt.show()
41 # privateBusNumber = []
42 # privateBusYear = []
43
44 # excursionBusNumber = []
45 # excursionBusYear = []
46
```

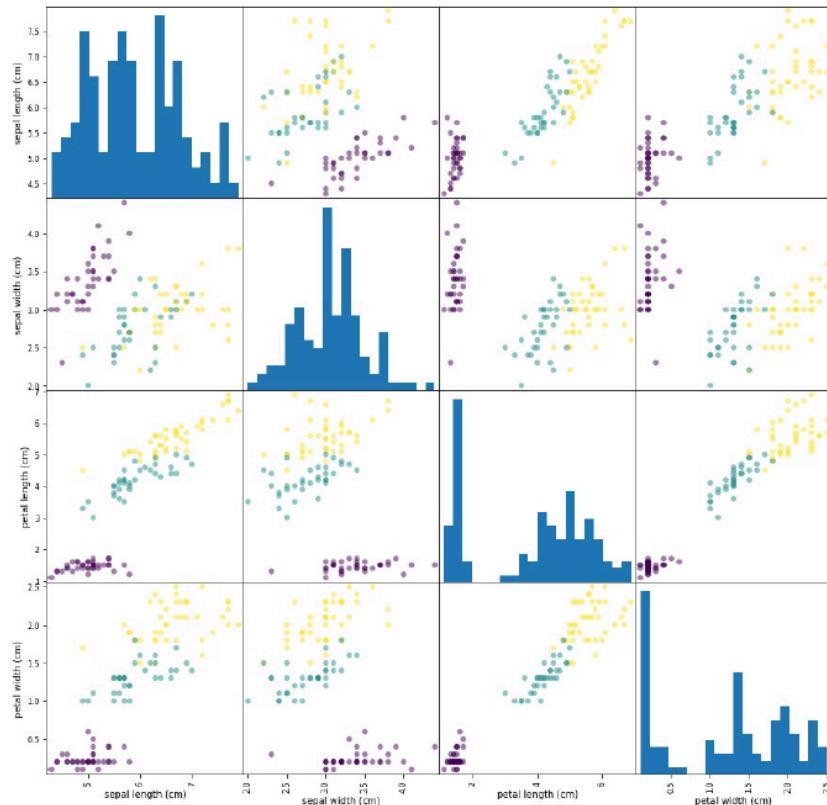


Q3

(Data Reading and Visualization, distribution)

Question 3:

The “iris” flower data set consists of measurements such as the length, width of the petals, and the length, width of the sepals, all measured in centimeters, associated with each iris flower. Get the data set “from sklearn.datasets import load_iris” and do a scatter plot as shown below. (Hint: you might need “from pandas.plotting import scatter_matrix”)



C: > University Stuff > Y2S2 > EE2211 > Tutut > tut2 > tut2q3.py > ...

```
1  from sklearn.datasets import load_iris
2
3  import pandas as pd
4
5  from pandas.plotting import scatter_matrix
6
7  import matplotlib.pyplot as plt
8
9  # Load the iris dataset
10 iris_dataset = load_iris()
11
12 print(iris_dataset)
13
```

sklearn build in dataset for trying shit out

```
2, 2, 2, 2,
), 'frame': None, 'target_names': array(['setosa', 'versicolor', 'virginica'], dtype='|<U10'), 'DESCR': '...
es: 150 (50 in each of three classes)\n      :Number of Attributes: 4 numeric, predictive attributes and the
```

- target name is specified here,
 - target name/species name should be the output of learning algo

```
• {'data': array([[5.1, 3.5, 1.4, 0.2],  
•                 [4.9, 3. , 1.4, 0.2],  
•                 [4.7, 3.2, 1.3, 0.2],  
•                 [4.6, 3.1, 1.5, 0.2],  
•                 [5. , 3.6, 1.4, 0.2],  
•                 [5.4, 3.9, 1.7, 0.4],
```

```
ansactions\n      on Information Theory, May 1972, 431-433.\n      - See also: 1988 MLC Proceedings, 54-64. Cheeseman et al's AUTOCLASS II\n', 'feature_names': ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)'], 'filename': 'iris.csv', 'data_modul
```

- feature is here
 - feature should be the input to learning algo

sklearn.model_selection.train_test_split

