

23/9/21

Class I (Online)

* What is ^{data} mining?

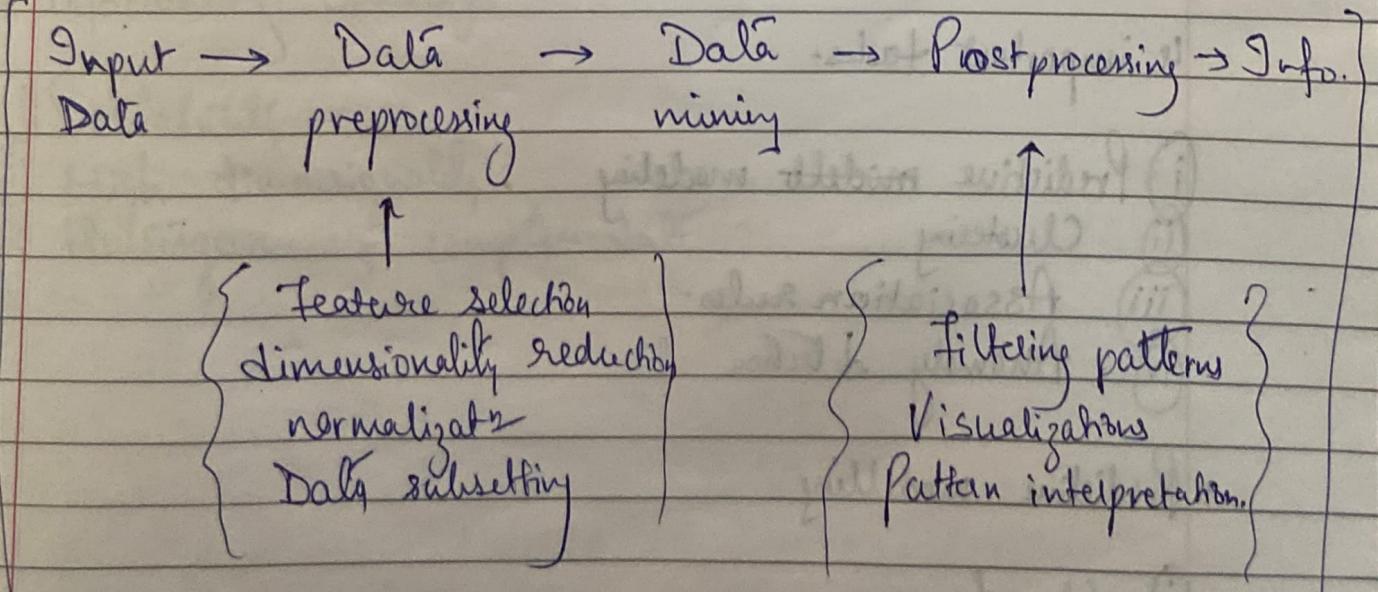
→ Technology that blends traditional methods with modern algorithms to process large amount of data.

→ Viewpoints:

- i) Commercial
- ii) Scientific
- iii)

→ Process of automatically retrieving useful information from large set of data.

→ Knowledge discovery in databases (KDD)



* Origin of DM.

- MLAT
- Statistics
- pattern recognition
- Database sys.
- Hypothesis
- Sampling

* Tasks

<p>Exploratory</p>	<p>independent variable</p>
<p>Predictive Tasks</p>	<p>Dependent variable (what is predicted)</p>

(i) Descriptive tasks → target variable value is predicted (correlations, trends, clusters, trajectories, anomalies)

Summarize relationships in data

↳ important tasks:

- i) Predictive modeling
- ii) Clustering
- iii) Association rules.
- iv) Anomaly detection !!

* Predictive modeling

- i) Classification
- ii) Regression.

① Classification

→ Target variable should take discrete values

→ Example in ppt.

② Regression

→ Predict value of given continuous valued variable based on the value of other variables, assuming linear or non-linear model of dependency.

→ Clustering
→ Association Rule Discovery
→ Dimensionality reduction

* Motivating challenges

- i) Scalability
- ii) High dimensionality
- iii) Heterogeneous & complex data

Online

Chpt - 2.

24 | 9/2021

$$\text{new value} = f(b) -$$

- Attributes of Obj
 - types of data
 - Data quality
 - Noise
 - Missing
 - Similarity of instance
 - Data preprocessing

→ Data (what is data?)

→ collection of data objects & their attributes

六

property or characteristic
of an object. A collection
of attributes describes an

→ Aktiv variable, fiktiv,
charakteristisch, dimension,

15

→ Attribute values: Numbers or symbols assigned to an attribute for a particular object.

→ Measurement of length
→ That day, we measure an attribute may not
match the attribute properties

→ Types of attributes

- ① Nominal ($=, \neq$) , any permutational value
 ② Ordinal ($<, >$) , an order preserving change of values
 ③ Interval (+, -) , new_value = $a * \text{old_value} + b$
 Ratio (*, /) , new_value = $a * \text{old_value}$.

→ Definition of attribute types is cumulative.

→ Disease of continuous attention

⇒ es gewisse Attribute

→ Types of dataset

\Rightarrow Dimensionalität

\Rightarrow Sparsity

→ C.
ige.

→ Record data = Data that consists of a collection of records each of which consists of fixed set of attributes.

Data Preprocessing

Some aggreq functors: min, max, sum, avg.

→ Data Aggregation =

→ Company XYZ

Year	Sales
2018	\$100,000
2019	\$200,000
2020	\$300,000
Total	\$600,000

Quarter	Sales
Q1	\$100,000
Q2	\$150,000
Q3	\$200,000
Q4	\$250,000

Year	Sales
2018	\$10,000
2019	\$50,000
2020	\$300,000
Total	\$360,000

→ (Combining 2 or more objects) attributes into a single object attribute.

→ (using storage, lesser processing time, lesser details).

→ allows use of more expensive data mining algorithms

→ allows change of scope | scale by providing high level view of data

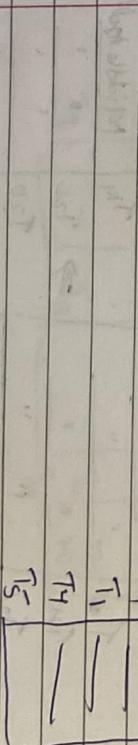
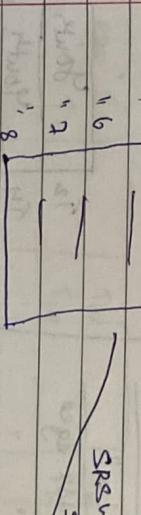
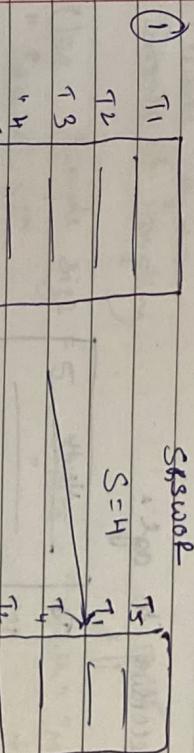
→ Statistical observation - aggregated data is more stable than that of individual objects | attributes. (less variability)

→ Disadvantages: Key details may be lost in aggregation process

* Sampling

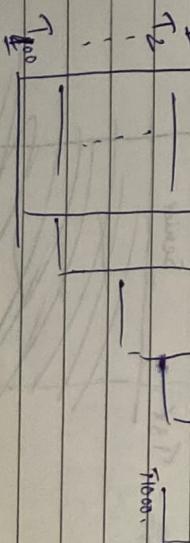
- (1) Simple random sample without replacement (SRSWR)
 (2) Simple random " with replacement (SRSWR).

Examples:



(3) Cluster sample.

→



Q) For the following data, for the attribute age = 13, 15, 16, 19,
20, 21, 22, 22, 23, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 40,

46, 52, 50

Sketch examples of each of the following sampling techniques

(i) Stratified Sampling:

according to age.

T_{10}	Young
T_{10}	"
T_{10}	"
T_{10}	Middle aged
T_{12}	Young
T_{12}	"
T_{12}	"
T_{12}	Middle aged
T_{20}	Young
T_{20}	"
T_{20}	"
T_{20}	Middle aged
T_{21}	Young
T_{21}	"
T_{21}	"
T_{21}	Middle aged
T_{22}	Young
T_{22}	"
T_{22}	"
T_{22}	Middle aged
T_{25}	Young
T_{25}	"
T_{25}	"
T_{25}	Middle aged
T_{30}	Young
T_{30}	"
T_{30}	"
T_{30}	Middle aged
T_{32}	Young
T_{32}	"
T_{32}	"
T_{32}	Middle aged
T_{35}	Young
T_{35}	"
T_{35}	"
T_{35}	Middle aged
T_{36}	Young
T_{36}	"
T_{36}	"
T_{36}	Middle aged
T_{40}	Young
T_{40}	"
T_{40}	"
T_{40}	Middle aged
T_{50}	Young
T_{50}	"
T_{50}	"
T_{50}	Middle aged
T_{52}	Young
T_{52}	"
T_{52}	"
T_{52}	Middle aged

- ① SRSWOR
- ② SRSWR
- ③ Cluster Sampling
- ④ Stratified Sampling

↑
12-25
26-50

Use Sample size = 5, Strata = "Young", "Middle Aged", "Senior" \rightarrow 60

↓

T_1	T_2	T_3	T_4	T_5
T_{10}	T_{10}	T_{10}	T_{10}	T_{10}
T_{12}	T_{12}	T_{12}	T_{12}	T_{12}
T_{20}	T_{20}	T_{20}	T_{20}	T_{20}
T_{21}	T_{21}	T_{21}	T_{21}	T_{21}
T_{22}	T_{22}	T_{22}	T_{22}	T_{22}
T_{25}	T_{25}	T_{25}	T_{25}	T_{25}
T_{30}	T_{30}	T_{30}	T_{30}	T_{30}
T_{32}	T_{32}	T_{32}	T_{32}	T_{32}
T_{35}	T_{35}	T_{35}	T_{35}	T_{35}
T_{36}	T_{36}	T_{36}	T_{36}	T_{36}
T_{40}	T_{40}	T_{40}	T_{40}	T_{40}
T_{50}	T_{50}	T_{50}	T_{50}	T_{50}
T_{52}	T_{52}	T_{52}	T_{52}	T_{52}

↓

13, 15, 16 22, 25, 25
16, 19, 20 25, 25, 30
20, 21, 22 3

Outline

30/9/21

- Q) Suppose that the data for analysis includes attribute 'age'. The age values for the data tuples are

(in increasing order): 13, 15, 16, 16, 19, 19, 20, 20, 21, 22, 22

36, 40, 45, 46, 46, 46, 46, 46

$$\begin{aligned} \text{Mean: } & 13, 15, 16 \Rightarrow 13+15+16/3 = 14.6 \approx 15 \\ & 16, 19, 20 \Rightarrow = 18 \\ & 20, 21, 22 = 21 \\ & 18, 18, 25 = 24 \\ & 25, 25, 25 = 27 \\ & 25, 25, 30 = 27 \\ & 25, 25, 33 = 33 \\ & 25, 33, 35 = 34 \\ & 35, 35, 36 = 35 \approx 35-3 \\ & 40, 45, 46 = 43 \approx 43-6 \\ & 46, 46, 46 = 46 \approx 46-6 \end{aligned}$$

- Q) Use Smoothing Avg bin means to smoother above data using bin depth of 3

Comment on effect of this techniques for given data.

$$\Rightarrow \text{Mean} = 13 + 15 + 16 + 16 + 19 + 20 + 20 + 21 + 22 + 22 + 25 + 25 + 25 + 30 + 33 + 33 + 33 + 35 + 35 + 36 + 40 + 40 + 40$$

$$13 + 16 + 16 + 16 + 19 + 20 + 20 + 20 + 21 + 22 + 22 + 25 + 25 + 25 + 30 + 33 + 33 + 33 + 35 + 35 + 36 + 40 + 40 + 40$$

$$28 \cdot 12 \Rightarrow 28$$

$$W = \frac{\text{max} - \text{min}}{\text{no. of bins}}$$

$$\text{No. of bins} = \frac{28-6}{3}$$

$$= 9.3$$

=

- C) What other methods are there for data smoothing?

10/21

* Data transformation:

- Data is transformed or consolidated into forms appropriate for mining.

Strategies

- ① Smoothing : works to remove noise from data.
 - Binning
 - Regression
 - Clustering
- ② Attribute construction: (feature construction)
 - Min - max normalization
 - 3 methods:
 - ① Min - max normalization
 - ② This performs linear transformation on original data
 - ③ Discretization
- ④ Normalization: attribute data are scaled so as to fall within smaller range such as -1.0 to 1.0 or 0.0 to 1.0
- ⑤ Discretization: raw values of numeric attribute are replaced by interval labels (e.g., 0-10, or categorical labels (e.g.: youth, adult, etc.))
 - ∴ The labels can be recursively organized into higher level concept resulting in concept hierarchy for numeric attribute.

(6)

- Concept hierarchy generation for nominal data -
- attributes such as street can be generalized to higher level concept like country or city
 - Delta transformation like normalization.

7

- Concept hierarchy generation for nominal data -
- attributes such as street can be generalized to higher level concept like country or city
 - Delta transformation like normalization.

8

- Concept hierarchy generation for nominal data -
- delta transformation like normalization.
 - It helps avoid dependence on choice of measurement units.
 - It gives all attributes equal weight.
 - Useful for classification algorithms & clustering algorithms
 - Many methods:
 - 3 methods:
 - ① Min - max normalization
 - ② This performs linear transformation on original data
 - ③ Discretization

9

- Concept hierarchy generation for nominal data -
- Suppose min_A and max_A are min & max values of an attribute -
 - Min - max normalization maps a value v_i of A to v'_i in the range [new-min, new-max] by computing

$$v'_i = \frac{v_i - \text{min}_A}{\text{new-max} - \text{new-min}}$$

∴ New-min & New-max

Q)

Income has min 12,000

has max 98,000

→ We would like to map income in the range

[0.0, 1.0]

★ V_i

What will 73,600 as income transform to?

Using min-max normalization

~~new-min $_A$ = 0.0~~

new-max $_A$ = 1.0

$$V_i' = \frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0.0) + 0$$

(3) Normalization by decimal scaling

→ Normalizes by moving decimal point of values of attribute A

→ No. of decimal points moved depends on max.

