

Lecture-1: Supervised and Unsupervised Machine Learning Algorithms.

- Classification and regression : supervised learning
- Clustering and association : unsupervised learning.
- between sup and unsup is called semisupervised learning.

Supervised Machine Learning:

- Practical learning uses supervised learning. (forms output file)
- input variable (x) and output variable (Y), use an algorithm to learn the mapping function from input to output.

$$Y = f(x) \Rightarrow \text{mapping function}$$

new input data (x) can predict output (Y) for data.

The process of an algorithm learning from the training dataset can be thought as a teacher supervising the learning process is called Supervised Learning.

* we know the correct ans, then algo makes prediction on training data.

Classification: output variable is a category.

Ex: red or blue ; 'disease' or 'no disease'

Regression: output variable is a real value.

Ex: dollars, weight

Example: of supervised ML:

recommendation, time series prediction.

① Linear regression
(regression problem)

② Random Forest
(regression ;
classification
problem)

③ SVM
(classification
problem)

Unsupervised Machine Learning:

UML is where only have input data (x) and no corresponding output variables.

Goal: underlying structure in the data, so it can learn more about the data.

There is no correct answer, no teacher. Algorithms are left to their own devices to discover and present to the interesting structure in the data. is called unsupervised learning.

UML further grouped into clustering and association problems.

Clustering: Discover the inherent groupings in the data,

Ex: grouping customers by purchasing behaviour.

सर्गो, अर्गोस सर्गो, लुगोस सर्गो @

Association :

Discover rules that describe large portions of data.

(* data ko sare rules ko dikhao)

Ex: agar x buy karta, agar y ko buy karta or agar z ko buy karta.

Example : ① k-means for clustering.

② Apriori " association.

① Semi-supervised ml

A large amount of input data (x) and only some of the data is labeled (Y) are called semi-supervised learning problems.

Example: photo archive, some of images are labeled and majority are ~~unabled~~ unlabeled.

- unlabeled data is cheap and easy to collect and store.
- many world fall into this area.
- use SSML to discover and learn the structure in the input variables.

① Supervised: All data is labeled and algorithms learn to predict the output from the input data.

② Unsupervised: All data is ~~unlabeled~~ unlabeled and algorithms learn to inherent structure from the input.

③ Semi-supervised: some data is labeled but most of data is unlabeled, mixture of supervised and unsupervised techniques are used.

Lecture : 2

KDD = Knowledge data discovery.

