

17. What is discrete attributes?

Ans:

A discrete attributes has a finite or countably infinite set of values; which may or may not be represented as integers.

18. What is Continuous attributes?

Ans:

A Continuous attribute has a numeric or Continuous attribute values.

Discrete: All the values are Categorical

The Weather Problem

Outlook	Temperature	Humidity	Wind	Play
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No

Table: Weather Table

1. প্রতিটি কলাম হল feature বা attribute.
2. উত্তরে চিহ্নিত 9 টি attribute আছে-

- Outlook
- Temperature
- Humidity
- Wind

3. Play কলাম হল class.

4. এই table হল একটি historical data

ମାତ୍ରାଣୀ feature ଏବଂ feature value ଆବଶ୍ୟକ,
ଯେଉଁ:

Outlook → Sunny
 → Overcast
 → Rain

ଉର୍ଦ୍ଧ୍ୱ

Temperature → Hot
 → Mild
 → Cool

ଉର୍ଦ୍ଧ୍ୱ

Humidity → High
 → Normal

ଉର୍ଦ୍ଧ୍ୱ

Wind → Weak
 → Strong

ଉର୍ଦ୍ଧ୍ୱ

ମାତ୍ରାଣୀ Permutation; Combination ବ୍ୟବହାର

$$3 \times 3 \times 2 \times 2 = 36 \text{ ଟା Combination}$$

ଅନୁସାରେ,

ମାତ୍ରାଣୀ Decision Tree ଏହା ସାହାଯ୍ୟରେ ଏହି Combination
ସୂଚକ rules ଆକାରରେ sort out କରାଯାଇ ଏହା,
ଅନ୍ତର୍ଗତ rules optimal କରିବା।

Machine Learning Models Rules

1. If Outlook = Sunny & Humidity = High then Play = No
2. If Outlook = Sunny & Humidity = Normal then Play = Yes
- ~~3. If Outlook = Sunny &~~
3. If Outlook = Overcast then Play = Yes
4. If Outlook = Rain & Wind = Strong then Play = No
5. If Outlook = Rain & Wind = Weak then Play = Yes

=> These rules generate a very simple model.
Temperature is a less important feature.

Table for Numerical value representation.

Table: Weather Data with some Numeric Attributes

Outlook	Temperature	Humidity	Wind	Play
Sunny	85	85	Weak	No
Sunny	80	90	Strong	No
Overcast	83	86	Weak	Yes
Rain	70	96	Weak	Yes
Rain	68	80	Weak	Yes
Rain	65	70	Strong	No
Overcast	64	65	Strong	Yes
Sunny	72	95	Weak	No
Sunny	69	70	Weak	Yes
Rain	75	80	Weak	Yes
Sunny	75	70	Strong	Yes
Overcast	72	90	Strong	Yes
Overcast	81	75	Weak	Yes
Rain	71	91	Strong	No

Example Rules for \Rightarrow

1. If Outlook = Sunny & Humidity > 75 then Play = No
2. If Outlook = Sunny & Humidity \leq 75 then Play = Yes
3. If Outlook = Overcast then Play = Yes
4. If Outlook = Rain & Wind = Strong then Play = No
5. If Outlook = Rain & Wind = Weak then Play = Yes.

Threshold point.

Humidity for separate weather,

Threshold Point for Cutpoint / Discrete Discrete

Define Concept.

Ans:

Concept is the thing to be learned.

Define Concept Description

Ans:

The output produced by a learning classifier.

Define Instances

Ans:

Things that are to be classified or associated or clustered.

Q Dataset কে matrix আকারে represent করা হয়,
 যেখানে একগাছি-খাটো instances আর অন্যগাছি-
 খাটো attributes

Q What is Flat file

Ans:

A flat file database is one that only contains a single table of data.

Q Define Attribute

Ans:

Attribute is a data field, representing a characteristic or feature of a data object.

Q What are the steps of data Preprocessing

Ans:

Data Preprocessing has four steps:

- Data cleaning: Missing values, noisy data.
- Data Integration: Redundancy, Correlation, instance duplication.
- Data Reduction: Dimensionality Reduction.
- Attribute subset selection.