→ A value V_i of A is normalized to V_i' by

computing $V_i' = \frac{V_i - \bar{A}}{\sigma_A}$ where $j \rightarrow$ smallest integer such that $|V_i'| < 1$

⇒ Value V_i of A is normalized to V_i' by computing $V_i' = \frac{V_i - \bar{A}}{\sigma_A}$

$$V_i' = \frac{V_i - \bar{A}}{\sigma_A}$$

~~(std. deviation of A)~~

Average (\bar{A}) = $\frac{1}{n} (V_1 + V_2 + V_3 + \dots)$

$$\sigma_A = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (V_i - \bar{V})^2}$$

Q)

Suppose mean and std. dev. of values for income are 54,000 and 16,000, what is the values for the income 73,600 using Z-score norm.

$$Z_{260} V_i' = \frac{73,600 - 54,000}{16,000} = 1.225$$

Q) Suppose that recorded values of attribute A range from -986 to 917. The max. absolute value of A is

from -986 to 917. The max. absolute value of A is 986.

To normalize by decimal scaling, we divide each value by 1800 (i.e. $j=3$) so that -986 normalizes to -0.986 and 917 normalizes to 0.917

Q] Normalize the following group of data using

① Min-max normalization by setting min=0 & max=1

② Z-score normalization.

$$\text{Data} = 300, 300, 400, 600, 1000$$

$$= \frac{300 - 300}{282.84} = -0.707$$

$$\downarrow \min(300) = 0 \quad \max(600)$$

$$= \frac{400 - 300}{282.84} = +0.353$$

$$\text{Q}_1 = \frac{-200}{1000 - 200} (1.0 - 0.0) + 0$$

$$= \frac{600 - 500}{282.84} = 0.353$$

* Q1 Discretization & Binarization

\Rightarrow Transform a continuous attribute into a categorical attribute (discretization)

\Rightarrow Both continuous & discrete attributes may need to be transformed into one or more binary attributes

(binarization)

$$\sigma_A = \sqrt{\frac{1}{5} ((200 - 500)^2 + (300 - 500)^2 + (400 - 500)^2 + (600 - 500)^2 + (1000 - 500)^2)}$$

$$= 282.84$$

\Rightarrow Technique for binarization of categorical attributes

\Rightarrow If there are n categorical values, then uniquely assign each original value to an integer m in the interval $[0, m-1]$.

\Rightarrow Convert each of these m integers to a binary number.

\Rightarrow Since $n = \lceil \log_2 m \rceil$ binary digits are needed to represent the integers. The binary nos. are represented in binary form.

(1) Categorical variable with 5 values
↳: 1, 2, 3, 4, 5

We require 3 binary variables n_1, n_2, n_3

Categorical value	Skewness	n_1	n_2	n_3	n_4
Walde	0.02	0.001			

awful	0	0	0
poor	1	0	0
OK	2	0	0
good	3	0	0
great	4	0	0
	14	12	10
	100%	83%	71%

⇒ ~~Association~~ convert a categorical attribute to fine
asymmetric binary attributes ~~using~~

categorical value $\text{integer } n_1, n_2, n_3, n_4, n_5$

awful	0	0	1	0	0	0
poor	1	0	1	0	0	0
OK	2	0	0	1	0	0
good	3	0	0	0	1	0
great	4	0	0	0	0	1

Median: Middle value of ordered set or the no. of values in this set.

- Q1) Suppose that data for analysis includes the attribute age.
 \Rightarrow Age values for the data tuples are
 (increasing order) 81

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 22, 25, 25, 25, 25,	30, 30, 33, 33,
25, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.	

↓
median.

Q2) What is the mean of the data?

Q3) What is the median?

Q4) What is the mode of the data?

Q5) Comment about data's modality (bimodal, trimodal, etc.)

Q6) What is the mid range of the data

Q7) Can you find roughly the first quartile (Q_1) & 3rd quartile (Q_3) of the data?

Q8) Give the five-number summary of these data.

② Mode : This data set has two values that occur with the same highest freq. So is therefore bimodal.

$$\text{③ Midrange} = \frac{15 + 25}{2} = 20$$

$$\text{④ Median} = \frac{13 + 20}{2} = 16.5$$

Avg of smallest & largest value in the dataset

$$\frac{1}{4}(N+1)^{\text{th}} \text{ term} \Rightarrow 28 \Rightarrow 13, 15, 16, 16, 19, 20, 20$$

$\frac{3}{4}(N+1)^{\text{th}}$ term

$$\frac{3}{4}(N+1)^{\text{th}} \text{ term} \Rightarrow 21 \Rightarrow 35, 36$$

first corr. to 2st. percentile is 20
third corr. " 25, percentile is 35 .

5 - Nos sum + \downarrow
min value, first Q, median, third Q, max value
 $15, 20, 25, 35, 40$

Q) Suppose a hospital tested the body fat data for 18 randomly selected adults with following results.

age	23	23	27	27	31	31	39	41	41	47	49	50	52	54	54
1.15	23	23	27	27	31	31	39	41	41	47	49	50	52	54	54
1.16	26.5	26.5	27.8	27.8	31.4	31.4	35.9	37.4	37.4	41.2	43.1	43.6	44.5	44.5	44.5
1.17	30.2	30.2	31.1	31.1	32.9	32.9	41.2	41.2	41.2	45.7	45.7	45.7	45.7	45.7	45.7
1.18	35.4	35.4	36.2	36.2	37.1	37.1	38.9	38.9	38.9	40.9	40.9	40.9	41.2	41.2	41.2

Calculate

- Mean, median, std. deviation of age & fat
- Normalize the two variables based on Z-score

$$\text{std. deviation (age)} \Rightarrow \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\text{Mean (age)} = \frac{36}{18} = 2.0 \quad \text{Median} = \frac{30+32}{2} = 31$$

$$\text{Mean (fat)} = \frac{51.8}{18} = 2.88 \quad \text{Median} = 32.9$$

Ans : 1.03 hours. In which 0.01 hr wholly lost in water. So, Ans

DOMS | Page No.

D.O.H.S	Page No.
----------------	----------

$$d_i = \sqrt{(x_i - \bar{x})^2}$$

23	-23.44	509.43
23	-23.44	509.43

29	-2.144	2.144
27	-1.944	1.944
27	-1.944	1.944

H7	0.56	0.31	0.44	-0.54	0.59
			41	51	52

149 2.50 6.50
150 2.50 6.50
151 2.50 6.50

5-2	5-36	30-91
54	7-56	54-15
7-56	54-15	

57	9.56	91.39
57	16.56	111.51

58 11.52 133.63
58 11.56 133.63
58 11.57 133.62

61 14.56 211.99

卷之三

1960 - 1961

P.B.D. - 100324T 8E 95 - 12

(R, S, K) 2

卷之三

200 w/ 25 new work at 3 min 202-23 231 (i)

1980-81 - 13

288-02

2000 2000 2000 2000

228-1 28-2
01

(i) Using data for age, answer the foll.

⇒ i) Use min-max norm. to transform value 85

for age onto the range [0.0, 1.0]

$$D_i' = \frac{(35 - 13)}{70 - 13} \times (1.0 - 0.0) + 0.0$$

$$= \frac{22}{57} \times 1$$

$$= 0.38$$

(ii) Use Z-score norm. to transform 35 for age
where std. dev. of age is 12.91

$$D_i' = \frac{35 - 30}{12.91}$$

$$= 0.386$$

(iii) Use normalization by decimal scaling to transform
value 85 for age.

$$D_i' = \frac{V}{10^j}$$

$$= \frac{35}{10^2} = 0.35$$

(iv) Comment on which method you would prefer to use
giving reasons as to why?

(i) What are the value ranges for foll. normalizing methods

new min new max

a) Min-max :- [0.0, 1.0]

b) Z-score :- $\frac{\text{Old min-mean}}{\text{std. devi}} , \frac{\text{old max-mean}}{\text{std. devi}}$

c) Norm. by decimal scaling :- (-1.0, 1.0)

11/10/21

* Similarity & dissimilarity measures (1) and (2)

- (1) Similarity measure
→ Numerical measure of how alike

1.0 to 0.0

∞ to ∞

* Euclidean Distance

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

$$P_1(3,2), P_3(3,1)$$

= 2.

$$\begin{aligned} d(P_1, P_2) &= \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \\ &= \sqrt{8} = 2\sqrt{2} = 2.828 \end{aligned}$$

$$P_1(0,2), P_2(0,0)$$

* Minkowski Distance

⇒ Generalization of Euclidean dist.

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

$$\boxed{\begin{aligned} x &= 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ y &= 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \end{aligned}}$$

$r=1 \rightarrow$ City block (Manhattan, taxicab, L₁ norm) dist

$r=d \rightarrow$ "Euclidean"

$r \rightarrow \infty \rightarrow$ Supremum, (L_∞ norm, L_∞ norm).

→ Max diff. w.r.t. any components of vector

$$\begin{aligned} d(P_2, P_3) &= \sqrt{[(2-3)^2 + (0-1)^2]} \\ &= \sqrt{1+1} = \sqrt{2} = 1.414 \end{aligned}$$

$$\text{L}_\infty \text{ norm: } P_3(3,1), P_4(5,1)$$

$$d(P_3, P_4) = \max \left[|(3-5)|, |(1-1)| \right]$$

$$d(P_1, P_3) = \max \left[|(3-0)|, |(2-1)| \right] = 3.$$

* Simple matching & Jaccard coeff.

SIMC = number of matches / number of attributes

$$= (f_{11} + f_{00}) / (f_{01} + f_{10} + f_{11} + f_{00})$$

$$\text{Jaccard} = (f_{11}) / (f_{01} + f_{10} + f_{11})$$

$$f_{01} = 2$$

$$f_{10} = 1$$

$$f_{00} = 7$$

$$f_{11} = 0$$

$$S_{MC} = \frac{(f_{11} + f_{00})}{(= 0 + 1)}$$

$$f_{21} + f_{10} + f_{11} + f_{20} = 2 + 1 + 0 + 1$$

$$T = \frac{f_1 + f_{10} + f_{11}}{1+1+0} = \frac{0}{2}, 0$$

~~18 | 10 | 21~~

⇒ ~~Document~~

Document team coach hockey basketball soccer

Doc 1 57803002

DOC 3

Doc 4

\Rightarrow let $x \in V$ lie the vector for comparison.

Using cosine similarity, $\text{cosine sim}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$

where $\| \cdot \|$ is the Euclidean norm of vector n .

defined by $\sqrt{x_1^2 + x_2^2 + \dots + x_p^2}$

Conceptually, it is the length of vector

$\|y\|$ is the Euclidean norm of vector y .

$$x = (5, 1, 0, 3, 1, 0, 1, 2, 1, 0, 1, 0, 1, 2, 0, 1, 0, 0)$$

سیمی = سیمی

$$\text{N.Y.} = 6500 + 0 + 6 + 0 + 2 + 0 + 0 + 2 + 0 + 0 = 6224$$

$$||y_1|| = \sqrt{1^2} (1+4+4) = \sqrt{13}$$

Cosine simil. doc1, doc4)

$$\Rightarrow \text{vec} \cdot \text{vec} = 2.$$

$$\|\text{vec}\| = \sqrt{2}$$

$$\|\text{vec}\| = \sqrt{2}.$$

$$\text{simil}(u, v) = 0.6748$$

* Dissimilarity for attributes of mixed type

Object	Test-1	Test-2	Test-3
id (nominal)	(nominal)	(decimal)	(measure)
1 code A	excellent	4.5	
2 code B	fail	2.2	
3 code C	good	6.4	
4 code A	excellent	2.8	

② The contribution of attribute f to the dissimilarity between objects i & j ($d_{ij}^{(f)}$) is completely depending on its type

$$① \text{ If } f \text{ is numeric, } d_{ij}^{(f)} = \sqrt{|x_{if} - x_{jf}|^2}$$

$$\text{max}_{1 \leq f \leq n} |x_{if} - \min_{1 \leq f \leq n} x_{if}|,$$

runs over all non-missing objects of attribute f.

⇒ Suppose the dataset contains 'p' attributes of mixed type, the dissimilarity d_{ij} b/w objects i & j is defined as

$$d_{ij} = \sum_{f=1}^p d_{ij}^{(f)}$$

② If f is nominal or binary, $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, otherwise $d_{ij}^{(f)} = 1$.

If f is ordinal, compute the ranks π_f if $z_{if} = \pi_{if} - 1$

$$\frac{\sum_{j=1}^n S_{ij}}{N_f - 1}$$

where indicator $S_{ij} = 1$ if $z_{if} < z_{jf}$ else 0 if either

③ If x_{if} or x_{jf} is missing (there is no measurement of attribute f for object i or object j)

or

if $x_{if} = x_{jf} = 0$, if the attribute f is quantitative binary, otherwise $d_{ij}^{(f)} = 1$

f1

f2

f3

f4

f5

f6

f7

f8

f9

f10

f11

f12

f13

f14

f15

f16

f17

f18

f19

f20

f21

f22

f23

f24

f25

f26

f27

f28

f29

f30

f31

f32

f33

f34

f35

f36

f37

f38

f39

f40

f41

f42

f43

f44

f45

f46

f47

f48

f49

f50

f51

f52

f53

f54

f55

f56

f57

f58

f59

f60

f61

f62

f63

f64

f65

f66

f67

f68

f69

f70

f71

f72

f73

f74

f75

f76

f77

f78

f79

f80

f81

f82

f83

f84

f85

f86

f87

f88

f89

f90

f91

f92

f93

f94

f95

f96

f97

f98

f99

f100

f101

f102

f103

f104

f105

f106

f107

f108

f109

f110

f111

f112

f113

f114

f115

f116

f117

f118

f119

f120

f121

f122

f123

f124

f125

f126

f127

f128

f129

f130

f131

f132

f133

f134

f135

f136

f137

f138

f139

f140

f141

f142

f143

f144

f145

f146

f147

f148

f149

f150

f151

f152

f153

f154

f155

f156

f157

f158

f159

f160

f161

f162

f163

f164

f165

f166

f167

f168

f169

f170

f171

f172

f173

f174

f175

f176

f177

f178

f179

f180

f181

f182

f183

f184

f185

f186

f187

f188

f189

f190

f191

f192

f193

f194

f195

f196

f197

f198

f199

f200

f201

f202

f203

f204

f205

f206

f207

f208

f209

f210

f211

f212

f213

f214

f215

f216

f217

f218

f219

f220

f221

f222

f223

f224

f225

f226

f227

f228

f229

f230

f231

f232

f233

f234

f235

f236

f237

f238

f239

f240

f241

f242

f243

f244

f245

f246

f247

f248

f249

f250

f251

f252

f253

f254

f255

f256

f257

f258

f259

f260

f261

f262

f263

f264

f265

f266

f267

f268

f269

f270

f271

f272

f273

f274

f275

f276

f277

f278

f279

f280

f281

f282

f283

f284

f285

f286

f287

Cosine sinil doc1, doc4)

$$\Rightarrow \text{cycle} \quad n \cdot y = 2.$$

$$||\eta_2|| = \sqrt{\mu_2}$$

$$11411 = \text{Fig.}$$

$$\sin(\alpha_{ij}) = 0.6748$$

* Disimilants for attributes of mixed type

* Disinhibiting for attributes of mixed type

	Test - 1	Test - 2	Test - 3
Code A	(Nominal)	(Decimal)	(Inches)
Code A	Excellent	4.5	0.115 (0.005)
Code B	Fair	2.2	0.055 (0.005)
Code C	Good	6.4	0.165 (0.005)
Code D	Excellent	2.8	0.075 (0.005)

⇒ Suppose the dataset contains ' p ' attributes of mixed types, the dimensionality of $\text{J}(i,j)$ b/w objects i, j is defined as

$$d_{ij} = \sum_{f=1}^F \alpha_{ij}^{(f)} d_{ij}^{(f)}$$

③ If f is ordered, compute the raw
heat ~~step~~ as if no number -

runs over all non-missing objects of attribute f .

max $r_{ij} - \min r_{ij}$,

(2) If f is nominal or binary, $d_{ij}^{(f)} = 0$ if $r_{ij} = r_{jf}$,
 otherwise $d_{ij}^{(f)} = 1$

sums over all non-missing objects of attribute f .

where

(3) If f is ordinal, compute the ranks $r_{if} \leq Z_{if} = r_{if} - 1$
 $\cdot \frac{N_f - 1}{2}$

① If f is numeric, $d_{ij}^{(f)} = \frac{x_{if} - \bar{x}_{if}}{\max_{in} x_{if} - \min_{in} x_{if}}$, then sum over all non-missing objects of attribute f .