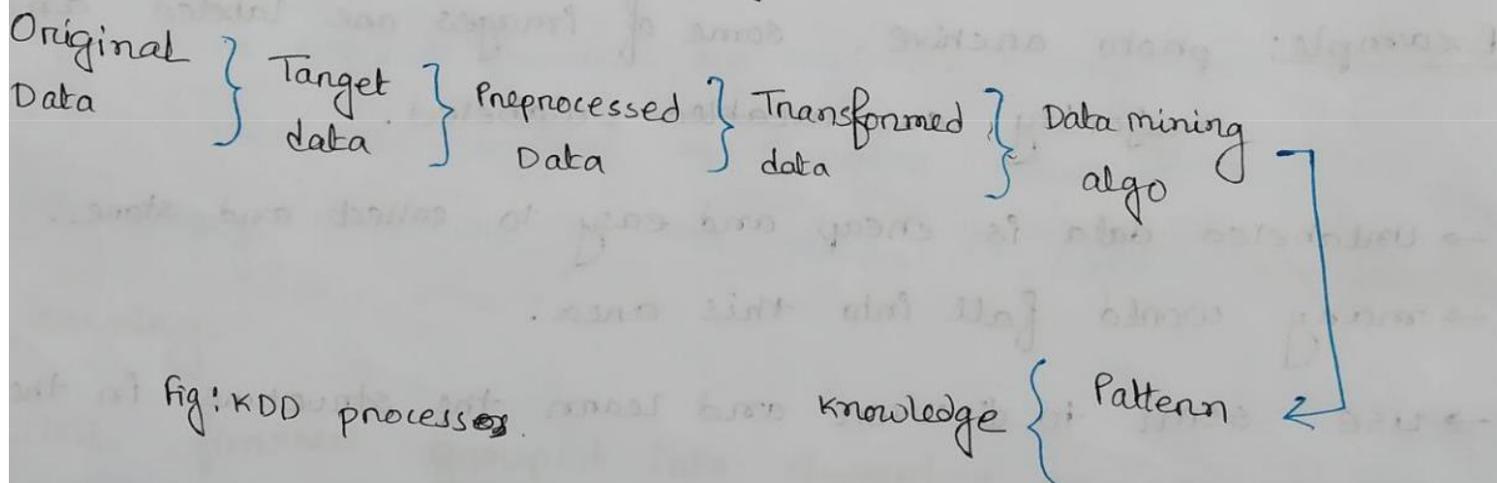


fig: KDD processes.

Lecture 2

Why we need to prepare data?

Data can be incomplete, inconsistent or noisy.

- 1 Data entry, transmission, data collection problem solve
- 2 duplicate data remove
- 3 incomplete missing data problem solve
- 4 avoid contradictions in data.

When data can't be trusted:

- * difficult to make knowledge from untrusted site.
- * Better chance to discover new knowledge or patterns

when data is clean.

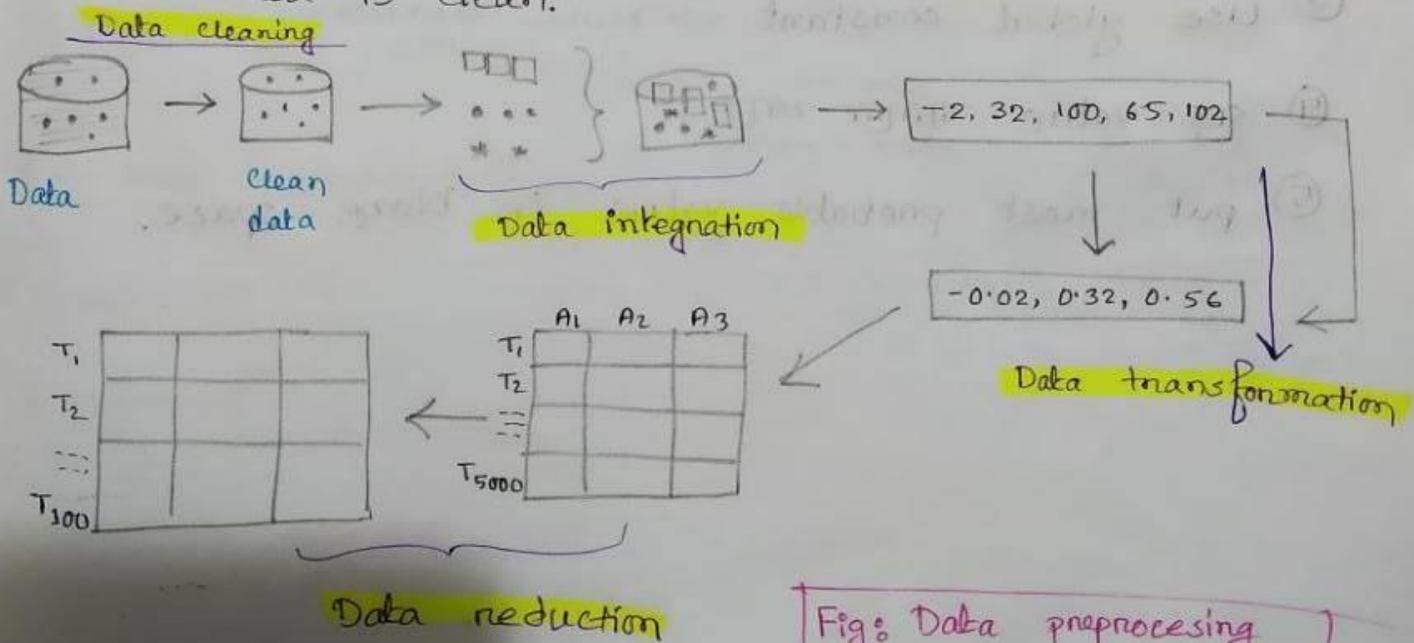


Fig: Data preprocessing

Data cleaning:

real world data is noisy, incomplete, inconsistent etc.

