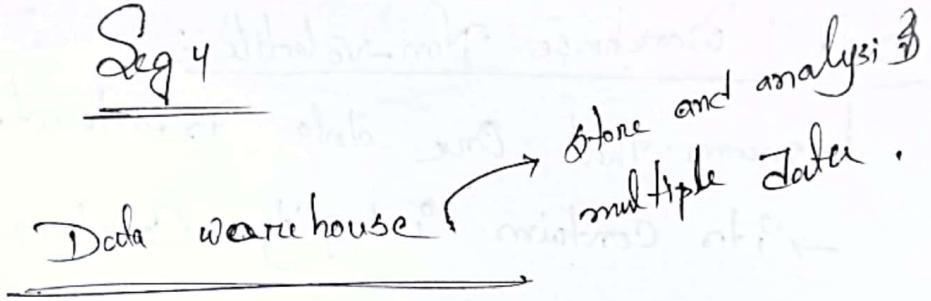


Sug 4



It is a subject oriented collection of data in support of management decision making process.

(Management decision make more obj of MIS)

Subject oriented Data Warehouse

→ organized around major

subject such as

Customer, product, Sales,

Focus are on decision
makers are obj of MIS

Data warehouse integrated

→ constructed by integrating multi-
ple, heterogeneous data sources

(Data integration, cleaning are
applied)

Data warehouse Time variant: time horizon → duration.
for this is significantly

→ historical data provide (5-10) years longer than that of operational system.

Data Warehouse Non-volatile:

↳ means that One data is entered. It's not changed.
 → It contains Integrity, Consistency.

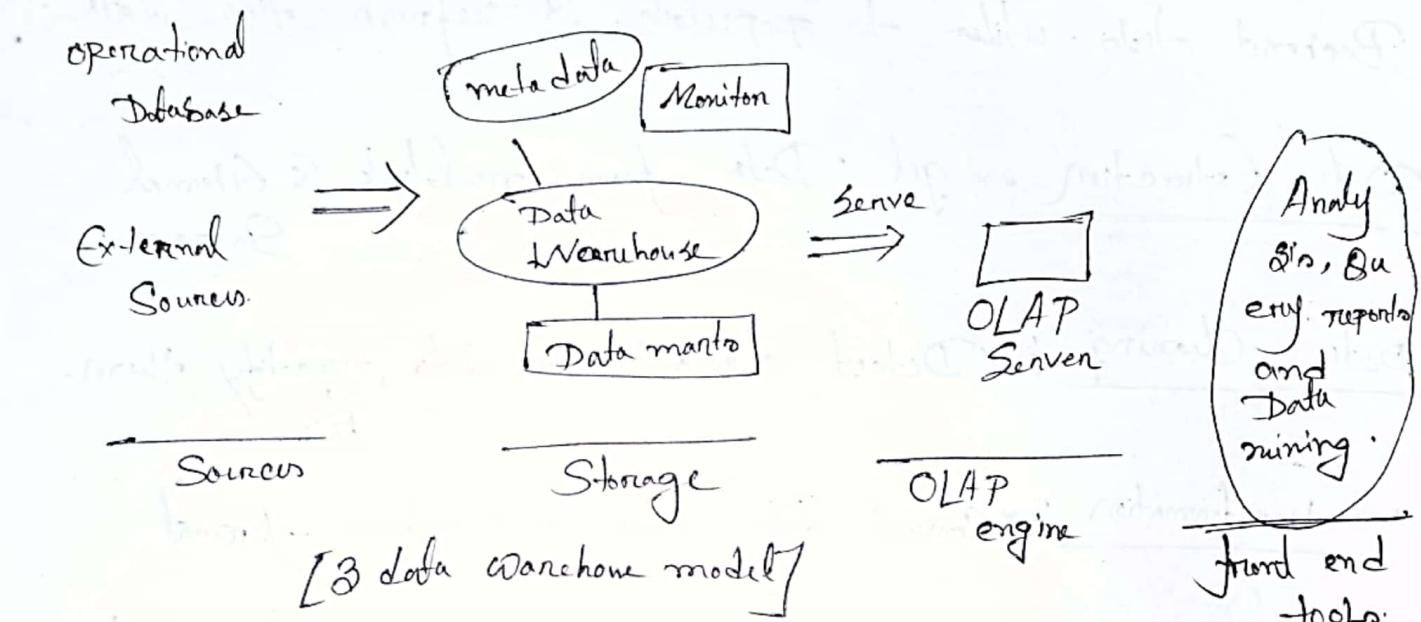
→ Database data daily update \Rightarrow ~~can't~~ \Rightarrow ~~can't~~ So

Compare decision makes better \Rightarrow

OLTP vs OLAP

| OLAP | OLTP |
|--|---|
| ① Online analytical processing | ① Online Transaction Processing |
| ② Consists of historical data | ② Consists of only current data |
| ③ Makes use of Data Warehouse | ③ It makes uses DBMS. |
| ④ Subject oriented used for Data mining, Decision making. | ④ Application oriented used for business tasks. |
| ⑤ It servers extract info for analysis | ⑤ Servers the purpose to insert/dele update info from database |
| ⑥ Only read and rarely write operation. | ⑥ Both read and write operation. |

Multi-tier Data warehousing Architecture



Data Mart: Smaller, focused version of a data warehouse designed for specific department → customer, item, sales analysis etc.

Virtual Data warehouse: Logical data warehouse that provides view of the physically storing.
 (budget friendly)
 [Or database known virtual data warehouse].