② If f is nominal or binary, $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, otherwise $d_{ij}^{(f)} = 1$.

③ If f is ordinal, compute the ranks $x_{if} \in Z_{if} = \{1, \dots, M_f\}$.

where indicator $i_{ijf} = 1$ if $x_{ijf} \neq x_{jif}$ and $i_{ijf} = 0$ if either $x_{ijf} = x_{jif}$ or x_{ijf} is missing (there is no measurement of attribute f for object i or object j)

② If $x_{ijf} = x_{jif} = 0$, & the attribute f is quantitative binary, otherwise $i_{ijf} = 1$

③ The contribution of attribute f to the dissimilarity below is $d_{ij}^{(f)}$ & is completely depending on its type

- ① If f is numerical, $d_{ij}^{(f)} = \frac{|x_{ijf} - x_{jif}|}{\max_n x_{nf} - \min_n x_{nf}}$, where runs over all non-missing objects of attribute f .
- ② If f is nominal or binary, $d_{ij}^{(f)} = 0$ if $x_{ijf} = x_{jif}$, otherwise $d_{ij}^{(f)} = 1$
- ③ If f is ordinal, compute the ranks $x_{ijf} \leq z_{ijf} = x_{if} - \frac{N_f - 1}{2}$ and $d_{ijf} = z_{ijf}$ as numerical.

Test 3	1	2	3	4
1	0			
2	0.55	0		
3	0.55	1	0	
4	+0.49	0.14	0.85	0

$$d(3,4) = \frac{64 - 22}{64 - 22}$$

$$d(3,4) = 64 - 22$$

$$(h_1,1) =$$

$$(h_1,2)$$

$$(h_1,3)$$

Dissimilarity

Test 1	1	2	3	4
1	0			
2	1	0		
3	0.55	0.55	0	
4	0	1	0.55	0

fair good excellent

$$0.13 \quad 0.71 \quad 0.986 \quad 0$$

↓ most similar

$$d(3,4) = \frac{1 \times 0.45 + 1 \times 0.55 + 1 \times 1}{3} = 0.65$$

$$d(2,1) = \frac{1 \times 0.55 + 1 \times 0.55 + 1 \times 1}{3} = 0.85$$

$$d(2,1) = \frac{1 - 0}{2}$$

$$(h_{1,1}) = \frac{1 - 0}{2}$$

$$(h_{1,2}) = \frac{2 - 0}{2} = 1$$

$$(h_{1,3}) = \frac{2 - 1}{2} = 1$$

$$d(h_1,3) = 0.786$$

Date

Signature

Correlation analysis

Date

⇒ Measures how strongly one attribute implies the other based on available data

⇒ Chi-Square stat: χ^2 : for only nominal data

⇒ Consider two attributes $A \& B$

Attribute $A = \{a_1, a_2, a_3, a_4, \dots, a_r\}$
 $B = \{b_1, b_2, b_3, \dots, b_s\}$

⇒ Data tuples described by $A \& B$ can be shown as Contingency table where

Let (A_i, B_j) denote joint event that attribute A takes on value a_i & attribute B takes

on value B_j

of value (A_i, B_j) Pearson's χ^2 statistic

is computed as

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

observed freq.

Calculated E_{ij}

		male	female	Total
Total	fiction	280 (40)	200 (30)	480
	non-fiction	50	100	150
	Total	330	1200	1530

Where O_{ij} is the factual count of joint events

(A_i, B_j) and e_{ij} is the expected freq. of (A_i, B_j) which can be computed as

$$e_{ij} = \frac{\text{Count}(A=a_i) \times \text{Count}(B=b_j)}{n}$$

where, $n \rightarrow$ no. of data tuples.

count ($A=a_i$) is the no. of tuples having values a_i for A, count ($B=b_j$) is the no. of tuples having b_j for B.

The χ^2 statistic tests the null hypothesis that A & B are independent i.e. no correlation b/w them. Test is based on significance level with $(r-1) \times (s-1)$ degrees of freedom.

Eg: Suppose that group of 1500 ppl was surveyed gender of each person was noted. Each person was polled to find out whether preferred type of reading material. (fiction or non-fiction)

⇒ The observed freq. count of each possible joint event is summarized in Contingency Table given b/w

$$e_{11} = \frac{\text{Count}(A=fiction) \times \text{Count}(B=fiction)}{n} = \frac{300 \times 480}{1500} = 96$$

$$e_{12} = \frac{\text{Count}(A=fiction) \times \text{Count}(B=non-fiction)}{n} = \frac{300 \times 450}{1500} = 60$$

$$e_{\text{eff}} = \text{cover}(\text{mail}) \times \text{cover}(\text{non-fishin})$$

$$= \frac{300 \times 1050}{1800} = 175$$

$$e_{22} = \text{cover (female)} \times \text{cover (fictitious)}$$

$$= 1200 \times \frac{16}{25} = 9840$$

$$c_1 = \sqrt{\frac{(250-90)^2 + (200-360)^2 + (80-210)^2}{940}}$$

$$\sum_{i=1}^n \overline{g_i g_i}$$

R Correlation coefficient for numeric data

For numeric attributes we can evaluate correlation b/w two attributes A & B by computing correlation coefficient (Pearson's product moment coeff.)

Where n is the number of tuples
 a_i & b_i are respective values of A_i & B in tuple i
 \bar{A}_B are respective mean values of A_i & B
 σ_A & σ_B are respective std. deviations of A_i & B

$\sum (a_i b_i)$ is the sum of AB base product
 $- 1 \leq r_A < +1$

If $r_{A,B} > 0$, A & B are positively correlated
when value of A ↑, B also ↑

∴ For 1 degree of freedom, χ^2 value needed to reject the hypothesis at 0.001 significance level is 16.828 (χ^2 distribution)

Since χ^2 is much higher than significance level, we reject ~~that~~ the hypothesis that gender $\{\leftarrow$ less preferred reading in independent, and it is highly correlated.

\Rightarrow Since F^2 is much higher than significance level, we reject ~~that~~ the hypothesis that gender $\{\leftarrow$ preferred reading is independent, and thus highly correlated.

* Covariance of numeric data

$$\text{Cov}(A, B) = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B}$$

Correlation and covariance are two similar measures for assessing how much 2 attributes change together. Consider two numeric attributes $A \in B$ and set of observations $\{(a_1, b_1), \dots, (a_n, b_n)\}$

\Rightarrow The mean values of $A \in B$ respectively are also known as expected values on $A \in B$ i.e $E(A) = \bar{A} = \sum_{i=1}^n a_i$

\Rightarrow If $A \in B$ are independent, $\text{Cov}(A, B) = E(A \cdot B) - \bar{A} \bar{B} = E(A) \cdot E(B) - \bar{A} \bar{B} = 0$

If $A \in B$ are not independent, then $\text{Cov}(A, B) = E(A \cdot B) - \bar{A} \bar{B}$

If $A \in B$ are positively correlated, then $\text{Cov}(A, B) > 0$

If $A \in B$ are negatively correlated, then $\text{Cov}(A, B) < 0$

If $A \in B$ are uncorrelated, then $\text{Cov}(A, B) = 0$

Consider the table given, which presents stock prices observed at five time points for AllElectronics of stocks as affected by game industry trends, will their prices rise or fall together?

Covariance b/w $A \in B$ is defined as

$$\begin{aligned} \text{Cov}(A, B) &= E((A - \bar{A})(B - \bar{B})) \\ &= \sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B}) \\ &\quad \text{Timepoints} \quad \text{AllElectronics} \quad \text{HighTech} \\ t_1 &\quad 6 \quad 20 \\ t_2 &\quad 5 \quad 10 \\ t_3 &\quad 4 \quad 14 \\ t_4 &\quad 3 \quad 5 \\ t_5 &\quad 2 \quad 5 \\ n. & \end{aligned}$$

Also, it can be shown that- $\text{Cov}(A, B) = E(A \cdot B) - \bar{A} \bar{B}$

for two attributes, $A \in B$ that tend to change together, if A is larger than \bar{A} , then B is likely to be larger than \bar{B} (covariance of $A \in B$ is positive)

Correlation

$$\text{E(Auselectronics)} = \frac{6+5+4+3+2}{5} = \$4$$

$$\text{E(High Tech)} = \frac{20+10+14+5+5}{5} = \$10$$

$$\text{E(Female)} = \frac{30+14+10}{5} = \$10$$

$$\text{E(Male)} = \frac{54}{5} = \$10.80$$

\Rightarrow Covariance of (Auselectronics, High Tech)

$$\text{Cov}(A, B) = E(AB) - \bar{A}\bar{B}$$

$$= 6 \times 20 + 5 \times 10 + 14 \times 20 + 5 + 2 \times 5$$

$$= \frac{5}{5} = 10$$

$$= 4 \times 10 + 8 \times 20 + 2 \times 5$$

$$= 80.2$$

$$= 43.2$$

$$= 7$$

Is there a relationship b/w gender & an individual's level of education they have obtained?

Two nominal attributes:

χ^2 test

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{\text{Count}(A=a_i) \times \text{Count}(B=b_j)}{n}$$

$$E_{11} = \frac{60 \times 20}{395} = 50.886$$

$$E_{12} = \frac{98 \times 20}{395} = 49.86$$

$$E_{21} = \frac{99 \times 20}{395} = 50.377$$

$$E_{22} = \frac{99 \times 19}{395} = 48.13$$

$$E_{14} = \frac{98 \times 19}{395} = 49.868$$

$$E_{24} = \frac{98 \times 19}{395} = 48.132$$

Covariance is positive, we can say that stock prices for both companies move together.

Q] A random sample of 395 ppl was surveyed and each person was asked to report the highest education level they obtained. The data that resulted from the survey is summarised in the foll. table

	High School	Bachelors	Masters	Ph.D.
Female	0.1 60	0.454	0.46	0.141
Male	0.2 40	0.244	0.253	0.57
Total	100	98	99	98

$$\chi^2 = \frac{(60 - 50.886)^2}{50.886} + \frac{(54 - 49.868)^2}{49.868}$$

$$+ \frac{(46 - 50.372)^2}{50.372} + \frac{(41 - 49.868)^2}{49.868}$$

$$+ \frac{(40 - 49.114)^2}{49.114} + \frac{(44 - 48.132)^2}{48.132}$$

$$+ \frac{(53 - 48.623)^2}{48.623} + \frac{(54 - 48.132)^2}{48.132}$$

$$= 1.632 + 0.342 + 0.380 + 1.576 + 1.691$$

$$+ 0.354 + 0.394 + 1.633$$

$$= 8.002 \quad (\approx 0.002)$$

Null hypothesis: Two attributes, gender & education level are independent.

$$\text{D.O.F.} = (r-1)(c-1) = \binom{4}{1}(4-1) = 12$$

$$= (1)(3) = 3.$$

Since, For 3 D.O.F at 5% significance level, the χ^2 value is 7.815.

The χ^2 obtained is much higher than significance level, we reject the null hypothesis that gender & education level are independent & are highly correlated, and has an impact on education level of an individual.

② Using the data for age & body fat.