```
sklearn.model_selection.train_test_split(*arrays, test_size=None, train_size=None, random_state=None, shuffle=True, stratify=None)
```

[source]

Split arrays or matrices into random train and test subsets.

Quick utility that wraps input validation, `next(ShuffleSplit().split(X, y))`, and application to input data into a single call for splitting (and optionally subsampling) data into a one-liner.

Read more in the [User Guide](#).

Parameters:	*arrays : sequence of indexables with same length / shape[0] Allowed inputs are lists, numpy arrays, scipy-sparse matrices or pandas dataframes.
	test_size : float or int, default=None If float, should be between 0.0 and 1.0 and represent the proportion of the dataset to include in the test split. If int, represents the absolute number of test samples. If None, the value is set to the complement of the train size. If <code>train_size</code> is also None, it will be set to 0.25.
	train_size : float or int, default=None If float, should be between 0.0 and 1.0 and represent the proportion of the dataset to include in the train split. If int, represents the absolute number of train samples. If None, the value is automatically set to the complement of the test size.
	random_state : int, RandomState instance or None, default=None Controls the shuffling applied to the data before applying the split. Pass an int for reproducible output across multiple function calls. See Glossary .
	shuffle : bool, default=True Whether or not to shuffle the data before splitting. If <code>shuffle=False</code> then <code>stratify</code> must be <code>None</code> .
	stratify : array-like, default=None If not <code>None</code> , data is split in a stratified fashion, using this as the class labels. Read more in the User Guide .
Returns:	splitting : list, length=2 * len(arrays) List containing train-test split of inputs. <i>New in version 0.16:</i> If the input is sparse, the output will be a <code>scipy.sparse.csr_matrix</code> . Else, output type is the same as the input type.

- random state is the seed
- https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html

https://pandas.pydata.org/docs/reference/api/pandas.plotting.scatter_matrix.html

Q5

(Missing Data)

Question 5:

The Pima Indians Diabetes Dataset involves predicting the onset of diabetes within 5 years in Pima Indians given medical details. Download the Pima-Indians-Diabetes data from

<https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv>.

It is a binary (2-class) classification problem. The number of observations for each class is not balanced. There are 768 observations with 8 input variables and 1 output variable. The variable names are as follows:

0. Number of times pregnant.
1. Plasma glucose concentration a 2 hours in an oral glucose tolerance test.
2. Diastolic blood pressure (mm Hg).
3. Triceps skinfold thickness (mm).
4. 2-Hour serum insulin (mu U/ml).
5. Body mass index (weight in kg/(height in m)^2).
6. Diabetes pedigree function.
7. Age (years).
8. Class variable (0 or 1).

- (a) Print the summary statistics of this data set.
- (b) Count the number of “0” entries in columns [1,2,3,4,5].
- (c) Replace these “0” values by “NaN”.

(Hint: you might need the “.describe ()” and “.replace (0, numpy.NaN)” functions “from pandas import read_csv”).

- {dataset}.describe() gives summary states, print out to read before processing data

Tutorial 3

Q5

`pdf(x, loc=0, scale=1)`

Probability density function.

`logpdf(x, loc=0, scale=1)`

Log of the probability density function.

`cdf(x, loc=0, scale=1)`

Cumulative distribution function.

The probability density function for `norm` is:

$$f(x) = \frac{\exp(-x^2/2)}{\sqrt{2\pi}}$$

for a real number x .

The probability density above is defined in the “standardized” form. To shift and/or scale the distribution use the `loc` and `scale` parameters. Specifically, `norm.pdf(x, loc, scale)` is identically equivalent to `norm.pdf(y) / scale` with $y = (x - \text{loc}) / \text{scale}$. Note that shifting the location of a distribution does not make it a “noncentral” distribution; noncentral generalizations of some distributions are available in separate classes.

- loc and scale is mean and standard deviation

Q10-Q11

Question 10: (Fill-in-blank)

The rank of the matrix $\begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$ is _BLANK_.

2

Question 11: (Fill-in-blank)

The rank of the matrix $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ is _BLANK_.

3

2, must do ero to ref

- find rank must do ref, for $> 2 \times 2$, cannot eye ball

Tutorial 4

- $X^T X$ is always symmetrical if $\det(X^T X) \neq 0$
- XX^T is always symmetrical if $\det(XX^T) \neq 0$

Q1

(Systems of Linear Equations)

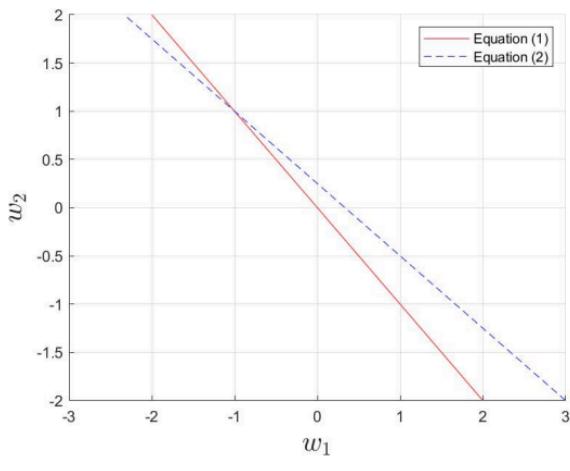
Question 1:

Given $\mathbf{X}\mathbf{w} = \mathbf{y}$ where $\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 3 & 4 \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$.

even
 $\det \mathbf{X} = 1 \neq 0$, invertible

- What kind of system is this? (even-, over- or under-determined?)
- Is \mathbf{X} invertible? Why?
- Solve for \mathbf{w} if it is solvable.

$\mathbf{w} = (-1, 1)$



Q2

(Systems of Linear Equations)

Question 2:

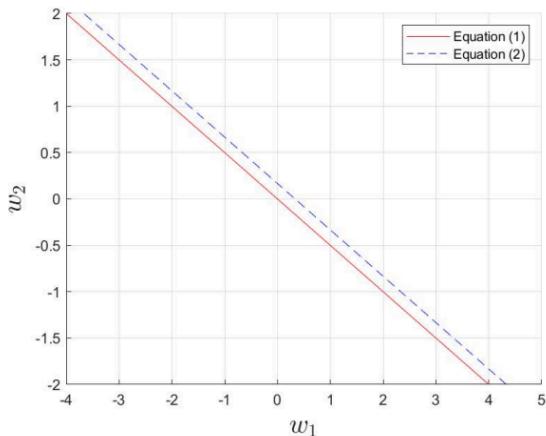
Given $\mathbf{X}\mathbf{w} = \mathbf{y}$ where $\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$.

even
 $\det \mathbf{X} = 0$, non invertible

- (a) What kind of system is this? (even-, over- or under-determined?)
- (b) Is \mathbf{X} invertible? Why?
- (c) Solve for \mathbf{w} if it is solvable.

not solvable? can use least square?

- Two parallel lines confirm no solution, there is no least square solution also



• Let $V \subseteq \mathbb{R}^n$ and $S = \{u_1, u_2, \dots, u_k\}$ be a basis for V . Then the orthogonal projection of any vector $w \in \mathbb{R}^n$ onto $V = A(A^T A)^{-1} A^T w$, where $A = (u_1, u_2, \dots, u_k) \in \mathbb{R}^{n \times k}$ $n \geq k$

↳ since S is basis, columns of A are linearly independent ($\text{Rank}(A) = k$, A is full ranked) so $A^T A$ is invertible

$\therefore u = (A^T A)^{-1} A^T w$ is the unique solution to $A^T A u = A^T w$

\therefore project of w onto $\text{Col}(A) = V = Au = A(A^T A)^{-1} A^T w \#$

- column must be independent so cannot use this
- if take out one column then can apply

Question 3:

Given $\mathbf{X}\mathbf{w} = \mathbf{y}$ where $\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 1 & -1 \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} 0 \\ 0.1 \\ 1 \end{bmatrix}$.

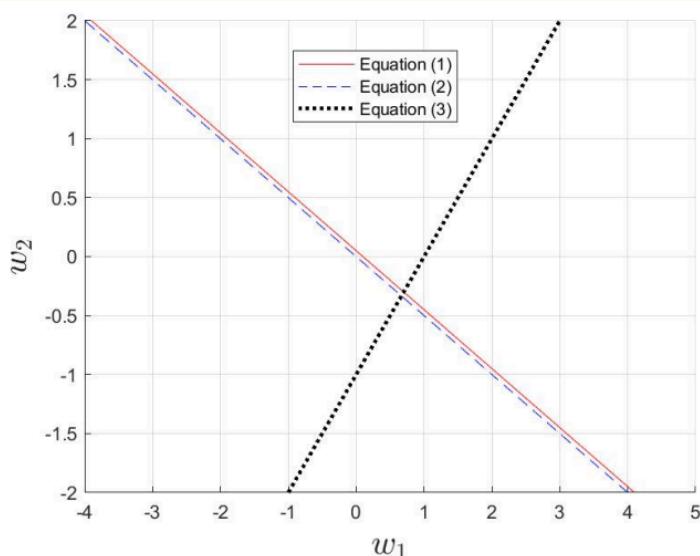
over

X is full rank, rank = 2, left inverse exist

$$\mathbf{w} = [0.68 \ 0.32]$$

- (a) What kind of system is this? (even-, over- or under-determined?)
- (b) Is \mathbf{X} invertible? Why?
- (c) Solve for \mathbf{w} if it is solvable.

- Two parallel line and one non parallel line intersecting at two points, the \mathbf{w} which is least square, gives the \mathbf{w} between ish between the two intersection points

**Q6****Question 6:**

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

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 1 & -1 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

- (a) What kind of system is this? (even-, over- or under-determined?)
- (b) Is \mathbf{X} invertible? Why?
- (c) Solve for \mathbf{w} if it is solvable.

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

over

after transpose its under determined

rank 2 = full rank, left inverse exist

$$\mathbf{w} = [0.06666667 \ 0.13333333 \ -0.33333333]$$

$$= [1/15, 2/15, -1/3]$$

- (a) This is an under-determined system (there are 3 unknowns with 2 equations).
- (b) \mathbf{X} is NOT invertible but $\mathbf{X}^T \mathbf{X}$ is. The determinant of $\mathbf{X}^T \mathbf{X} = 6 \times 21 - 9 \times 9 = 45$.
- (c) A constrained solution (exact) is given by

$\hat{\mathbf{w}}^T = (\mathbf{X}\mathbf{a})^T$ (The 3-dimensional vector \mathbf{w} can be constrained by projecting \mathbf{X} onto a 2-dimensional vector \mathbf{a})

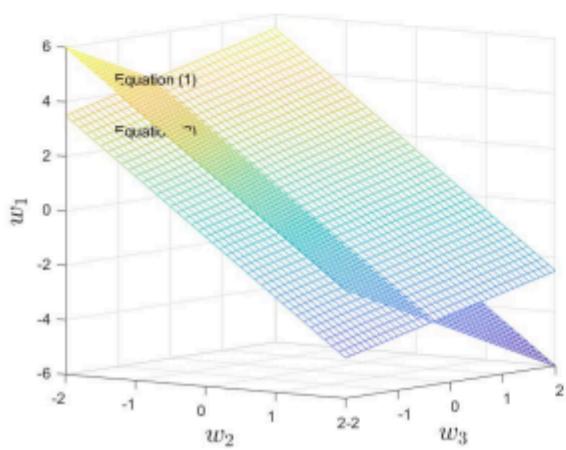
$$= \mathbf{a}^T \mathbf{X}^T$$

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

$$= [0 \quad 1] \begin{bmatrix} 0.4667 & -0.2 \\ -0.2 & 0.1333 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & -1 \end{bmatrix}$$

$$= [0.0667 \quad 0.1333 \quad -0.3333].$$

Note: $\dim(\mathbf{X})$ is 3×2 , $\dim(\mathbf{a})$ is 2×1 , estimation is done/constrained on/to the lower dimension of (3×2) and then projected back to the higher dimension 3.



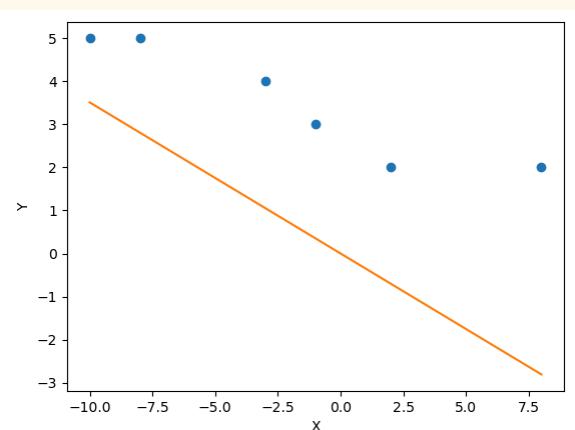
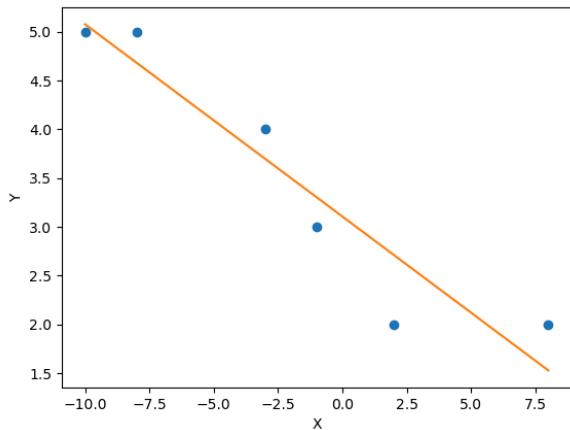
Tutorial 5

Q1

