

Machine Learning with Python

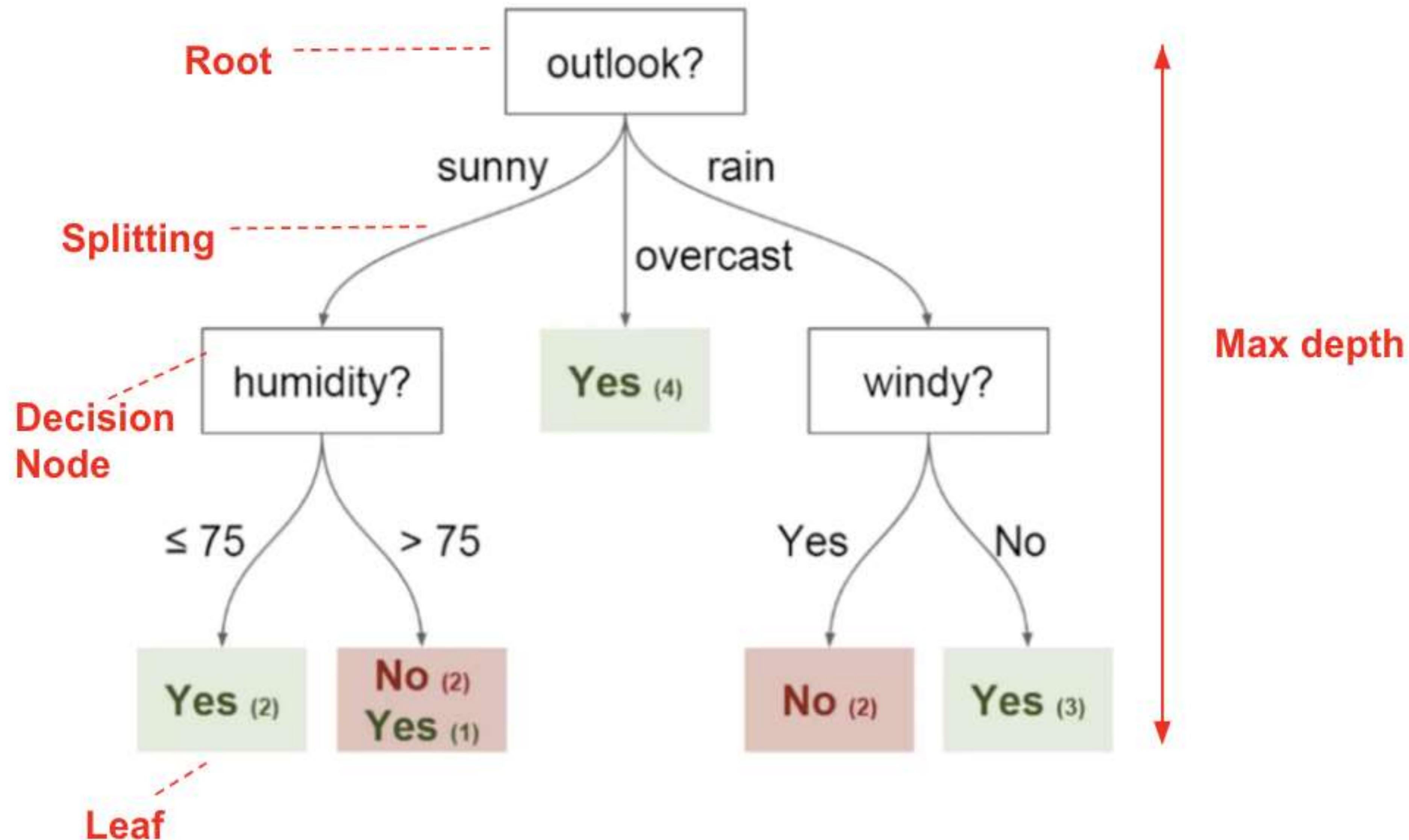
Classification and Regression Trees

Arghya Ray

Decision Tree

- A decision tree is a popular classification method that results in a flow-chart like tree structure where each node denotes a test on an attribute value and each branch represents an outcome of the test. The tree leaves represent the classes.
- Decision tree is a model that is both predictive and descriptive.
- **Advantages:**
 - Decision tree approach is widely used since it is efficient and can deal with both continuous and categorical variables.
 - The decision tree approach is able to deal with missing values in the training data and can tolerate some errors in data.
 - The decision tree approach is perhaps the best if each attribute takes only a small number of possible values.
- **Disadvantages:**
 - Decision trees are less appropriate for tasks where the task is to predict values of a continuous variable like share price or interest rate.
 - Decision trees can lead to a large number of errors if the number of training examples per class is small.
 - The complexity of a decision tree increases as the number of attributes increases.
- **Measuring the quality of a decision tree** is an interesting problem altogether. **Classification accuracy** determined using test data is obviously a good measure but other measures like, **average cost** and **worst case cost** of classifying an object may be used.

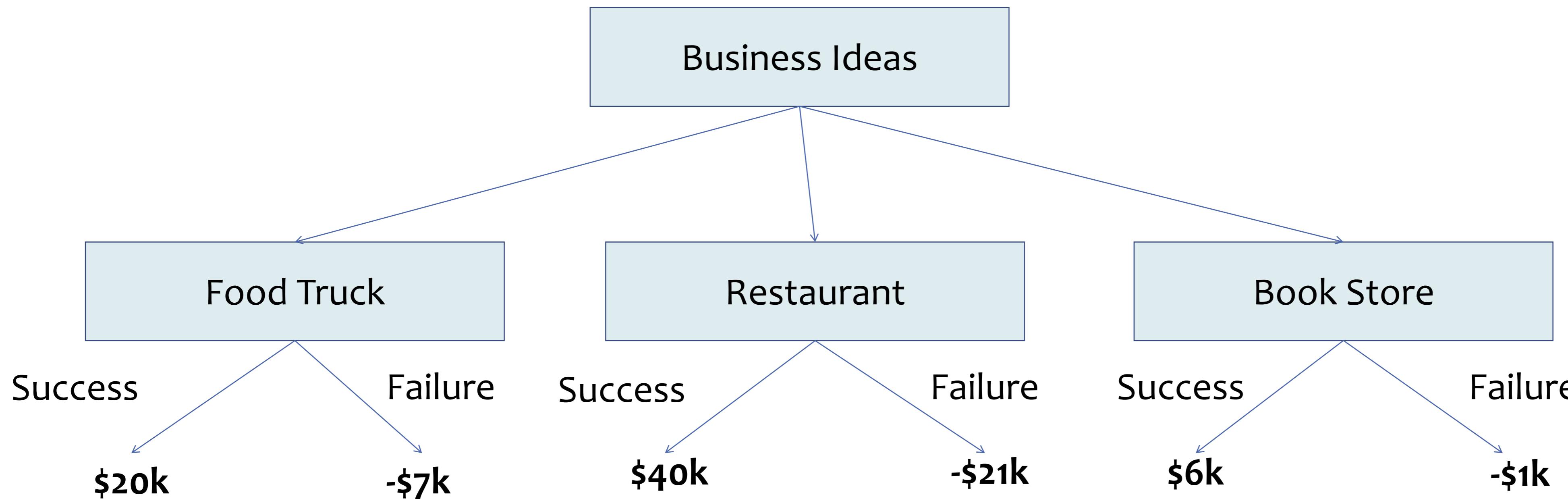
Decision Tree Diagram



1. A decision tree is an approach to analysis that can help you make decisions.

Suppose for example you need to decide whether to invest a certain amount of money in one of the three business projects: a food-truck business, a restaurant, or a bookstore based on the data given below.

	Business Success Percentage		Business Value Changes	
Business	Success Rate	Failure Rate	Gain (USD)	Loss (USD)
Food Truck	60%	40%	20000	-7000
Restaurant	52%	48%	40000	-21000
Bookstore	50%	50%	6000	-1000



- In these cases, ***the expected value*** calculated based on all possible outcomes helps in figuring out the business decision making.
- Expected Value for the food truck business = $(60\% \text{ of USD } 20000) + (40\% \text{ of USD } (-7000)) = \text{USD } 9200$.
- Expected Value of restaurant business = $(52\% \text{ of USD } 40000) + (48\% \text{ of USD } (-21000)) = \text{USD } 10720$.
- Expected Value of bookstore business = $(50\% \text{ of USD } 6000) + (50\% \text{ of USD } (-1000)) = \text{USD } 2500$
- Here the expected value reflects the average gain from investing in the business. Based on the above hypothetical figures, the results reflect that if you attempt to invest in a businesses say Food Truck business several times (under the same circumstances each time), your average profit will be USD 9200 per business.

2. Decision trees can also be used to visualize classification rules.

Classification and Regression Trees

Goal: Classify or predict an outcome based on a set of predictors.

The output is a set of **rules**

Example:

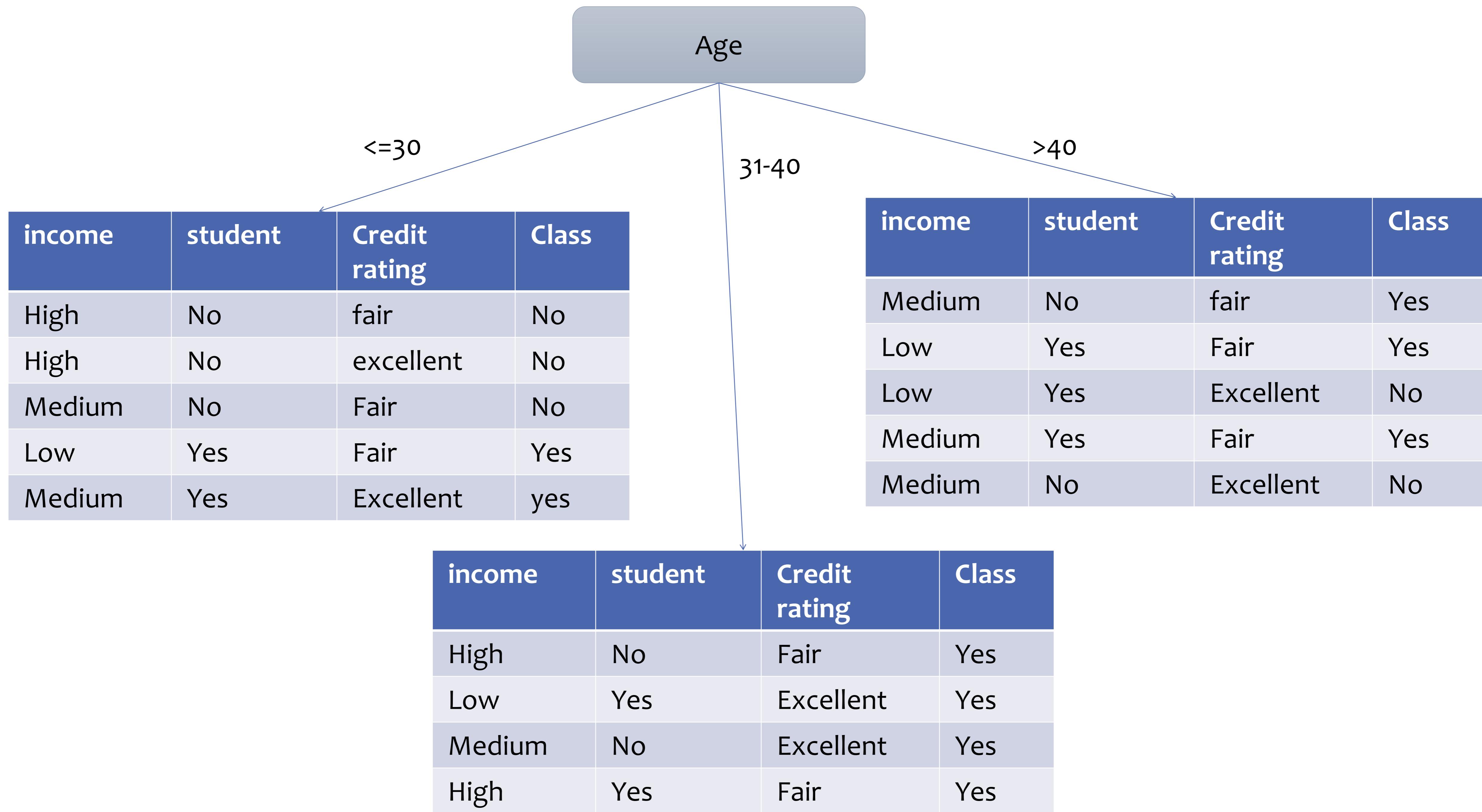
- Goal: classify a record as “will accept credit card offer” or “will not accept”
- Rule might be “IF (Income > 92.5) AND (Education < 1.5) AND (Family <= 2.5) THEN Class = 0 (non-acceptor)”
- **Recursive partitioning:** Repeatedly split the records into two parts so as to achieve maximum homogeneity within the new parts

Recursive partitioning steps:

- Pick one of the predictor variables, x_i
- Pick a value of x_i , say s_i , that divides the training data into two (not necessarily equal) portions
- Measure how “pure” or homogeneous each of the resulting portions are
- “Pure” = containing records of mostly one class
- Algorithm tries with different variables (x) and different values of x_i , i.e., s_i to maximize purity in a split
- After you get a “maximum purity” split, repeat the process for a second split, and so on

Forming a tree from the given example

RID	Age	Income	Student	Credit rating	Class (buys computer)
1	<=30	High	No	Fair	No
2	<=30	High	No	Excellent	No
3	31-40	High	No	Fair	Yes
4	>40	Medium	No	Fair	Yes
5	>40	Low	Yes	Fair	Yes
6	>40	Low	Yes	Excellent	No
7	31-40	Low	Yes	Excellent	Yes
8	<=30	Medium	No	Fair	No
9	<=30	Low	Yes	Fair	Yes
10	>40	Medium	Yes	Excellent	Yes
11	<=30	Medium	Yes	Excellent	Yes
12	30-40	Medium	No	Excellent	Yes
13	30-40	High	Yes	Fair	Yes
14	>40	Medium	No	Excellent	No