Calculate the correlation coefficient (Pearson's product moment coeff.). Are these two attributes positively or negatively correlated? Compute their covariance

$$r_{A,B} = \frac{\text{Cov}(A, B)}{\sqrt{A} \sqrt{B}}$$

$$\text{Cov}(A, B) = \sum_{i=1}^{n-1} (a_i - \bar{A})(b_i - \bar{B})$$

$$n = 9.$$

$$\bar{A}(\text{Age}) = 36.22$$

$$\bar{B}(\text{body fat}) = 22.44$$

$$\text{Cov}(A, B) = \frac{1}{9} \left[(23 - 36.22)(9.5 - 22.44) + (23 - 36.22)(26.5 - 22.44) + (36.22)(9.8 - 22.44) + (36.22)(31.4 - 22.44) \right]$$

$$+ (41 - 36.22)(25.9 - 22.44) + (47 - 36.22)(27.4 - 22.44)$$

$$+ (49 - 36.22)(24.2 - 22.44) + (53 - 36.22)(23.2 - 22.44)$$

$$= \frac{1}{9} (175.03 - 19.7072 + 86.206524 + 24.0948 + 15.1048 + 50.2348 + 57.9988 + 116.548)$$

$$\sigma_A = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}} = \sqrt{\frac{1}{9} \times 9} + (47 - 36.22)^2 + (41 - 36.22)^2 + (39 - 36.22)^2 + (30 - 36.22)^2 + (31 - 36.22)^2 + (29 - 36.22)^2 + (27 - 36.22)^2 + (25 - 36.22)^2$$

$$= \sqrt{\frac{1}{9} \times 9 \cdot 24} = 1.013$$

$$\sigma_B = \sqrt{\frac{1}{9} \times 9} = 1$$

$$\rho_{A,B} = \frac{\text{Cov}(A,B)}{\sigma_A \sigma_B}$$

$$\text{Cov}(A,B) = \frac{1}{9} \left[(23 - 36.22)(9.5 - 22.74) + (23 - 36.22)(26.5 - 22.74) \right.$$

$$+ (27 - 36.22)(9.5 - 22.74) + (27 - 36.22)(25.9 - 22.74) \\ + (31 - 36.22)(31.4 - 22.74) + (41 - 36.22)(25.9 - 22.74) \\ + (47 - 36.22)(27.4 - 22.74) + (49 - 36.22)(27.2 - 22.74) \\ + (50 - 36.22)(31.2 - 22.74) \right]$$

$$\rho_{A,B} = \frac{6.3 \cdot 51}{10 \cdot 6.43 \times 8.401} \approx 0.7163$$

Since, $\rho_{A,B} = 0.7163$ which is positive, the attributes age and body fat are positively correlated. As age increases, body fat also increases.

$$= 63.51$$

\Rightarrow

$$\sigma_A = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}} = \sqrt{\frac{1}{9} \times 1019.5556}$$

$$\sigma_B = \sqrt{\frac{1}{9} \times 113.28} = 10.6733$$

Offive

9/11/2014 Histogram Visualization Techniques

- * Visualizing small no. of attributes:

(i) Stem & leaf plot
→ Shows data into distribution of 1-D integer / continuous

→ Insights

data.

- Split the values into groups:
→ Each group contains those values that are the same except for the last digit
- Each group becomes the stem, while the last digit of a group are the leaves.

Age attribute of employee table :

22, 22, 23, 23, 23, 25, 28, 28, 30, 31, 31, 32, 35, 40, 40, 41,

42, 45, 45, 48, 50, 52, 54, 55

2 : 2, 2, 3, 3, 3, 5, 8, 8 → leaves
3 : 0, 1, 1, 2, 5
4 : 0, 0, 1, 2, 5, 8

Stem ↗

→ By plotting stem & leaf vertically & leaf horizontally
• visual representation of data distribution

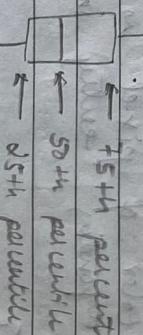
① Histogram:

- Continuous attribute
- Range of values is divided into bins of equal width usually represented by a bar. Area of bar & no. of objects / values in each bin is counted.

→ Box Bar plot is constructed, wherein each bin is represented by a box bar. Area of bar & no. of values / objects that fall into the corr. range.

* Box plot.

Min + 1.5 * IQR → upper tail (90th percentile)



- Used to illustrate linear correlation
- Each data object is plotted as a point in the plane using values of two attributes as x & y coordinate
- Attribute are either integers / real valued.

Other techniques

Star plots \rightarrow Similar to parallel co-ords; axes radiate from central pt.; line connecting values of obj in polygon

~~Class off faces~~
 \rightarrow Class off faces; approach associates each attribute with class or a point.

Star plot

~~Face~~ \rightarrow Face \rightarrow Shape of face \rightarrow Size of face
~~Face~~ \rightarrow Forehead \rightarrow Relative size of forehead
~~Face~~ \rightarrow Jaw \rightarrow Shape of jaw

- * Contour plot:
 - \rightarrow Used for 3-D data of which two attributes specify position in a plane
 - \rightarrow Third attribute (continuous attribute)
 - Eg: Temp | elevation, breaks the plane into regions with roughly the same attribute values.

Matrix plot

- array of pixels where each pixel is characterized by the color & brightness
- \rightarrow A matrix is a rectangular array of values
- \rightarrow A matrix can be visualized as an image by associating each entry of data matrix with a pixel in the image
- \rightarrow Brightness or color of pixel is determined by the corresponding entry in matrix.

Parallel 1D-ordinates

- \Rightarrow Have one co-ord. axis for each attribute; the different ones are parallel to one another instead of 1's (perpendicular)
- \Rightarrow Object is represented as line instead of point.

10.11.2021 * SIAP 3 Multidimensional data analysis

OLAP 33

~~Q&P~~ QMP

卷之三

* Represent this data as multidimensional array

→ We create a table by describing with only two attributes i.e.

petal length & petal width = low, medium, high.
then count number of flowers that have the
a particular combination of petal width, petal length
and species (class)

→ **pellet width** - $\begin{cases} \text{low} & - [0, 0.75], \\ \text{medium} & - [0.75, 1.75], \\ \text{high} & - [1.75, \infty) \end{cases}$

total length - { low - [0, 2.5)
 medium - [2.5, 5)
 high - [5, ∞) }

→ rated length within species type count

low	low	glossy	H6	+ ¹⁵ / ₁₄	Petal width	low	1.10?
low	medium	glossy	2	+ ¹ / ₉	med.	0	1.3
medium	low	glossy	2	+ ¹ / ₁₄	high	0	2.2
medium	medium	vesicular	H3	+ ¹³ / ₉	high	2	2.2
medium	high	"	3	+ ¹ / ₁₄			
"	high	Virginia	3				
high	medium	vesicular	"				
"	"	Wingville	3				
high	high	vesicular	2				
					Petal width		
				low	med.	high	
				low	0	0	
				med.	0	0	
				high	0	0	
				0	3	4	

DOMS	Page No.
/	/

DOMS	Page No.
/	/

15/11/2020

- Q] Suppose that a datawarehouse consists of 3D time, doctor, patient & two measures count & charge. where charge is the fee a doctor charges a patient for a visit

a) Enumerate 3 classes of schema that are popularly used for modelling data warehouse

b) Draw a star schema diagram for the above data warehouse.

c) Starting with the base cubeid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2020.

d) To obtain the same list, write a SQL query assuming the data is stored in RDB with 4 schema - fee (day, month, year, doctor, hospital, patient, charge, count, charge)

Ex: Star, snowflake, fact constellation schema

constellation fact table

Time dimension table



fact table

doctor

dimension table

② Draw the snowflake schema for this information:

MARCH 1960

(6) Starting with these carboid f students, course
starts in [unclear]: What specific OLR

the average grade of CS course for each operation should one perform in order to last

② Big Data. Student. Each dimension has 5 levels (including ALL).

a student L

Major L
Shaw L

How many cells will this culture contain including the dead and open cells?

Student dimension
false

Student - id
course - id
Semester - id
instructor - id

McBrown Cover Aug - grade.

Want out? You're welcome with this

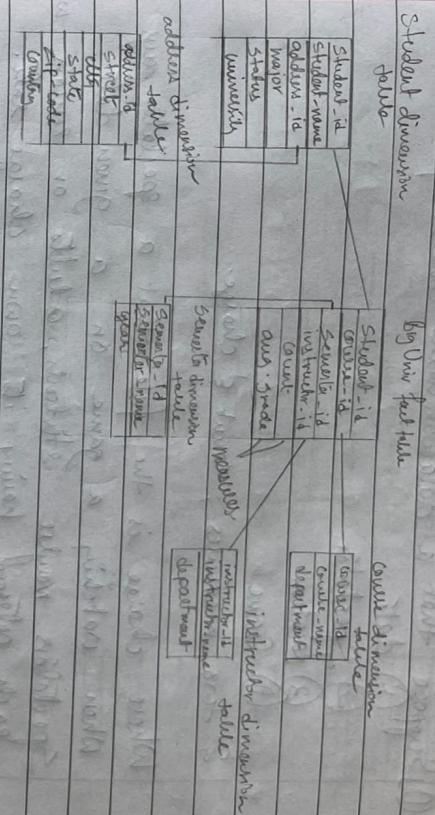
10/13 100% (100%) 100% (100%) 100% (100%) 100% (100%) 100% (100%)

(*Ex. 204*) *reinforcement* (*Ex. 205*) *reinforcement*

most common species along streams

TANAKA, YOSHIO AND KAZUO WATANABE: SHELF CLOUDS

S61m



5

Roll up on student from student-id to dept-
in Dice up on course ; dept = 'CS' & course = "Engg Univ"

c

$$\text{Number of levels} = 5$$

$$\therefore \text{Ans. of numbers} = (1+1)^4 = 6^{14}$$

\Rightarrow $\left(\begin{matrix} A \\ B \end{matrix}\right)$

No. of Culexoids = (Level) ^{No. of} dimensions.

17/11/20

Soln

Spectator dimension table

- Q) Suppose that a database consists of 4 dimensions
- (1) date
 - (2) Spectator
 - (3) location
 - (4) game

§ two measures ; count & charge.

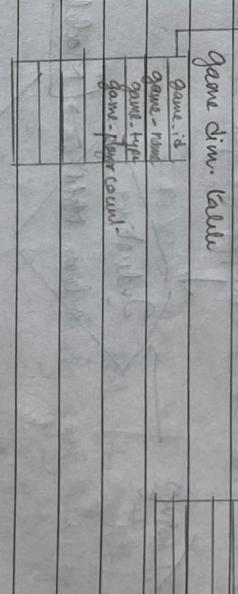
Where charge is the fee that a spectator pays

When watching a game on a given date

Spectators include students, adults or seniors with each category having its own charge rate.

- Q) Draw a star schema diagram for the database.

b) Starting with the base fact cuboid [date, spectator, location, game], what specific OLAP operations should one perform in order to list the total charge paid by student spectators at GM place in the year 2016



- Q) cuboid [date, spectator, location, game]

Q) Roll up on date from date-id to year.
 Q) Roll up on location from location-id to location-name.
 Q) Roll up on spectator from spectator-id to spectator_category = student, year = 2016 & location-name = GM place.

- Q) Roll up on game from game-id to ALL.

End of
Unit 5

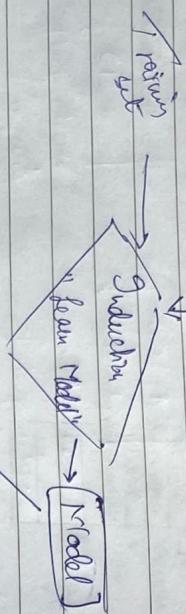
18/11/2021

Classification

→ Given a collection of records (training set); each

* General approach for building classifier model

i) Learning rule



ii) If D_t contains records that belong to same class to each other then t is leaf node labelled as y_t .

* Hoeffding algorithm
→ Let D_t be the set of tr. re. that learn a node t .

→ Gen. procedure:

i) If D_t contains records that belong to same class to each other then t is leaf node labelled as y_t .

ii) If D_t contains records " " " more than one class, use an attribute test to split data into smaller subsets. Recursively apply to procedure to each subset.

* Design issues -

- ① How should splitting training records be split?
- ② How should splitting procedure stop?

* Classification tech.

- Base classifier
- Ensemble

(1) Naive Bayes

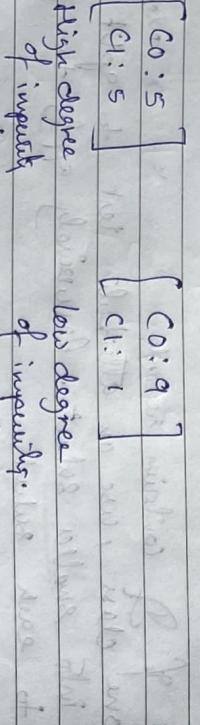
14/12/2021

* How to determine best split?

① Greedy approach

→ Nodes with pure class distribution are preferred.

→ Need a measure of node impurity.



High degree of impurity
↓
→ equally divided.

→ 3 measures of node impurity.

② Gini index

$$\text{Gini Index} = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

$p_i(t) \rightarrow \text{freq. of class } i \text{ at node } t$

c in total
no. of classes

③ Entropy

$$\text{Entropy} = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

$$\downarrow$$

$$N_{111} \quad N_{112}$$

$$\downarrow$$

$$N_{21} \quad N_{22}$$

$$\text{Classification error} = 1 - \min [p_i(t)]$$

* finding the best split

① Gini impurity measure (P) before splitting
→ Gini impurity measure of each child node

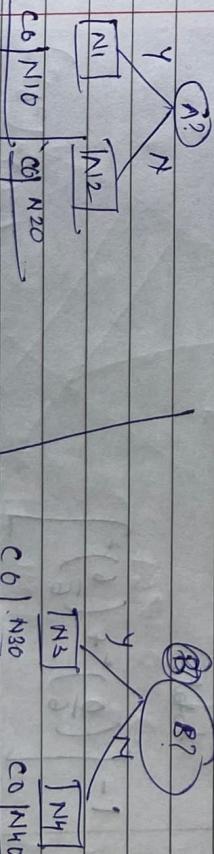
→ M is weighted impurity of the child node.

② Close attribute test condition that produces higher gain

$$\text{Gain} = P - M$$

or equivalently, lowest impurity measure after splitting.

$$\text{Before splitting} : \frac{C_0 | N_{100}}{C_1 | N_{01}} \rightarrow P$$



N_{111}

\curvearrowleft

N_{112}

\curvearrowleft

N_{21}

\curvearrowleft

N_{22}

$$\text{Gain} = P - M_1 \quad \text{vs} \quad P - M_2$$

\therefore

Max $\left[\textcircled{1}, \textcircled{2} \right]$
in chosen for splitting.