- some attribute has no value
- random error
- irrelevant records

data cleaning
etc.

Data cleaning steps:

- 1 fill in missing values
- 2 pull out noisy data
- 3 correct inconsistencies
- 4 remove irrelevant data

Solving missing data problem:

- 1 ignore the record with missing value
- 2 fill the missing value manually.
- 3 Use global constant \Rightarrow NULL, unknown
- 4 use value from other records
- 5 put most probable value in blank space.

plate

Smoother Noisy Data 3 methods

① Binning method.

Avg Boundaries

② Clustering

③ Regression

Binning method technique

- ① Data sort kro
- ② 3 bin ka 3 avg value prso. when limit is given = 3
- ③ " " ka avg
- ④ 3 bin ka 3 avg value avg pr replace kro.
- ⑤ in cast of float, keep the nearest int value.
- ⑥ Then put the value according to ID in data table.

Data Transformation Nonnormalization

① Min max method →
$$v' = \frac{v - \min}{\max - \min}$$
 when range is [0, 1]

② Z-score nonnormalization →
$$v' = \frac{(v - \text{mean})}{\text{Stdv (Standard Deviation)}}$$

① $(\text{mean} - x_i)^2$

$$\frac{\sum (\text{mean} - x_i)^2}{N} = \text{variance}$$

$$\text{Stdv} = \sqrt{\text{variance}}$$

Data transformation: Discretization

Methods

Nominal = values from unordered list

Ordinal = " " ordered list

Numeric = real numbers.

① Binning

② Histogram Analysis

③ Clustering

④ Decision-tree

⑤ Correlation

① Discretization use for: continuous attribute reduce or

* some DM algo allow only categorical data/ Attribute, cannot handle continuous attribute value.

** Discretization can also used to generate high level concepts.

age continuous value convert high level concept

convert to discretization or

Example: young, middle age, old → category or age class.

(*** continuous value to discrete convert)

| Age |
|-----|
| 25 |
| 38 |
| 51 |
| 75 |
| 86 |

low = 25 - 40

MA = 41 - 60

old = 60 <

| Age |
|-----|
| Low |
| Low |
| MA |
| old |
| old |

Data reduction

Often Data become too large, reduce data for improve performance.

- ① dimensionally reduce ~~size~~ size
- ② Discretization
- ③ Regression / clustering reduction.

Chap = 9

- ① continuous value ~~array~~ dynamically define new discrete value with range / intervals.
- ② missing = assign probable value (possible)
= " most common value.

① Lecture = 3

Confusion Matrix

TP = Dengue α α ,
test positive } actual result.

FN = Dengue α β ,
test negative

TN = Dengue β β ,
test negative } actual result

| | | Predicted Class | |
|----------------|-----------------|-------------------------|--------------------------------|
| | | test | |
| Original Class | Dengue α | 1 Yes α | 0 No β |
| | Dengue β | 1 Yes α TP ++ | 0 No β FN +- TN -- |

FP = Dengue β α ,
test positive.

| | | |
|---|----|----|
| | 1 | 0 |
| 1 | TP | FN |
| 0 | FP | TN |

① Classification

find group, rule
शर्त,
(Supervise)

① Clustering

find group, rules अनशर्त,
(unsupervise)

Confusion Matrix

① Accuracy = $\frac{TP + TN}{ALL}$

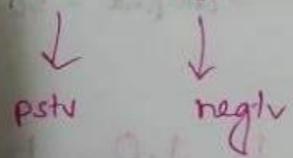
① Error rate = 1 - accuracy

or = $\frac{FN + FP}{ALL}$

① Sensitivity = $\frac{TP}{P}$ (TP rate)

① Specificity = $\frac{TN}{N}$

① bayes - rate should be = 60 - 40



P = Actual true
N = Negative

other classification evaluation metrics.

Precision: original 70% positive cases, machine says 95% for predict cases.

