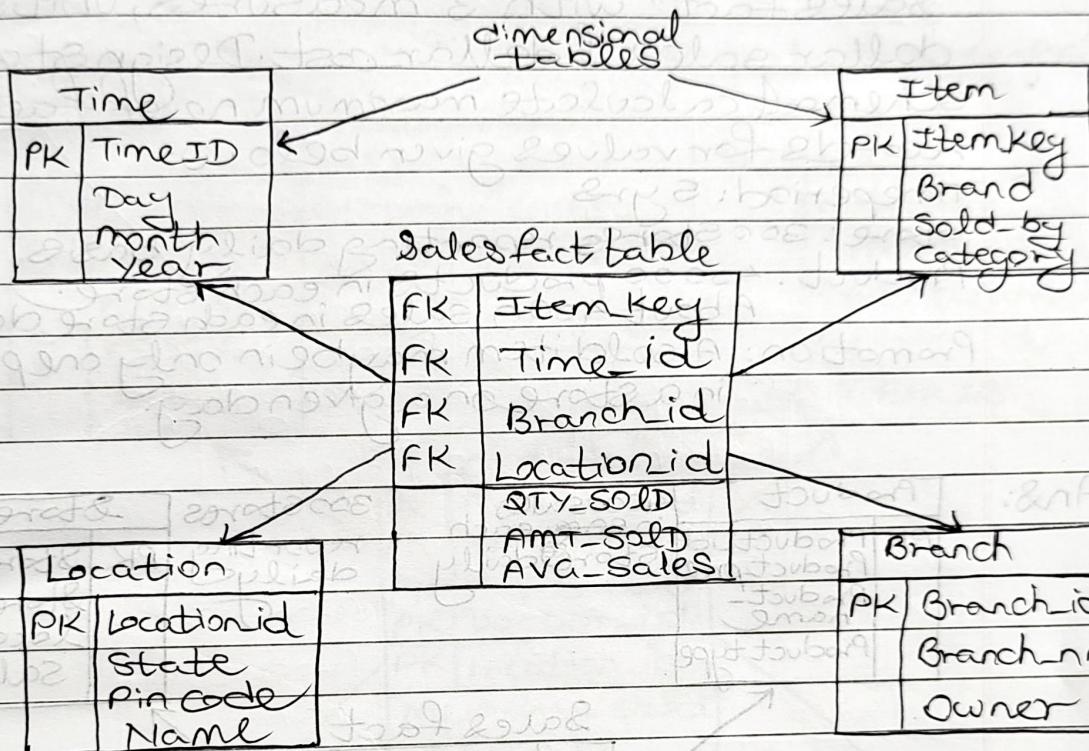


* Data warehousing & Mining *

18.7.2024

• Fact & dimensional table:

⇒ Star schema:



- Q1) Diff between er-modelling & dimensional modelling?
- Q2) Data warehouse architecture
- Q3) Star schema & Snowflake schema

Q1) for a supermarket chain, consider the following dimensions,

i) product, ii) store, iii) time, iv) promotion.

The schema consists a central facttable "Sales fact" with 3 measures, unitsales, dollarsales & dollarcost. Design star schema to calculate maximum no. of facttable records for values given below:

Time period: 5 yrs

Store: 300 stores reporting daily basis

Product: 4000 products in each store.
About 400 sales in each store daily.

Promotion: A sold item may be in only one promotion in a store on a given day.

Ans:

Product		4000 sales in every each store daily	300 stores reporting daily	Store					
PK	Product-id	Product-price	Product-name	Product-type	PK	Store-id	Store-emp	Location	Sales

Sales fact.

FK	Product-id
FK	Store-id
FK	Time-id
FK	Promo-id
	Unit-sales
	dollar-sales
	dollar-cost

Time	
PK	Time-id
	Day
	month
	Year

Time: 5 yrs

1 promotion/day

Promotion	
PK	Promo-id
	Promo-cost
	Promo-brand
	Duration

Time period = $5 \text{ yrs} \times 365 \text{ days} = 1825 \text{ days}$
 maximum no. of rows = $1825 \times 300 \times 4000 \times 1$
 (records) = 2 billion record will be stored in fact table

- Q2) Draw star schema for hospital management system.