* Gain

~~choose in step~~

$$\textcircled{2} \quad 1 - \sum_{i=0}^5 p_i(i)^2$$

$$1 - \left[\left(\frac{1}{6} \right)^2 + \left(\frac{5}{6} \right)^2 \right]$$

$$\text{for 2-class prob } (P, 1-P):$$

$$G(\text{IN}) = 1 - P^2 - (1-P)^2$$

$$= 2P(1-P)$$

(1)

C1	1
C2	5

$$\textcircled{2} \quad \frac{C1}{C2} \sqrt{\frac{3}{5}}$$

$$1 - \left[\left(\frac{2}{6} \right)^2 + \left(\frac{4}{6} \right)^2 \right].$$

$$= 0.444.$$

$$G(\text{IN}) = 0.278 \rightarrow \text{Not uniformly distributed}$$

6 seconds in
total:

$$1 - \left[\left(\frac{0}{6} \right)^2 + \left(\frac{6}{6} \right)^2 \right]$$

$$1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right]$$

$$\Rightarrow 0.500 \rightarrow \text{Equal distribution}$$

$$20 \rightarrow \text{equal}$$

* GINI for single node

$$\textcircled{1} \frac{C_1}{C_2} \begin{vmatrix} 0 \\ 6 \end{vmatrix}$$

* GINI for collection of nodes.

when a node P is split into k parts

GINI split *

$$\textcircled{2} \frac{C_1}{C_2} \begin{vmatrix} 0.486 \\ 5 \\ 12 \end{vmatrix}$$

Parson - 5

$$\textcircled{3} \frac{C_1}{C_2} \begin{vmatrix} 0.361 \\ 5 \\ 12 \end{vmatrix}$$

Parson - 5

$$\therefore 1 - \left[\left(\frac{4}{12} \right)^2 + \left(\frac{5}{12} \right)^2 \right]$$

$$\textcircled{4} \frac{C_1}{C_2} \begin{vmatrix} 0.486 \\ 5 \\ 12 \end{vmatrix}$$

Multivariate split, Two-way split

* GINI for categorical attributes

$$\text{Gain} = 0.486 - 0.361$$

Weighted GINI :- $W_1 N_1 + W_2 N_2$

$$\Rightarrow \frac{6}{12} * 0.486 + \frac{6}{12} * 0.361$$

$$\Rightarrow 0.361$$

of least GINI index to be selected to split
giving one split

Decision split

spouse working

Outcome, 0.375, 0, 0.248, 0.219

Weighted gain :- $\frac{4}{20} * 0.375 + \frac{8}{20} * 0.248 + \frac{8}{20} * 0.219$

$$= 0.163 \cdot \underline{\underline{0.1622}}$$

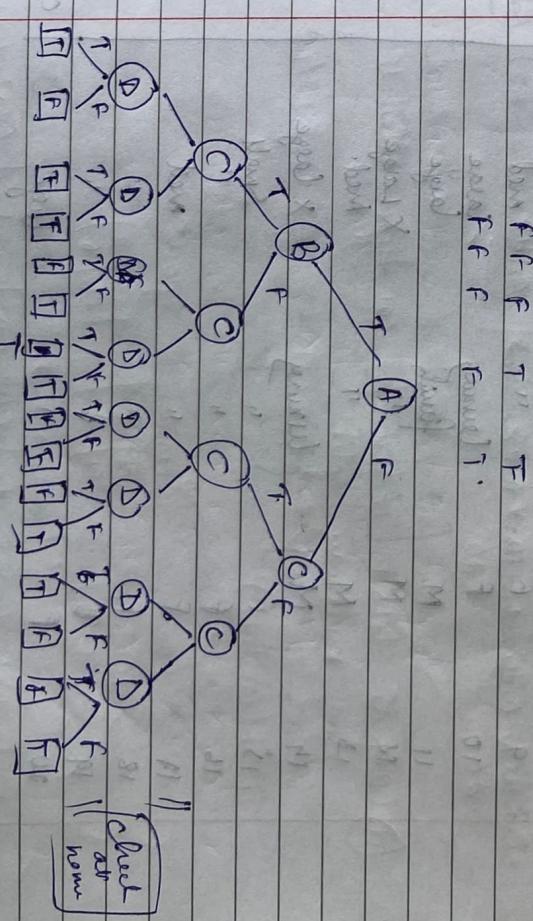
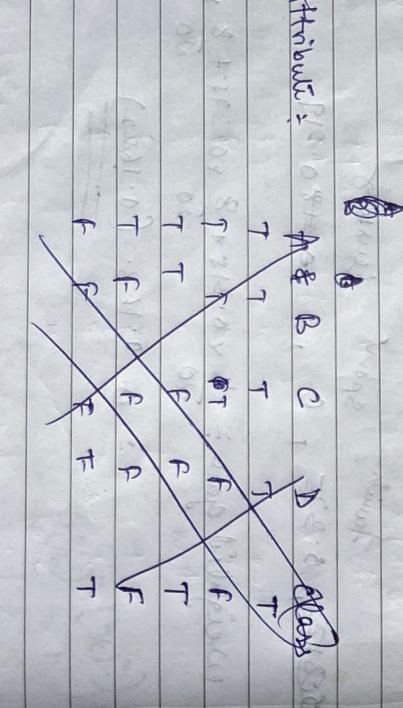
$$\textcircled{5} \begin{array}{|c|c|} \hline C_1 & N_1 \\ \hline 5 & 2 \\ \hline C_2 & N_2 \\ \hline 1 & 4 \\ \hline \end{array} \quad 1 - \left(\frac{2}{6} \right)^2 \cdot \left(\frac{4}{6} \right)^2$$

$$\Rightarrow 0.486 \Rightarrow 0.444$$

* Continuous attribute GINI index method

- Sort attribute on values.
- Linearly scan these values each time updating the count matrix & computing Gini
- Choose split position that has least GINI index.

- ~~2011/2012~~
- Draw a full decision tree for the parity function of 4 boolean attributes A, B, C, & D. Assume even parity
- Even no. of 1's \rightarrow 1
Odd no. of 1's \rightarrow 0
- $$= \epsilon[A, B, C, D]$$



in field

- Q) Consider the training examples showing below for binary classifier

Customer ID	Gender	Car type	Shirt size	Class
1	M	Fewish	Small	C0
2	M	Sports	med	"
3	M	"	"	"
4	M	"	Large	"
5	M	"	X Large	"
6	M	"	X Large	"
7	F	"	Small	"
8	F	"	"	"
9	F	"	med	"
10	F	Luxury	large	C0
11	M	familiy	large	C1
12	M	"	X Large	"
13	M	"	med.	"
14	M	Luxury	X Large	"
15	F	"	Small	"
16	F	"	"	"
17	F	"	med	"
18	F	"	"	"
19	F	"	"	"
20	F	"	"	"

- Q) Compute GINI index for small collection of training examples

$$GINI(t) = 1 - \sum_{i=0}^{C-1} \left[p_i(t) \right]^2$$

\downarrow
fraction of
data records
belonging to class
priori.

$$GINI(t) = 1 - \left[\left(\frac{10}{20} \right)^2 + \left(\frac{10}{20} \right)^2 \right]$$

$$= 0.5$$

- Q) Compute GINI index for customer id attribute

$$\text{Cust. Id} = 1 \quad [C0 : 1] \quad [C1 : 8]$$

$$\Rightarrow 1 - \left[\left(\frac{1}{9} \right)^2 + \left(\frac{8}{9} \right)^2 \right] = 0$$

∴ Overall minimum id is GINI index = 0.

DOMS	Page No.
/	/

Lesson entropy, pure class

DOMS	Page No.
/	/

3 instances	↓	↓	↓	↓
6 GINI index	↓	↓	↓	↓
Entropy	↓	↓	↓	↓
Misclassification rate	↓	↓	↓	↓
Date (class prior)	↓	↓	↓	↓

$$\Rightarrow 4 \text{ records} - + \text{ class} \\ 5 " - - , "$$

Q Consider the training examples shown in the table below for luxury classification problem.

Instance	a ₁	a ₂	a ₃	target class
1	T	T	1.0	+
2	F	T	6.0	+
3	X	F	5.0	Wor
4	F	F	4.0	+
5	T	F	3.0	-
6	F	F	0.8	-
7	F	T	7.6	+
8	T	F	0.8	-
9	F	T	5.0	-

Q What are the info. gains of a₁, a₂ related to this training examples.

	a ₁	+	-	a ₂	+	-
T	3	10 × 8	3	T	2	3
F	1	4	0	F	2	2

$$\therefore \text{Entropy} = -\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) = 0.8112$$

$$\text{Entropy}_{a_1(-)} = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{1}{4} \log_2 \left(\frac{4}{5} \right) = 0.7219$$

$$\text{Entropy}_{a_2(+)} = \overbrace{\frac{5}{9}}^{0.941} \times 0.7219 + \overbrace{\frac{4}{9}}^{0.059} \times 0.8112$$

- Q What is the entropy of this collection of training example w.r.t positive class
fraction w.r.t first column to last
- $$\text{Entropy} = -\sum_{i=1}^{120} p(a_i) \log_2 p(a_i) \cdot P(i|+)$$
- Where
- C = number of classes $\neq 0$ $\log_2 0 = 0$ for loss function calculations.

$$\text{Entropy gain} = 0.991 - 0.7615$$

$$\text{for } a_1 = 0.2295$$

$$\log_2 \left(\frac{1}{4} \right) = \log_2 \left(\frac{5}{9} \right) \\ 0.3010$$

$$\text{Entropy} = -\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{3} \right)$$

$$G_A(+)$$

$$= \frac{1}{2}$$

$$\text{Entropy} = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{3} \right)$$

$$G_B(-)$$

$$= 0.9709.$$

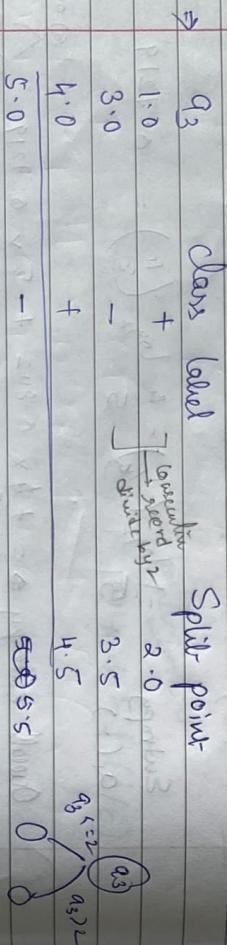
$$\text{Overall entropy} = \frac{4}{9} \times 1 + \frac{5}{9} \times 0.9709$$

$$= 0.9829$$

$$\therefore \text{Info. gain} = 0.9911 - 0.9829$$

$$= 0.0072$$

Q) For a_3 which is a continuous attribute, compute info. gain for every possible split.



$$\Rightarrow \frac{1}{9} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{1} \right) \right] + \frac{8}{9} \left[-\frac{3}{8} \log_2 \left(\frac{3}{8} \right) - \frac{5}{8} \log_2 \left(\frac{5}{8} \right) \right]$$

$$\Rightarrow q_3(+)$$

$$-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - \frac{0}{1} \log_2 \left(\frac{0}{1} \right)$$

$$= 0.8484.$$

\Rightarrow q_3 class label Split point

1.0 + 3.5 2.0
3.0 - 3.5 2.0
4.0 + 4.5 4.5
5.0 - 5.5 5.5

\Rightarrow q_3 class label Split point

1.0 + 4.5 4.5
2.0 - 4.5 4.5
3.0 + 4.5 4.5
4.0 - 4.5 4.5
5.0 - 4.5 4.5

Point of split	Entropy	Inf. gain
3.5	0.9889	0.9911 - 0.9889
4.5	0.9855	0.9911 - 0.9855
5.5	0.9828	0.9911 - 0.9828
4.0	0.9839	0.9911 - 0.9839
5.0	0.9928	0.9911 - 0.9928
6.0	0.0183	0.9911 - 0.0183
7.0	0.8889	0.9911 - 0.8889
8.0	0.1022	0.9911 - 0.1022

Q) What is the best split among $a_1, a_2 \& a_3$ according to info. gain?

$\Rightarrow a_3$

postulate

- * Algorithm for Decision tree induction.
- * Problem with large no. of partitions.

* Gain ratio

Gain $R = \frac{\text{Gain}_{\text{split}}}{\text{Split Info}}$

$$\text{Split Info} = - \sum_{i=1}^k n_i^{\text{left}} \log_2 \frac{n_i^{\text{left}}}{n}$$

(1)

$$C_1 \quad 0$$

$$C_2 \quad 6 \Rightarrow P(C_1) = 0/6 = 0, \quad P(C_2) = 6/6 = 1$$

$$\therefore \text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

$$C_1 \quad 1$$

$$C_2 \quad 5 \Rightarrow P(C_1) = 1/6, \quad P(C_2) = 5/6$$

$$C_1 \quad 4$$

$$C_2 \quad 5 \Rightarrow \text{Error} = 1 - \min\left(\frac{1}{6}, \frac{5}{6}\right)$$

$$\text{Split Info}_{1.52} = 0.42 \quad \text{Gain}_{a_3} = 0.94$$

$$-\frac{4}{6} \log_2 \left(\frac{4}{6}\right) - \frac{16}{20} \log_2 \left(\frac{16}{20}\right) = \frac{8}{20} \times \log_2 \left(\frac{8}{20}\right)$$

$$-\frac{8}{20} \log_2 \left(\frac{8}{20}\right) - \frac{4}{20} \log_2 \left(\frac{4}{20}\right) = \frac{12}{20} \log_2 \left(\frac{12}{20}\right) - \frac{8}{20} \log_2 \left(\frac{8}{20}\right) = 6.7219 - 20.9409 = 1.5219$$

* Classification error.

\Rightarrow Class. error at node t.

$$[\text{Error}(t) = 1 - \max_i [P_i(t)]]$$

\downarrow
max freq. belonging to particular class

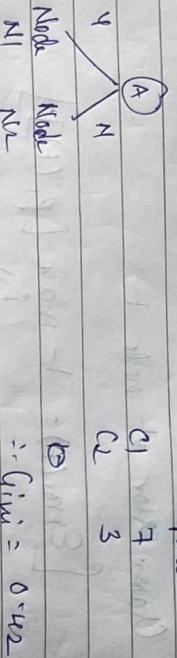
\Rightarrow Max. of $1 - \frac{1}{c}$ when records are equally distributed.

among all classes.

\Rightarrow

Misclassification Error vs. Jarden

Gini
Parent



$$Gini(N_1) = 0 \quad Gini(N_2) = 0.489$$

(weighted)

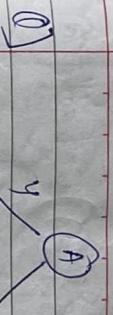
$$C_1 \quad \begin{cases} 3 \\ 7 \end{cases} \quad Gini(\text{Children}) = \frac{3}{10} \times 0.343 + \frac{7}{10} \times 0.42$$

$$\overline{P(C_1)} = \frac{7}{10} \quad P(C_2) = \frac{3}{10}$$

$$\therefore \text{Error} = 1 - \left(\frac{7}{10}, \frac{3}{10} \right)$$

Q] For the following dataset, generate the decision tree

Ref	Age	Income	Student	Gender
1	Youth	High	No	Fair
2	Youth	Low	No	Excl.
3	Middle-Aged	High	No	Fair
4	Senior	Medium	No	Excl.
5	Senior	Low	Yes	Excl.
6	Senior	Low-Med	Yes	Excl.
7	Middle-Aged	Med-Low	Yes	Excl.
8	Youth	Low-Medium	No	Fair
9	"	Med-Low	Yes	"
10	Senior	Med	Yes	"
11	Youth	Med.	Yes	Excl.
12	Middle-Aged	Med-Low	No	Excl.
13	"	High	Yes	Fair
14	Senior	Medium	No	Excl.