```
# Question 1
print("Question 1")

X = np.array([[-10], [-8], [-3], [-1], [2], [8]])
Y = np.array([[5], [5], [4], [3], [2], [2]])
w = helper.linearRegressionWithBias(X, Y, printResult=False, printFeature=1)
print(w)

X = np.array([[-10], [-8], [-3], [-1], [2], [8]])
Y = np.array([[5], [5], [4], [3], [2], [2]])
w = helper.linearRegressionWithoutBias(X, Y, printResult=False, printFeature=0)
print(w)
```



```
W =
[[ 3.10550459]
[-0.19724771]]
```

```
W =
[[ -0.35123967]]
```

bias term allows the regression to move away from the origin and reduce error

Q2

even-determined

$$X = \begin{bmatrix} 1 & 0 & 1 \\ 2 & -1 & 1 \\ 1 & 1 & 5 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

- even determined with full rank can get answer by finding the inverse

```
Question 2
rank(X) = 3
W =
[[ 0.33333333]
[-0.66666667]
[ 0.66666667]]
MSE = 1.472540356412555e-29
Ytest1_nobias =
[[3.66666667]]
Ytest2_nobias =
[[-3.66666667]]
rank(X) = 3
W =
[[-1.5]
[-1. ]
[ 3. ]
[-0.25]]
MSE = 34.895833333333336
Ytest1_bias =
[[3.5]]
Ytest2_bias =
[[12.75]]
```

$$\hat{\omega} = X^\top (X X^\top)^{-1} y$$

$$= \begin{bmatrix} -0.1429 \\ 0.5238 \\ -0.4762 \\ 0.6190 \end{bmatrix}$$

- should use right inverse

$$y = X \hat{\omega} = \begin{bmatrix} 1 & 0 & 1 \\ 2 & -1 & 1 \\ 1 & 1 & 5 \end{bmatrix} \begin{bmatrix} -0.1429 \\ 0.5238 \\ -0.4762 \\ 0.6190 \end{bmatrix}$$

$$= \begin{bmatrix} 3.3333 \\ -2.6190 \end{bmatrix}$$

```

# this is wrong, should use right inverse
W_bias = helper.linearRegressionWithBias(
    X, Y, printResult=True, printFeature=None)

# add in bias
Xtest1 = np.array([[1, -1, 2, 8]])
Xtest2 = np.array([[1, 1, 5, -1]])

Ytest1_bias = helper.testData(Xtest1, W_bias, printResult=False)
print("Ytest1_bias =\n", Ytest1_bias)
Ytest2_bias = helper.testData(Xtest2, W_bias, printResult=False)
print("Ytest2_bias =\n", Ytest2_bias)