Q. Draw block diagram of Pattern mining Process.

Ans

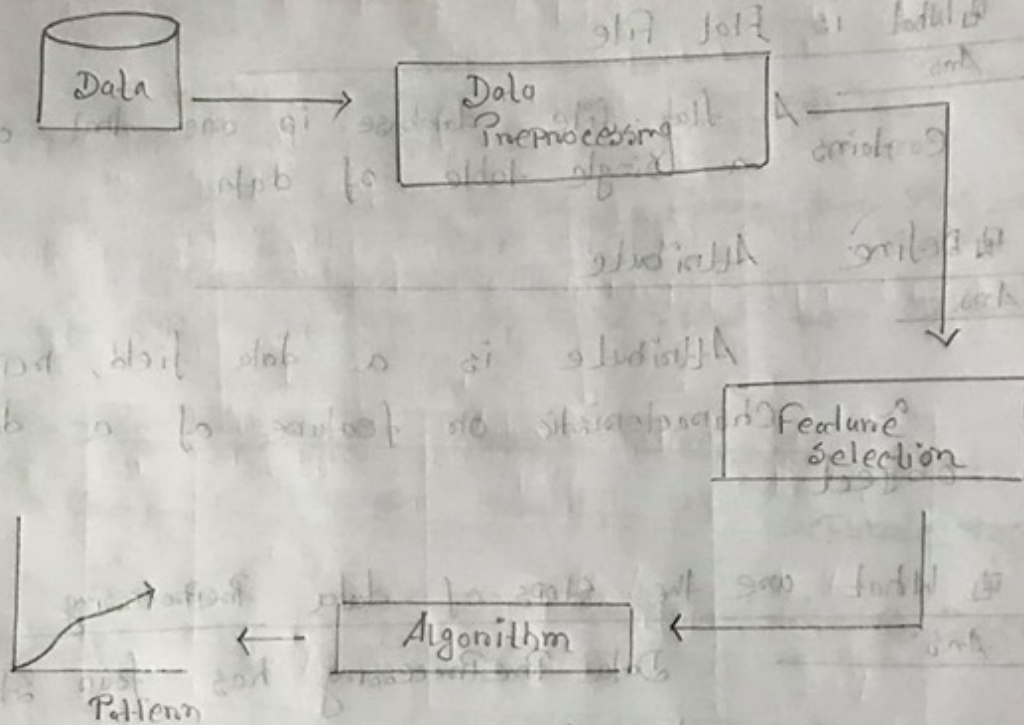


Fig: Data Mining Process.

Q. What is ARFF

Ans:

ARFF stands for "Attribute Relation File Format".

It's mainly an ASCII text file that describes a list of instances sharing a set of attributes.

Q Write the Work Process of KNN algorithm

Ans:

KNN is a classification algorithm. It is supervised. It takes a bunch of labelled points & uses them to learn how to label other points. To label a new point it looks at the labelled points closest to the new point & has those neighbors vote.

Q Write the KNN algorithm.

Ans:

Input: $D = \{x_1, \dots, x_n\}$

Input: $D = \{x_1, \dots, x_i, \dots, x_n\}$

Output: KNN Classifier, KNN

Method:

a. Find $x \in D$, that identify the k nearest neighbours. Regardless of class label, c_i .

b. Out of these instances, $X = \{x_1, x_2, \dots, x_k\}$ identify the numbers of instances, k_i that belong to class c_i . Where $i = 1, 2, n$; $\sum_i k_i = k$.

c. Assign x_{test} to the class c_i with the maximum number of k_i of instances.

ii. D data learning data
 iii. X data Nearest neighbours of test data

iv. KNN is dataset for training and testing.

- i. Training set
- ii. Validation set
- iii. Test set.

v. Dataset for training and testing.

	Training set	Validation set	Test set
1.	80 %	10 %	10 %
2.	70 %	15 %	15 %
3.	70 %	x	30 %

Training set

- a. Model for training and testing
- b. Classifier construct

Validation set

- a. Hyper-parameters tune for accuracy
- b. Accuracy increase,
- c. Validation Accuracy maximize
- d.

a. Test set actual accuracy হলো,

কি KNN এর distance বোঝে বসিয়ে দেয়, তার distance
- মতোই হবে বসিয়ে দেয়।

a. Euclidean Distance

b. Manhattan Distance

Euclidean Distance:

(x, y) একটি Point ; (a, b) আর একটি Point,
এদের মধ্যকার Euclidean Distance হলো

$$(x, y), (a, b) = \sqrt{(x-a)^2 + (y-b)^2}$$

Manhattan Distance:

(x, y) একটি Point ; (a, b) আর একটি Point,
এদের মধ্যকার Manhattan Distance হলো

$$(x, y), (a, b) = |x-a| + |y-b|$$

The following dataset (Table 1) with two features (acid durability & strength) to classify whether a special tissue paper is good or not.

A_1	A_2	Class
7	7	Bad
7	4	Bad
3	4	Good
1	4	Good

Now, the factory produces a new tissue paper (x_{new}) that pass laboratory test with $A_1 = 3$, $A_2 = 7$. So classify the x_{new} using k-Nearest neighbor (kNN) classifier, where $k=3$ & distance function is

Distance [Spring 19]

Ans:

$$|d - b| + |a - x| = (d, a), (b, x)$$

A_1	A_2	Distance
7	7	$\sqrt{(7-3)^2 + (7-7)^2} = \sqrt{16} = 4$
7	4	$\sqrt{(7-3)^2 + (4-7)^2} = \sqrt{25} = 5$
3	4	$\sqrt{(3-3)^2 + (4-7)^2} = \sqrt{9} = 3$
1	4	$\sqrt{(1-3)^2 + (4-7)^2} = \sqrt{13} = 3.60$

Number of instances	Distance	Nearest N	Majority class	Majority voting
1	4	Yes	Bad	
2	5	No	X	
3	3	Yes	Good	Good
4	3.60	Yes	Good	

So, X_{new} belongs to class "Good"

Q Explain the difference & similarity between classification & regression.

Ans.

Differences between Classification vs Regression

Classification	Regression
a. Group the output into a class.	a. Predict the output value using training data.
b. If it is discrete/ categorical variable then it is classification Problem.	b. If it is real numbers/ continuous then it is regression Problem.
c. Data is labelled into one or more classes.	c. Prediction of a quantity.

Similarities between Classification & Regression

- Both are supervised learning
- Develop Predictive model based on both input & output data

- | | |
|------------------|-------------------------|
| * Classification | → Supervised Learning |
| * Regression | → Supervised Learning |
| * Clustering | → Unsupervised Learning |

Q Write the disadvantage of KNN -

Ans.

KNN is a lazy learners. It does not learn anything from training data. It simply uses the training data itself for classification.

Q Write the drawback of KNN

Ans.

Its complexity is $O(KN)^2$.
Where

N = Total number of neighbours

K = Number of nearest neighbours.

It takes too much time. For new test data it runs from the beginning.