Enterprise Data warehouse: Provide constant data source for the entire data organization.

→ Integrate from multiple sources

→ Supports enterprise-wide data analysis

Ex: Your sales company Sales, and custom data analyze to All regions etc.

Backend Tools

Backend tools used to populate & refresh tier data.

Data Extraction → get Data from multiple & external Sources.

Data Cleaning : → Detect errors in the data, rectify them.

" Transformation " → Convert data from → Warehouse format

" Landing " : Sort, Summarize, Compute views, clean & integrate

" refresh " : Update from the data warehouse

Meta data Repository

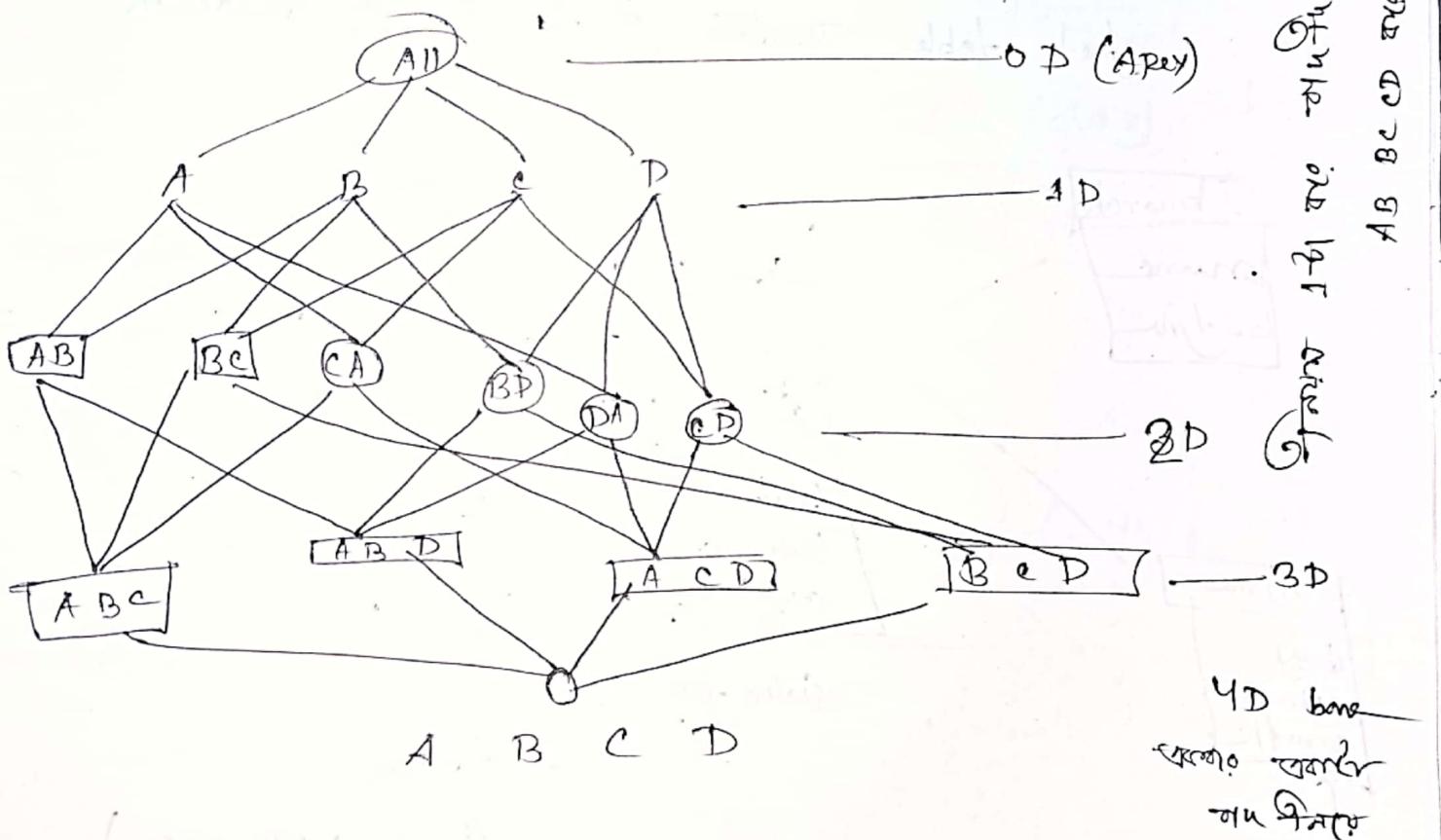
It provides detailed info about the structure, definition, and attributes of the data within a system.

It includes details like data source, type, relationship, understand and manage more effectively.

Data cubes

Multidimensional data model views data in the form of cube

- Literature, an n-d base called a base cuboid
- Top most 0-D cuboid hold highest level summarization
called grey Cuboid.