# after adding bias, its a wider matrix so should solve with right inverse
W_bias_right = helper.rightInverse(helper.paddingOfOnes(X)) @ Y
print("W_bias_right =\n", W_bias_right)

```

You, 1 second ago • Uncommitted changes

```

Ytest1_bias = helper.testData(Xtest1, W_bias_right, printResult=False)
print("Ytest1_bias =\n", Ytest1_bias)
Ytest2_bias = helper.testData(Xtest2, W_bias_right, printResult=False)
print("Ytest2_bias =\n", Ytest2_bias)

```

```

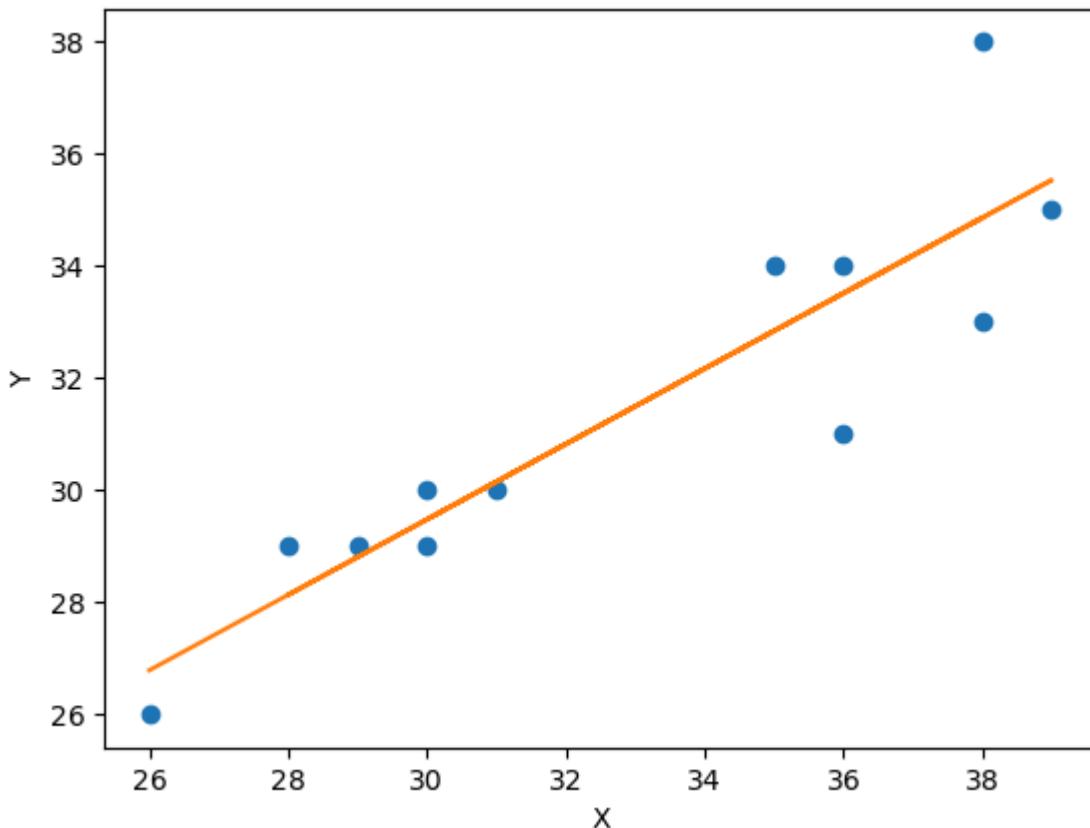
W_bias_right =
[[[-0.14285714]
 [ 0.52380952]
 [-0.47619048]
 [ 0.61904762]]]
Ytest1_bias =
[[3.33333333333]]
Ytest2_bias =
[[-2.61904762]]

```

The wide matrix has a left inverse as the 3×3 is already full rank, so the dim of row is 3, after padding ones, dimension remains

- hence can find left inverse, but it's not the most accurate regression as there is no interpretation of least squares error in this case

Q3



```

rank(X) = 2
W =
[[9.3
[0.67272727]]
MSE = 1.9530303030303038
Ytest1 =
[[29.48181818]]
Ytest2 =
[[12.66363636]]

```

always include bias if not specified

- predicted that 30 student buy around 30 books, one per student
- predicted that 5 student buy around 12 books, 2 per student
- 5 student is out of the range of regression, not appropriate to use to same regression
 - not sure whether the regression still fits, cannot extrapolate
 - like gae1000

Q4

Semester	Students	Books
1	36	31
2	26	20
3	35	34
4	39	35
5	26	20
6	30	30
7	31	30
8	38	38
9	36	34
10	38	33
11	26	20
12	26	20

Without removing duplicates VS removing the duplicates

```

X = np.array([[36], [26], [35], [39], [26], [30],
|           | [31], [38], [36], [38], [26], [26]])
Y = np.array([[31], [20], [34], [35], [20], [30],
|           | [30], [38], [34], [33], [20], [20]])

# xw = Y regression
w = helper.linearRegressionWithBias(X, Y, printResult=True, printFeature=1)

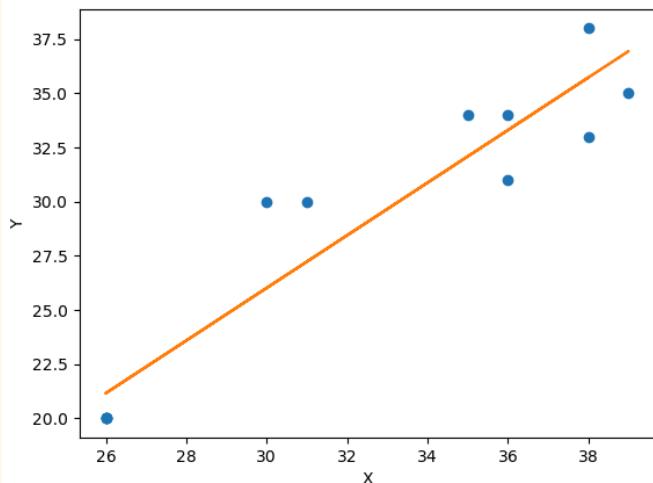
# predict
Xtest1 = np.array([[1, 30]])
Ytest1 = helper.testData(Xtest1, w, printResult=False)
print("Ytest1 =\n", Ytest1)

```

```

rank(X) = 2
W =
[-10.41257051]
[ 1.21434327]
MSE = 4.562181036798282
Ytest1 =
[[26.01772764]]

```



```

# remove duplicates
print("X =\n", X)
print("Y =\n", Y)
combined = np.hstack((X, Y))
print("combined =\n", combined)
new_array = [tuple(row) for row in combined]
print("new_array =\n", new_array)
uniques = np.unique(new_array, axis=0)
print("uniques =\n", uniques)
X_unique = uniques[:, 0].reshape(-1, 1)
Y_unique = uniques[:, 1].reshape(-1, 1)
# In this context, -1 is used as a placeholder for "figure out what the dimension should be". It means that the length in that dimension is inferred from the length of the array and remaining dimensions.
# So, .reshape(-1, 1) means to reshape the array to have one column and as many rows as necessary to accommodate the original data.
print("X_unique =\n", X_unique)
print("Y_unique =\n", Y_unique)