Ans

Medicines		Patients	
PK	medicine_id	Patient-wt	
	name	Patient-name	
	stock	Patient_id	
	bill	Address	

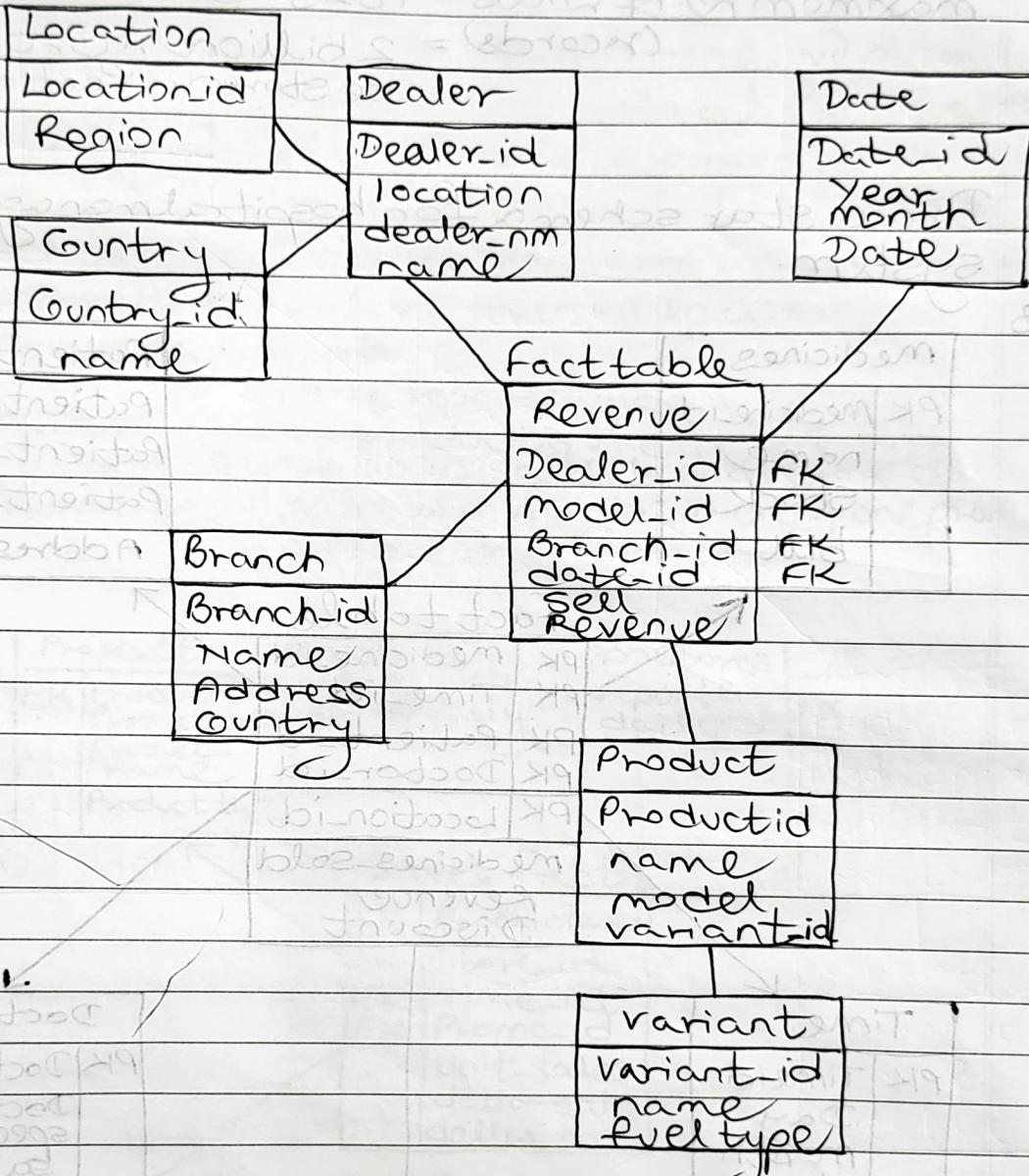
Fact table			
PK	medicine_id	Time_id	Patient_id
PK	Time_id	Patient_id	Doctor_id
PK	Patient_id	Doctor_id	Location_id
PK	Doctor_id	Location_id	medicines_sold
			Revenue
			Discount