① Basic algo. (Greedy algo.)
→ Tree is constructed in top-down recursive divide & conquor manner.

Complexity of algo: $O(n \times 10 \times \log 10)$
where, $n \rightarrow$ no. of attributes describing tuples in D and $10!$ is no. of training tuples in D .

Class: lungs - complete

1	yes
2	no
3	yes
4	no
5	yes
6	no
7	yes
8	no
9	yes
10	no
11	yes
12	no
13	yes
14	no
15	yes
16	no
17	yes
18	no
19	yes
20	no
21	yes
22	no
23	yes
24	no
25	yes
26	no
27	yes
28	no
29	yes
30	no
31	yes
32	no
33	yes
34	no
35	yes
36	no
37	yes
38	no
39	yes
40	no
41	yes
42	no
43	yes
44	no
45	yes
46	no
47	yes
48	no
49	yes
50	no

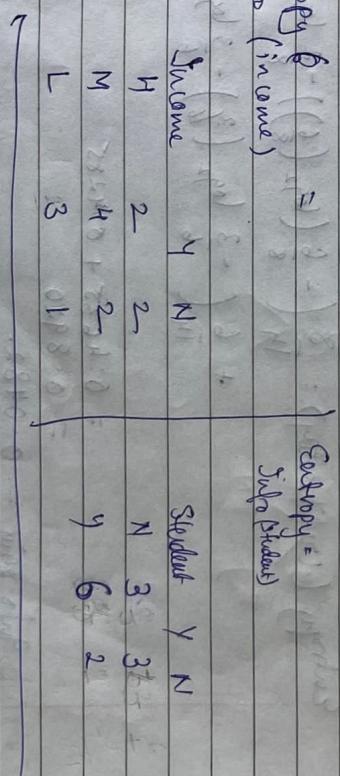
Use entropy & info. gain

(d) Entropy of entire system

$$g = -\frac{1}{4\pi G}$$

$$\therefore \text{Entropy of dataset} = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

extending credit \Rightarrow Credit rating $\neq N$
(Credit rating) fear 6 2



$$\therefore \text{gain (age)} = 0.9402 - 0.694 = 0.246.$$

$$\text{Info (age)} = \frac{5}{14} \left(-2 \log_2 \frac{2}{5} - 3 \log_2 \frac{3}{5} \right)$$

$$\text{Entropy} = \frac{5}{14} \left(\frac{-2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \left(-\frac{1}{4} \log_2 \frac{1}{4} \right) + \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} \right)$$

$$= 0.694$$

212

= 0.9402

卷之二

1

1

2

۲

23-4

ω

100

2

July)

1

100

100

$$\text{Entropy (income)} = \frac{5}{14} \left(-2 \log_2 \left(\frac{2}{4} \right) - 2 \log_2 \left(\frac{2}{4} \right) \right)$$

$$+ \frac{6}{14} \left(-\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) \right)$$

$$+ \frac{4}{14} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right)$$

$$= 0.6912 + *$$

$$= 0.9110$$

$$\therefore \text{Gdp. gain} : 0.9462 - 0.9110 = 0.0292$$

$$\text{Entropy (student)} = \frac{6}{14} \left(-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right)$$

$$+ \frac{8}{14} \left(-\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right)$$

\Rightarrow NII \Rightarrow Gdp. gain on attributes Income, student, credit rating to determine to which one to use

$$\text{Entropy (Gdp. gain)} = \frac{8}{14} \left(-\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right)$$

$$+ \frac{6}{14} \left(-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right)$$

Income Y N

	High	Med.	Low
High	0	1	0
Med.	1	0	1
Low	0	1	0

Credit rating Y N

	fair	good
fair	1	2
good	1	1

$$= 0.4635 + 0.44285$$

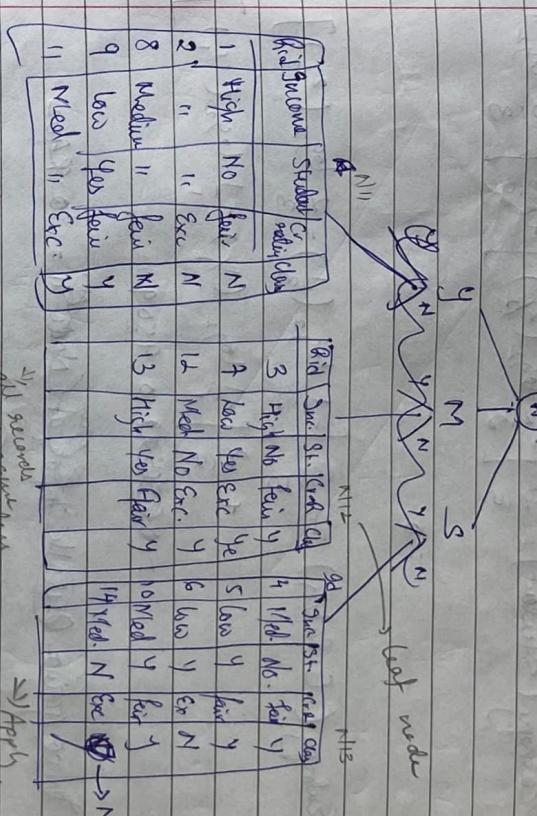
$$\therefore \text{Gdp. gain} = 0.0482$$

$$\text{Entropy (income)} \Rightarrow \frac{2}{5} \left(-2 \log_2 \left(\frac{2}{5} \right) \right) + \frac{1}{5} \left(-\frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right)$$

$$+ \frac{2}{5} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right)$$

$$= 0.4482 \quad 0.4$$

$$\text{Gdp. gain} =$$



$$\text{Entropy (Student)} = \frac{2}{5} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] + \frac{3}{5} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right]$$

M13 → category of income, student, credit rating.

Income	Y	N	Student	Y	N
Med.	12	81	Y	2	1
low	1	N	2	0	

$$\text{Entropy}(\text{C.C.Ratio}) = \frac{3}{5} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right]$$

$$\text{Water grain: } \frac{\text{Volume of water}}{\text{Volume of soil}} = 0.9402$$

= 0.95050509

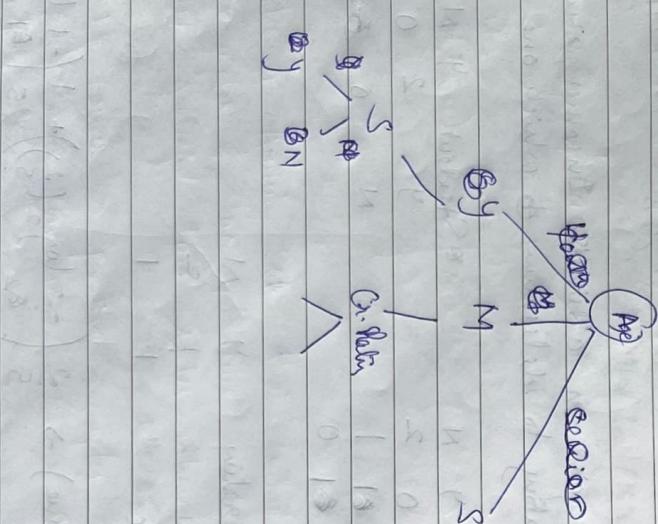
$$\text{Entropy } E_{\text{Gibbs}} = \frac{3}{5} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right)$$

$$+ \frac{2}{5} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right)$$

$$= 0.9509 \cdot 0.0001 = 9.509 \cdot 10^{-6}$$

$$\text{Entropy (Student)} = \frac{3}{5} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \right)$$

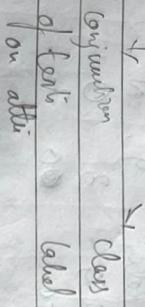
$$\text{Entropy (red. satn)} = \frac{3}{5} \left(-\frac{3}{3} \log_2 \left(\frac{3}{3}\right) \right)$$



Rule-based classifier:

→ Classify records using a collection of "if .. then .." rules.

Rule : (Condition) \rightarrow Y



e.g.: (Blood type = warm) \wedge (lay eggs) = Yes \rightarrow Birds.

(Tawalli know C50) \wedge (Mifawed

* Rule coverage & accuracy

⇒ Only correct antecedent consequent

↓
antecedent
consequent

$$R = 2 \cdot \sum_{i=1}^k f_i \log_2 \left(\frac{f_i}{e_i} \right)$$

where
 $k \rightarrow$ no. of classes

$f_i \rightarrow$ observed freq. of class i examples

that are covered by the rule.

- Q) Consider a training set that contains 60 positive examples & 100 negative examples. Suppose we have generated two rules where rule

R1 covers 50 positive examples & 5 negative examples

R2 covers 60 " " 0 " "

Rule R1:

Can you compute accuracy of these two rules?

$$\Rightarrow \text{Total} = 160$$

$$= 100 + 60$$

∴ Rule accuracy:

$$\text{Accuracy}(R_1) = \frac{50}{160} = 90.9\%$$

$$\therefore (x_2) = \frac{2}{e_2} = 100 \text{ %}$$

① Likelihood ratio statistic

R has a Chi-square distribution with $k-1$ degrees of freedom

rule R1 covers 55 examples, $e_1 = n_1$ of example expected freq. of the class (e_1) = no. of example covered \times total freq.

$$= 55 \times \frac{60}{160} = 20.625$$

$$\text{Expected freq. of } e_-(-) \text{ class} = 55 \times \frac{100}{165} = 34.375$$

Assignment 2

$$\therefore \text{Likelihood ratio for } R_1 = 2 \times \left[50 \times \log_2 \left(\frac{50}{20.625} \right) + 50 \times \log_2 \left(\frac{50}{34.375} \right) \right]$$

$$R(R_2) \div$$

$$P_{\text{err}} = 1 - 0.9993 = 0.0007$$

$$\text{Total} \Rightarrow 500 (+ve + -ve) = 100 \rightarrow 400$$

$$\text{Accuracy}(R_1) = \frac{4}{5} = 80\%.$$

$$\text{Accuracy}(R_2) = \frac{30}{40} = 75\%.$$

$$\text{Accuracy}(R_3) = \frac{100}{190} = 52.631\%.$$

$$\therefore R(R_2) = 2 \times \left[2 \times \log_2 \left(\frac{2}{0.25} \right) + 0 \right]$$

Likelihood

$$R = 2 \sum_{i=1}^k f_i \log_2 \left(\frac{f_i}{e_i} \right)$$

$$R(R_1) \rightarrow e_+ = 25 \times \frac{100}{500}$$

$$= 1.$$

$$e_- = 5 \times \frac{400}{500} = 4.$$

$$\therefore R(R_1) = 2 \times \left[2 \times \log_2 \left(\frac{4}{1} \right) + 50 \times \log_2 \left(\frac{1}{4} \right) \right]$$

Expt

$$R(R_2) \rightarrow e_+ = 50 \times \frac{100}{500} = 8$$

$$e_- = 100 \times \frac{400}{500} = 32$$

(d) $R(R_3)$: A \rightarrow + (Correct +ve & 1-ve examples)

A₂: B \rightarrow + (" 30 +ve & 10 -ve ")

A₃: C \rightarrow + (" 100 +ve & 90 -ve examples")

$$R(R_2) = 2 \left[30 \times \log_2 \left(\frac{30}{8} \right) + 10 \times \log_2 \left(\frac{10}{32} \right) \right]$$

$$= 80.8819.$$

$$R(R_3) =$$

$$e^+ = 190 \times \frac{100}{500} = 38$$

$$e^- = 190 \times \frac{400}{500} = 152.$$

$$= 2 \left[100 \times \log_2 \left(\frac{100}{38} \right) + 90 \times \log_2 \left(\frac{90}{152} \right) \right]$$

$$\approx 143.0923$$

where $n \rightarrow$ no. of examples covered by the rule.

$f^+ \rightarrow$ no. of positive examples covered by the rule.

$k \rightarrow$ no. of classes.

(III) m-estimate measure

→ Given as

$m\text{-estimate } f^+ + K p^+$, $p^+ \rightarrow$ prior probability of the class

$n+k$

$\boxed{Q} \rightarrow$ Problem solved on 29th Nov

(IV) Laplace

$$f^+(R_1) \rightarrow 50$$

$$n = 55$$

$$k = 2.$$

$$\therefore \text{Laplace}(R_1) = \frac{51}{57} (53+1) = 0.8947.$$

$$f^+(R_2) = \frac{51}{57} (53+2)$$

$$n = 2$$

$$k = 2$$

④ Laplace measure

→ Evaluation metric that takes into account

rule coverage

→ Given by formula:

$$\text{Laplace} = \frac{f^+ + 1}{n+k}$$

officer
6/12/2023

$$R(R_2) = 2 \left[30 \times \log_2 \left(\frac{30}{8} \right) + 10 \times \log_2 \left(\frac{10}{32} \right) \right]$$

$$= 80.8879.$$

$$R(R_3) = e_+ = 190 \times \frac{100}{300} = 38$$

$$e_- = 190 \times \frac{400}{300} = 152.$$

$$= 2 \left[100 \times \log_2 \left(\frac{100}{38} \right) + 90 \times \log_2 \left(\frac{90}{152} \right) \right]$$

$$= 143.0923$$

IV Laplace measure

→ Evaluation metric that takes into account rule coverage

→ Given by formula:

$$\text{Laplace} = \frac{f_+ + 1}{n+k}$$

where $n \rightarrow$ no. of examples covered by the rule.

$f_+ \rightarrow$ no. of positive examples covered by the rule.