# xw = Y regression
w_unique = helper.linearRegressionWithBias(X_unique, Y_unique, printResult=True, printFeature=1)

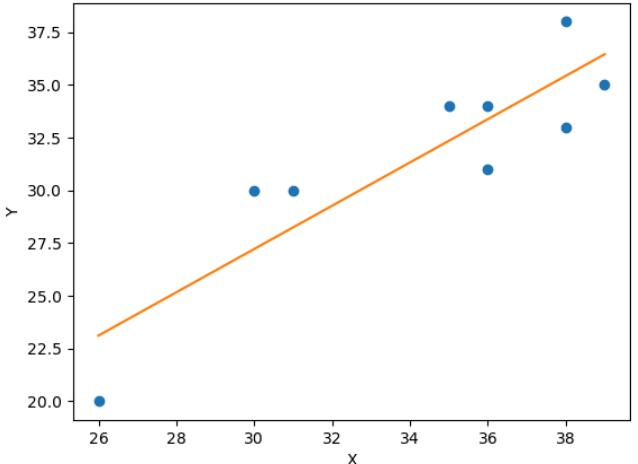
# predict
Ytest1_unique = helper.testData(Xtest1, w_unique, printResult=False)
print("Ytest1_unique =\n", Ytest1_unique)

```

```

rank(X) = 2
W =
[[ -3.55844156]
 [ 1.02597403]]
MSE = 4.8773448773448775
Ytest1_unique =
[[27.22077922]]

```



- can see that the effect (26, 20) has on the regression is less

Combine so it's easier to compare

```

# Combining the plots for easier comparison
plt.plot(X, Y, 'o', color='blue')

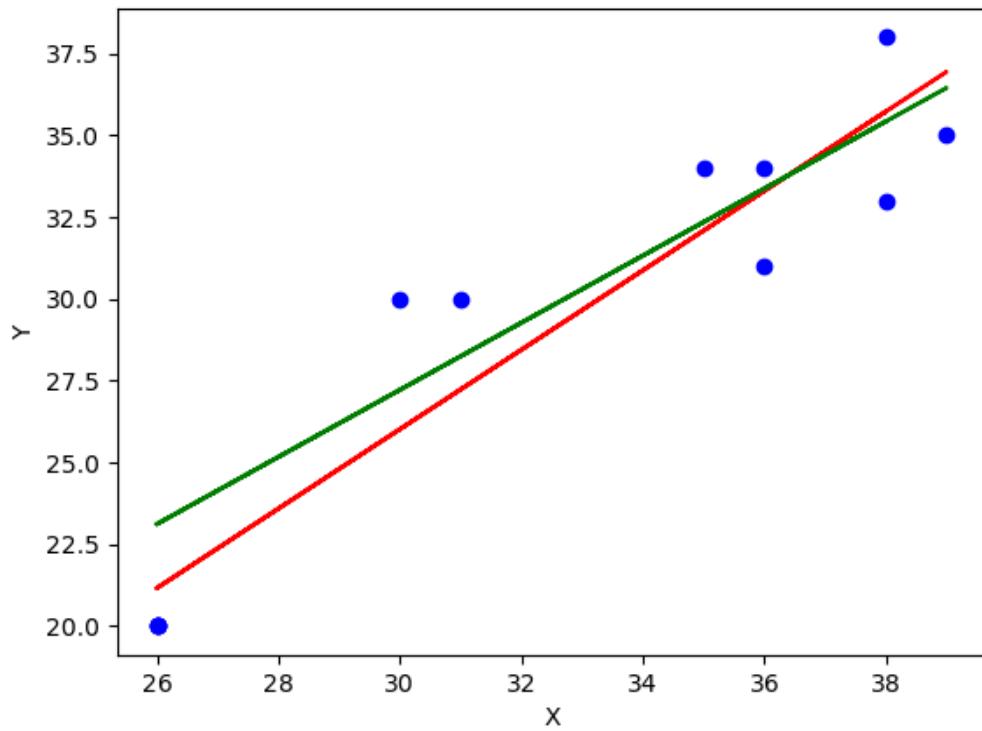
# plot the graph for the unique data, X needs to be padded for the offset
plt.plot(X, np.hstack((np.ones((X.shape[0], 1)), X)) @ w, '-', color='red')

# plot the graph for the unique data
plt.plot(X, np.hstack((np.ones((X.shape[0], 1)), X)) @ w_unique, '--', color='green')

plt.xlabel('X')
plt.ylabel('Y')

plt.show()

```



- more optimistic after purging before 37

Q5

```
# Question 5 #####
print("Question 5 #####")
exp_df = pd.read_csv("GovernmentExpenditureonEducation.csv")
exp_df.info()

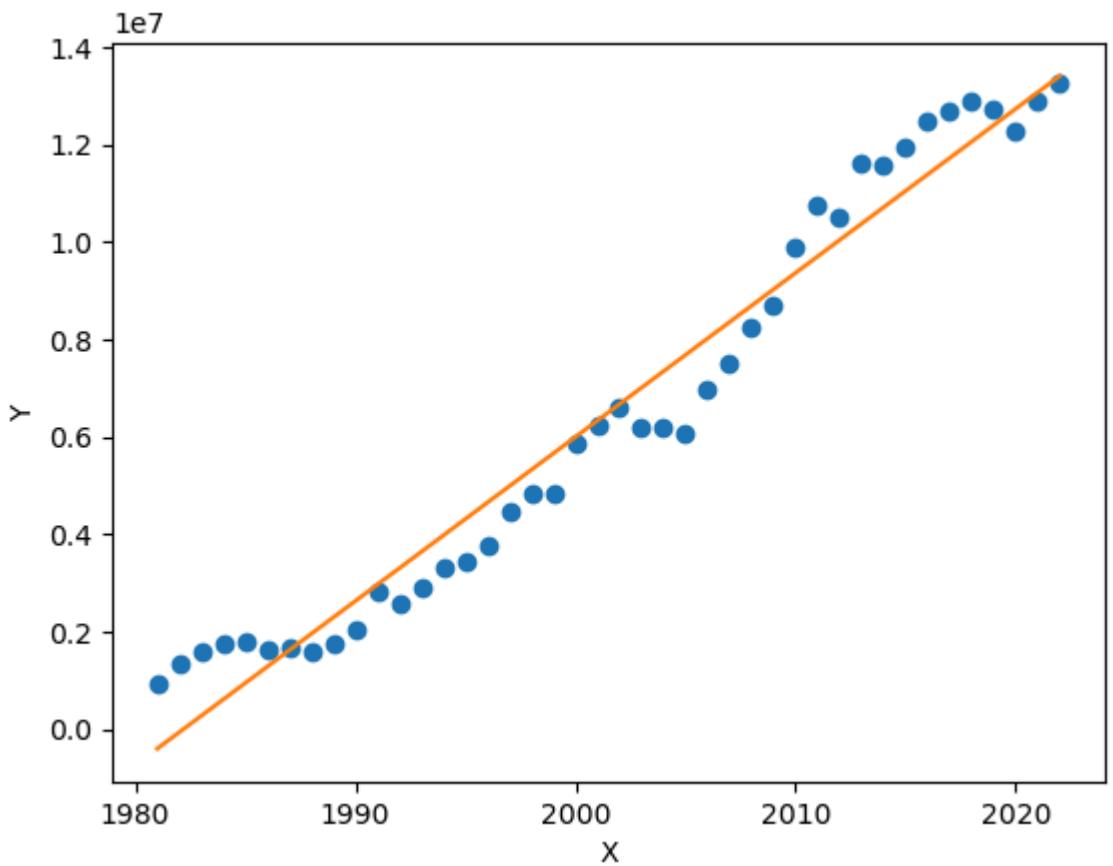
X = np.array(exp_df['year']).reshape(-1, 1)
Y = np.array(exp_df['total_expenditure_on_education']).reshape(-1, 1)
print("X =\n", X)
print("Y =\n", Y)