Time	
PK	Timeid
	Day
	month
	Year

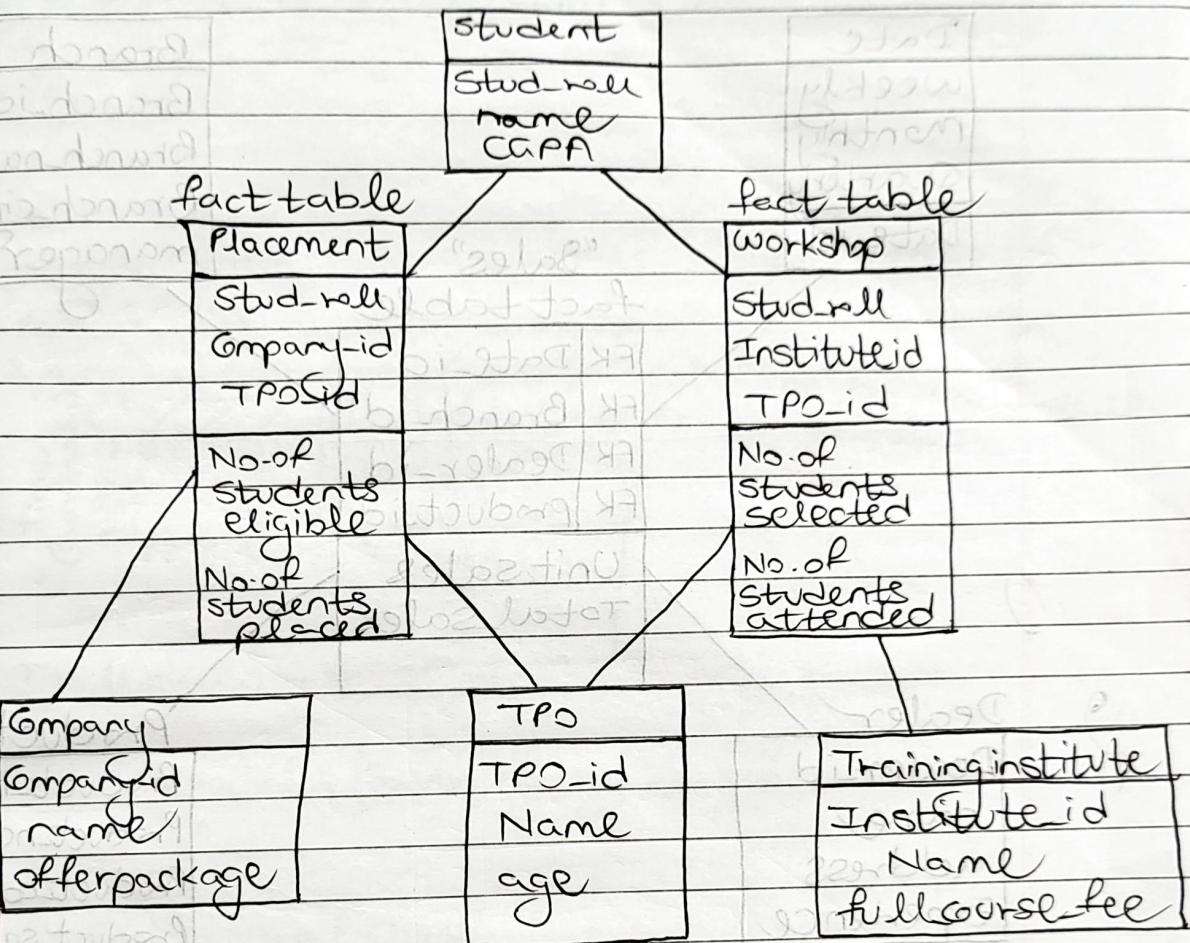
Doctor	
PK	Doctor_id
	Doctor_name
	specialization
	salary

Location	
PK	Location_id
	State
	pincode
	Name

2) Snowflake schema:-



* Fact Constellation in Data warehouse modelling:-



Problem Statement #2

All electronics may create a sales data warehouse in order to keep records of the store's sales with respect to dimensions Date, Dealer, branch & Product. These dimensions are allow to store to keep track of things like sales of items & items were sold.