Measuring Impurity

- Gini Index (measure of impurity)

- Gini Index for rectangle A containing m cases

$$I(A) = 1 - \sum_{k=1}^m p_k^2$$

p = proportion of cases in rectangle A that belong to class k

- $I(A) = 0$ when all cases belong to same class (most pure)

- Entropy (measure of impurity)

$$\text{entropy}(A) = - \sum_{k=1}^m p_k \log_2(p_k)$$

p = proportion of cases (out of m) in rectangle A that belong to class k

Entropy ranges between 0 (most pure) and $\log_2(m)$ (equal representation of classes)

Using the principle of ‘Information entropy’ build a ‘decision tree’ using the training data given below. Divide the ‘credit rating’ attribute into ranges as follows: (0, 1.6], (1.6,1.7], (1.7,1.8], (1.8,1.9], (1.9,2.0], (2.0,5.0]

Sr. No.	Profession	Credit rating	Class
1	Business	1.6	Buys only laptop
2	Service	2.0	Buys laptop with CD Writer
3	Business	1.9	Buys laptop with printer
4	Business	1.88	Buys laptop with printer
5	Business	1.70	Buys only laptop
6	Service	1.85	Buys laptop with printer
7	Business	1.60	Buys only laptop
8	Service	1.70	Buys only laptop
9	Service	2.20	Buys laptop with CD writer
10	Service	2.10	Buys laptop with CD writer
11	Business	1.80	Buys laptop with printer
12	Service	1.95	Buys laptop with printer
13	Business	1.90	Buys laptop with printer
14	Business	1.80	Buys laptop with printer
15	Business	1.75	Buys laptop with printer

Profession

Business

Service

Sr. No.	Credit rating	Class
1	1.6	Buys only laptop
2	1.9	Buys laptop with printer
3	1.88	Buys laptop with printer
4	1.70	Buys only laptop
5	1.60	Buys only laptop
6	1.80	Buys laptop with printer
7	1.90	Buys laptop with printer
8	1.80	Buys laptop with printer
9	1.75	Buys laptop with printer

Sr. No.	Credit rating	Class
1	2.0	Buys laptop with CD Writer
2	1.85	Buys laptop with printer
3	1.70	Buys only laptop
4	2.20	Buys laptop with CD writer
5	2.10	Buys laptop with CD writer
6	1.95	Buys laptop with printer

Credit Rating

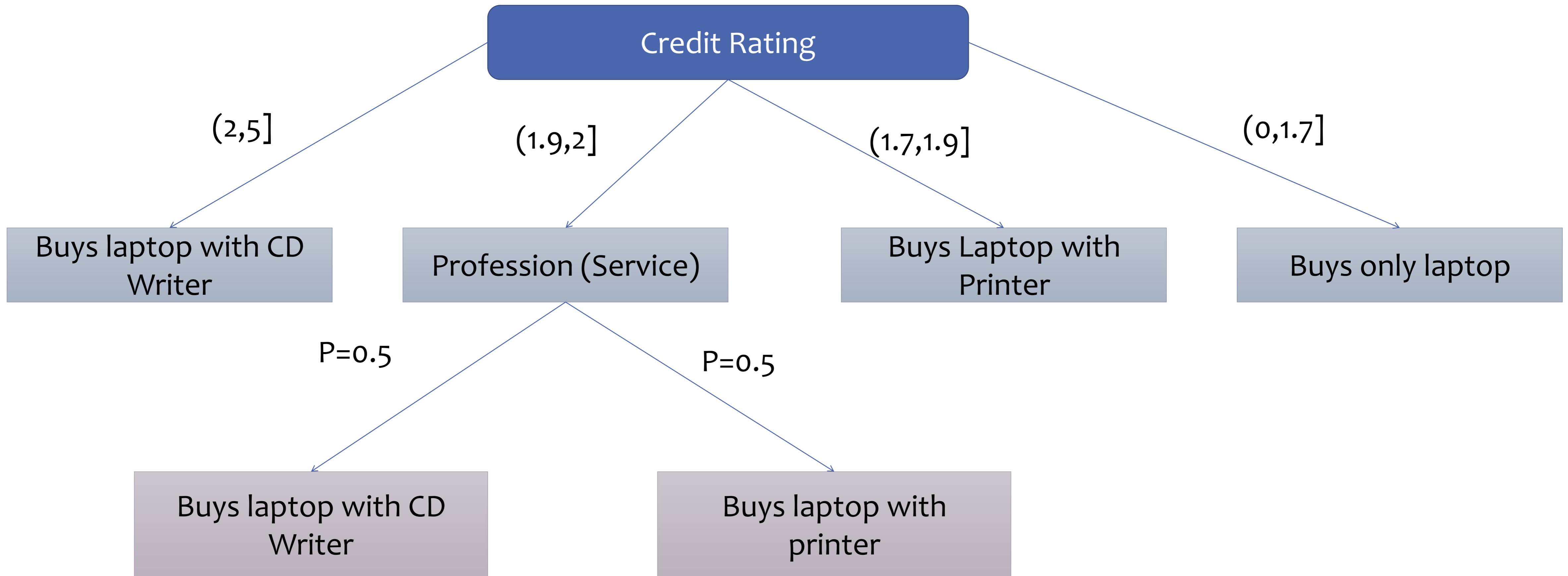
(0,1.6]

(1.6,1.7]

Sr. No.	Profession	Class
1	Business	Buys only Laptop
2	Business	Buys only Laptop

Sr. No.	Profession	Class
1	Business	Buys only Laptop
2	Service	Buys only Laptop

- Initially there are 3 classes: Buys only laptop, buys laptop with CD writer, buys laptop with printer
- Initial Overall Entropy (E_0) = $-\sum_{i=1}^3 p_i \log_3 p_i = -\left[\frac{4}{15} \log_3 \frac{4}{15} + \frac{3}{15} \log_3 \frac{3}{15} + \frac{8}{15} \log_3 \frac{8}{15}\right] = 0.918$
- Based on Profession : 9 Business, 6 Service
- $\text{Entropy}(\text{Profession}) = \frac{9}{15} \text{Entropy}(\text{business}) + \frac{6}{15} \text{Entropy}(\text{service}) = \frac{9}{15} \left(-\frac{3}{9} \log_3 \frac{3}{9} - \frac{6}{9} \log_3 \frac{6}{9}\right) + \frac{6}{15} \left(-\frac{1}{6} \log_3 \frac{1}{6} - \frac{2}{6} \log_3 \frac{2}{6}\right)$



The content of the slides are prepared from different textbooks.

References:

- Data Mining and Predictive Analytics, By Daniel T. Larose. Copyright 2015 John Wiley & Sons, Inc.
- Predictive Analytics for Dummies, By Anasse Bari, Mohamed Chaouchi, & Tommy Jung, Copyright 2016, John Wiley & Sons, Inc.
- Introduction to Data Mining with Case Studies, By G.K. Gupta. Copyright 2014 by PHI Learning Private Limited.



—
Thank you..

Machine Learning with Python

Session 14: Agglomerative Clustering

Arghya Ray

Measuring Distance Between Clusters:

- **Single Linkage**
 - Minimum Distance (Cluster A to Cluster B)
 - Distance between two clusters is the distance between the pair of records A_i and B_j that are closest.
- **Complete Linkage**
 - Maximum Distance (Cluster A to Cluster B)
 - Distance between two clusters is the distance between the pair of records A_i and B_j that are farthest from each other
- **Average Linkage**
 - Distance between two clusters is the average of all possible pair-wise distances
- **Centroid**
 - Distance between two clusters is the distance between the two cluster centroids.
 - Centroid is the vector of variable averages for all records in a cluster $(x_1, y_1, z_1) \quad (x_2, y_2, z_2) \quad (x_3, y_3, z_3)$
$$C_1 = ((x_1+x_2+x_3)/3, (y_1+y_2+y_3)/3, (z_1+z_2+z_3)/3)$$

Q. Consider the following three clusters, each with four members:

Cluster 1:	$\{(1,5), (2,4), (3,3), (2,1)\}$	Centroid- $((1+2+3+2)/4, (5+4+3+1)/4)$ (2,3.25)
Cluster 2:	$\{(5,4), (6,6), (7,5), (8,8)\}$	
Cluster 3:	$\{(4,1), (3,0), (5,1), (6,2)\}$	

Distance Between	Single Link	Complete-Link	Centroid	Average-Link
Cluster 1 & 2	2.24	9.22	5.15	5.43
Cluster 2 & 3	2.24	9.43	5.15	5.38
Cluster 3 & 1	1.41	5.83	3.36	3.76

- Whichever distance algorithm is applied, two nearby clusters can be merged if an agglomerative approach is being used.
- It has been reported that the **complete link algorithm** generally produces compact and more useful clusters.
- The **single link algorithm** tends to suffer from chaining effects and elongated clusters in some situations. However, the single link algorithm is found to be effective in some applications.
- Both the complete link algorithm and single link algorithm can suffer from the presence of outliers.

Agglomerative Method:

The basic idea of the agglomerative method is to start out with n-clusters for ‘n’ data points and keep on combining points.

Steps involved:

1. Allocate each point to a cluster of its own. Thus we start with n clusters for n objects.
2. Create a distance matrix by computing distances between all pairs of clusters using one of the distance measuring methods (e.g. single link metric or complete link metric). Sort these distances in ascending order.
3. Find the two clusters that have the smallest distance between them.
4. Remove the pair of objects and merge them. When you are merging, take the average value of the two cluster distances.
5. If there is only one cluster left, then stop.
6. Compute all distances from the new cluster and update the distance matrix after the merger and go to step 3.

Use agglomerative clustering method for clustering the data (use centroid method for calculating distance between clusters).

Student	Age	Mark1	Mark2	Mark3	
S1	18	73	75	57	$S1 \rightarrow (18, 73, 75, 57)$
S2	18	79	85	75	$S2 \rightarrow (18, 79, 85, 75)$
S3	23	70	70	52	
S4	20	55	55	55	
S5	22	85	86	87	
S6	19	91	90	89	
S7	20	70	65	60	
S8	21	53	56	59	
S9	19	82	82	60	
S10	47	75	76	77	

S1 and S2 → 34

S2 and S3 → 52

S1 and S3 → 18

S2 and S4 → 76

Step 1 and 2: Allocate each point to a cluster and compute the distance matrix using the centroid method. The distance matrix is symmetric.

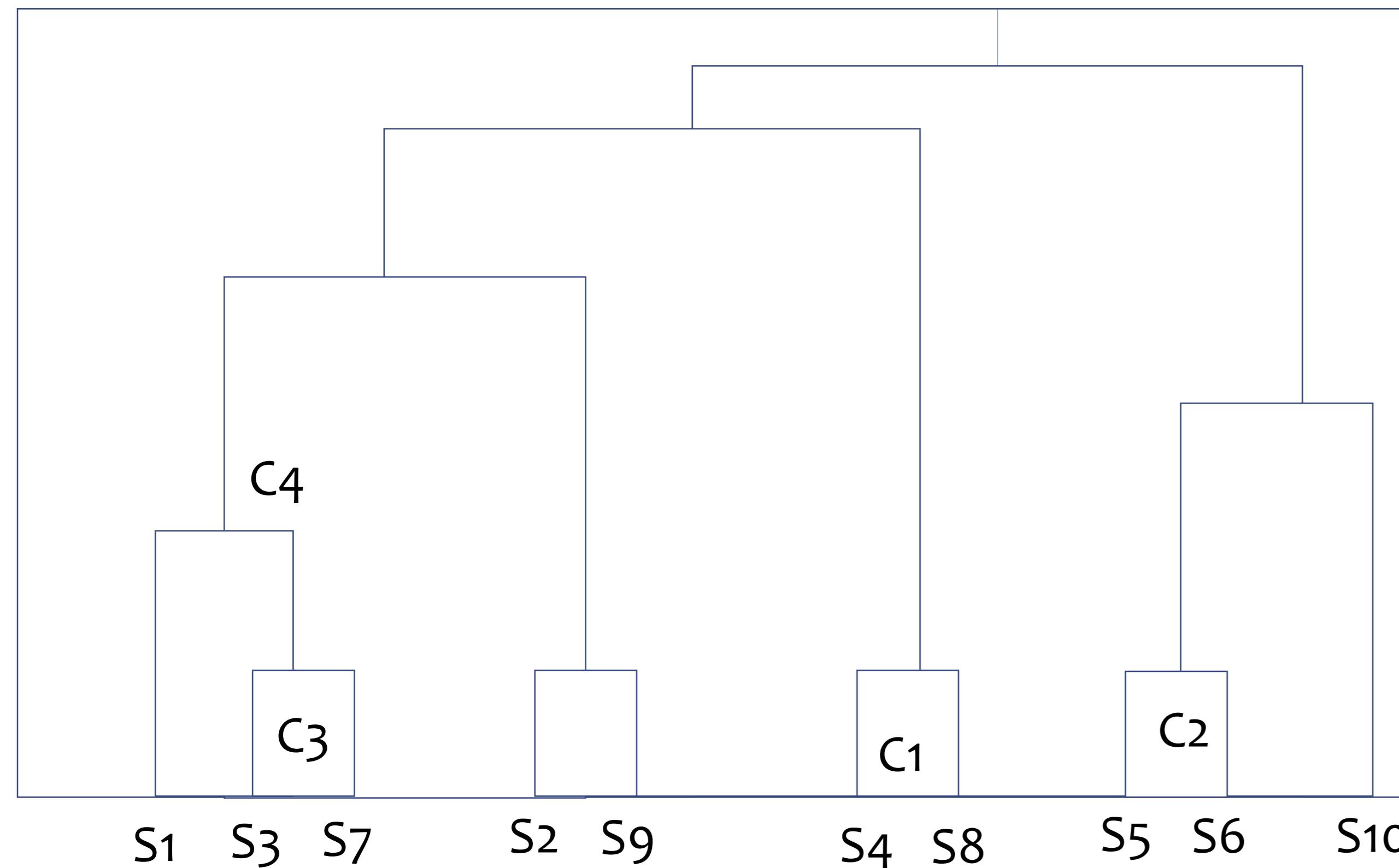
	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10
S1	0									
S2	34	0								
S3	18	52	0							
S4	42	76	36	0						
S5	57	23	67	95	0					
S6	66	32	82	106	15	0				
S7	18	46	16	30	65	76	0			
S8	44	74	40	8	91	104	28	0		
S9	20	22	36	60	37	46	30	58	0	
S10	52	44	60	90	55	70	60	86	58	0

Step 3 and 4: The smallest distance is 8 between objects S4 and S8. We combine this to cluster (C1) and put it where S4 was.