Data cube (Lattice of Cuboids)

fact

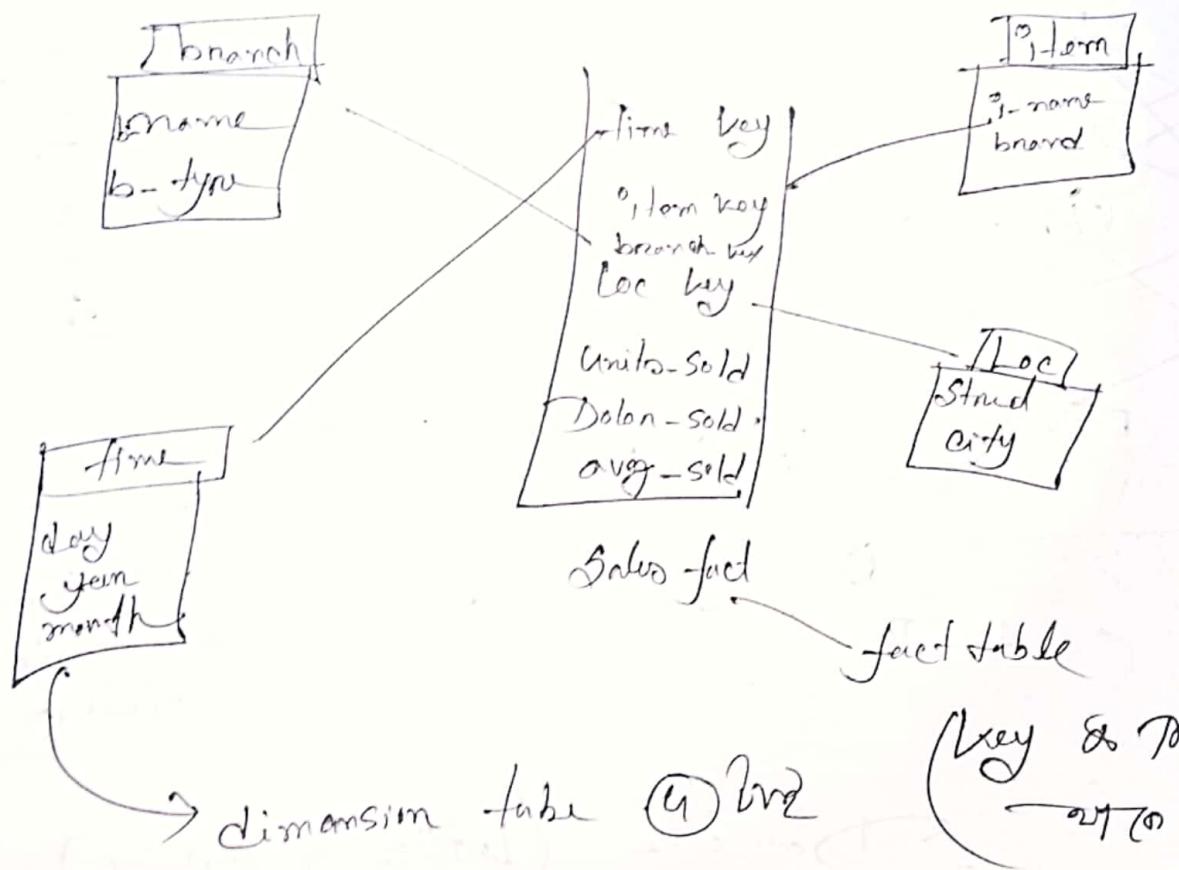
Dimension → Order Sale, Product name on time (day wise)

fact-table → Contains measure (Dolan + Sold)

key of each of the related dimension tables

Store Schema

fact-table middle of the set of dimension

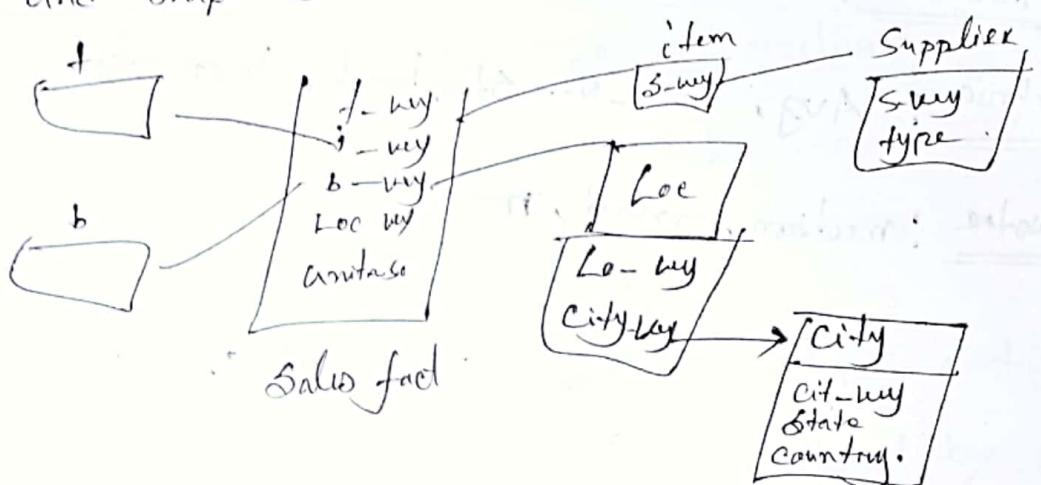


(Key & Measure)
→ (① ②)