Example: original 70% actually heart diseases cases, but machine predict 95% actually heart disease cases.

machine says positive cases = ~~FN~~ TP + FP.

Precision = $\frac{TP}{TP + FP}$

on / m

$\xrightarrow{0}$ original 70
 \xrightarrow{F} 95%, machine says, where HD has or not, but machine says yes.

Recall: machine says predict cases, originally 50% cases are positive cases.

~~70%~~ 50% originally samsung use cases, but machine predict 30% cases.

m / on

recall = $\frac{TP}{TP + FN}$

\xrightarrow{F} machine predict actual users = 30
 $\xrightarrow{0}$ original = 50

Perfect score for both is 1.0, but there is often a

a trade-off between precision and recall.

F measure : (F_1 or F-score)

harmonic mean of precision and recall,

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{Overall correct rate} = \frac{TP + FN}{TP + FP + TN + FN}$$

Association Analysis

Now, find, after buying Bread, how many prob to buy milks.

| TID | Items |
|-----|----------------------------|
| 1 | Bread, milk |
| 2 | Bread, Diapers, Beer, Egg |
| 3 | milk, Diapers, Beer, Cola |
| 4 | Bread, milk, Diapers, Beer |
| 5 | Bread, milk, Diapers, Cola |

$$\text{Support}_{x \rightarrow y} = \frac{|x \cup y|}{N}$$

$$\text{Confidence}_{x \rightarrow y} = \frac{|x \cup y|}{\sigma(x)}$$

* $x \cup y$ = x aur y ke items.

* N = total num of transition

* $\sigma(x)$ = total num of x

| TID | Bread | milk | Diaper | Beer | Egg | Cola |
|-----|-------|------|--------|------|-----|------|
| 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 2 | 1 | 0 | 1 | 1 | 1 | 0 |
| 3 | 0 | 1 | 1 | 1 | 0 | 1 |
| 4 | 1 | 1 | 1 | 1 | 0 | 0 |
| 5 | 1 | 1 | 1 | 0 | 0 | 1 |

0/1 binary representation

Q: for Bread \rightarrow Milk.

$$\text{Support} : \frac{3}{5} ; \quad \text{confidence} = \frac{3}{4}$$

* Total num of combination : $R = 2^d - 2^{d+1} + 1$

brute force approach for mining s & c for every possible rule.

d = num of item.

B, D \rightarrow M

D, M \rightarrow B etc.

E, D \rightarrow M

Threshold = 50%

①

| TID | Items |
|-----|---------|
| 001 | 1 3 4 |
| 002 | 2 3 5 |
| 003 | 1 2 3 5 |
| 004 | 2 5 |

from item purchase info

②

| Item | sup |
|------|-----|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {4} | ① |
| {5} | 3 |

each item name/oid

item support

③

| Item | sup |
|------|-----|
| {1} | 2 |
| {2} | 3 |
| {3} | 3 |
| {5} | 3 |

maximum product sell = 3
threshold = 50%.
So, we not accept
(3 x 50%) = 1.5
So, item 4 not.

⑤

| Itemset | sup |
|---------|-----|
| {1, 2} | ① |
| {1, 3} | 2 |
| {1, 5} | ① |
| {2, 3} | 2 |
| {2, 5} | 3 |
| {3, 5} | 2 |

set of product info

④

| Itemset |
|---------|
| {1, 2} |
| {1, 3} |
| {1, 5} |
| {2, 3} |
| {2, 5} |
| {3, 5} |

item of combination

this combination also available in table 1

⑥

| Item | sup |
|------|-----|
| 1, 3 | 2 |
| 2, 3 | 2 |
| 2, 5 | 3 |
| 3, 5 | 2 |