$k \rightarrow$ no. of classes

(III) m-estimate measure

→ Given as

$m\text{-estimate } f_+ + k p_+ \rightarrow p_+ \rightarrow \text{prior probability of the class}$

$n+k$

Q) → Problem solved on 29th Mar

IV Laplace

$$f_+(R_1) \rightarrow 50$$

$$n = 55$$

$$k \rightarrow 2.$$

$$\therefore \text{Laplace}(R_1) = \frac{51}{57} (53+1) = 0.8947.$$

$$\frac{57}{57} (53+2)$$

$$f_+(R_2) = 2 \Rightarrow \text{Laplace}(R_2) = \frac{2+1}{2+2} = \frac{3}{4} = 0.75$$

$$n = 2$$

$$k \rightarrow 2$$

M-estimate

$$\hat{p}_+ = 0.2$$

$$M\text{-est}(\mathcal{R}_1) = \frac{50 + 8 \times 0.2}{54}$$

$$= 0.8842$$

$$mest(\mathcal{R}_2) = \frac{2 + 2 \times 0.2}{4}$$

$$\begin{aligned} p_1 &\rightarrow 50 & p_2 &\rightarrow 2 \\ n_1 &\rightarrow 5 & n_{p_2} &\rightarrow 0 \end{aligned}$$

$$\therefore \text{FOL}'_3 = 50 \left(\log_2 \frac{50}{55} - \log_2 \frac{60}{160} \right)$$

$$= 2.83004$$

\mathcal{R}_2 's info. gain

→ Takes into account the support count of the rule.

→ Support count → No. of the examples covered by the rule.

Suppose rule $\mathcal{R}_2: A \rightarrow +$ covers p_0 positive examples & n_0 negative examples.

After adding a new conjunct B , the extended rule $\mathcal{R}_1: A \wedge B \rightarrow +$ covers p_1 positive exs.

∴ \mathcal{R}_1 negative examples.

∴ \mathcal{R}_1 info gain is given by

$$= p_1 \times \left(\log_2 \frac{p_1}{p_1 + n_1} - \log_2 \frac{p_0}{p_0 + n_0} \right)$$

base prob:

$$p_0 \rightarrow 60$$

$$n_0 \rightarrow 100$$

DOMS	Page No.
1	1

DOMS	Page No.
1	1

Date

Page No.

Date

Page No.

DOMS	Page No.
/	/

DOMS	Page No.
/	/

$$\text{Laplace}(R_2) = \frac{30+1}{40+2}$$

$$= \frac{31}{42}$$

$$= 0.7380$$

$$B$$

$$\text{Faplace}(R_3) = \frac{100+1}{190+2} \approx 0.5260$$

$$\text{GOLIS}(R_3) = 100 \left(\log_2 \frac{30}{40} - \log_2 \frac{100}{190} \right)$$

$$= 57.42067$$

$$= 139.592 \quad B$$

M - estimation

$\hat{p}_+ = 0.1$.

$\hat{p}_- = 0.2$.

$\text{M}-\text{est}(R_1) = \frac{4+2\times 0.2}{5+2}$

(19)

$$= 0.6285$$

(19)

$$\text{M}-\text{est}(R_2) = \frac{30+2\times 0.2}{40+2}$$

(19)

$$= 0.7238$$

(19)

$$\text{M}-\text{est}(R_3) = \frac{100+2\times 0.2}{190+2}$$

(19)

$$= 0.7380$$

(19)

- (a) Given dataset illustrates the coverage of classifier under R_1 , $R_2 \leq R_3$. Determine which is the best rule according to
- (i) Likelihood ratio statistic
 - (ii) Laplace measure.
 - (iii) M-estimation.
- (b) Rule accuracy after R_1 has been discarded where none of the examples covered by R_1 are discarded.
- (c) Rule " R_1 " " " R_1 " " " R_1 " " " R_1 " only true examples covered by R_1 are discarded.
- (d) Rule " R_1 " " " R_1 " " " R_1 " " " R_1 " both true & -ve examples covered by R_1 are discarded.

FOILS: W - P - 3 SNR - 0.001

SNR

0.001

0.01

0.1

1

10

100

1000

10000

100000

1000000

10000000

100000000

1000000000

10000000000

100000000000

1000000000000

10000000000000

$$P_0 = 100 ; P_1 = 4 ; P_2 = 30 ; P_3 = 100$$

$$n_0 = n_{100} ; n_1 = 1 ; n_2 = 10 ; n_3 = 90$$

FOLIS(R_1) =



3 + n

W

10

100

1000

10000

100000

1000000

10000000

100000000

1000000000

10000000000

100000000000

1000000000000

10000000000000

$\therefore 8$

r_3 when one egg is discarded

+ + + + + + + + + +

1
+
-
1
.
—
—
1
1

1 t. —
— 1
1 300. DOL f

- 1 -

—

1924-25 1925-26

Ward 4 - 29 = 100% (100% of the 29 houses)

$$10 \cdot \frac{1}{2} = 5$$

R1

No. of positive = 12
No. of negative = 3

R₂

11 11
11 11

P3 when home of Eng. of R1 are discussed.

七
八
九
十

five = 5

1
三
二

R3 when both true & false legs are discarded.

$$= \text{ant}$$

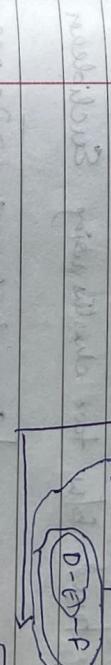
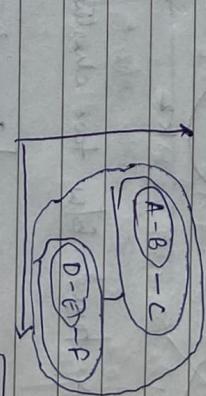
→ Computer measure

K-means clustering

Clustering : process of dividing the data into groups or clusters depending on relationships that exists in a data.

Types of clustering :

- ① Exclusive clustering
 - Hard clustering
 - Data points fall exclusively belong to one cluster.



→ Eg : K-means clustering.

- ② Overlapping clustering

- Soft clustering
- Data point item can be belong to multiple clusters.
- Eg : fuzzy | C-means clustering.

→ K-Means clustering :

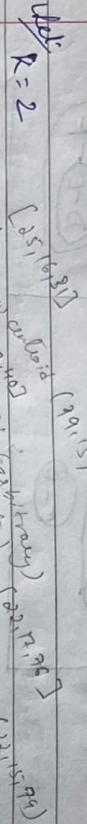
- Iterative algo.
- Partition dataset into k-clusters.

Mall Customer Dataset

Cluster	Customer ID	Age	Income	Spending Score (1-100)
K1	1	19	15	39
	2	21	15	81
K1	3	26	16	6
	4	23	16	77
K1	5	31	17	40
K2	6	22	17	76

B, C, D, E, F

k = no. of clusters which has to be known apriori



$k_1 = [19, 15, 19]$

$k_2 = [21, 15, 19]$

Similarly we have find two clusters using Euclidean distance.

$$k_1, \textcircled{3} = \sqrt{(20-19)^2 + (16-15)^2 + (16-39)^2} = 33.03$$

$$[20, 16, 6] \Rightarrow k_1, \textcircled{3} = 33.03.$$

place where distance is min

∴ New centroid

$$k_1, \textcircled{5} \rightarrow 21.7255$$

$$k_2, \textcircled{5} \rightarrow 14.6044$$

$$k_1 = 19.5, 15.2, 22.5$$

$$k_2 = 14.6 - 4.582$$

∴ New centroid

$$[31, 19, 40]$$

$k_1 = [19, 15, 39]$

$k_2 = [21, 15, 81]$

∴ Update centroid : Take average of $\frac{\textcircled{1} + \textcircled{3}}{2} = 19.5, 15.2, 22.5$

New:

$$\left[\frac{21+19}{2}, \frac{15+39}{2}, \frac{22.5+81}{2} \right]$$

$$k_1, \textcircled{6} \rightarrow 24.55, 45.11$$

$$k_2, \textcircled{6} \rightarrow 3.605$$

k_1	19	15	22	1.	19
k_2	21	15	81	2	19

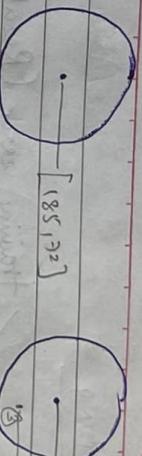
$$\left[\frac{21+19}{2}, \frac{15+39}{2}, \frac{22.5+81}{2} \right]$$

DOMS	Page No.
/	/
Date	/ /

Algorithm:

- ① Start
- ② Decide on the no. of clusters (k)
- ③ Find the centroids
- ④ Take each record & calculate the Euclidean distance from all centroids
- ⑤ Select the min. dist.
- ⑥ Update cluster centroid
- ⑦ Update record with cluster no.
- ⑧ Stop after executing all records.

$$k_1, k_2 \rightarrow 620, 808$$



Height	Weight	Cluster
185	82	k_1
170	56	k_2
168	60	
179	68	
182	72	
188	77	k_1
180	71	k_1
180	70	k_1
183	84	k_2
180	88	k_2
180	64	k_1
177	70	k_1

\checkmark Use K-means

$K=2$

overfitting & MODEL OVERFITTING

Error \leftarrow generalization error.

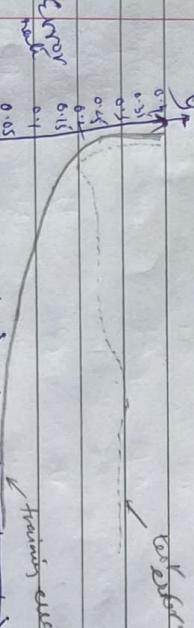
Training error (Resubstitution error) \downarrow
 appears as overfitting.

Training error = no. of misclassification errors committed on training records. (Training records)

Generalization error = expected error of model on previously unseen records (test data set)

Good classifying model should have low training error as well as low generalization error. Model that fits the data too well is known as model overfitting.

Training \leq test error rule -



number of nodes in tree

Model overfitting: Smaller the tree size, training \leq test errors are large. This situation is called model underfitting.

When the tree becomes too large, its test error will begin to increase even though the training error will continue to decrease. This phenomenon is

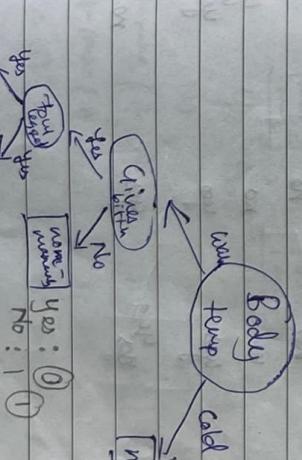
known as model overfitting.

i) Overfitting due to the presence of noise.

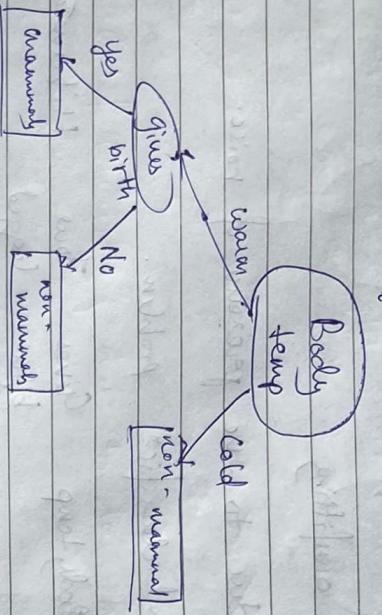
Eg: Mammal classification problem.

Name	Body temp	Gives birth	Hibernates	Class label
porcupine	warm	yes	yes	yes
cat	"	yes	yes	no
dog	"	yes	no	yes
whale	"	no	no	no
Salamander	cold	no	yes	yes
Komodo dragon	cold	no	yes	no
pygmy marmoset	"	no	yes	no
salmon	"	no	no	no
eagle	warm	no	no	no
puppy	cold	yes	no	no
cow	"	no	no	no

Decision tree - A model with 3 nodes
if Body temp is warm then
if Cries then
if non-mammal then



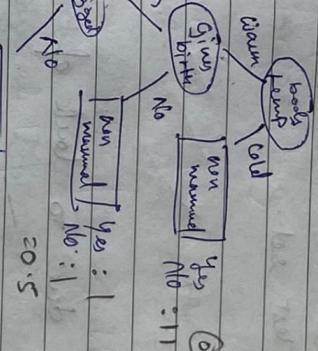
Decision tree for M_2 :



Eg: Test set for classifying mammals.

name	body temp	give birth	fur	hibernates	Class
Human	warm	yes	no	no	non-mammal
Pigeon	warm	no	no	no	non-mammal
Elephant	warm	yes	yes	no	non-mammal
Cheetah	cold	yes	no	no	non-mammal
Snowy owl	cold	yes	no	no	non-mammal
Fulmar	cold	no	yes	no	non-mammal
Penguin	cold	no	no	no	non-mammal
Seal	cold	no	no	no	non-mammal
Dolphin	warm	yes	no	no	non-mammal
Sperm whale	warm	no	yes	yes	non-mammal
Auklet	cold	no	yes	yes	non-mammal
Albatross	cold	yes	yes	no	non-mammal

Model M1



Model M2

Model M1

Training error rate for T_L

Training error rate.

$$Eg(T_L) = \frac{1}{24} [0.25 + 3 \cdot 0.33 + 3 \cdot 0.25 + 4 \cdot 0.25]$$

$$\Rightarrow 0.1666.$$

↓ lower after
decreasing error
before it is

$$T_R = 0.25$$

$$T_L = 0.1666.$$

* Incorporating model complexity

→ chance for model overfitting increases as model becomes more complex

→ prefers simpler models

→ Occam's Razor or principle of parsimony.

Occam's Razor: Given two models with same generalization error, simpler model is preferred over the more complex model.

\Rightarrow it is intuitive.

* Pessimistic error estimate

Generalization error = training error + penalty term for model complexity

(lower the pessimistic error, better it is)

↓ let $n(t)$ be no. of training records classified by node t & $e(t)$ be no. of misclassified records

* Pessimistic error estimate of a decision tree T ,

$$Eg(T) = \sum_{i=1}^k [e(t_i) + \Omega(t_i)]$$

$$= e(T) + \Omega(T)$$

$$\text{where } k \rightarrow \text{no. of leaf nodes}$$

$e(t) \rightarrow$ overall training error of the decision tree

$\Omega(t_i) \rightarrow$ penalty term associated with each node t_i

→ for the binary decision trees, $T_L \leq T_R$ if penalty sum = 0.5

Pessimistic error estimate for T_L

$$eg(T_L) = 0.167 + \frac{7 \times 0.5}{24} = 0.167 + 0.3328$$

$$eg(T_R) = 0.25 + \frac{7 \times 0.5}{24} = 0.3333.$$

② Min. description length principle or MDL.

→ Overall cost of decision tree would be cost of encoding the model + cost of encoding unlabelled records.

$$\rightarrow \text{cost(model, data)} = \text{cost(model)} + \text{cost(data | model)}$$

According to MDL, we should seek a model that minimizes the overall cost funcⁿs.