→ dimension table ④ ⑤ ⑥

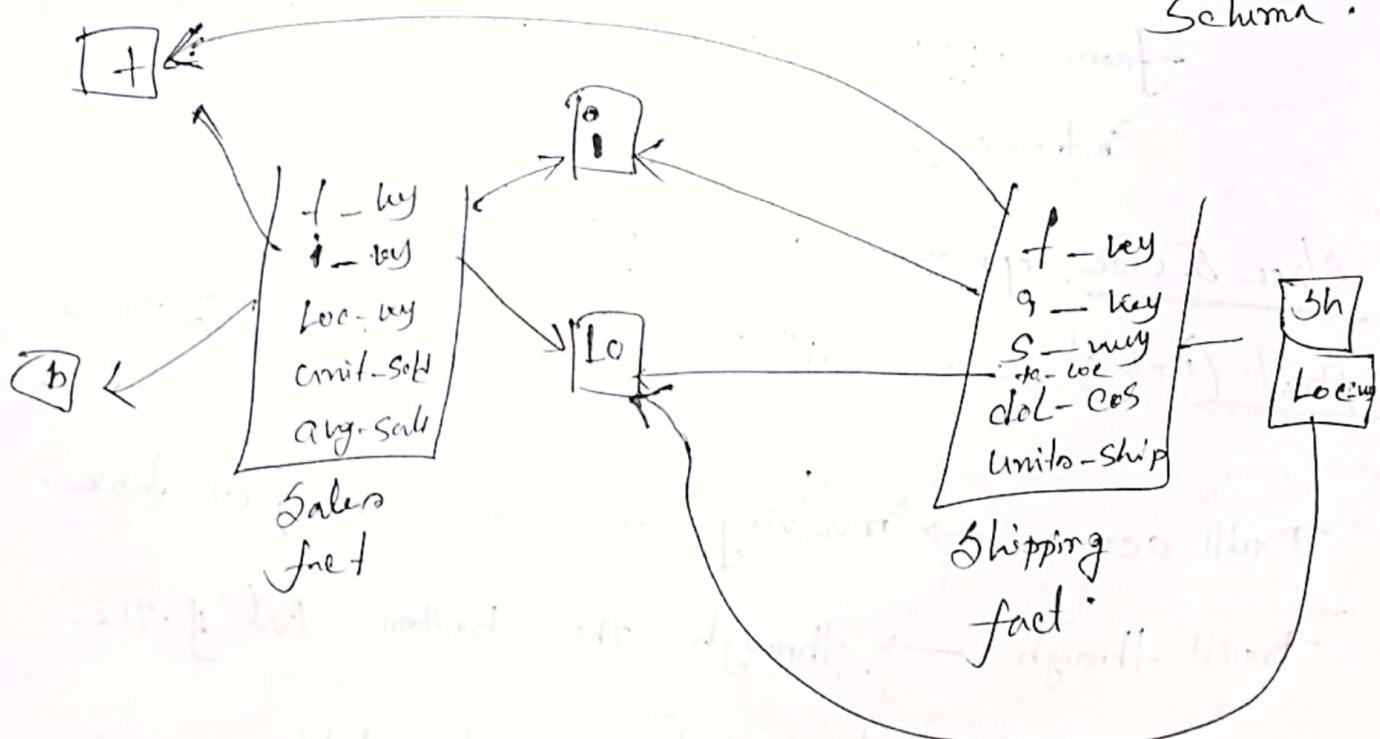
Snowflake Schema

Star Schema is normalized into a set of small dimension tables and shape similar to snowflake.



fact - constellation / galaxy Schema

multiple fact tables share dimension tables, called galaxy schema.



Data Cube Measure

Distributive: Count, Sum, min, max

Algebraic: Avg, Min-N - Standard deviation.

Holistic: median, mod, rank

Typical OLAP Operation

roll up → Summarize data ← dimension reduction.

Drill down (zoom down) → reverse of roll up.

- from higher lvl summary to lower lvl summary
Introducing new dimension.

Slice & dice: project and select

Pivot (rotate) → visualization, 3d to sever. of 2D plans

Drill across → involving more than one fact table

Drill through → through the bottom lvl of the cube to its back-end tables.

Design of Data Warehouse

Data friendly warehouse

Design of Data Warehouse

4 views regarding the design:

Top Down view: Allow selection of the relevant info necessary for the data warehouse to match current & future needs
~~and analytical~~
~~multidimensional~~
~~should~~

Data Source view: the info captured by stored and managed by op system. It's modeled by CASE and entity relationship model. tools

to provide details understanding and data are fully integrated

Data Warehouse View: Consists of fact & Dimension table represents info that stored in the warehouse

Design
After, Breakline on glory
Selma
for analysis

Business Query View: See the prospective of Data from end user

Detailed analysis
Financial Dashboard
ensure the supports of the analytical and responding needs the user.

Design process

① Top Down, bottom up approaches are a combination of both.

- Start with overall system design and planning
- Starts with prototype & experiments (rapid)