# xw = Y regression
W = helper.linearRegressionWithBias(X, Y, printResult=True, printFeature=1)

x_test = helper.paddingOfOnes(np.array([[2021]]))

y_predict = helper.testData(x_test, W, printResult=False)

print("y_predict =\n", y_predict)      You, 1 second ago • Uncommitted changes
```



```
y_predict =
[[13075210.44420362]]
```

Answer:

The predicted educational expenditure in year 2021 is 12102904.270643068.

- diff from answer probably cuz diff dataset

Q6

```
+ CategoryInfo          : InvalidArgument: (:) [Invoke-WebRequest], ParameterBindingException
+ FullyQualifiedErrorId : MissingArgument,Microsoft.PowerShell.Commands.InvokeWebRequestCommand
PS C:\University Stuff\Y2S2\EE2211\Tutut\tut5> curl -o ./wine.csv https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv
PS C:\University Stuff\Y2S2\EE2211\Tutut\tut5> ls
Directory: C:\University Stuff\Y2S2\EE2211\Tutut\tut5

Mode                LastWriteTime         Length Name
----                -----        ----  --
-a----   2/20/2024 12:43 PM           4325 tut5.py
-a----   2/19/2024 10:11 PM      193462 Tutorial15-Questions.pdf
-a----   2/20/2024 12:54 PM          84199 wine.csv

Ask Copilot or type / for commands >
```

- downloading files using curl

```

# Question 6 #####
print("Question 6 ######")
# wine_df = pd.read_csv("wine.csv", sep=';')
# print(wine_df)

wine_df = pd.read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv", sep=";")
wine_df.info()

```

- get data without downloading

```

wine_df = pd.read_csv("wine.csv", sep=';')
print(wine_df)

X = np.array(wine_df.drop('quality', axis=1)) # drop the quality column, the other columns are the features
Y = np.array(wine_df['quality'].tolist()).reshape(-1, 1) # the quality column is the output

print("X\n", X)
print("Y\n", Y)

Xtrain = X[:1500]
Ytrain = Y[:1500]
print("Xtrain\n", Xtrain)
print("Ytrain\n", Ytrain)

Xtest = X[1500:]
Ytest = Y[1500:]
print("Xtest\n", Xtest)
print("Ytest\n", Ytest)

# Training
w = helper.linearRegressionWithBias(Xtrain, Ytrain, printResult=True, printFeature=None)

# Prediction Test
Ytest_pred = helper.testData(helper.paddingOfOnes(Xtest), w, printResult=False)
print("Ytest_pred\n", Ytest_pred)

MSE = mean_squared_error(Ytest_pred, Ytest)
print("MSE = ", MSE)

```

```

W =
[[ 2.22330327e+01]
 [ 2.68702621e-02]
 [-1.12838019e+00]
 [-2.06141685e-01]
 [ 1.22000584e-02]
 [-1.77718503e+00]
 [ 4.29357454e-03]
 [-3.18953315e-03]
 [-1.81795124e+01]
 [-3.98142390e-01]
 [ 8.92474793e-01]
 [ 2.77147239e-01]]
MSE =  0.4216923253120808

```

MSE on training data

```

MSE =  0.3435263811803818

```

MSE on testing data, lower than 0.2 better regression

Q8

Question 8:

MCQ: There could be more than one answer.

Suppose $f(\mathbf{x})$ is a scalar function of d variables where \mathbf{x} is a $d \times 1$ vector. Then, without taking data points into consideration, the outcome of differentiation of $f(\mathbf{x})$ w.r.t. \mathbf{x} is

- a) a scalar
- b) a $d \times 1$ vector
- c) a $d \times d$ matrix
- d) a $d \times d \times d$ tensor
- e) None of the above

Answer: b) $\frac{df}{d\mathbf{x}} = \begin{bmatrix} a \\ a \\ \vdots \\ a \end{bmatrix} \left\{ \text{row} = d \right. \right)$

$$f(\mathbf{x}) = \frac{ax}{\uparrow \text{constant}} + \frac{b}{\downarrow \text{independent of } x} ?$$

$$f(\mathbf{x}) = \left[\begin{array}{c} ax_1 + b \\ ax_2 + b \\ \vdots \\ ax_d + b \end{array} \right] \left\{ \text{row} = d \right. \right)$$

Q9

```
X = np.array([[3, -1, 0], [5, 1, 2], [9, -1, 3], [-6, 7, 2], [3, -2, 0]])
Y = np.array([[1, -1], [-1, 0], [1, 2], [0, 3], [1, -2]])
```

```
# Training
w = helper.linearRegressionWithBias(X, Y, printResult=True, printFeature=None)

Xtest = helper.paddingOfOnes(np.array([[8, 0, 2]]))
Ytest = helper.testData(Xtest, w, printResult=False)
print("Ytest\n", Ytest)
```