Date
Weekly
Monthly
Quarterly
Date_id

Jan
Feb
Mar
Apr
May

Branch
Branch_id
Branch_name
Branch_city
manager

"Sales" fact table

FK	Date_id
FK	Branch_id
FK	Dealer_id
FK	productid
	Unitsales
	Total_sales

Dealer
Dealerid
name
address
experience

Product
Productid
Productname
Productcost
Productsales

Variant
Variantid
name
fueltype

- OLAP operations:-

Roll up (drill-up) : summarize data

Drill-down (roll down)

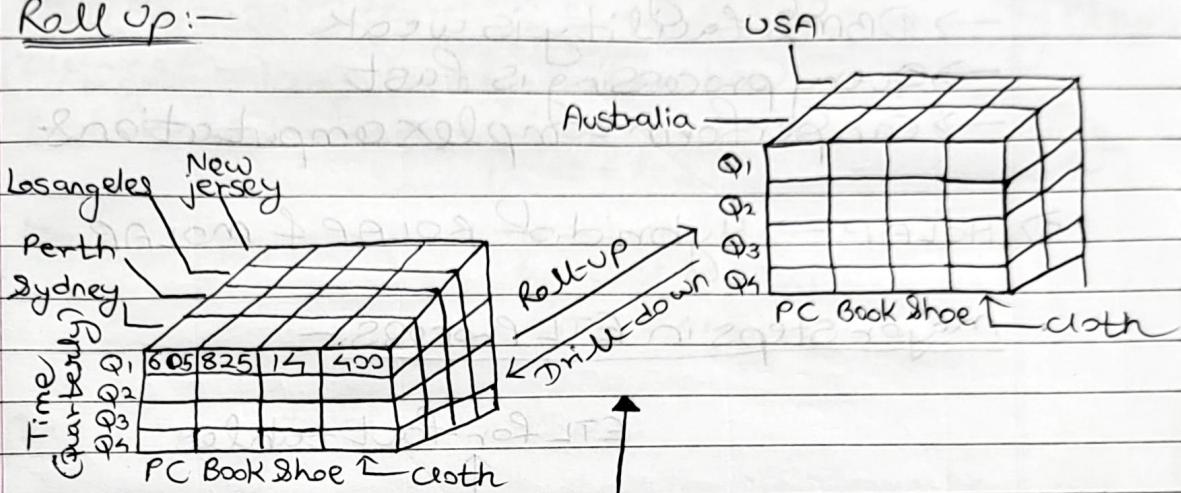
Slice: select

Dice: project

Pivot (rotate) : reorient the cube

Q) Difference between OLAP & OLTP?

I Roll up:-



II Drill down:-

* OLAP models/Servers:

1) ROLAP:

- Relational OLAP
- Highly scalable
- Relap tool analyze
- large amt of data
- ROLAP tools store & analyze highly volatile & changeable data
- Poor query performance
- DBMS feature is strong
- Requires expertise to use OLAP