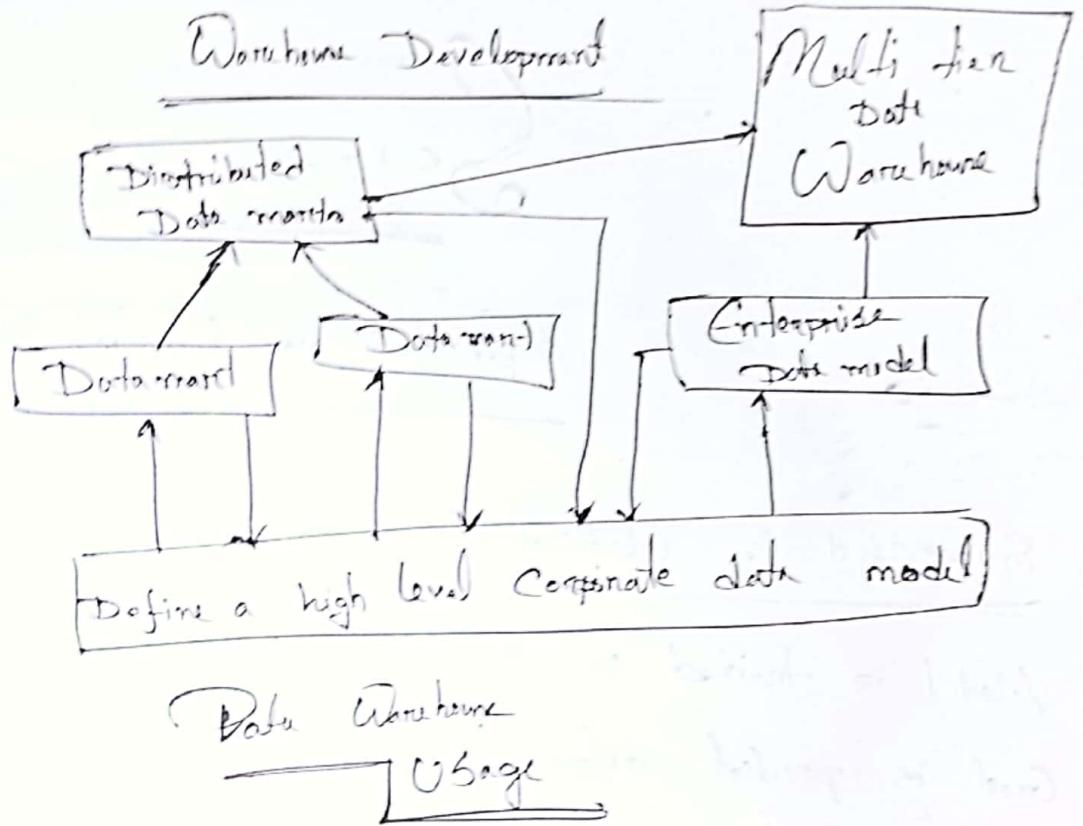
② From Dev eng stand of view:

Waterfall: Structured and systematic analysis at each step before proceeding to the next.

Spiral: Increased fun system, Short turn around time

③ Typical Data Warehouse design process:

- choose a business process to model → vendor, supplier
- choose a atomic lvl data of bus. proc → individual transaction, individual daily snapshot
- choose a dimension that will apply to each fact table record → typical dimension are item, customer, supplier
- choose a measure that will populate each fact table → units sold, quantity sold



3 kinds of warehouse application.

① Information processing

- Supports querying, basic statistical analysis, reporting tables, charts and graphs.

② Analytical processing

- Multidimensional analysis of Data warehouse data.
- Supports basic OLAP operations, drilling, pivoting

③ Data mining

- Knowledge discovery from hidden patterns.

- Performing classification and prediction
- Presenting the mining results using visualization tools

Seg - 5

Supervised vs Unsupervised Learning

Supervised is classification

Model is trained on labeled data.

Goal is predict output for new input.

Ex - Linear Regression, Logistic Regression, SVM, NN

unsupervised (clustering) → grouping similar data.

where model works with unlabeled data.

Algorithm tries to find patterns, structures or relationships.

Ex → K-means clustering, Autoencoder.

→ Classification is discrete categories (yes or no)

→ Numeric prediction is continuous numeric value

Classification

Numeric Prediction

→ Accuracy, recall, precision

→ Mean Squared, R² score

→ predict new inputs
output

→ Predict unknown or (loan approval)
missing values.

Classification 2 step

① Model Construction (Training)

- Build a model using a labeled ~~no~~ training set.
 - * Model represents as classification rules, dt or mathematical formulae.
 - * Data consists of samples with predefined class labels.

② Model usage (Testing) → prediction

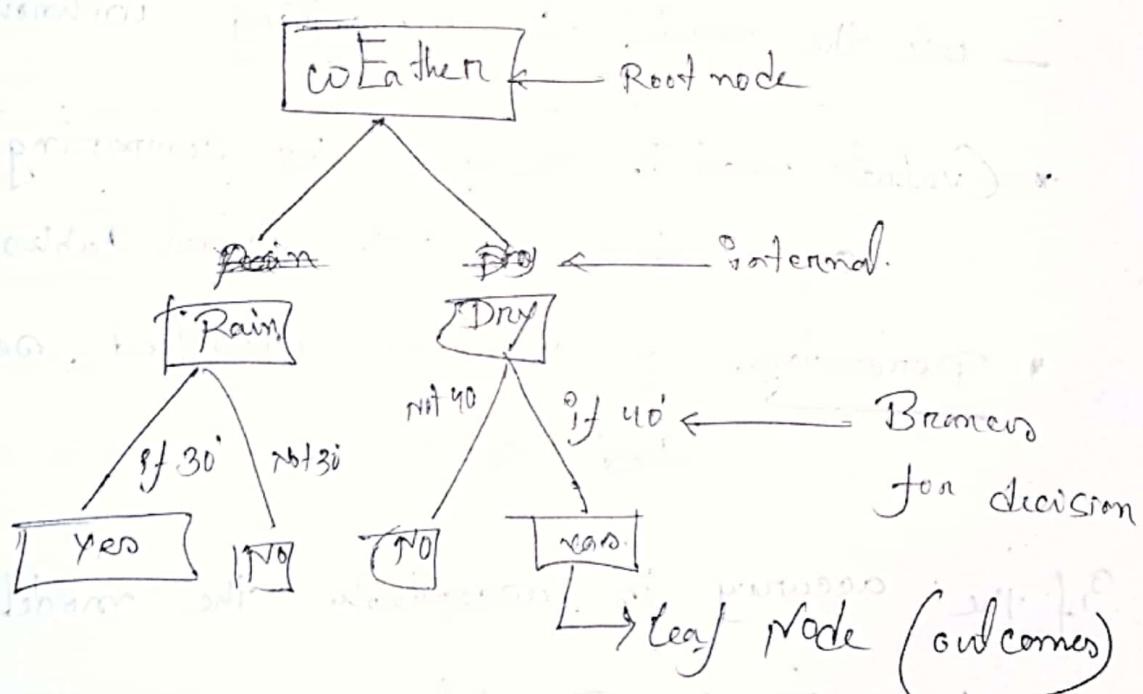
- use the model for classifying unknown data.
- * Evaluate model's accuracy by comparing prediction on a test set with actual labels.
- * Percentage of correctly classified samples in the test set.