	S1	S2	S3	C1	S5	S6	S7	S9	S10
S1	0								
S2	34	0							
S3	18	52	0						
C1	41	75	38	0					
S5	57	23	67	95	0				
S6	66	32	82	106	15	0			
S7	18	46	16	29	65	76	0		
S9	20	22	36	59	37	46	30	0	
S10	52	44	60	88	55	70	60	58	0

Step 5 and 6: The smallest distance is now 15 between objects S5 and S6. We combine this to cluster (C2) and put it where S5 was.

	S1	S2	S3	C1	C2	S7	S9	S10
S1	0							
S2	34	0						
S3	18	52	0					
C1	41	75	38	0				
C2	61.5	27.5	74.5	97.5	0			
S7	18	46	16	29	69.5	0		
S9	20	22	36	59	41.5	30	0	
S10	52	44	60	88	62.5	60	58	0



	S1	S2	S3	C1	C2	S9	S10
S1	0						
S2	34	0					
C3	15	49	0				
C1	41	75	30	0			
C2	61.5	27.5	71.5	97.5	0		
S9	20	22	33	59	41.5	0	
S10	52	44	60	88	62.5	58	0

Divisive Hierarchical Method:

The basic idea of the divisive method is that it starts with the whole dataset as one cluster and then proceeds to recursively divide the cluster into two sub-clusters and continues until each cluster has only one object or some other stopping criterion has been reached. There are two types of divisive methods:

- **Monothetic:** It splits a cluster using only one attribute at a time.
- **Polythetic:** It splits a cluster using all attributes together.

Steps involved in a polythetic divisive method:

1. Decide a method of measuring the distance between two objects. Also decide a threshold distance.
2. Create a distance matrix by computing distance between all pairs of objects within the cluster. Sort these distances in ascending order.
3. Find the two objects that have the largest distance between them. They are most dissimilar.
4. If the distance between the two objects is smaller than the pre-specified threshold and there is no other cluster that needs to be divided then stop, otherwise continue.
5. Use the pair of objects as seeds of a K-means method to create two new clusters
6. If there is only one object in each cluster, then stop otherwise continue with Step 2.

Use divisive clustering method for clustering the data (use centroid method for calculating distance between clusters).

Student	Age	Mark1	Mark2	Mark3
S1	18	73	75	57
S2	18	79	85	75
S3	23	70	70	52
S4	20	55	55	55
S5	22	85	86	87
S6	19	91	90	89
S7	20	70	65	60
S8	21	53	56	59
S9	19	82	82	60
S10	47	75	76	77

Step 1 and 2: Allocate each point to a cluster and compute the distance matrix using the centroid method. The distance matrix is symmetric.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10
S1	0									
S2	34	0								
S3	18	52	0							
S4	42	76	36	0						
S5	57	23	67	95	0					
S6	66	32	82	106	15	0				
S7	18	46	16	30	65	76	0			
S8	44	74	40	8	91	104	28	0		
S9	20	22	36	60	37	46	30	115	0	
S10	52	44	60	90	55	70	60	98	99	0

The largest distance is 115 between the objects S8 and S9. They becomes the seed for two new clusters. K-Means is used to split the group into two clusters.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10
S8	44	74	40	8	91	104	28	0	115	98
S9	20	22	36	60	37	46	30	115	0	99

Cluster C1: S4, S7, S8, S10

Cluster C2: S1, S2, S3, S5, S6, S6, S9

Cluster C1: S₄, S₇, S₈, S₁₀

Cluster C2: S₁, S₂, S₃, S₅, S₆, S₆, S₉

If the stopping criteria is not met, we can follow the previous steps and divide these two clusters again one by one.

For Cluster C2:

	S ₁	S ₂	S ₃	S ₅	S ₆	S ₉
S ₁	0					
S ₂	34	0				
S ₃	18	52	0			
S ₅	57	23	67	0		
S ₆	66	32	82	15	0	
S ₉	20	22	36	37	46	0

The largest distance is 82 between the objects S₃ and S₆. They becomes the seed for to new clusters.

For Cluster C1:

	S ₄	S ₇	S ₈	S ₁₀
S ₄	0			
S ₇	30	0		
S ₈	8	28	0	
S ₁₀	90	60	98	0

The largest distance is 98 between the objects S₈ and S₁₀. They becomes the seed for to new clusters.

This continues until one of the stopping criteria is met.

The content of the slides are prepared from different textbooks.

References:

- Data Mining and Predictive Analytics, By Daniel T. Larose. Copyright 2015 John Wiley & Sons, Inc.
- Predictive Analytics for Dummies, By Anasse Bari, Mohamed Chaouchi, & Tommy Jung, Copyright 2016, John Wiley & Sons, Inc.
- Introduction to Data Mining with Case Studies, By G.K. Gupta. Copyright 2014 by PHI Learning Private Limited.



—
Thank you..

Machine Learning with Python

Session 12: Cluster Analysis

Arghya Ray

Identifying Similarities in Data

- Large amount of information are constantly being generated, organized, analyzed and stored.
- Identifying important patterns, associations and groupings of similar data can be helpful for customers as well as organizations.
- **Data Clustering** can help us make sense of the huge amount of data by discovering hidden groupings of similar items.
- Clustering can help in analysis of online social networks.
- Clustering can help to distinguish between different items. E.g. Fresh vegetables are more similar to each other than frozen items.
- Clustering is useful in **market segmentation** by partitioning the target market data into groups such as customer who share the same interests or those with common needs. Identifying clusters of similar items can help develop a marketing strategy that addresses the needs of specific clusters.
- Data Clustering can help **to identify, learn, or predict the nature of new data items**— especially how new data can be linked with making predictions. For e.g., in pattern recognition analyzing patterns in the data (such as buying patterns in particular regions or age groups) can help to develop predictive analytics to predict the nature of future data items that can fit well with established patterns.
- Clustering can help in **dividing the e-mail dataset into spam and non-spam messages**.
- Data Clustering is also helpful in **image segmentation** for analyzing the image more easily.
- Data Clustering can help **in information retrieval from a collection of data**, mainly documents (using tf-idf concepts).
- On the other hand, finding important **association rules in a dataset** of customer transactions helps a company to maximize revenue by deciding which products should be on sale, how to position products in the store's aisles, and how and when to offer promotional pricing.

Types of Cluster Analysis Methods:

- **Partitional Methods:** Partitional methods obtain a single level partition of objects. These methods are usually based on a greedy heuristics that are used to obtain a local optimum solution. Given n objects, these methods make $k \leq n$ clusters.
 - **K-means:** Each of the K-clusters is represented by the mean of the objects inside each cluster.
 - **Density-Based:** It is based on the assumption that clusters have high density collection of data of arbitrary shape that are separated by a large space of low density data (which is assumed to be the noise).
 - **Expectation-Maximization:** The EM method assigns objects to different clusters with certain probabilities in an attempt to maximize the expectation (or likelihood) of assignment.
- **Hierarchical Methods:** Hierarchical methods obtain a nested partition of the objects resulting in a tree of clusters.
 - **Agglomerative:** Start with each object in an individual cluster and then try to merge similar clusters into larger and larger clusters.
 - **Divisive:** Start with one cluster and then split into smaller and smaller clusters.
- **Grid Based method:** The object space rather than the data is divided into a grid. Grid partitioning is based on characteristics of the data and such methods can deal with non-parametric data more easily.
- **Model based method:** A model is assumed based on a probability distribution. Essentially the algorithm tries to build clusters with a high level of similarity with them and a low level of similarity between them. It tries to minimise the squared-error function.

Requirements of Clustering in Data Mining

The following points throw light on why clustering is required in data mining -

Scalability – We need highly scalable clustering algorithms to deal with large databases.

Ability to deal with different kinds of attributes – Algorithms should be capable to be applied on any kind of data such as interval-based (numerical) data, categorical, and binary data.

Discovery of clusters with attribute shape – The clustering algorithm should be capable of detecting clusters of arbitrary shape. They should not be bounded to only distance measures that tend to find spherical cluster of small sizes.

High dimensionality – The clustering algorithm should not only be able to handle low-dimensional data but also the high dimensional space.

Ability to deal with noisy data – Databases contain noisy, missing or erroneous data. Some algorithms are sensitive to such data and may lead to poor quality clusters.

Interpretability – The clustering results should be interpretable, comprehensible, and usable.

Some important concepts:

- **Measuring Distance**

- Between records: Distance between each record in a cluster.

- Between clusters: Distance between each cluster.

- **Distance Between Two Records:** Euclidian distance is most popular.

$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

- **Normalizing:**

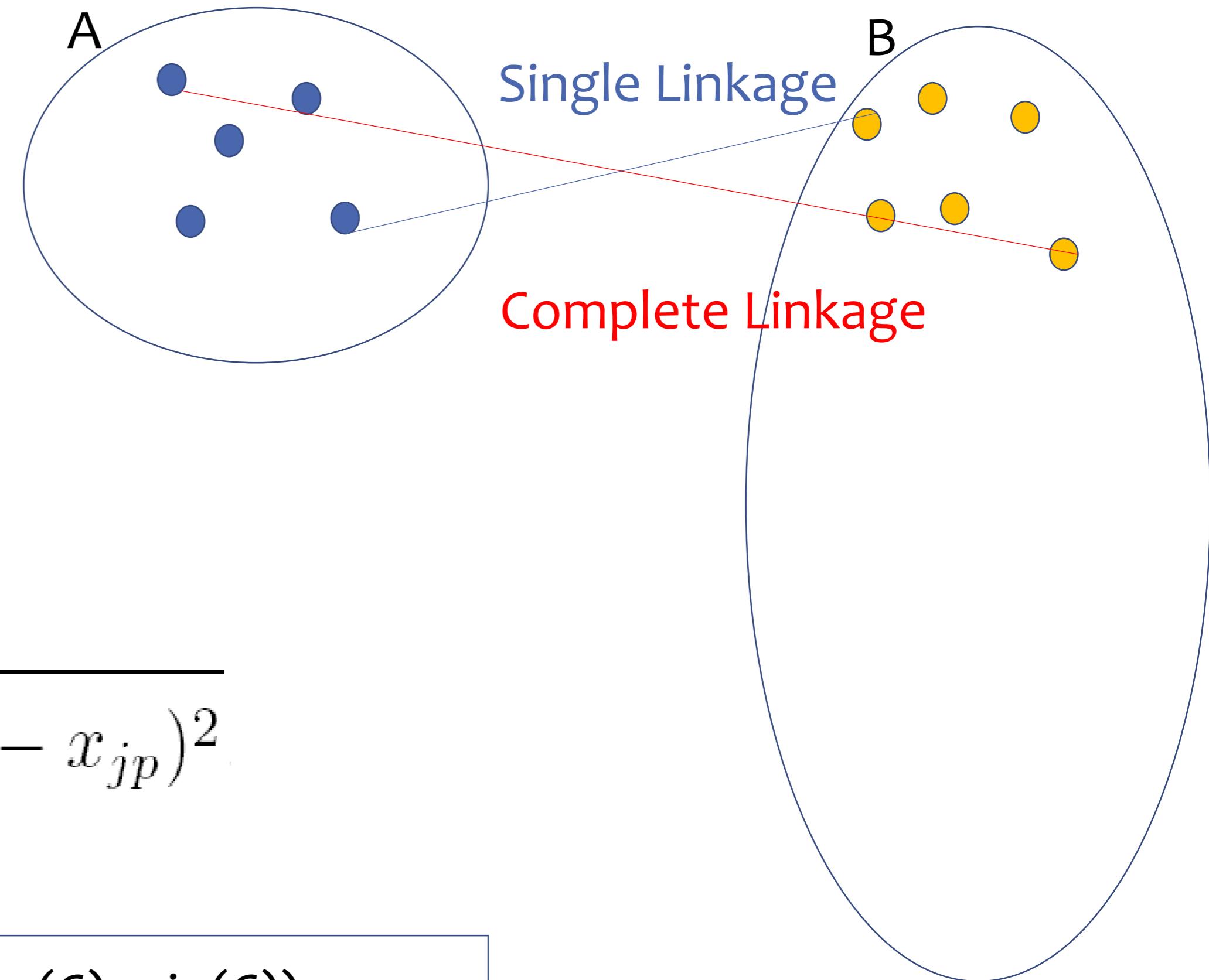
$$(x - \text{min}(C)) / (\text{max}(C) - \text{min}(C))$$

- **Problem:** Raw distance measures are highly influenced by scale of measurements

- **Solution:** Normalize (standardize) the data first.

- Subtract mean, divide by std. deviation. (Also called z-scores).

- Example: For 22 utilities, Avg. sales = 8,914; Std. dev. = 3,550. Hence, Normalized score: $(9,077 - 8,914) / 3,550 = 0.046$



Measuring Distance Between Clusters:

- **Single Linkage**
 - Minimum Distance (Cluster A to Cluster B)
 - Distance between two clusters is the distance between the pair of records A_i and B_j that are closest.
- **Complete Linkage**
 - Maximum Distance (Cluster A to Cluster B)
 - Distance between two clusters is the distance between the pair of records A_i and B_j that are farthest from each other
- **Average Linkage**
 - Distance between two clusters is the average of all possible pair-wise distances
- **Centroid**
 - Distance between two clusters is the distance between the two cluster centroids.
 - Centroid is the vector of variable averages for all records in a cluster

K-Means Clustering:

- The K-means method may be described as follows:
 1. Select the number of clusters. Let this be k .
 2. Pick k seeds as centroids of the k clusters. The seeds may be picked randomly unless the user has some insight about the data.
 3. Compute the Euclidean distance of each object in the dataset from each of the centroids.
 4. Allocate each object to the cluster it is nearest to based on the distances computed in the previous step.
 5. Compute the centroids of the clusters by computing the means of the attribute values of the objects in each cluster.
 6. Check if the stopping criteria has been met. If yes, stop. If not, go to step 3.
- The ***k-means method*** uses the Euclidean distance method, which appears to work well with compact clusters.
- If the Manhattan distance is used the method is called ***k-median method***. This method may be less sensitive to outliers.
- K-means Algorithm: Choosing k and Initial Partitioning
 - Choose k based on the how results will be used. E.g., “How many market segments do we want?”
 - Also experiment with slightly different k ’s.
 - Initial partition into clusters can be random, or based on domain knowledge. If random partition, repeat the process with different random partitions
- For clustering to be effective all attributes should be converted to a similar scale unless you want to give more weight to some attributes that are relatively large in scale.