remove, 1, 5 > item

⑦

| Itemset |
|---------|
| 2 3 5 |

Here, 1, 3, 5 is not available in table 1.
(2, 3, 5) available

⑧

| Item | sup |
|---------|-----|
| {2 3 5} | 2 |

Here (2, 3, 5) combination occur 2 times

in table 1

Final stage

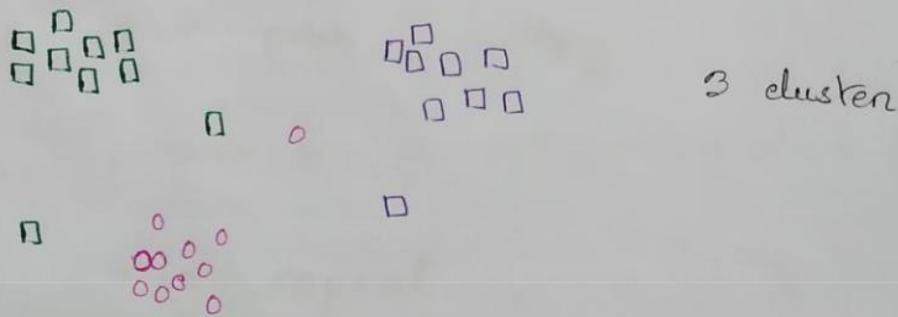
K-mean clustering

Clustering: Clustering is a task of grouping set of object in such way that objects in the same group are more similar to each other.

* same group को cluster कहते हैं। प्रत्येक cluster के objects के characteristics same, बराबर के होते हैं।

* distance measure को grouping करते हैं।

Example:



Type of clustering:

① Hierarchical Algorithm: Hierarchy cluster / ranking cluster build करते हैं। Clustering fall into two types

① bottom-up

उन्हें → छोटे समूह

② top-down

एक → कई

② Partitional clustering: determine cluster at once.

⇒ k-means and derivatives

⇒ fuzzy c-means clustering.

⇒ OT clustering.

① Common distance measure: distance first similarity measure

cluster shape influence

① The Euclidean distance $\Rightarrow d = \sqrt{\sum |x_i - y_i|^2}$

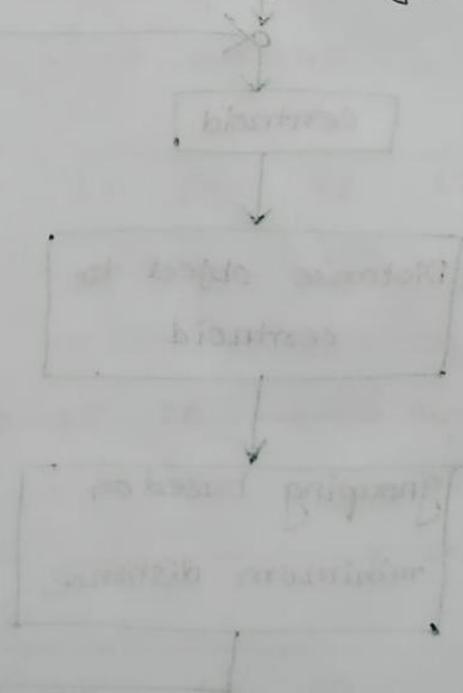
② The Manhattan Distance $\Rightarrow d = \sum |x_i - y_i|$

③ maximum norm $\Rightarrow d = \max |x_i - y_i|$

④ Mahalanobis distance

⑤ Inner product space

⑥ Hamming distance.