```
rank(X) = 4
W =
[[ 1.14668974 -0.95997404]
[-0.630463 -0.33427088]
[-1.10601471 -0.24426655]
[ 1.3595846  1.77953267]]
MSE = 0.254110774556469
Ytest
 [[-1.17784509 -0.07507572]]
```

Tutorial 6

Q7

Question 7

MCQ: there could be more than one answer. Given three samples of two-dimensional data points $\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 3 & 3 \end{bmatrix}$

with corresponding target vector $\mathbf{y} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$. Suppose you want to use a full third-order polynomial model to fit these data. Which of the following is/are true?

- a) The polynomials model has 10 parameters to learn
- b) The polynomial learning system is an under-determined one
- c) The learning of the polynomial model has infinite number of solutions
- d) The input matrix \mathbf{X} has linearly dependent samples
- e) None of the above

Answer: a, b, c, d

Q8

Question 8

MCQ: there could be more than one answer. Which of the following is/are true?

- a) The polynomial model can be used to solve problems with nonlinear decision boundary.
- b) The ridge regression cannot be applied to multi-target regression.
- c) The solution for learning feature \mathbf{X} with target \mathbf{y} based on linear ridge regression can be written as $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$ for $\lambda > 0$. As λ increases, $\hat{\mathbf{w}}^T \hat{\mathbf{w}}$ decreases.
- d) If there are four data samples with two input features each, the full second-order polynomial model is an over-determined system.

Answer: a, c

PY

Three balls are drawn from three urns sequentially, one ball from each urn. The first urn contains 1 blue and 7 red balls, the second urn contains 2 blue and 6 red balls, and the third urn contains 3 red and 5 green balls. Find the probability that 2 red balls are chosen. (3 marks)

Currently Selected: B

A 270/1024

B 270/512

C 226/64

D 226/512

E None of the rest.

This question is related to the understanding of linear systems and partial derivatives. Which of the following statements below is correct?

Currently Selected: A

A $\begin{bmatrix} 1 & 4 \\ 2 & 7 \\ -3 & 11 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -2.5 \\ 4 \end{bmatrix}$ has no exact solution but an approximated solution is available using the left inverse.

B If $f(\mathbf{x})$ is a vector-valued function of size $p \times 1$ and \mathbf{x} is an $m \times 1$ vector, then differentiation of $f(\mathbf{x})$ with respect to \mathbf{x} is an $m \times p$ matrix.

C A linear function needs to satisfy the properties of homogeneity only.

D In over-determined linear systems, the number of parameters is greater than the number of unknown equations.

E None of the other options.

Which of the following task is likely to be achieved via supervised learning?

Currently Selected: A

A Using historical data for weather forecast.

B Grouping together users with similar viewing patterns in order to recommend similar content.

C Grouping a number of oranges by their size.

D None of the rest.

Causality is a deterministic relationship; suppose we know A and B have causal relation, if A occurs, B is for sure to take place.

Currently Selected: A

A True

B False

- False
- [Causality is a Statistical Relationship](#)

A set of linear equations is written as $\mathbf{w}^T \mathbf{X} = \mathbf{y}^T$ where $\mathbf{X} \in \mathcal{R}^{3 \times 2}$ and $\mathbf{y} \in \mathcal{R}^{2 \times 1}$. How many simultaneous equations are there in this set of equations?

Currently Selected: B

A 1

B 2

C 3

D 5

E 4

A discrete random variable takes a finite number of values, while a continuous random variable can only take infinite number of values. (2 marks)

Currently Selected: A

A True

B False

- False: Discrete can be infinite also

This question is related to the understanding of modelling assumptions. $f(x) = 5x - 3$ is a linear function. (2 marks)

Currently Selected: B

A True

B False

A machine learning algorithm takes the temperature as one of its input features. The temperature is measured in Celsius. Please select the correct option.

Currently Selected: E

A The temperature in Celsius is considered as interval data.

B We can calculate the mean and standard deviation of temperature.

C The temperature in Celsius is considered as ratio data.

D None of the rest.

E (a), and (b)

F (a), and (c)

A machine learning algorithm takes the letter grade of students as one of its input features. The letter grade can take any element from {A+, A, A-, B+, B, B-, C+, C, D+, D, F}, subject to some distribution curving. For example, A+ corresponds to 95%, A corresponds to 85%, and A- corresponds to 80%. Which of the following statements is/are true?

Currently Selected: F

A The letter grade is an example of nominal variable.

B The letter grade is an example of ordinal variable.

C The letter grade is an example of interval variable.

D The letter grade is an example of discrete variable.

E (a), and (c)

F (b), and (d)

G (b), and (c)

A person draws 2 cards from a deck of 52 cards, one after another without replacing the previous card back. What is the probability of drawing two Queens in a row?

Currently Selected: B

A 4/52

B 1/221

C 3/51

D 2/52

One key step in Data Cleaning is to check the missing features of data samples. When we have insufficient number of training samples in our dataset, we may consider removing the examples with missing features.

Currently Selected: B

A True

B False

Suppose the random variable X has a probability mass function (pmf) given in the table below.

X	1	2	3	4	5
Pr[X]	0.1	(BLANK1)	0.2	0.4	(BLANK2)

We also know that the expected value of X is 3.5.

1) What is the probability of $\text{Pr}[X=5]?$ ① 0.2 (2 Marks)

2) What is the probability of $\text{Pr}[X \leq 2]?$ ② 0.2 (2 Marks)

Answers 1 - 2

1. 0.2

2. 0.2

1. Range - Min:0.19999 Max:0.20001
2. Range - Min:0.19999 Max:0.20001

We have a collection of 1,000 images from three classes: cat, bird, and dog. With these images, we would like to train an image classifier that categorizes an input image into one of the three classes.

To ensure the 1,000 images are of good quality for training the classifier, we ask a well-trained human inspector to go through all the images, to label the images and remove noisy ones. Eventually, we removed 200 images of low quality suggested by the inspector, and use the remaining 800 images to train the classifier; the 800 images comprise 200 cat images, 300 bird images, and 300 dog images.

Please select the correct option.

Currently Selected: A

- A The human inspection process can be considered as a data cleaning step.

If we are to use one-hot encoding for the labels of the three classes, we can set:

- B Cat = [1 1 1]
Dog = [0 1 0]
Bird = [0 0 1]

- C The image classification conducted here is an unsupervised-learning task.

- D If we keep the 200 noisy images (suggested by the human inspector), we will end up having more training images and hence a better-performed image classifier.

- E (a) and (b)

- F (a), (b), and (d)

- G None of others is correct.

You are given a collection of 5 training data points of two features (x_1, x_2) and their target output (y) which are packed as follows:

$$\text{Feature matrix: } X = \begin{bmatrix} 1 & 2 \\ 0 & 6 \\ 1 & 0 \\ 0 & 5 \\ 1 & 7 \end{bmatrix}, \text{ Target output: } y = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}.$$

Predict the output (up to 4 decimal places) of $(x_1, x_2) = (1, 3)$ using the linear regression model. (4 marks)

1) What is the mean of squared error of the estimated model? 1. 1.3887 (up to 4 decimal places, 2 marks)

2) The prediction for y is 2. 2.6603 (up to 4 decimal places, 2 marks).

Answers 1 - 2

1. 1.3887

2. 2.6603

1. Range - Min:1.3886 Max:1.3888
2. Range - Min:2.9999 Max:3.0001

- type in wrongly lmao, get 3