1 6 3 7 4 9

Step 1: K= 2

Step 2: C1: 1 6 3 C2: 7 4 9
Mean 3.333 6.667

Step 3: C1: 1 6 3
M1: 2.3333 2.667 0.3333
M2: 5.667 0.667 3.667

C2: 7 4 9
M2: 1.667 2.667 2.333
M1: 4.337 0.7 5.7

Step 4: 1 → C1; 6 → C2; 3 → C1; 7 → C2 ; 4 → C1 ; 9 → C2

Step 5: C1: 1 3 4 C2: 6 7 9
Mean 2.667 7.333

Step 6: C1: 1 3 4 C2: 6 7 9
M1: 1.667 0.333 1.33 M1: 3.333 4.333 6.333
M2: 6.333 4.333 3.333 M2: 1.333 0.333 1.67

Step 7: 1 → C1; 3 → C1; 4 → C1

6 → C2; 7 → C2; 9 → C2

Q1.

Use k-mean clustering to divide the following set of numbers into two clusters.

1 2 3 4 5 6 7 8 9 10

Q2.

Use k-median method to form clusters of students based on the data given below:

Student	Age	Mark1	Mark2	Mark3
S1	18	73	75	57
S2	18	79	85	75
S3	23	70	70	52
S4	20	55	55	55
S5	22	85	86	87
S6	19	91	90	89
S7	20	70	65	60
S8	21	53	56	59
S9	19	82	82	60
S10	47	75	76	77

K=3

Steps 1 and 2: Let the three seeds be the first three students as shown.

Student	Age	Mark1	Mark2	Mark3
S1	18	73	75	57
S2	18	79	85	75
S3	23	70	70	52

Step 3 and 4: Compute the distances using the four attributes and using the sum of absolute differences (k-Median method)

	Age	Mark1	Mark2	Mark3	Distance from Clusters			Allocation to the nearest Cluster
C1	18	73	75	57	From C1 From C2 From C3			C1 C2 C3
C2	18	79	85	75				
C3	23	70	70	52				
S1	18	73	75	57	0	34	18	C1
S2	18	79	85	75	34	0	52	C2
S3	23	70	70	52	18	52	0	C3
S4	20	55	55	55	42	76	36	C3
S5	22	85	86	87	57	23	67	C2
S6	19	91	90	89	66	32	82	C2
S7	20	70	65	60	18	46	16	C3
S8	21	53	56	59	44	74	40	C3
S9	19	82	82	60	20	22	36	C1
S10	47	75	76	77	52	44	60	C2

Step 5: The new cluster means of clusters are given in the table below:

Manhattan Distance-> $|x_1-x_2| + |y_1-y_2| + |z_1-z_2| + |w_1-w_2|$

Student	Age	Mark1	Mark2	Mark3
C1	18.5	77.5	78.5	58.5
C2	26.5	82.5	84.3	82.0
C3	21	61.5	61.5	65.5

Step 3 and 4: Using this new cluster means compute the distances of each object to each of the means and allocate to nearest cluster.

	Age	Mark1	Mark2	Mark3	Distance from Clusters			Allocation to the nearest Cluster
C1	18.5	77.5	78.5	58.5	From C1	From C2	From C3	
C2	26.5	82.5	84.3	82.0				
C3	21	61.5	61.5	65.5				
S1	18	73	75	57	10	52.3	28	C1
S2	18	79	85	75	25	19.8	62	C2
S3	23	70	70	52	27	60.3	23	C3
S4	20	55	55	55	51	90.3	16	C3
S5	22	85	86	87	47	13.8	79	C2
S6	19	91	90	89	56	28.8	92	C2
S7	20	70	65	60	24	60.3	16	C3
S8	21	53	56	59	50	86.3	17	C3
S9	19	82	82	60	10	32.3	46	C1
S10	47	75	76	77	52	41.3	74	C2

Step 6: The clusters have not changed and hence we can stop.

Therefore, **Cluster 1:** S1, S9. **Cluster 2:** S2, S5, S6, S10. **Cluster 3:** S3, S4, S7, S8.

Cluster	C1	C2	C3
C1	5.9	26.5	23.3
C2	29.5	14.3	42.6
C3	23.9	41.0	10.7

Within cluster and between cluster distances

The content of the slides are prepared from different textbooks.

References:

- Data Mining and Predictive Analytics, By Daniel T. Larose. Copyright 2015 John Wiley & Sons, Inc.
- Predictive Analytics for Dummies, By Anasse Bari, Mohamed Chaouchi, & Tommy Jung, Copyright 2016, John Wiley & Sons, Inc.
- Introduction to Data Mining with Case Studies, By G.K. Gupta. Copyright 2014 by PHI Learning Private Limited.



—
Thank you..

Machine Learning with Python

Session 12: Cluster Analysis

Arghya Ray

Rational for Measuring Cluster Goodness

- What are the optimal number of cluster to identify?
- How do I measure weather one set of clusters is preferable to another?
- The **silhouette** method and the **pseudo-F statistic** will help us address these questions by measuring cluster goodness.

Concepts Measures Should Address

- Cluster separation represents how distant the clusters are from each other
- Cluster cohesion refers to how tightly related the records within the individual clusters are
- Good measures should incorporate both as do the silhouette and pseudo-F statistic
- However, the sum of squares error (SSE) only accounts for cluster cohesion and is monotonically decreasing with increasing numbers of clusters.

Measuring Cluster Goodness:

The Silhouette Method

For each data value i the silhouette is used to gauge how good the cluster assignment is for that point:

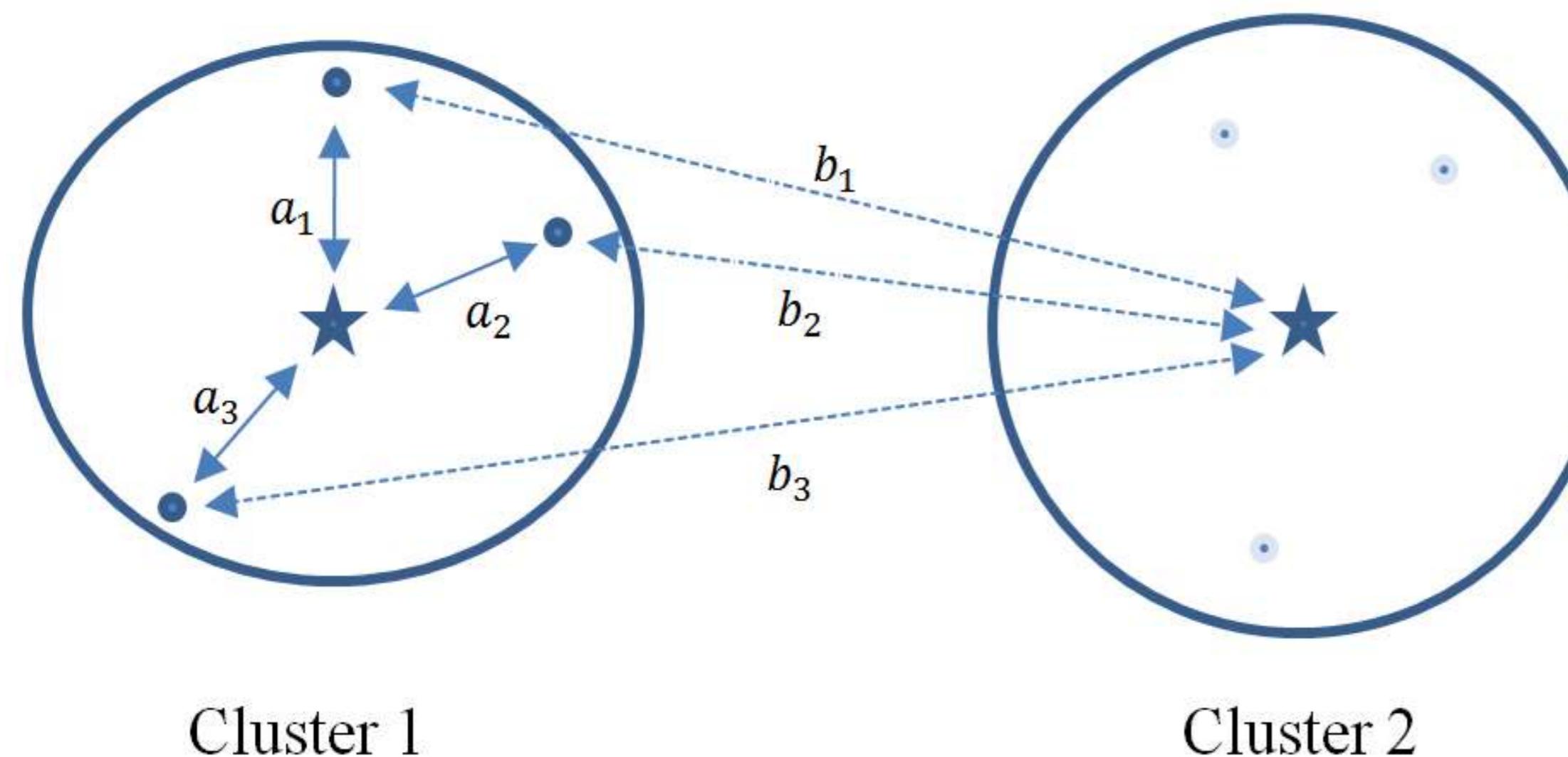
$$\text{Silhouette}_i = s_i = \frac{b_i - a_i}{\max(b_i, a_i)}$$

where a_i is the distance between the data value and its cluster center and represents cohesion

and b_i is the distance between the data value and the next closest cluster center and represents separation

Silhouette Accounts for Separation & Cohesion

Each data value in Cluster 1 has its values of a_i and b_i represented by solid and dotted lines, respectively



$b_i > a_i$ for each data value, thus each data value's silhouette is positive, indicating the data are not misclassified

Measuring Cluster Goodness:

The Silhouette Method contd...

- A positive value indicates that the assignment is good, with higher values better than lower values.
- A value close to zero is considered to be weak since the observation could have been assigned to the next cluster with little negative consequence.
- A negative value is considered to be misclassified since assignment to the next closest cluster would have been better.

The Average Silhouette Value

The average silhouette value over all records yields a measure of how well the cluster solution fits. A thumbnail interpretation, meant as a guide only:

- 0.5 or better provides good evidence of the reality of the clusters in the data
- 0.25 – 0.5 provides some evidence of the reality of the clusters in the data.
- Less than 0.25 provides scant evidence of cluster reality

Silhouette Example (cont.)

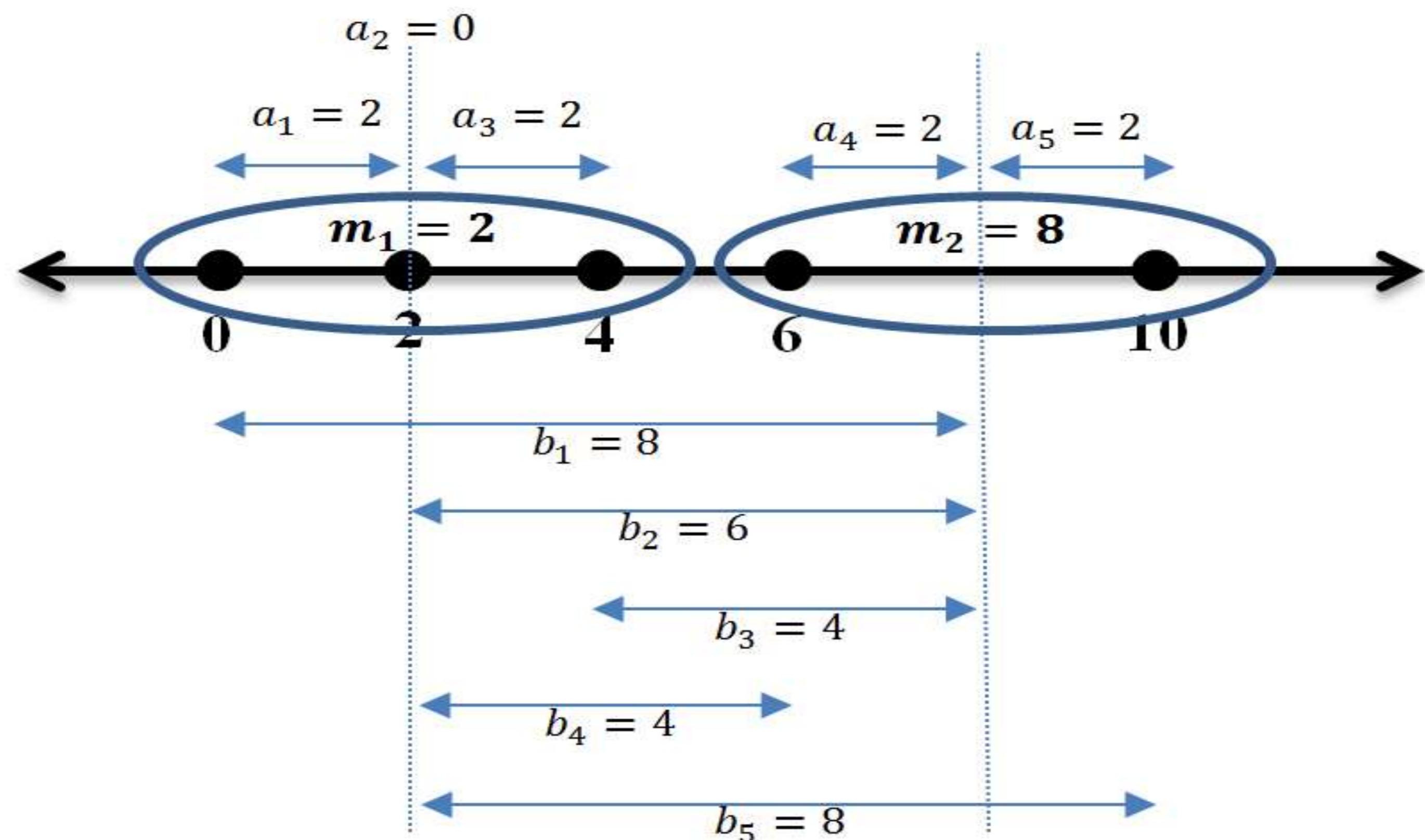
- Apply k-means clustering to the following data set:

$$x_1 = 0 \quad x_2 = 2 \quad x_3 = 4 \quad x_4 = 6 \quad x_5 = 10$$

- The first three data values are assigned to Cluster 1 and the last two to Cluster 2
- Center for Cluster 1 is $m_1 = 2$ and for Cluster 2 is $m_2 = 8$
- Values for a_i are distance between x_i and its cluster center; values for b_i are distance between x_i and the other cluster center

Silhouette Example (cont.)