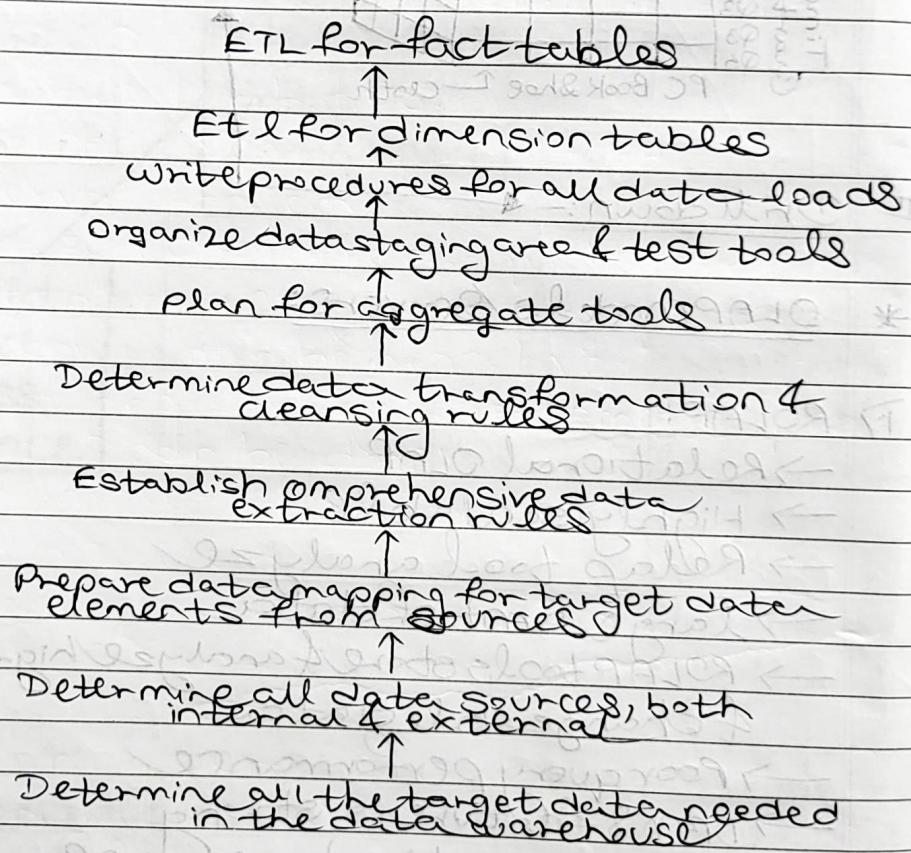
2) MOLAP:

- multidimensional OLAP
- MOLAP are not capable of containing detailed data

- It is fast.
- very easy to use
- DBMS facility is weak
- Query processing is fast
- can perform complex computations.

3) HOLAP → Hybrid of ROLAP & MOLAP

- major steps in ETL process:-



* Data mining - Unit 2:

IMP: KDD Process: Iterative sequence of following steps

```

    (Cleaning & integration) Clean   (Selection)
    Database → Data Transformation
    Database → Data Warehouse [Flatfiles] → Data mining
  
```

Knowledge ← evaluation & presentation
(patterns)

- Data mining task primitives:



↓
Knowledge type
to be mined:
characterization,
discrimination,
association,

Task relevant

data
database or

data warehouse

name database

tables or data

warehouse cubes

Conditions for data

~~selection~~

relevant attributes

or dimensions

Date grouping criteria

- What is an attribute?

Types:-

(i) Nominal: Symbols or names of things. Each value represents some kind of category, code or state & so nominal attributes:
eg - hair color, marital status, occupation
Non-quantitative

(ii) Binary attributes:

Symmetric: both of its states are equally valuable & carry same weight.
eg - gender

Unsymmetric: Opposite of symmetric

(iii) Ordinal attributes:

possible values that have an order but magnitude between successive values is not known. eg - Grade: A+, A, A-, B-

(iv) Numeric attributes:

Quantitative

→ interval scaled:

measured, ordered, equidistant, doesn't have any meaningful zero, interval can be negative

grouping, sorting & arithmetic operations can be done

example - thermometer

→ ratio scaled:

measured in form of numbers.

It has rank & order
equidistant

meaning zero

no negative values

example: distance travelled, weight, age

grouping, sorting, arithmetic

arithmetic op can be performed

mean, median, mode

Statistical description of data:

mean: center of set of data

$$\bar{x} = \frac{\sum x_i}{N}$$

Five number summary:

Q1	minimum	1 st quantile	median	2 nd quantile	max
	10	25th	50th	75th	59
		↓	↓	↓	
		25	33	36	

data points: 10 11 12 25 ~~28~~ 27 31 ~~33~~ 34
34 35 36 43 50 59

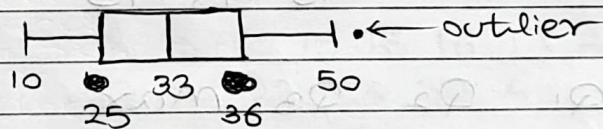
1st quantile = 25, 2nd quantile = 36

$$IQR = 36 - 25 = 11$$

outliers: Data value $< Q_1 - 1.5(IQR) < Q_3 + 1.5(IQR)$

Data value $> Q_3 + 1.5(IQR) > 52.5$

$$= 59$$



Q2) Dataset: 6 7 8 9 10 11 12

$$\sigma^2 = 4, \text{SD} = \sqrt{4} = 2$$

Q3) Data points: 5 10 13 15 16 16 20 20 21

22 22 ~~22~~ 25 25 25 ~~25~~ 30 33 33 35

35 35 35 36 40 45 46 52 70 85

$$\Rightarrow \text{mean} = 30.68$$

$$\text{median} = 25$$

$$Q_1 = 25^{\text{th}} \text{ percentile} = \frac{25(30)}{100} = 7.50$$

$$= 20 + 0.50(20 - 20) = 20$$

~~median = 80th~~

$$Q_3 = 75\text{th percentile} \neq \frac{75(30)}{100} = 22.5$$

$$= \frac{35+36}{2} = 35.5$$

Boxplot:

Min	Q1	Q2	Q3	Max
5	20	25	35	85

$$IQR = Q_3 - Q_1 = 35 - 20 = 15$$

Outliers: Data value $< Q_1 - 1.5(IQR) < 2.5$
 Data value $> Q_3 + 1.5(IQR) > 57.5$

Q5) Find median, mean mode for data:
 Show boxplot

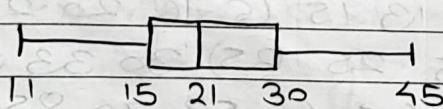
Ans Data points = 11 13 13 15 15 15 16 19 20

20 21 21 22 23 24 30 40
 45 45 45

Min	Q1	Q2	Q3	Max
11	15	21	30	45

$$IQR = 30 - 15 = 15$$

Outliers: Data point $< 15 - 1.5(IQR) = 7.5$
 Data point $> 30 + 1.5(IQR) = 52.5$



cab car cap

* Binning → dealing with noisy data

Q17 4, 8, 15, 21, 21, 24, 25, 28, 34 make 3 equal bins

Ans Bin 1: 4 8 15, Avg = 9

Bin 2: 15 21 21, Avg = 22

Bin 3: 24 25 28, Avg = 29

decade

(i) Smooth by Bin mean

Bin 1: 9 9 9

Bin 2: 22 22 22

Bin 3: 29 29 29

(ii) Smooth by Bin median:

Bin 1: 8 8 8

Bin 2: 21 21 21

Bin 3: 25 25 25

(iii) Smooth by Bin boundaries:

Bin 1: 4 4 9 5

Bin 2: 15 21 21

Bin 3: 24 24 28

Q27 Data points: 11 13 13 15 15 16 19 20 20 20 21
 21 22 23 24 30 40 45 45 45 71 72 73 75

make 4 equal bins

Ans Bin 1: 11 13 13 15 15 16, Avg = 13.80

Bin 2: 19 20 20 20 21 21, Avg = 20.10

Bin 3: 22 23 24 30 40 45, Avg = 30.60

Bin 4: 45 45 71 72 73 75, Avg = 63.50

(i) Smooth by Bin mean:

Bin 1: 13.8 13.8 13.8 13.8 13.8 13.8

2: 20.1 20.1 20.1 20.1 20.1 20.1

3: 30.6 30.6 30.6 30.6 30.6 30.6

4: 63.5 63.5 63.5 63.5 63.5 63.5

(ii) Smoothing by Bin median:

Bin 1: 14 14 14 14 14 14

2: 20 20 20 20 20 20

3: 27 27 27 27 27 27

4: 71.5 71.5 71.5 71.5 71.5 71.5

(iii) Smooth by Bin boundary:

Bin 1: 11 9 11 11 16 16

2: 19 19 19 21 21 21

3: 22 22 22 22 25 25

4: 45 45 75 75 75 75

Bayes Theorem:-

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

* ID3 - Decision Tree Classification

$$\begin{aligned} P = & \text{yes} = 9 \\ n = & \text{no} = 5 \end{aligned} \Rightarrow 14$$

$$\begin{aligned} I(p, n) &= -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n} \\ &= 0.940 \end{aligned}$$

Compute entropy - outlook

outlook	p	n	$I(p, n)$
Sunny	2	3	0.971
Overcast	5	0	0
Rain	3	2	0.971

$$\begin{aligned} E(\text{outlook}) &= \frac{5}{14} I(2, 3) + \frac{5}{14} \cdot I(5, 0) + \frac{5}{14} \cdot I(3, 2) \\ &= 0.694 \end{aligned}$$

$$\text{gain}(\text{outlook}) = I(p, n) - E(\text{outlook}) \\ = 0.940 - 0.694 = 0.246$$

compute entropy - temperature

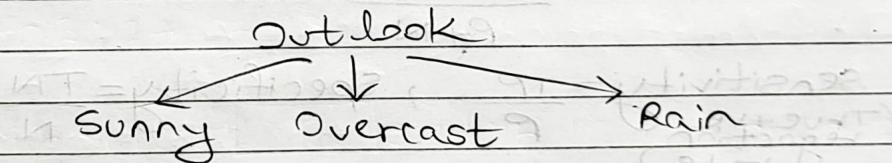
temperature	p	n	$I(p, n)$
hot	2	2	1
cold	3	1	0.81
mild	4	2	0.918

$$E(\text{temp}) = \frac{4}{14} I(2, 2) + \frac{5}{14} I(3, 1) + \frac{6}{14} I(4, 2) \\ = 0.91$$

$$\text{gain}(\text{temp}) = I(p, n) - E(\text{temp}) = 0.940 - 0.91 \\ = 0.03$$

$$\text{gain}(\text{humidity}) = 0.940 - 0.786 = 0.154$$

$$\text{gain}(\text{wind}) = 0.048$$



- Using If-Then Classification:

The If part of the rule is called antecedent or precondition.
The part is called rule consequent.

Assessment of a rule:

$\text{N}_{\text{covers}} = \# \text{ of tuples covered by } R$

$\text{N}_{\text{correct}} = \# \text{ of tuples correctly classified by } R$

$\text{Coverage}(R) = \text{N}_{\text{covers}} / |\text{D}|$, D = training data set

$\text{accuracy}(R) = \text{N}_{\text{correct}} / \text{N}_{\text{covers}}$.

X Q17 R1 which covers 2/14 tuples correctly classified
tuple 2. $\rightarrow P_0 \cdot 0 - C_0 \cdot 0 =$

Ans average = $\frac{2}{14} = 0.14$

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

• Confusion matrix:

		predicted value		Type I error
		N	P	
N	True Negative(TN)	False Positive(FP)		
	False Negative(FN)	True Positive(TP)		

↑ Type 2 error

Accuracy = $\frac{TP + TN}{P + N}$

Error rate = 1 - accuracy
 $= \frac{FP + FN}{P + N}$

Sensitivity = $\frac{TP}{P}$, Specificity = $\frac{TN}{N}$
 (True positive rate)

Precision = $\frac{TP}{TP + FP}$, Recall = $\frac{TP}{TP + FN}$

measure of exactness. What % of tuples labelled as +ve & are actually +ve?

Recall is measure of completeness. What % of positive tuples are labelled as +ve?

23-08-2024

(4)

60 50 10

105 100

- Imp questions:

Q1) Total types in dataset is 165. Among which 60 was known classified as 'No' & 105 are classified as 'Yes'. After applying a model, 50 were predicted 'No' & 100 were predicted 'Yes'.

Design confusion matrix, accuracy, error rate, precision & recall.