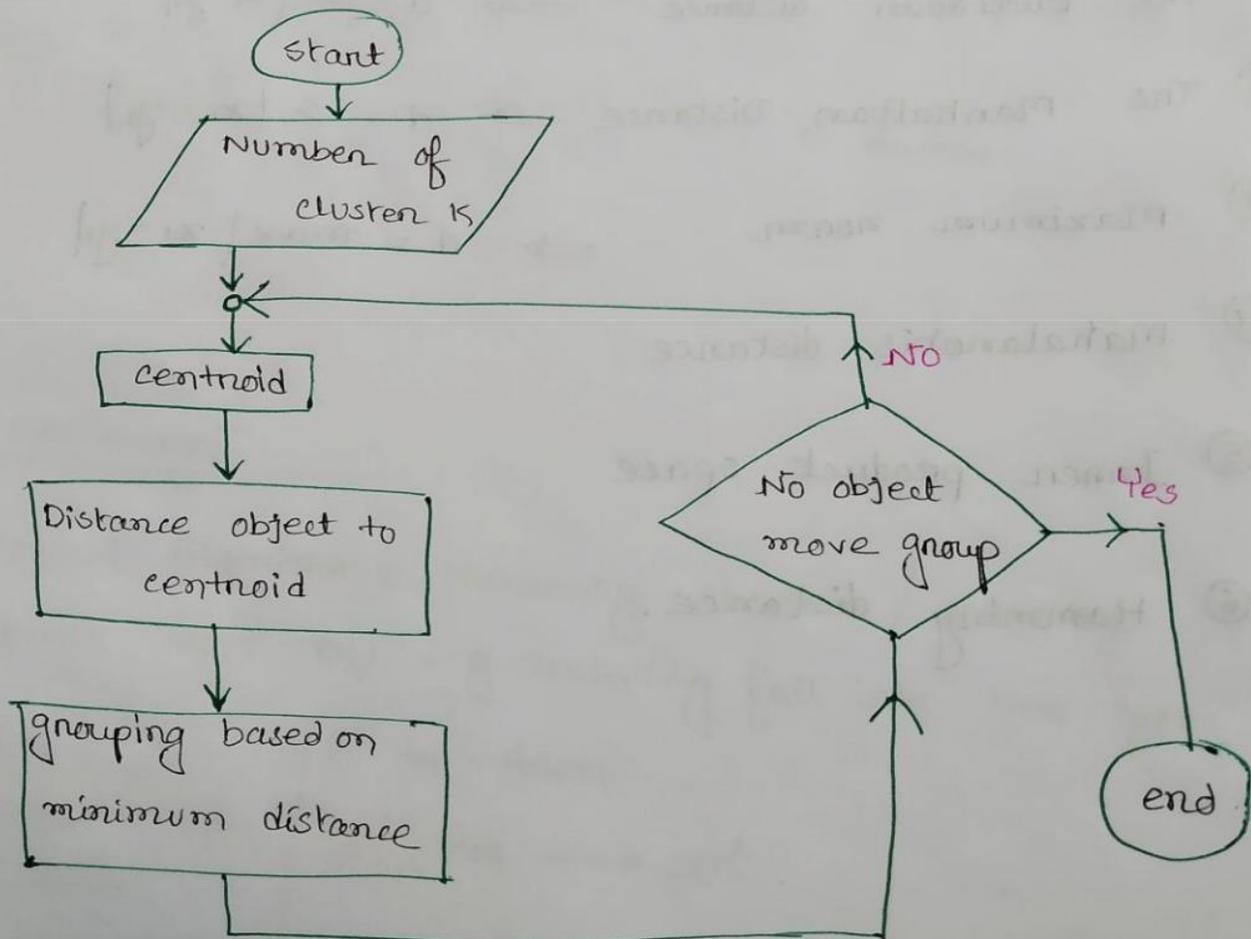


K-mean clustering

n = object number, K = partition, $n > K$.

एकसमय 2 से n तक गिने g_1, g_2 को, then, सबसे
close object को cluster में शामिल, distance measure
करें, next g_1, g_2 इन group object को average.

centroid बदलने चंग शुरू, जब तक g_1, g_2
बदलने 2 बार same रहते हों clustering done.



1 5 9 15 20 25 35 50 70

Here $k = 2$.

Solution: let, centroid 1 = 1, centroid 2 = 20.