Distances between the data values and cluster centers:



Silhouette Example (cont.)

Calculations for individual data values:

x_i	a_i	b_i	$\text{Max}(a_i, b_i)$
-------	-------	-------	------------------------

$$\text{Silhouette}_i = s_i = \frac{b_i - a_i}{\text{Max}(b_i, a_i)}$$

0	2	8	8	$\frac{8-2}{8} = 0.75$
2	0	6	6	$\frac{6-0}{6} = 1.00$
4	2	4	4	$\frac{4-2}{4} = 0.50$
6	2	4	4	$\frac{4-2}{4} = 0.50$
10	2	8	8	$\frac{8-2}{8} = 0.75$

Mean Silhouette = 0.7

The pseudo-*F* Statistic

Let:

k be number of clusters

$\sum n_i = N$ be total sample size

x_{ij} refer to the j^{th} data value in the i^{th} cluster

m_i refer to cluster center (centroid) of the i^{th} cluster

M represent the grand mean of all the data

and $Distance(a, b) = \sqrt{\sum(a_i - b_i)^2}$

The pseudo-*F* Statistic (cont.)

Then the *sum of squares between* the clusters is:

$$SSB = \sum_{i=1}^k n_i \cdot Distance^2(m_i, M)$$

And the *sum of squares error*, or the *sum of squares within* the clusters is:

$$SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} Distance^2(x_{ij}, m_i)$$

And the *pseudo-F* statistic is:

$$F = \frac{MSB}{MSE} = \frac{SSB/k - 1}{SSE/N - k}$$

The pseudo- F Statistic (cont.)

- The hypotheses being tested are:

H_0 : There are no clusters in the data.

H_a : There are k clusters in the data.

- Reject H_0 for sufficiently small p-value where:

$$p\text{-value} = P(F_{k-1,n-k}) > \text{Pseudo F value}$$

- The pseudo- F statistic rejects the null hypothesis too easily.

The pseudo-*F* Statistic (cont.)

The pseudo-*F* statistic should not be used to determine the presence of clusters but can be used to select the optimal number of clusters as follows:

1. Use a clustering algorithm to develop a clustering solution for a variety of values of k .
2. Calculate the pseudo-*F* statistic and p-value for each candidate, and select the candidate with the smallest p-value as the best clustering solution.

Pseudo- F Statistic Example

- Apply k -means clustering to the following data set:

$$x_1 = 0 \quad x_2 = 2 \quad x_3 = 4 \quad x_4 = 6 \quad x_5 = 10$$

- The first three data values are assigned to Cluster 1 and the last two to Cluster 2
- Center for Cluster 1 is $m_1 = 2$ and for Cluster 2 is $m_2 = 8$
- $n_1 = 3$ and $n_2 = 2$ data values, and $N = 5$, the grand mean is $M = 4.4$. And, because we are in one dimension, $Distance(m_i, M) = |m_i - M|$

Pseudo- F Statistic Example (cont.)

Then

$$\begin{aligned}SSB &= \sum_{i=1}^k n_i \cdot Distance^2(m_i, M) \\&= 3 \cdot (2 - 4.4)^2 + 2 \cdot (8 - 4.4)^2 = 43.2\end{aligned}$$

And

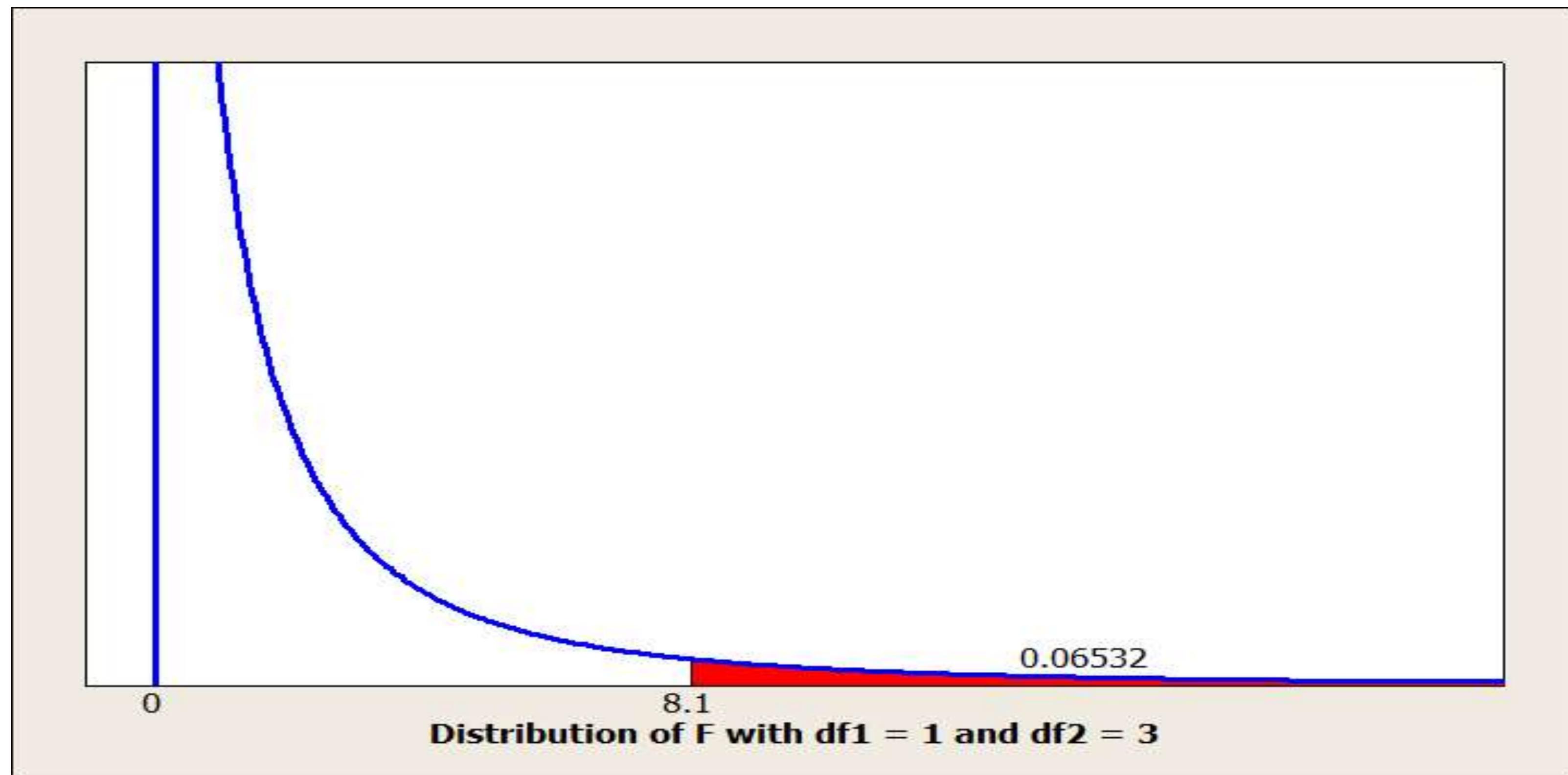
$$\begin{aligned}SSE &= \sum_{i=1}^k \sum_{j=1}^{n_i} Distance^2(x_{ij}, m_i) \\&= (0 - 2)^2 + (2 - 2)^2 + (4 - 2)^2 + (6 - 8)^2 + (10 - 8)^2 = 16\end{aligned}$$

And

$$F = \frac{MSB}{MSE} = \frac{SSB/k - 1}{SSE/N - k} = \frac{43.2/1}{16/3} = \frac{43.2}{5.33} = 8.1$$

Pseudo-*F* Statistic Example (cont.)

Distribution of the F statistic shows that p-value of 0.06532 does not indicate strong evidence of clusters:



Cluster Validation

- As with any other data mining modeling technique, cluster analysis should be subject to cross-validation to ensure the clusters are real
- A simple graphical and statistical approach with the goal of confirming the clusters found in the test data match those in the training data is:
 1. Apply cluster analysis to training data
 2. Apply cluster analysis to test data
 3. Use graphics and statistics to confirm the clusters in the training data match those in the test data

The content of the slides are prepared from different textbooks.

References:

- Data Mining and Predictive Analytics, By Daniel T. Larose. Copyright 2015 John Wiley & Sons, Inc.
- Predictive Analytics for Dummies, By Anasse Bari, Mohamed Chaouchi, & Tommy Jung, Copyright 2016, John Wiley & Sons, Inc.
- Introduction to Data Mining with Case Studies, By G.K. Gupta. Copyright 2014 by PHI Learning Private Limited.



—
Thank you..

Machine Learning with Python

Session 11: Measures of Proximity, Components of
Machine Learning

Arghya Ray

Introduction:

The term proximity between two objects is a function of the proximity between the corresponding attributes of the two objects. Proximity measures refer to the **Measures of Similarity and Dissimilarity**.

Similarity and Dissimilarity are important because they are used by a number of data mining techniques, such as clustering, nearest neighbour classification, and anomaly detection.

What is Similarity?

- It is a numerical measure of the degree to which the two objects are alike.
 - Higher for pair of objects that are more alike.
 - Usually non-negative and between 0 & 1.
- 0 ~ No Similarity, 1 ~ Complete Similarity

What is Dissimilarity?

- It is a numerical measure of the degree to which the two objects are different.
- Lower for pair of objects that are more similar.
- Range 0 to infinity.

Transformation Function

It is a function used to convert similarity to dissimilarity and vice versa, or to transform a proximity measure to fall into a particular range. For instance:

$$s' = (s - \min(s)) / (\max(s) - \min(s))$$

range

where,

s' = new transformed proximity measure value,

s = current proximity measure value,

$\min(s)$ = minimum of proximity measure values,

$\max(s)$ = maximum of proximity measure values

This transformation function is just one example from all the available options out there.

Similarity and Dissimilarity between Simple Attributes

The proximity of objects with a number of attributes is usually defined by combining the proximities of individual attributes, so, we first discuss proximity between objects having a single attribute.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{if } x \neq y \end{cases}$	$s = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$
Ordinal	$d = x - y /(n - 1)$ (values mapped to integers 0 to $n-1$, where n is the number of values)	$s = 1 - d$
Interval or Ratio	$d = x - y $	$s = -d, s = \frac{1}{1+d}, s = e^{-d},$ $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

To understand it better, let us go through some examples.

Consider objects described by one **nominal** attribute. How to compare similarity of two objects like this? Nominal attributes only tell us about the distinctness of objects. Hence, in this case similarity is defined as 1 if attribute values match, and 0 otherwise and oppositely defined would be dissimilarity.

For objects with a single **ordinal** attribute, the situation is more complicated because information about order needs to be taken into account. Consider an attribute that measures the quality of a product, on the scale {poor, fair, OK, good, wonderful}. We have 3 products P₁, P₂, & P₃ with quality as wonderful, good, & OK respectively. In order to compare **ordinal** quantities, they are mapped to successive integers. In this case, if the scale is mapped to {0, 1, 2, 3, 4} respectively. Then, dissimilarity(P₁, P₂) = 4 - 3 = 1.

For **interval or ratio** attributes, the natural measure of dissimilarity between two objects is the absolute difference of their values. For example, we might compare our current weight and our weight a year ago by saying “I am ten pounds heavier.”