If the accuracy is acceptable → the model is used.
to classify new data.

If test set used for model selection \Rightarrow 1 is
Called a validation set. (test set)

Decision tree

It is a flowchart structure where each internal node represents a test on an attribute. each branch represents the outcome of a test, each node represents a class label (classification) or value (regression).



DT Example

| Day | Weather | Temp | Humidity | Wind | Play |
|-----|---------|------|----------|------|------|
| 1 | S | H | H | w | N |
| 2 | S | H | H | s | N |
| 3 | C | H | H | w | Y |
| 4 | R | M | H | w | Y |
| 5 | R | C | N | w | Y |
| 6 | R | C | N | s | N |
| 7 | C | C | N | s | Y |
| 8 | S | M | H | w | N |
| 9 | S | C | N | w | N |
| 10 | R | M | N | w | Y |
| 11 | S | M | N | s | Y |
| 12 | C | M | A | s | Y |
| 13 | C | H | N | w | Y |
| 14 | R | M | H | s | N |

- (1) $\text{Gain}_{\text{Temp}}$
(2) $\text{Gain}_{\text{Humidity}}$

Step 1
Weather

Calculate Information Gain

- Entropy of entire dataset

here $y = 9$, $N = 14$ $\rightarrow \text{Total} = 14$

$$\text{Entropy} = -J \log_2 J$$

$$S(f_9, f_5) = - \frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94.$$

- entropy of all attributes

$$\text{Sunday entropy } (+2, -3) = - \frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

$$\text{Cloudy } (+4, 0) = - \frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$\text{Rain } (+8, -2) = - \frac{8}{14} \log_2 \frac{8}{14} - \frac{2}{14} \log_2 \frac{2}{14} = 0.97$$

whole tree size(m) $\frac{21 \text{ leaf}}{\text{Total}}$ first

Information gain = Entropy(whole data) - $\frac{5}{14} \text{ent}(S) - \frac{9}{14} \text{ent}(C) \frac{5}{14} \text{ent}(R)$

weather

$$= \cancel{0.94} - (0.97 - 0 = 0.97)$$

$$= 0.94 - \frac{5}{14} \times 0.97 - \frac{9}{14} \cdot 0 - \frac{5}{14} \times 0.97$$

$$= 0.2471$$

② calculate Ig for Temp

$$\{H, S, C\} = 0.94 \leftarrow \text{whole Ent. } (\text{Same, always})$$

$$\text{Ent of Hot } \{+2, -2\} = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1.0$$

$$\text{,, " mild } \{+4, -2\} = -\frac{4}{8} \log_2 \frac{4}{8} - \frac{2}{6} \log_2 \frac{2}{6} = 0.92$$

$$\text{,, " cold } \{+3, -1\} = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.81$$

$$\text{Ig for Temp} = \text{whole ent} - \frac{4}{14} E(H) - \frac{6}{14} E(M) \\ - \frac{4}{14} E(C)$$

$$= 0.94 - \frac{4}{14} \times 1.0 - \frac{6}{14} \times 0.92$$

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

$$= 0.029$$

③ Some term Humanity

$$\text{whole } \epsilon = 0.94.$$

$$\textcircled{e} \text{ High (+3,-4)} = 0.98$$

$$\textcircled{e} \text{ Normal (-6,-1)} = 0.59$$

$$Ig \text{ for Humanity} = \epsilon \text{ of whole} - \frac{7}{14} \epsilon(H) - \frac{7}{14} \epsilon(N)$$
$$= 0.15$$

④ wind

$$\text{whole } \epsilon = 0.94$$

$$\textcircled{e} \text{ Strong (+3,-3)} = 1.0$$

$$\epsilon \text{ } \cancel{\text{Normal}} \text{ } \cancel{(-6,-2)} = \cancel{-} 0.81 \\ \text{weak.}$$

$$Ig \text{ of wind} = \text{whole } \epsilon - \frac{6}{14} - \frac{8}{14} * \epsilon(w)$$
$$= 0.94 - \frac{6}{14} \times 1.0 - \frac{8}{14} \times 0.81$$
$$= 0.0478$$