Ans Predicted

		N	Y
		50	10
N	5	105	50
	100	5	100

$$\text{Accuracy} = \frac{100 + 50}{165} = 0.90 = 90\%$$

$$\text{error rate} = 1 - \text{accuracy} = 1 - 0.90 = 10\%$$

$$\text{precision} = \frac{100}{105} = 0.90 = 90\%$$

$$= 110$$

$$\text{recall} = \frac{100}{110} = 90\% = 95\%$$

- * Decision tree problems:

Q1) [Gatesmashers]

Ans (i) IG of weather

Entropy of entire dataset

$$\Rightarrow S[+9, -5] = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

entropy of attributes:

$$\text{Sunny: } S[+2, -3] = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

$$\text{Cloudy: } S[+4, 0] = 0$$

$$\text{Rain: } S[+3, -2] = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$$

Information gain = Entropy (whole data)

$$-\frac{5}{14} \text{Ent}(S) - \frac{4}{14} \text{Ent}(C)$$

$$-\frac{5}{14} \text{Ent}(R)$$

$$= 0.246$$

(ii) IG of Temperature:

Entropy of entire dataset

$$\mathcal{S}[+9, 5] = 0.94$$

Entropy of all attributes:

$$\text{Hot: } \mathcal{S}[+2, -2] = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 1.0$$

$$\text{Mild: } \mathcal{S}[+4, -2] = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.91$$

$$\text{Cold: } \mathcal{S}[+3, -1] = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{1}{5} \log_2 \frac{1}{5} = 0.81$$

$$\text{Information gain} = 0.94 - \frac{4}{14} \times 1 - \frac{6}{14} \times 0.91$$

$$- \frac{4}{14} \times 0.81 = 0.029$$

(iii) IG of humidity:

Entropy of entire dataset

$$\mathcal{S}[+9, 5] = 0.94$$

Entropy of all attributes:

$$\text{High: } \mathcal{S}[+3, -4] = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.98$$

$$\text{Normal: } \mathcal{S}[+6, -1] = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.59$$

$$\text{Gain} = 0.94 - \frac{7}{14} \times 0.98 - \frac{7}{14} \times 0.59 = 0.15$$

(iv) IG of wind:

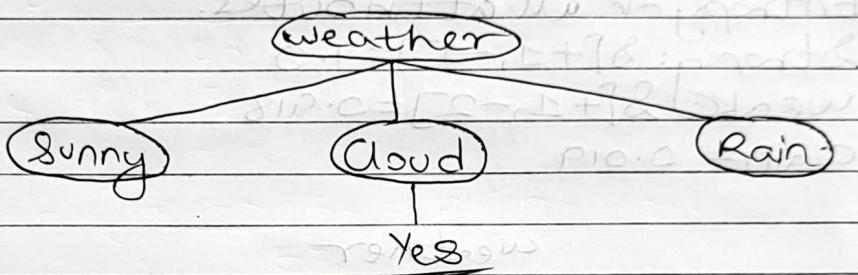
Entropy of entire dataset = 0.94

Entropy of all attributes:

$$\text{Strong: } S[+3, -3] = \frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1.0$$

$$\text{Normal: } S[16, -2] = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.81$$

$$I\text{-Gain} = 0.94 - \frac{6}{15} \times 1.0 - \frac{8}{15} \times 0.81 = 0.0478$$



Calculating for sunny:

(i) IG of temperature:

$$\text{Entropy of sunny} \Rightarrow S[+2, -3] = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

Entropy of all attributes:

$$\text{Hot: } S[+2, -2] = 0$$

$$\text{mild: } S[+1, -1] = 1.0$$

$$\text{Cold: } S[+1, -0] = 0$$

$$\begin{aligned} I\text{-Gain of temp} &= 0.97 - \frac{2}{5} \times 0 - \frac{2}{5} \times 1 - \frac{1}{5} \times 0 \\ &= 0.57 \end{aligned}$$

(ii) IA of humidity:

$$\text{Entropy of sunny} = 0.97$$

Entropy of all attributes:

$$\text{High: } S[+1, -3] \approx$$

$$\text{Normal: } S[+2, -0] = 0$$

$$\text{Gain} = 0.97$$

(iii) IA of wind:

$$\text{Entropy of wind} = 0.97$$

Entropy of all attributes:

$$\text{Strong: } S[+1, -1] = 1.0$$

$$\text{weak: } S[+1, -2] = 0.918$$

$$\text{Gain} = 0.019$$