```

X = np.array([[1, 2], [0, 6], [1, 0], [0, 5], [1, 7]])
Y = np.array([[1], [2], [3], [4], [5]])

w = helper.linearRegressionWithBias(X, Y, printResult=True, printFeature=None)

xtest = np.array([[1, 3]])
xtest = helper.paddingOfOnes(xtest)
ytest = helper.testData(xtest, w, True)

```

```

rank(X) = 3
W =
[[1.13207547]
[0.8490566 ]
[0.33962264]]
MSE = 1.388679245283019
X @ W = Y =
[[3.]]

```

The values of feature x and their corresponding values of target y are shown in the table below.

x	3	4	5	6	7
y	5	4	3	2	1

Find the least square regression line $y = a + bx$ and then estimate the value of y when $x = 8$.

Currently Selected: A

A $y = 0$

B $y = -1$

C $y = +1$

D $y = 8$

E None of the above

```

X = np.array([[3], [4], [5], [6], [7]])
Y = np.array([[5], [4], [3], [2], [1]])

w = helper.linearRegressionWithBias(X, Y, printResult=True, printFeature=None)

xtest = np.array([[8]])
xtest = helper.paddingOfOnes(xtest)
ytest = helper.testData(xtest, w, True)

```

```

rank(X) = 2
W =
[[ 8.]
[-1.]]
MSE = 4.733165431326071e-30
X @ W = Y =
[[4.4408921e-15]]

```

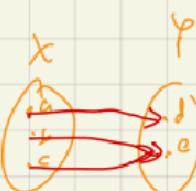
Functions

Functions : (roughly speaking) data of $(X, Y, x \mapsto f(x))$
 where X, Y are sets and $x \mapsto f(x)$ assigns $x \in X$ to some $f(x) \in Y$.

X : domain
 Y : target / codomain

Notation $f: X \rightarrow Y$

image diagram: e.g. $f: \mathbb{R} \rightarrow \mathbb{R}$ { e.g. $f(x) = x + 1$



- function is defined with three stuff
 - input space X
 - output space Y
 - and the transformation

Right-Inverse: $L F =$

Left-Inverse: $F L =$

DEFINITIONS

Let $f: X \rightarrow Y$ be a function.

- f is injective (one-one) if for all $x_1, x_2 \in X$, f has the property that $f(x_1) = f(x_2) \Rightarrow x_1 = x_2$ implies

- f is surjective (onto) if for all $y \in Y$, there exist some $x \in X$ s.t. $f(x) = y$

- Injective means the function
 - $x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$, if input different, output must be different
 - $f(x_1) = f(x_2) \Rightarrow x_1 = x_2$, if output same, input must be same

- Contrapositive example:
 - if rain \Rightarrow cloud
 - not contrapositive, converse: no rain \Rightarrow no cloud
 - this is not the same statement, no rain could be cloud or no cloud
 - if apply contrapositive to (no rain \Rightarrow no cloud): cloud \Rightarrow rain which is confirm not the same statement
 - contrapositive: no cloud \Rightarrow no rain
 - this is the same statement, if rain must have cloud, if no cloud must be no rain
 - contrapositive is done by negation and flipping \Rightarrow
- for linear function:
 - injective $\Leftrightarrow (f(x) = 0 \Rightarrow x = 0)$
 - showing \Leftarrow
 - assume $(f(x) = 0 \Rightarrow x = 0)$ is true
 - want to show: $f(x_1) = f(x_2) \Rightarrow x_1 = x_2$
 - assume $f(x_1) = f(x_2)$
 - $f(x_1) - f(x_2) = 0$
 - by linearity of f , $f(x_1 - x_2) = 0$
 - by hypothesis, $x_1 - x_2 = 0$
 - $x_1 = x_2$
 - have shown $f(x_1) = f(x_2) \Rightarrow x_1 = x_2$
 - showing \Rightarrow
 - assume injective: $f(x_1) = f(x_2) \Rightarrow x_1 = x_2$
 - want to show: $f(x) = 0 \Rightarrow x = 0$
 - assume $f(x) = 0$
 -
 - have shown $f(x) = 0 \Rightarrow x = 0$
 - Example: $f(x) = mx$ for $m \neq 0$
 - $f(x_1) = f(x_2)$
 - $mx_1 = mx_2$
 - $x_1 = x_2$
 - \Rightarrow injective
- Surjective means the function can output all y in Y as define in the set
 - surjectivity of a function depends on the Y that is defined
- Example: $f(x) = x^2$, $X = \mathbb{R}$, $Y = \mathbb{R}$
 - f is not injective as x^2 maps many to one
 - f is not surjective as f cannot map to negative numbers
 - f does not have inverse, nor right nor left inverse
- Example: $f(x) = x^2$, $X = \mathbb{R}$, $Y = \text{non negative } \mathbb{R}$
 - f is not injective as x^2 maps many to one
 - f is surjective as f can map to all values in Y
 - f has a right inverse (means apply right inverse to input first)
 - sqrt() then square
 -

A2 Dump

<https://stackoverflow.com/questions/50132322/how-does-multiple-target-ridge-regression-work-in-scikit-learn>

Admin

- Introduction and Preliminaries (Xinchao)
 - Introduction
 - Data Engineering
 - Introduction to Probability and Statistics
- Fundamental Machine Learning Algorithms I (Yueming)
 - Systems of linear equations
 - Least squares, Linear regression
 - Ridge regression, Polynomial regression
- Fundamental Machine Learning Algorithms II (Yueming)
 - Over-fitting, bias/variance trade-off
 - Optimization, Gradient descent
 - Decision Trees, Random Forest
- Performance and More Algorithms (Xinchao)
 - Performance Issues
 - K-means Clustering
 - Neural Networks

- Schedule
 - **12 Weeks Lectures**, starting from Week 1
 - **12 Weeks Tutorials**, starting from Week 2
 - **2 Programming Tutorials** (optional and highly recommended)
 - Week 1 – 2, Friday
 - Right after the lecture (i.e., 2 to 3 PM)
 - **1 Mid-term Quiz** (using ExamSoft)
 - Tentatively held offline on 9 March 2024 (Recess Week)
 - Content up to Week 5 (inclusive)
 - **1 Briefing Session on ExamSoft**
 - Tentative on Week 3, exact time to be confirmed with CIT Staff
 - **1 Final Exam** (using ExamSoft)
 - Held on 3-May-2024
 - **3 Assignments**
 - Assignment 1: released on Week 4, due on Week 6 (tentatively)
 - Assignment 2: released on Week 6, due on Week 9 (tentatively)
 - Assignment 3: released on Week 9, due on Week 13 (tentatively)

- midterm is 16 march

- 3 Assignments (36%) + Tutorial Attendance (4%)
 - 1 Mid-term (30%)
 - 1 Final Exam (30%)
-
- Held online:
 - Lectures
 - Held offline (in classrooms):
 - Tutorials
-
- Videos of lectures are made available after lectures.
 - midterm is open book