distance from 1: 0 4 8 14 19 24 34 49 69

" " 20: 19 15 11 5 0 5 15 30 50

Now, (1) group 1: 1 5 9 \longrightarrow avg = 5

(20) group 2: 15 20 25 35 50 70 \longrightarrow avg = 35

distance from 5: 4 0 4 10 15 20 30 45 65

" " 35: 34 30 26 20 15 10 0 15 35

(5) group: 1 5 9 15 \longrightarrow avg = 8

(35) group: 20 25 35 50 70 \longrightarrow avg = 40

d.f 8: 7 3 1 7 12 17 27 42 62

d.f 40: 30 35 31 25 20 15 5 10 30

(8) group: ~~1~~ ~~5~~ ~~9~~ ~~15~~ 1 5 9 15 20 \longrightarrow avg = 10

(40) group: 25 35 50 70 \longrightarrow avg = 45

d.f 10 = 9 5 1 5 10 15 25 40 60

d.f 45 = 44 40 36 20 25 20 5 5 25

1 5 9 15 20 25 35 50 70

$$(10) \text{ group} = 1 \ 5 \ 9 \ 15 \ 20 \ 25 \rightarrow 12.5$$

$$(45) \text{ group} = 35 \ 50 \ 70. \rightarrow 51.66$$

d.f 13: 12 8 4 2 7 12 22 37 57

d.f 52: 51 47 43 37 32 27 17 2 18

$$(13) \text{ group} = 1 \ 5 \ 9 \ 15 \ 20 \ 25$$

$$(52) \text{ group} = 35 \ 50 \ 70.$$

as repeat 2 time so, group / clustering done.

$$\text{Ans: } [1 \ 5 \ 9 \ 15 \ 20 \ 25]$$

$$[35 \ 50 \ 70]$$

A

Example of 2D data (x, y) .

$$\text{Euclidean Distance} = \sqrt{(x-x_1)^2 + (y-y_1)^2}$$

given that,

$$K=2$$

$$m_1 = 1.0, 1.0$$

$$m_2 = 5.0, 7.0$$

| indv | v_1 (x) | v_2 (y) |
|------|-----------|-----------|
| 1 | 1.0 | 1.0 |
| 2 | 1.5 | 2.0 |
| 3 | 3.0 | 4.0 |
| 4 | 5.0 | 7.0 |
| 5 | 3.5 | 5.0 |
| 6 | 4.5 | 5.0 |
| 7 | 3.5 | 4.5 |

Now, calculate distance,

for m_1 :

$$1 = \sqrt{(1.0-1.0)^2 + (1.0-1.0)^2} = 0$$

$$2 = \sqrt{(1.0-1.5)^2 + (1.0-2.0)^2} = 1.12$$

$$3 = \sqrt{(1.0-3.0)^2 + (1.0-4.0)^2} = 3.60$$

$$4 = \sqrt{(1.0-5.0)^2 + (1.0-7.0)^2} = 7.21$$

$$5 = \sqrt{(1.0-3.5)^2 + (1.0-5.0)^2} = 4.72$$

$$6 = \sqrt{(1.0-4.5)^2 + (1.0-5.0)^2} = 5.31$$

$$7 = \sqrt{(1.0-3.5)^2 + (1.0-4.5)^2} = 4.30$$

for m2:

$$1 = \sqrt{(5.0 - 1.0)^2 + (7.0 - 1.0)^2} = 7.21$$

$$2 = \sqrt{(5.0 - 1.5)^2 + (7.0 - 2.0)^2} = 6.10$$

$$3 = \sqrt{(5.0 - 3.0)^2 + (7.0 - 4.0)^2} = 3.61$$

$$4 = \sqrt{(5.0 - 5.0)^2 + (7.0 - 7.0)^2} = 0$$

$$5 = \sqrt{(5.0 - 3.5)^2 + (7.0 - 5.0)^2} = 2.5$$

$$6 = \sqrt{(5.0 - 4.5)^2 + (7.0 - 5.0)^2} = 2.06$$

$$7 = \sqrt{(5.0 - 3.5)^2 + (7.0 - 4.5)^2} = 2.92$$

Now

| | cen 1 | cen 2 |
|---|-------------------|-------|
| 1 | 0 | 7.21 |
| 2 | 0 1.12 | 6.10 |
| 3 | 3.61 | 3.61 |
| 4 | 7.21 | 0 |
| 5 | 4.72 | 2.5 |
| 6 | 5.31 | 2.06 |
| 7 | 4.30 | 2.92 |

group 1 {

} group-2

Weakness of K-mean:

- ① Data points are grouped into clusters.
- ② cluster, k must be determined before, cluster depends on initial assignment.
- ③ never know the real cluster.
- ④ different initial condition or gives different cluster.

Application of K-mean:

① efficient & fast. time = $O(tkn)$

n = number of object

k = " " cluster

t = " " iteration.