$$\text{Grain } (S, \text{ weather}) = 0.246 \rightarrow \text{so root node.}$$

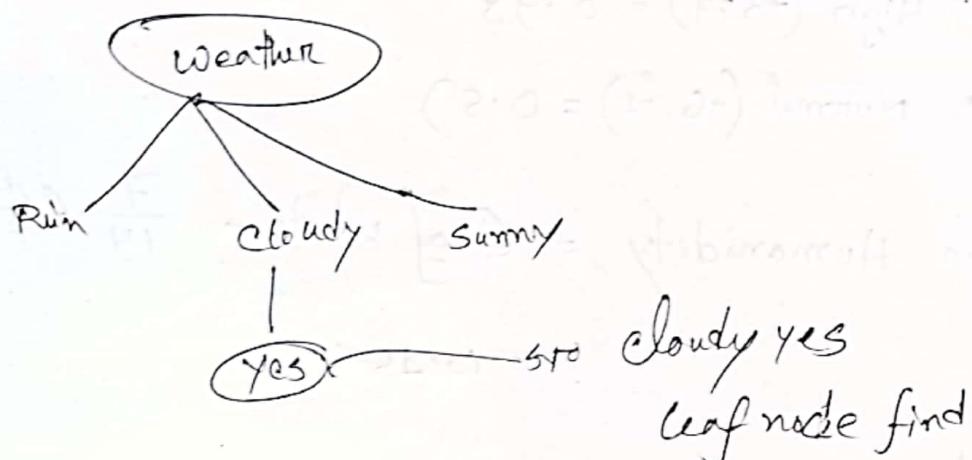
$$\text{, } (S, \text{ Temp}) = 0.029$$

$$\text{, } (S, \text{ Humanity}) = 0.15$$

$$\text{, } (S, \text{ wind}) = 0.0478$$

So weather node root.

Cloudy, Rain, Sunny



So start Rain and Sunny go Ig ck.

Sunny

Rain Ig

(Weather Root so next 3 in)

Temperature

$$\text{Ent for sunny } (+2, -3) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.97$$

$$\text{Ent (H)} \quad (-0, -2) = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log \frac{2}{2} = 0$$

$$\text{Ent (M)} \quad (+1, -1) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

$$\text{Ent (C)} \quad (+1, -0) = -\frac{1}{1} \log \frac{1}{1} - \frac{1}{1} \log \frac{1}{1} = 0$$

$$\begin{aligned}\text{Ig of Temp} &= \text{whole} - \frac{2}{5} \text{Ent (H)} - \frac{2}{5} \text{Ent (M)} - \frac{1}{5} \text{Ent (C)} \\ &= 0.57\end{aligned}$$

Humidity

$$\textcircled{e} \text{ whole sunny} = 0.97$$

$$e(H) = 0$$

$$e(N) = 0$$

$$Ig = 0.97 - 0 - 0 = 0.97$$

wind

$$e \text{ whole sunny} = 0.97$$

$$e(W) = 1$$

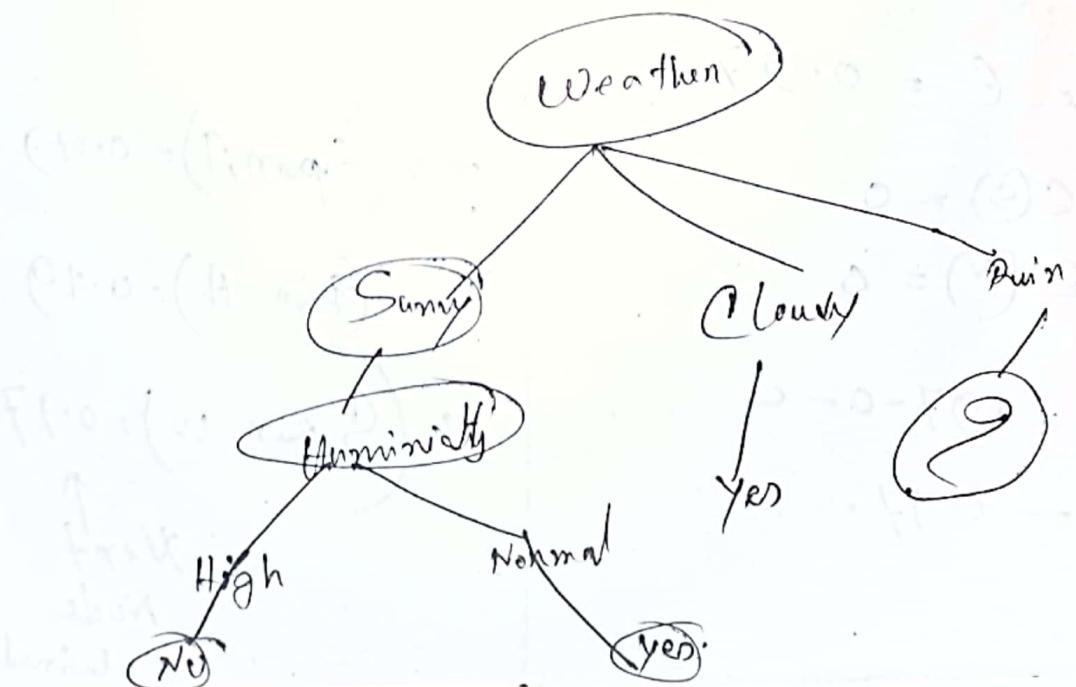
$$e(w) = 0.918$$

$$Ig(w) = 0.97 - 1 \times \frac{2}{5} + \frac{3}{5} \times 0.918 \\ = 0.019$$

$$G_e(S_{\text{sunny}}, T) = 0.957$$

$$G_e(S_{\text{sunny}}, H) = 0.97$$

$$G_e(S_{\text{sunny}}, w) = 0.019$$



for Rain

② Temp = $(+3 - 2) = 0.97$

$\epsilon(H) = 0$

$\epsilon(N) = 0.918$

$\epsilon(C) = 1.0$

$$Dg \text{ of Temp} = 0.97 - \frac{0}{5} \cdot 0 - \frac{3}{5} \cdot 0.918 - \frac{2}{5} \cdot 1 \\ = 0.019$$