Dissimilarities between Data Objects

Euclidean Distance

The Euclidean distance, d , between two points, x and y , in one, two, three, or higher-dimensional space, is given by the following formula:

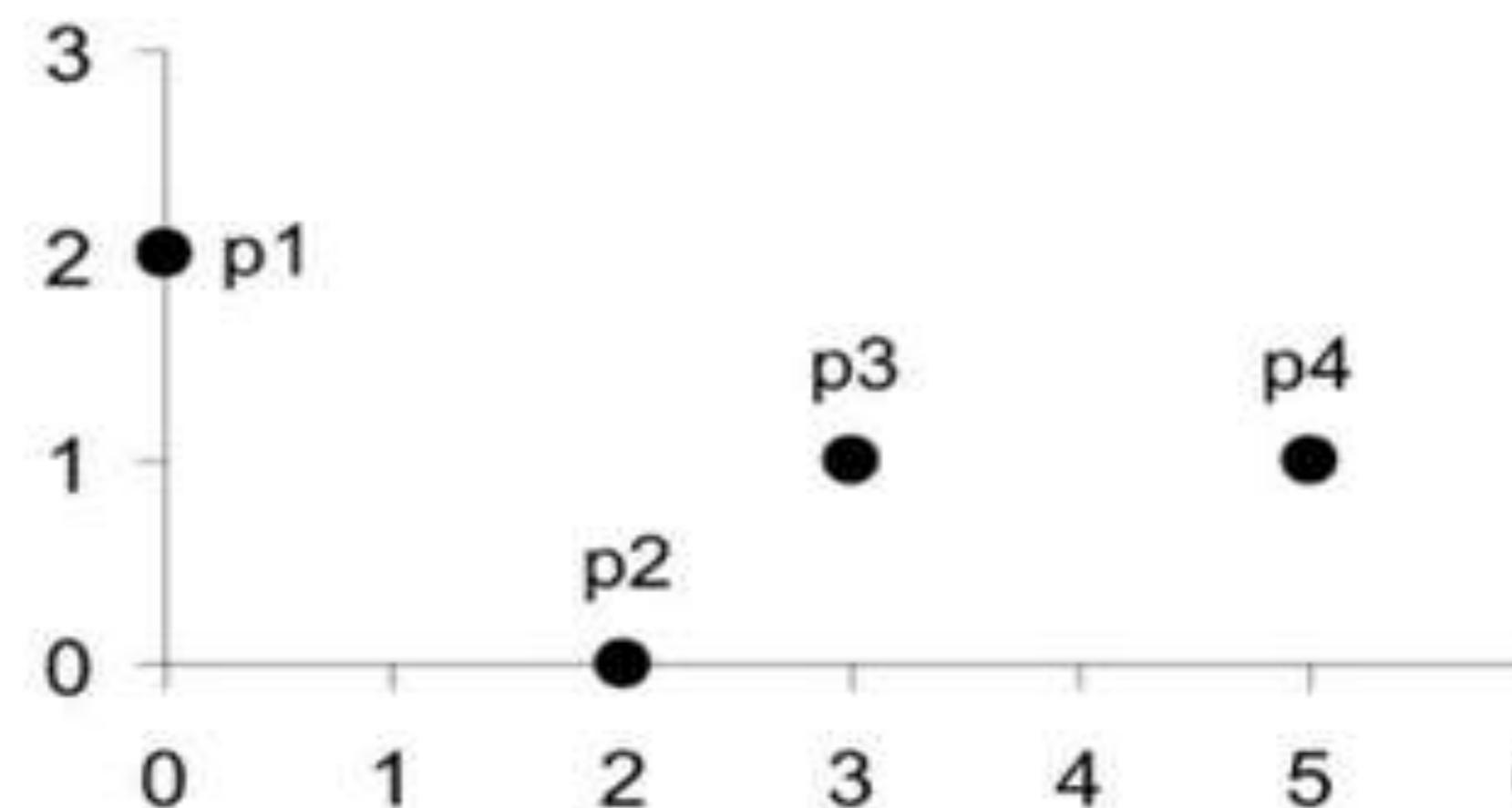
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2},$$

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

where n is the number of dimensions, and $x(k)$ and $y(k)$ are respectively, the k th attributes (components) of x and y .

Dissimilarities between Data Objects

Example:



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

Dissimilarities between Data Objects

Minkowski Distance

It is the generalisation of Euclidean distance. It is given by the following formula:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

where r is a parameter. The following are the three most common examples of Minkowski distances.

Dissimilarities between Data Objects

→ $r = 1$. City block (Manhattan, taxicab, $L1$ norm) distance. A common example is the **Hamming distance**, which is the number of bits that are different between two objects that have only binary attributes, i.e., between two binary vectors.

$$\begin{array}{r} 1\ 1\ 1\ 0\ 0\ 0\ 1 \\ 1\ 1\ 0\ 1\ 1\ 0 \\ \hline 0\ 0\ 1\ 1\ 1\ 1 \end{array} = 5$$

→ $r = 2$. Euclidean distance($L2$ norm).

$$\begin{array}{r} \text{Barun} \\ \text{Beran} \\ \hline 0+1+0+1+0 \end{array} = 2$$

→ $r = \text{infinity}$. Supremum ($L(\max)$, or $L(\infty)$ norm) distance. This is the maximum difference between any attribute of the objects. This is defined by the following formula:

$$d(\mathbf{x}, \mathbf{y}) = \lim_{r \rightarrow \infty} \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}.$$

Dissimilarities between Data Objects

Example:

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

L ∞	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

Distance Matrix

Distances, such as the Euclidean distance, have some well-known properties. If $d(x, y)$ is the distance between two points, x and y , then the following properties hold.

Positivity

- a) $d(x, y) > 0$ for all x and y ,
- b) $d(x, y) = 0$ only if $x = y$

2. Symmetry

$$d(x, y) = d(y, x) \text{ for all } x \text{ and } y$$

3. Triangle Inequality

$$d(x, z) \leq d(x, y) + d(y, z) \text{ for all points } x, y \text{ and } z$$

The measures that satisfy all three properties are called **metrics**.

Similarities between Data Objects

For similarities, the triangle inequality typically does not hold, but symmetry and positivity typically do. To be explicit, if $s(x, y)$ is the similarity between points x and y , then the typical properties of similarities are the following:

$$s(x, y) = 1 \text{ only if } x = y. (0 \leq s \leq 1)$$

$$s(x, y) = s(y, x) \text{ for all } x \text{ and } y. (\text{Symmetry})$$

There is no general analog of the triangle inequality for similarity measure.

Similarity Measures for Binary Data are called **similarity coefficients** and typically have values between 0 and 1. The comparison between two binary objects is done using the following four quantities:

f_{00} = the number of attributes where x is 0 and y is 0

f_{01} = the number of attributes where x is 0 and y is 1

f_{10} = the number of attributes where x is 1 and y is 0

f_{11} = the number of attributes where x is 1 and y is 1

Simple Matching Coefficient

It is defined as follows:

$$SMC = \frac{\text{number of matching attribute values}}{\text{number of attributes}} = \frac{f_{11} + f_{00}}{f_{01} + f_{10} + f_{11} + f_{00}}.$$

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

Jaccard Coefficient

It is defined as follows:

$$J = \frac{\text{number of matching presences}}{\text{number of attributes not involved in 00 matches}} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}.$$

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

An example comparing these two similarity methods:

$$\mathbf{x} = (1, 0, 0, 0, 0, 0, 0, 0, 0)$$

$$\mathbf{y} = (0, 0, 0, 0, 0, 1, 0, 0, 1)$$

$f_{01} = 2$ the number of attributes where \mathbf{x} was 0 and \mathbf{y} was 1

$f_{10} = 1$ the number of attributes where \mathbf{x} was 1 and \mathbf{y} was 0

$f_{00} = 7$ the number of attributes where \mathbf{x} was 0 and \mathbf{y} was 0

$f_{11} = 0$ the number of attributes where \mathbf{x} was 1 and \mathbf{y} was 1

$$SMC = \frac{f_{11} + f_{00}}{f_{01} + f_{10} + f_{11} + f_{00}} = \frac{0+7}{2+1+0+7} = 0.7$$

$$J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{0}{2+1+0} = 0$$

Q. Using the given snapshot of a movie recommender system, find the Jaccard's coefficient and Matching coefficient between (a) User 1 and User 3 (b) User 2 and User 3.

	User 1	User 2	User 3
There will be blood	1	0	0
Gravity	0	1	1
X-Men	0	0	1
Inception	0	1	1
Jurassic Park	1	0	0
Avengers: End Game	1	1	1

Between User 1 & User 3:

Jaccard's Coefficient= $1/6$

Matching coefficient = $1/6$

Between User 2 and User 3:

Jaccard's Coefficient= $3/4$

Matching coefficient = $5/6$

Cosine Similarity

Documents are often represented as vectors, where each attribute represents the frequency with which a particular term(word) occurs in the document. The **cosine similarity**, is one of the most common measure of document similarity. If x and y are two document vectors, then

$$\cos(x, y) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|},$$

	cos	sin	tan
D1 →	0	1	2
D2 →	1	2	3
D3 →	2	1	1

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

where $.$ indicates *dot product* and $\|x\|$ defines the length of vector x .

An example of **cosine similarity** measure is as follows:

$$\mathbf{x} = (3, 2, 0, 5, 0, 0, 2, 0, 0)$$

$$\mathbf{y} = (1, 0, 0, 0, 0, 0, 1, 0, 2)$$

$$\mathbf{x} \cdot \mathbf{y} = 3 * 1 + 2 * 0 + 0 * 0 + 5 * 0 + 0 * 0 + 0 * 0 + 0 * 0 + 2 * 1 + 0 * 0 + 0 * 2 = 5$$

$$\|\mathbf{x}\| = \sqrt{3 * 3 + 2 * 2 + 0 * 0 + 5 * 5 + 0 * 0 + 0 * 0 + 0 * 0 + 2 * 2 + 0 * 0 + 0 * 0} = 6.48$$

$$\|\mathbf{y}\| = \sqrt{1 * 1 + 0 * 0 + 0 * 0 + 0 * 0 + 0 * 0 + 0 * 0 + 0 * 0 + 1 * 1 + 0 * 0 + 2 * 2} = 2.24$$

$$\cos(\mathbf{x}, \mathbf{y}) = 0.31$$

Correlation

It is a measure of the linear relationship between the attributes of the objects having either binary or continuous variables. **Correlation** between two objects x and y is defined as follows:

$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\text{covariance}(\mathbf{x}, \mathbf{y})}{\text{standard_deviation}(\mathbf{x}) * \text{standard_deviation}(\mathbf{y})} = \frac{s_{xy}}{s_x s_y},$$

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

where the notations used are defined in standard as:

$$\text{covariance}(\mathbf{x}, \mathbf{y}) = s_{xy} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$\text{standard_deviation}(\mathbf{x}) = s_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2}$$

$$\text{standard_deviation}(\mathbf{y}) = s_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \text{ is the mean of } \mathbf{x}$$

$$\bar{y} = \frac{1}{n} \sum_{k=1}^n y_k \text{ is the mean of } \mathbf{y}$$

Introduction to Data Mining — Pang-Ning Tan, Michael Steinbach, Vipin Kumar

Correlation Example: <https://www.statisticshowto.com/probability-and-statistics/correlation-coefficient-formula/>

The Chi-square goodness of fit test:

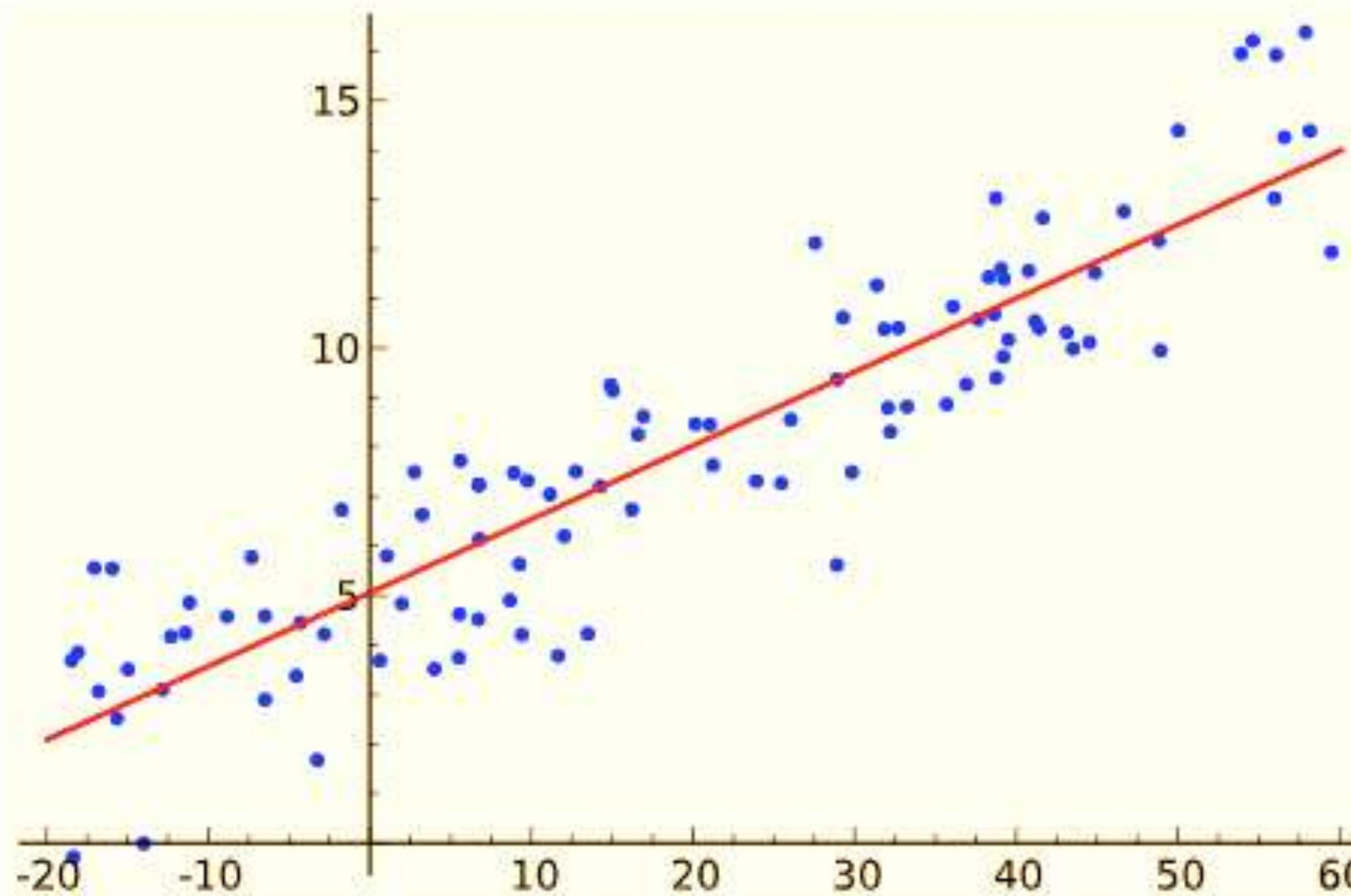
What is the goodness of fit?

A goodness-of-fit is a statistical technique. It is applied to measure “**how well the actual(observed) data points fit into a Machine Learning model**”. It summarizes the divergence between actual observed data points and expected data points in context to a statistical or Machine Learning model.