* How overfitting is handled in decision tree induction.

① Pruning (Early stopping rule)

→ Tree-growing algo. is halted before generating a fully grown tree. that perfectly fits an entire training data.

Eg: Stop expanding leaf nodes when gain in impurity measure falls below certain threshold.

→ Advantage: Avoids generating overly complex sub-tree that overfit the traini data.

→ Disadvantage: Choosing the right threshold value

② Post pruning

→ The decision tree is initially grown to its maximum size

→ Tree pruning step is applied to trim it fully grown tree in a bottom up fashion.

→ Following steps done: By applying a subtree with

① A set new, leaf node whose class label

is determined from majority class of records affiliated with subtree.

② Most freq. used branch of tree subtree

* Rule based classifier

→ Rule induced RUPPER algo

2 class prob: choose one as the and other as -ve

→ Learn rules for positive class

→ Negative class will be default class.

Eg: Classes ~~are~~ ^{are}

Class Y₂

Multiclass prob:

→ Order classes acc. to increasing class prevalence
 $\rightarrow y_1, y_2, \dots, y_c$; y_1 is the least freq. class

→ draw rule set for smallest class first, treat rest as negative

→ Repeat with next smallest class as the class

→ Growing a rule:

→ Start from empty rule.

→ Add conjuncts as long as they improve FOIL's info. gain

→ Stop when rule no longer covers negative egs.
 \rightarrow Prune the rule immediately using incremental reduction

→ Measure of purity $\Rightarrow V = \frac{p-n}{p+n}$

P: no. of true class covered by rule in validation set.

n: no. of -ve rays " "

" "

→ Pruning method: delete any fixed sequence of conditions that maximize P.

→ End fit advantage

Nearest-neighbour classifier:

~~Eager learner~~ Lazy learner,

→ Generalized model from training data set is constructed using the model, the class of test dataset is predicted by decision trees, rule based classifier

→ Training dataset is stored

→ On querying, similarity is calculated to predict the class of test data.

→ Using the model, the class of test dataset is predicted

→ Using the model, the class of test dataset is predicted

→ Using the model, the class of test dataset is predicted

→ Using the model, the class of test dataset is predicted

→ Using the model, the class of test dataset is predicted

→ Using the model, the class of test dataset is predicted

$C_{ij} =$

Species	dist
Setosa	0.608
Setosa	0.704
Virginica	0.701
Setosa	0.36
Setosa	0.22
Virginica	0.82
Virginica	0.22
Virginica	0.94
Virginica	0.1
Virginica	0.89
Virginica	0.1
Virginica	0.89
Virginica	0.7
Virginica	1.36
Virginica	0.60
Virginica	1.25
Virginica	0.95

Quarry
Sapalaeopter Specie, width species.

5.1 ~~3.1~~ ?

K-1 - 2:15

Step 1: Distance from our eg. to each of friend eg.

\Rightarrow Euclidean distance: $\sqrt{(x-a)^2 + (y-b)^2}$

$$Dist(\text{sepal length}, \text{sepal width}) = \sqrt{(5.2 - 5.0)^2 + (3.4 - 3.9)^2} = 0.698$$

→ Calculate for all ~~less~~ records

k=1, last sq. would be assigned the species *lobosa*.
 k=2, " " " "
 k=3, " " " "
 ;
 k=5, majority voting scheme is used. Last sq. is assigned to the cloud *lobosa*.

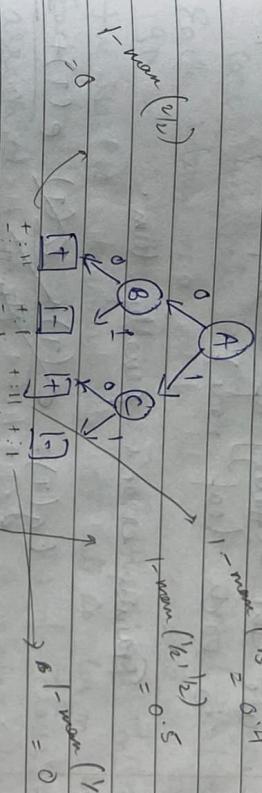
Terada Terada +

$$\text{Classification error } (\epsilon) = 1 - \max_i [P(i|t)]$$

$$\rightarrow \frac{2}{10} (0) + \frac{2}{10} (5.6) + \frac{5}{10} (5.4) + \frac{1}{10} (0)$$

13 = 0.3

Consider decision tree given below, compute generalization over the rest of tree using noisy optimistic approach.



Training A & B C Class Just now

10

1

七

1

1

18

30

1

1

1

b) Compute the generalization error rate of the tree using pessimistic approach (for simplicity, we use strategy of adding a factor of 0.5 to each leaf node).

Training error + penalty term.

$$\text{Eg } C(T) = \underline{e(T)} + \underline{\alpha(T)} = \underline{e(T)} + \frac{\alpha(T)}{N_T}$$

$e(T)$

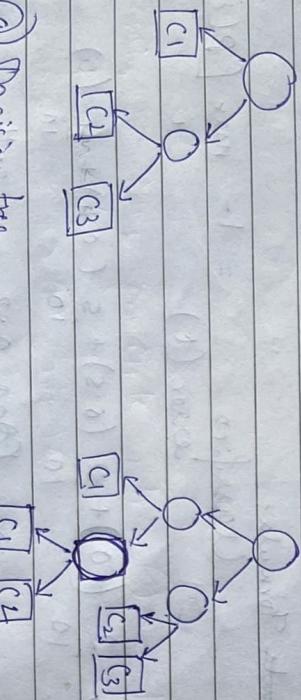
$$\text{Eg } C(T) = 0.3 + \frac{0.5 \times 0.5}{10}$$

$$= 0.35 + 0.05 = 0.4$$

c) Consider the decision shown below, assume they are generated from a dataset that contains 16 binary attributes & 3 classes

$$16 \text{ binary attributes} \leq C_1, C_2, C_3.$$

Which decision tree is better according to the MDL principle.



$$\rightarrow \text{cost}(\text{tree}, \text{data}) = \text{cost}(\text{tree}) + \text{cost}(\text{data}|\text{tree})$$

in attribute $\rightarrow \log_2 m$ bits
 $16 \text{ attributes} \rightarrow \log_2 16 = 4$ bits

k classes $\Rightarrow \lceil \log_2 k \rceil$ bits

3 classes $\Rightarrow \lceil \log_2 3 \rceil = 2$ bits

- a) Decision tree with 7 errors.
 b) Decision tree with 4 errors.

The cost of each misclassified error is $\log_2 n$ bits.

- a) Compute the total description length of each decision tree according to min. description length principle. The total description length of a tree is given by $\text{cost}(\text{tree}, \text{data}) = \text{cost}(\text{tree}) + \text{cost}(\text{data}|\text{tree})$. Each internal node of the tree is encoded by the ID of splitting attribute. If there are m attributes, the cost of splitting attribute is \log_m bits. Each leaf is encoded using ID of class; it is associated with. If there are k classes, the cost of encoding a class is $\log_2 k$ bits.
- $\text{cost}(\text{tree})$ is the cost of encoding all the nodes in the tree. To simplify, we compute you can assume that the total cost of the tree is obtained by adding up the cost of encoding each internal node of each leaf node.
 - $\text{cost}(\text{data}|\text{tree})$ is encoded using classification errors the tree commits on the dataset training set. Each error is encoded by $\log n$ bits where n is the total no. of training instances.

Fulleren nodes = $O(\text{circle})$
 Leaf nodes = \square

DOMS	Page No.
Date / /	Page No. / /

∴ Overall cost of decision tree (a) is $\Theta(n)$

$$2 \times 4 + 3 \times 2 + 7 \times \log_2 n$$

$$14 + 7 \log_2 n$$

∴ Overall cost of decision tree (b) is $\Theta(n)$

$$2 \times 4 + 5 \times 2 + 4 \times \log_2 n$$

$$\Rightarrow 16 + 10 + 4 \times \log_2 n$$

$$\Rightarrow 26 + 4 \log_2 n$$

Height (cm)	Weight (kg)	Class
167	51	underweight
182	62	normal
196	69	overweight
193	64	overweight
192	65	overweight
174	56	underweight
169	55	normal
173	57	normal
170	55	normal
170	57	normal

Test	170	170	57	?
Set C				

$$k=5$$

$$k=5$$

Predict class of above record.

$$n=2$$

$$(a) 14 + 7 \log_2 (2) = 21$$

$$(b) 26 + 4 \log_2 (2) = 30$$

∴ Cost of decision tree (b) is $\Theta(n)$

Q) Implement k-nearest neighbour classifier on dataset

- According to NID principle, tree (a) is better than
 (b) if $n < 16$
 worse than (b) if $n > 16$

Unit 4: Association analysis

→ Helps find hidden relationship among data.

e.g.: Customer Purchase Data.

Market Basket Data:

- 1) {Bread, milk}
- 2) {Bread, diapers, beer, eggs}
- 3) {Milk, diapers, beer, cola}
- 4) {Bread, milk, diapers, beer}
- 5) {Bread, milk, diapers, cola}

* Association rule:

Eg.: {Diaper} \rightarrow {Beer}.

Let $T = \{i_1, i_2, i_3, \dots, i_n\}$ be set of all items in market basket data.

$I = \{milk, bread, beer, diapers, cola, eggs\}$.

Ans

→ A collection of one or more items in a transaction is called an itemset.

If there are k -items in item set, it is called k -item set.

Eg.: {milk, bread, eggs} \rightarrow 3 item set

{beer, eggs} \rightarrow 2 item set

Support count: This refers to a no. of transactions that contain a particular item set.

Mathematically, represented using:

$$\text{Support count } S(X) = \left\{ |t_i| : X \subseteq t_i, t_i \in T \right\}$$

Support count {Beer, diapers, milk} = ?

Association rule: Implication expression of the form $X \rightarrow Y$ where $X \nsubseteq Y$ are disjoint sets i.e. $(X \cap Y) = \emptyset$

Strength of association rules:

Measured using two metrics: Support & Confidence.

Support: Determine how often a rule is applicable to a given dataset.

$$\text{Support} : S(X \rightarrow Y) = \frac{|XY|}{N}$$

Confidence: Determines how frequently items in Y appear in transactions that contain X .

$$\text{Confidence } C(X \rightarrow Y) = \frac{|XY|}{|X|}$$

e.g.: Consider the rule:

Wilk, Diaper \rightarrow {Beer}.

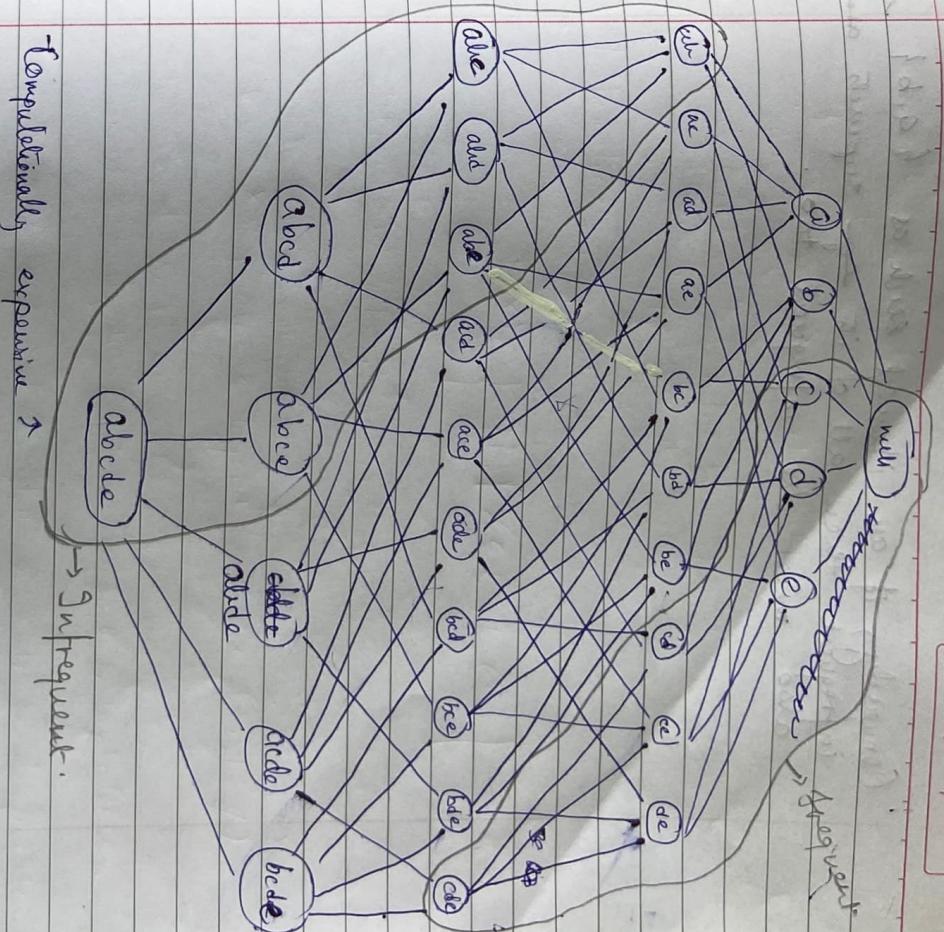
$$S = 2/5 \quad C = 2/3$$

* formulation of association rule mining problem:

\rightarrow Given a set of transaction T, find all the rules having support \geq minimum support and confidence \geq min. confidence where, minimum support \leq minimum confidence are corresponding support \leq confidence thresholds.

Algorithm:

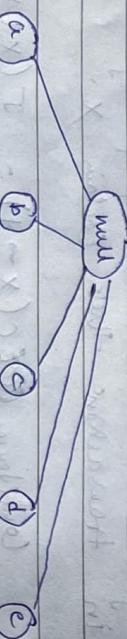
- ① frequent item set generation \Rightarrow to some
- ② Rule generation \rightarrow all high confidence rules are extracted from frequent item set.



- ① Reduce the number of candidate itemsets
- ② Reduce " " " comparison.

Lattice structure \Rightarrow

$$\text{e.g. } T = \{a, b, c, d, e\}$$



\Rightarrow If an item set is frequent then all of its subsets must also be frequent.

\Rightarrow E.g.: {c, d, e} \Rightarrow {c}, {d}, {e}, {c, d}, {c, e}, {d, e}, {c, d, e}.

Conversely) if an itemset such as {a, b} is infrequent then all of its supersets are also must be infrequent too.