② machine learning & data mining @ use.

③ wave \rightarrow k category @ use, wave, speech data @ use.

④ color plot & graphical display @ use.

Lecture = 5 Classification

Classification process: 3 steps

- ① model construction (learning) pre defined class & record
target, to determine target attribute value,
value of target attribute = class labels,
Training set = learning & use it.
- ② model Evaluation: (Accuracy) estimate accuracy rate
based on test set.
- ③ model use: (classification) model is used to
classify unseen instances. Predict value.

Classification method: ① Decision tree

- ② Bayesian classification
- ③ K-Nearest neighbor
- ④ Neural networks
- ⑤ SVM
- ⑥ Genetic Algorithm.

Naive Bayesian Learning

Naive bayes play important role on prediction.

$$P(H|x) = \frac{P(x|H) \cdot P(H)}{P(x)}$$

$P(A|B)$ = prob of A, given B.

B is known.

$$P(A \wedge B) = P(A|B) \cdot P(B). \quad \boxed{****}$$

Bayes's rule = $P(A \wedge B) = P(A|B) \cdot P(B)$

$$P(B \wedge A) = P(B|A) \cdot P(A).$$

but, $P(A \wedge B) = P(B \wedge A).$

$$\therefore P(A|B) \cdot P(B) = P(B|A) \cdot P(A).$$

$$\Rightarrow P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}.$$

So, $P(H|E) = \frac{P(E|H) \cdot P(H)}{P(E)}$

$$P(A|B) = P(A) \quad ; \quad P(B|A) = P(B)$$

two constraints are logically equal.

$$P(A|B) = \frac{P(A \wedge B)}{P(B)} = P(A)$$

$$P(A \wedge B) = \cancel{P(A)} P(A|B) \cdot P(B) = P(A) \cdot P(B)$$

$$P(C_i|E) = \frac{P(C_i) \cdot P(E|C_i)}{P(E)}$$

math

Problem: Too many possible combination.

Estimate = Estimated based on training data.

= compute mean & standard deviation

Smoothing: small training set, \bar{x} estimate smooth \bar{x}

$$P(\text{Yes}) = \frac{9}{14} \begin{array}{l} \text{--- yes} \\ \text{--- all} \end{array}$$

$$P(\text{sunny} | \text{Yes}) = \frac{2}{9} \begin{array}{l} \text{--- sunny} \rightarrow \text{yes} \\ \text{--- total yes.} \end{array}$$

$$P(\text{humidity} | \text{No}) = P(\text{high} | \text{No}) = \frac{4}{5}$$

4 for high, 5 for actual no.

| outlook | P | N |
|----------|-----|-----|
| Sunny | 2/9 | 3/5 |
| Overcast | 4/9 | 0 |
| rain | 3/9 | 2/5 |

| Humid. | P | N |
|--------|-----|-----|
| high | 3/9 | 4/5 |
| normal | 6/9 | 1/5 |

| Temp | P | N |
|------|-----|-----|
| hot | 2/9 | 2/5 |
| mild | 4/9 | 2/5 |
| cool | 3/9 | 1/5 |

| wind | P | N |
|-------|-----|-----|
| true | 3/9 | 3/5 |
| false | 6/9 | 2/5 |

given,

$$x = \langle \text{sunny, mild, high, true} \rangle$$

$$P(x | \text{No}) = \frac{3}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{3}{5} = \frac{72}{625}$$

$$P(x | \text{Yes}) = \frac{2}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{3}{9} = \frac{72}{6561}$$

Now,

$$P(x | \text{No}) \cdot P(\text{No}) = \frac{72}{625} \times \frac{5}{14} = 0.041$$

$$P(x | \text{Yes}) \cdot P(\text{Yes}) = \frac{72}{6561} \times \frac{9}{14} = 0.007$$

convert in probabilities,

$$P(\text{No} | x) = \frac{0.04}{0.04 + 0.007} = 0.85$$

$$P(\text{Yes} | x) = \frac{0.007}{0.04 + 0.007} = 0.15$$

Answer: Not play. classified as No.