Assessment of divergence between the observed data points and model-predicted data points is critical to understand, a decision made on poorly fitting models might be badly misleading. A seasoned practitioner must examine the fitment of actual and model-predicted data points.

Why do we test Goodness of fit?

Goodness-of-fit tests are statistical tests to determine whether a set of actual observed values match those predicted by the model. Goodness-of-fit tests are frequently applied in business decision making. For example, if we check linear regression function. The goodness-of-fit test here will compare the actual observed values to the predicted values.



The Chi-square test for a goodness-of-fit test is

$$\chi^2_c = \sum \frac{(O_i - E_i)^2}{E_i}$$

Where:

O_i = an observed count for bin*i*

E_i = an expected count for bin*i*, asserted by the null hypothesis.

What are the most common goodness of fit tests?

Broadly, the goodness of fit test categorization can be done based on the distribution of the predict and variable of the dataset.

- The chi-square
- Kolmogorov-Smirnov
- Anderson-Darling

The expected frequency is calculated by:

$$E_i = \left(F(Y_u) - F(Y_l) \right) N$$

where:

F = the cumulative distribution function for the probability distribution being tested.

Y_u = the upper limit for class*i*,

Y_l = the lower limit for class*i*, and

N = the sample size

The Chi-Square Goodness of Fit Test

Chi-square goodness of fit test is conducted when the predicted variable in the dataset is categorical. It is applied to determine whether sample data are consistent with a hypothesized distribution.

Chi-Square test can be applied when the distribution has the following characteristics:

- The sampling method is random.
- Predicted variables are categorical.
- The expected value of the number of sample observations at each level of the variable is at least 5. It requires a sufficient sample size for the chi-square approximation to be valid.

Merits of the Chi-square Test

- A distribution-free test. It can be used in any type of population distribution.
- It is widely applicable not only in social sciences but in business research as well.
- It can be easy to calculate and to conclude.
- The Chi-Square test provides an additive property. This allows the researcher to add the result of independence to related samples.
- This test is based on the observed frequency and not on parameters like mean, and standard deviation.

Until now we have defined and understood both similarity and dissimilarity measures amongst data objects. Now, let's discuss the issues faced in proximity calculations.

Issues in Proximity Calculation

- how to handle the case in which attributes have different scales and/or are correlated,
- how to calculate proximity between objects that are composed of different types of attributes, e.g., quantitative and qualitative, and
- how to handle proximity calculation when attributes have different weights i.e., when not all attributes contribute equally to the proximity of objects.

Selecting the Right Proximity Measure

The following are a few general observations that may be helpful. First, ***the type of proximity measure should fit the type of data.*** For many types of dense, continuous data, metric distance measures such as Euclidean distance are often used.

Proximity between continuous attributes is most often expressed in terms of differences, and distance measures provide a well-defined way of combining these differences into an overall proximity measure.

For sparse data, which often consists of asymmetric attributes, we typically employ similarity measures that ignore 0-0 matches. Conceptually, this reflects the fact that, for a pair of complex objects, similarity depends on the number of characteristics they both share, rather than the number of characteristics they both lack. For such type of data, Cosine Similarity or Jaccard Coefficient can be used.

Main Components of Machine Learning Algorithm:

1) Feature Extraction + Domain knowledge

First and foremost we really need to understand what type of data we are dealing with and what eventually we want to get out of it. Essentially we need to understand how and what features need to be extracted from the data. For instance assume we want to build a software that distinguishes between male and female names. All the names in text can be thought of as our raw data while our features could be number of vowels in the name, length, first & last character, etc of the name.

2) Feature Selection

In many scenarios we end up with a lot of features at our disposal. We might want to select a subset of those based on the resources and computation power we have. In this step we select a few of those influential features and separate them from the not-so-influential features. There are many ways to do this, information gain, gain ratio, correlation etc.

3) Choice of Algorithm

There are wide range of algorithms from which we can choose based on whether we are trying to do prediction, classification or clustering. We can also choose between linear and non-linear algorithms. Naive Bayes, Support Vector Machines, Decision Trees, k-Means Clustering are some common algorithms used.

4) Training

In this step we tune our algorithm based on the data we already have. This data is called training set as it is used to train our algorithm. This is the part where our machine or software learn and improve with experience.

5) Choice of Metrics/Evaluation Criteria

Here we decide our evaluation criteria for our algorithm. Essentially we come up with metrics to evaluate our results.

Commonly used measures of performance are precision, recall, f1-measure, robustness, specificity-sensitivity, error rate etc.

6) Testing

Lastly, we test how our machine learning algorithm performs on an unseen set of test cases. One way to do this, is to partition the data into training and testing set. The training set is used in step 4 while the test set is then used in this step.

Techniques such as cross-validation and leave-one-out can be used to deal with scenarios where we do not have enough data.

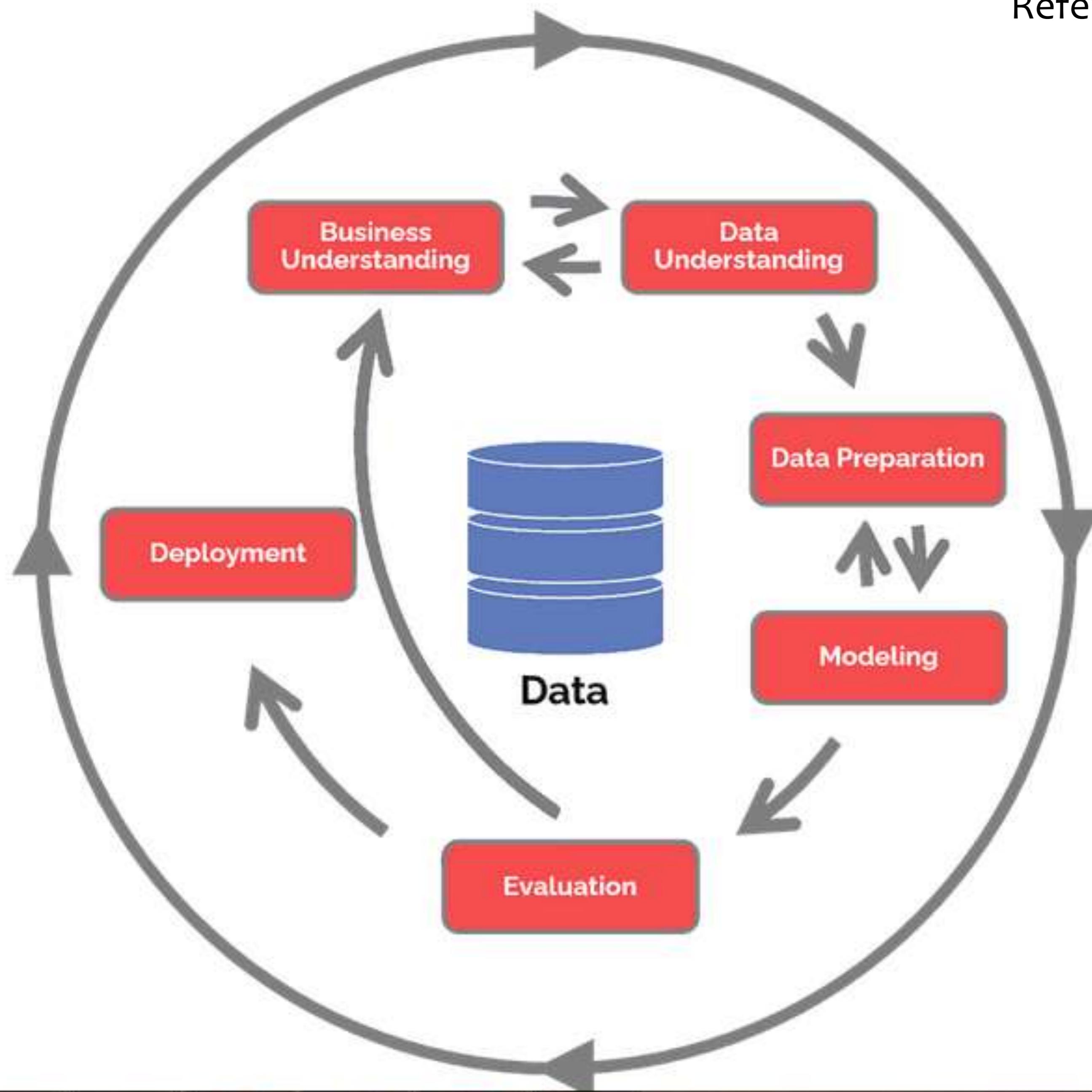
(Reference: <https://www.linkedin.com/pulse/20140822073217-180198720-6-components-of-a-machine-learning-algorithm>)

Main steps involved for End-to-End Machine Learning Project: (Chapter 2 of textbook)

1. Look at the big picture
2. Get the data
3. Discover and visualize the data to gain insights
4. Prepare the data for Machine Learning algorithms
5. Select a model and train it
6. Fine tune your model
7. Present your solution
8. Launch, monitor and maintain your system

CRISP DM (The CRoss Industry Standard Process for Data Mining)

Reference: <https://www.datascience-pm.com/crisp-dm-2/>



The content of the slides are prepared from different textbooks.

References:

Proximity Measures:

- <https://towardsdatascience.com/measures-of-proximity-in-data-mining-machine-learning-e9baaed1aafb>

Chi-Square Goodness of Fit readings:

- <https://www.mygreatlearning.com/blog/understanding-goodness-of-fit-test/>
- <https://machinelearningmastery.com/chi-squared-test-for-machine-learning/>
- <https://towardsdatascience.com/machine-learning-chi-square-test-in-evaluating-predictions-486404dd5bc>
- <https://towardsdatascience.com/chi-square-test-for-feature-selection-in-machine-learning-206b1fob8223>
- <https://www.analyticsvidhya.com/blog/2019/11/what-is-chi-square-test-how-it-works/> (stepwise calculation)
- <https://medium.com/wenyi-yan/a-simple-explanation-to-understand-chi-square-test-1814fa261499> (step wise calculation simple)

Correlation:

https://en.wikipedia.org/wiki/Pearson_correlation_coefficient

Extra Read:

<https://towardsdatascience.com/types-of-data-sets-in-data-science-data-mining-machine-learning-eb47c80af7a>



—
Thank you..

Machine Learning with Python

Measuring Performance of Classifiers

Arghya Ray

Why Evaluate?

- Multiple methods are available to classify or predict
- For each method, multiple choices are available for settings
- To choose best model, need to assess each model's performance



Reference: <https://www.kaggle.com/usenecoder/performance-metrics-for-classification-problems>

Accuracy Measures (Classification)

Misclassification error

- Error = classifying a record as belonging to one class when it belongs to another class.
- Error rate = percent of misclassified records out of the total records in the validation data

Naïve Rule

Naïve rule: classify all records as belonging to the most prevalent class

- Often used as benchmark: we hope to do better than that
- Exception: when goal is to identify high-value but rare outcomes, we may do well by doing worse than the naïve rule (see “lift” – later)

Separation of Records

“High separation of records” means that using predictor variables attains low error

“Low separation of records” means that using predictor variables does not improve much on naïve rule

Confusion Matrix

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Total Value = $(TP+FP+FN+TN)$

Accuracy = $(TP+TN)/\text{Total values}$

$1 - \text{Accuracy} = (FP+FN)/\text{Total Values}$
=Error rate

Confusion Matrix

Classification Confusion Matrix		
		Predicted Class
Actual Class	1	0
1	201	85
0	25	2689

201 1's correctly classified as “1”

85 1's incorrectly classified as “0”

25 0's incorrectly classified as “1”

2689 0's correctly classified as “0”

$$\begin{aligned} \text{Accuracy} &= (201+2689)/(201+85+25+2689) \\ &= 0.96 \end{aligned}$$

$$\text{Error rate} = 1-0.96=0.04$$

Error Rate

Classification Confusion Matrix		
		Predicted Class
Actual Class	1	0
1	201	85
0	25	2689

Overall error rate = $(25+85)/3000 = 3.67\%$

Accuracy = $1 - \text{err} = (201+2689)/3000 = 96.33\%$

If multiple classes, error rate is:

$(\text{sum of misclassified records})/(\text{total records})$

Cutoff for classification

Most DM algorithms classify via a 2-step process:

For each record,

1. Compute **probability of belonging to class “1”**
 2. Compare to cutoff value, and classify accordingly
-
- Default cutoff value is 0.50
 - If ≥ 0.50 , classify as “1”
 - If < 0.50 , classify as “0”
 - Can use different cutoff values
 - Typically, error rate is lowest for cutoff = 0.50

Cutoff Table

Actual Class	Prob. of "1"	Actual Class	Prob. of "1"
1	0.996	1	0.506
1	0.988	0	0.471
1	0.984	0	0.337
1	0.980	1	0.218
1	0.948	0	0.199
1	0.889	0	0.149
1	0.848	0	0.048
0	0.762	0	0.038
1	0.707	0	0.025
1	0.681	0	0.022
1	0.656	0	0.016
0	0.622	0	0.004

- If cutoff is 0.50: eleven records are actually in class “1”
- If cutoff is 0.80: seven records are actually in class “1”

Confusion Matrix for Different Cutoffs

Cut off Prob.Val. for Success (Updatable)

0.25

Classification Confusion Matrix		
	Predicted Class	
Actual Class	owner	non-owner
owner	11	1
non-owner	4	8

Cut off Prob.Val. for Success (Updatable)

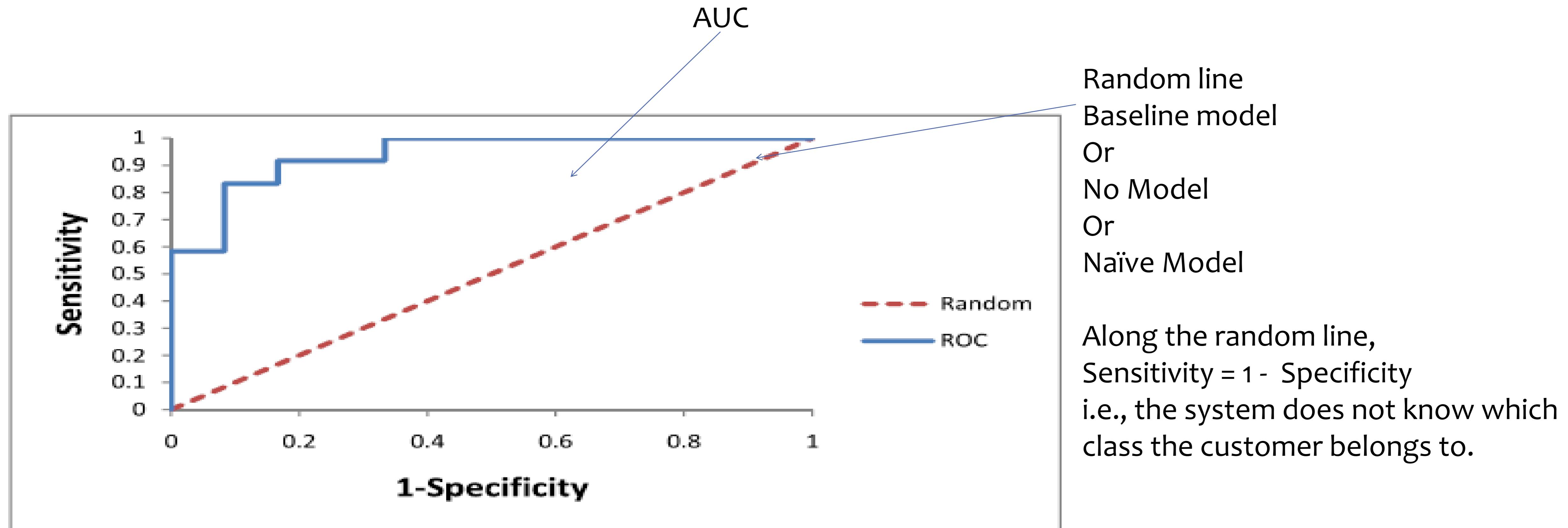
0.75

Classification Confusion Matrix		
	Predicted Class	
Actual Class	owner	non-owner
owner	7	5
non-owner	1	11

Other performance measures.

		Predicted Class		F1-score/ F-score = $2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{(TN + FP)}$
	Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$	Recall

Receiver Operating Characteristic curve (ROC Curve)



A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied.

ROC Curve

Compare performance of DM model to “no model, pick randomly”

Measures ability of DM model to identify the important class, relative to its average prevalence

Charts give explicit assessment of results over a large number of cutoffs

Asymmetric Costs

Misclassification Costs May Differ

The cost of making a misclassification error may be higher for one class than the other(s)

Looked at another way, the benefit of making a correct classification may be higher for one class than the other(s)

Example – Response to Promotional Offer

Suppose we send an offer to 1000 people, with 1% average response rate
 (“1” = response, “0” = nonresponse)

- “Naïve rule” (classify everyone as “0”) has error rate of 1% (seems good)
- Using DM we can correctly classify eight 1’s as 1’s
It comes at the cost of misclassifying twenty 0’s as 1’s and two 0’s as 1’s.

The Confusion Matrix

	Predict as 1	Predict as 0
Actual 1	8	2
Actual 0	20	970

Error rate = $(2+20) = 2.2\%$ (higher than naïve rate)

Introducing Costs & Benefits

Suppose:

- Profit from a “1” is \$10
- Cost of sending offer is \$1

Then:

- Under naïve rule, all are classified as “0”, so no offers are sent: no cost, no profit
- Under DM predictions, 28 offers are sent.
 - 8 respond with profit of \$10 each
 - 20 fail to respond, cost \$1 each
 - 972 receive nothing (no cost, no profit)
- Net profit = \$60

Profit Matrix

	Predict as 1	Predict as 0
Actual 1	\$80	0
Actual 0	(-\$20)	0

Generalize to Cost Ratio

Sometimes actual costs and benefits are hard to estimate

- Need to express everything in terms of costs (i.e., cost of misclassification per record)
- Goal is to minimize the average cost per record

A good practical substitute for individual costs is the **ratio** of misclassification costs (e.g., “misclassifying fraudulent firms is 5 times worse than misclassifying solvent firms”)

Minimizing Cost Ratio

q_1 = cost of misclassifying an actual “1”,

q_0 = cost of misclassifying an actual “0”

Minimizing the **cost ratio** q_1/q_0 is identical to
minimizing the average cost per record

Software* may provide option for user to specify cost ratio

*Currently unavailable in XLMiner

Note: Opportunity costs

- As we see, best to convert everything to costs, as opposed to a mix of costs and benefits
- E.g., instead of “benefit from sale” refer to “opportunity cost of lost sale”
- Leads to same decisions, but referring only to costs allows greater applicability

Cost Matrix (inc. opportunity costs)

	Predict as 1	Predict as 0
Actual 1	\$8	\$20
Actual 0	\$20	\$0

Recall original confusion matrix (profit from a “1” = \$10, cost of sending offer = \$1):

	Predict as 1	Predict as 0
Actual 1	8	2
Actual 0	20	970

Multiple Classes

For m classes, confusion matrix has m rows and m columns

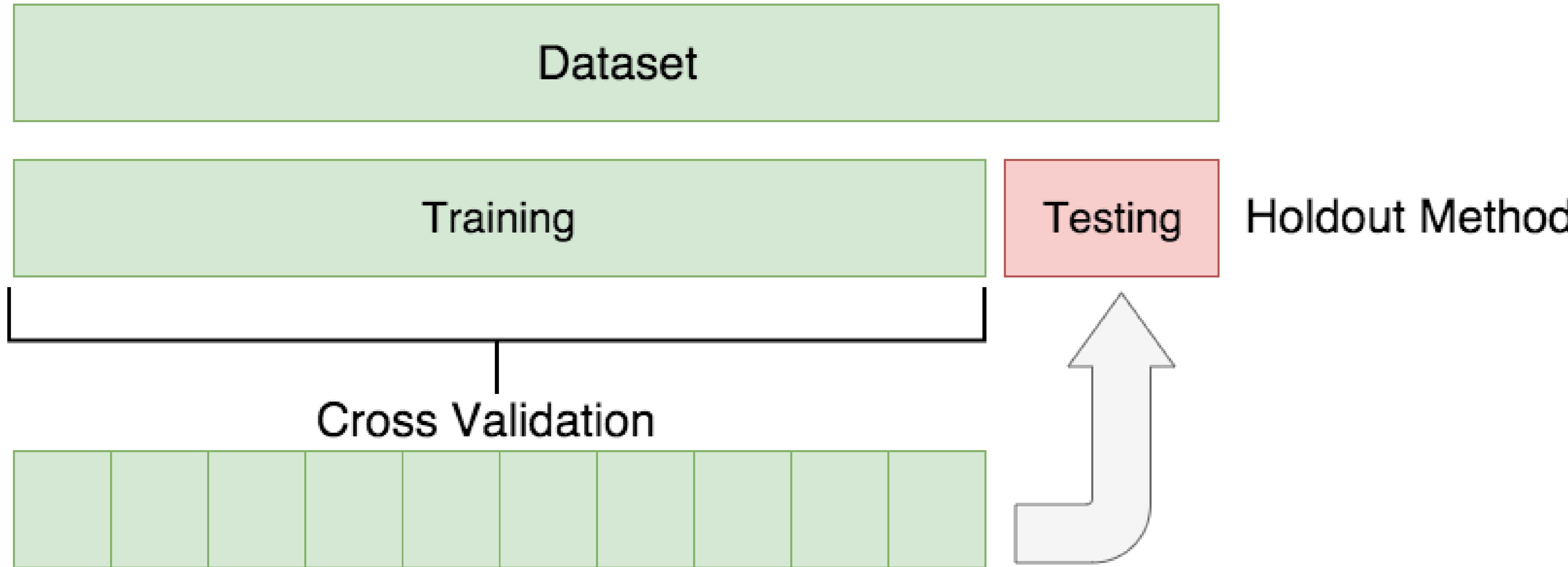
- Theoretically, there are $m(m-1)$ misclassification costs, since any case could be misclassified in $m-1$ ways
- Practically too many to work with
- In decision-making context, though, such complexity rarely arises – one class is usually of primary interest

Confusion Matrix for Multi-class problems

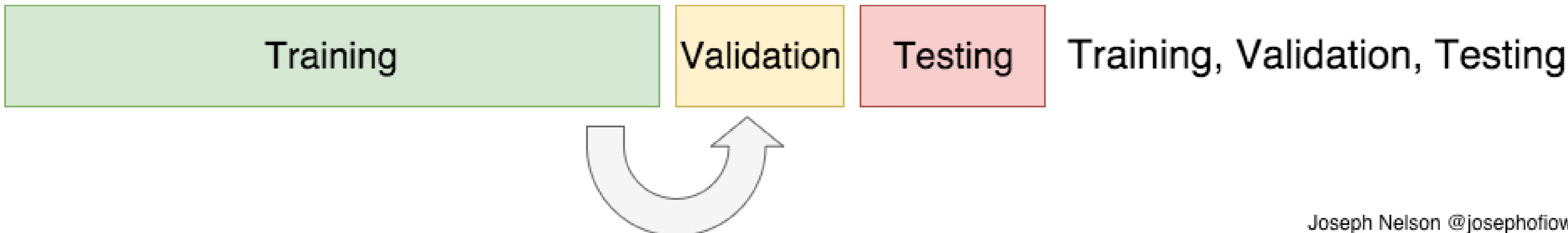
[https://www.analyticsvidhya.com/blog/2021/06/confusion-matrix-for-multi-class-classification/#:~:text=The%20confusion%20matrix%20is%20a,and%20False%20Negative\(FN\).](https://www.analyticsvidhya.com/blog/2021/06/confusion-matrix-for-multi-class-classification/#:~:text=The%20confusion%20matrix%20is%20a,and%20False%20Negative(FN).)

<https://towardsdatascience.com/confusion-matrix-for-your-multi-class-machine-learning-model-ff9aa3bf7826>

Dividing the dataset into training and testing sets



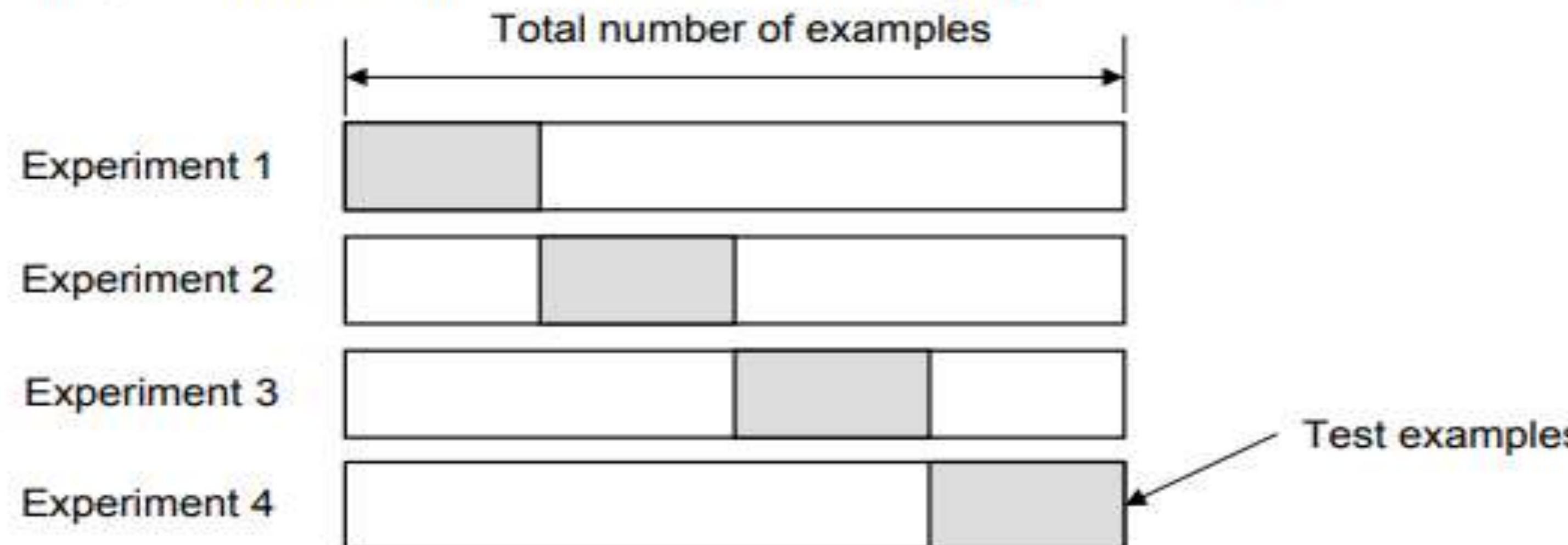
Data Permitting:



k-Fold Cross Validation

- **Create a K-fold partition of the dataset**

- For each of K experiments, use K-1 folds for training and a different fold for testing
 - This procedure is illustrated in the following figure for K=4



- **K-Fold Cross validation is similar to Random Subsampling**

- The advantage of K-Fold Cross validation is that all the examples in the dataset are eventually used for both training and testing

- **As before, the true error is estimated as the average error rate on test examples**

$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

Summary

- Evaluation metrics are important for comparing across DM models, for choosing the right configuration of a specific DM model, and for comparing to the baseline
- Major metrics: confusion matrix, error rate, predictive error
- Other metrics when
 - one class is more important
 - asymmetric costs
- When important class is rare, use oversampling
- In all cases, metrics computed from validation data

The content of the slides are prepared from different textbooks.

References:

- Data Mining and Predictive Analytics, By Daniel T. Larose. Copyright 2015 John Wiley & Sons, Inc.
- Predictive Analytics for Dummies, By Anasse Bari, Mohamed Chaouchi, & Tommy Jung, Copyright 2016, John Wiley & Sons, Inc.
- Introduction to Data Mining with Case Studies, By G.K. Gupta. Copyright 2014 by PHI Learning Private Limited.



—
Thank you..