Humidity

whole $\epsilon = 0.97$

$\epsilon(H) = 1$

$\epsilon(N) = 0.918$

$$Dg(\text{Hum}) = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.918 \\ = 0.019$$

Wind

whole $\epsilon = 0.97$

$\epsilon(S) = 0$

$\epsilon(N) = 0$

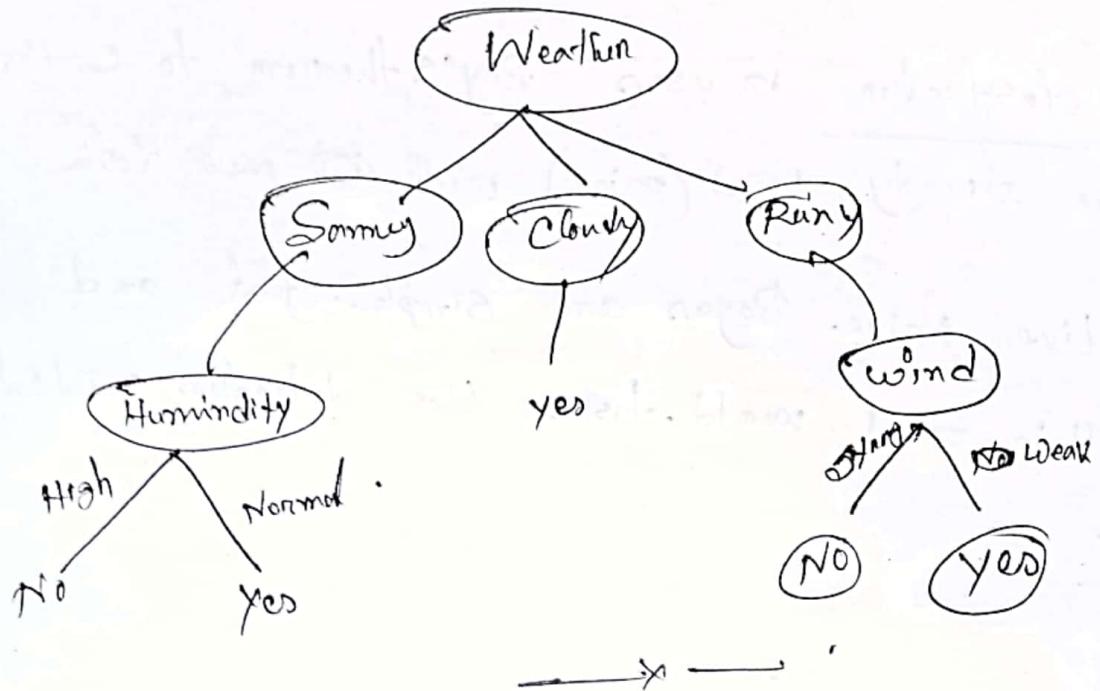
$$Dg(w) = 0.97 - 0 - 0 \\ - 0.97 \cdot$$

gain(S from T) = 0.019

" (S from H) = 0.019

gain(S from w) = 0.97

↑
Next
Node
Wind



Measure impurity

$$\text{Entropy} = \sum_j -P_j \log_2 P_j$$

$$\text{Gini Index} = 1 - \sum P_j^2$$

$$\text{Classification Error} = 1 - \max\{P_j\}$$

Bayesian classification uses Bayes' theorem to combine what we already know (prior) with new data.

Methods like Naive Bayes are simple, fast, and work well in real world tasks like detection, medical diagnosis.

Bayes theorem

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

$$\rightarrow = P(x|H) \times P(H) / P(x)$$

Let x be unknown.

Naive \rightarrow variables are independent

Perform classification tasks.

feature data will be independent

→ to make decision tasks simple Naive

Naive Bayes \rightarrow classify them

Naïve Bayes

| <u>Day</u> | <u>Outlook</u> | <u>Temp</u> | <u>Humidity</u> | <u>Wind</u> | <u>Play</u> |
|-----------------|----------------|--------------------------|-----------------|-------------|-------------|
| D ₁ | S | H | H | S | N |
| D ₂ | S | H | H | W | N |
| D ₃ | O | H | H | W | Y |
| D ₄ | P | M | H | W | Y |
| D ₅ | P | C | H | S | N |
| D ₆ | P | C | H | W | N |
| D ₇ | O | M | H | W | N |
| D ₈ | S | M | H | W | Y |
| D ₉ | S | M | N | S | Y |
| D ₁₀ | P | M | N | S | Y |
| D ₁₁ | S | M | N | W | Y |
| D ₁₂ | O | M | H | - | N |
| D ₁₃ | O | H | R | - | - |
| D ₁₄ | P | [O=s, T=c, H=M, W=W] P=? | | | |

Probability play? Yes? No?

Ans

$$\text{play(Yes)} = \frac{9}{14} = 0.64$$

$$\text{play(No)} = \frac{5}{14} = 0.36$$

| <u>Outlook</u> | <u>Yes</u> | <u>No</u> |
|----------------|------------|-----------|
| Sunny | 2/9 | 3/5 |
| overcast | 4/9 | 0/5 |
| rainy | 3/9 | 2/5 |

| <u>Temp</u> | <u>Yes</u> | | <u>No</u> | | <u>Humidity</u> |
|-------------|------------|-----------|------------|-----------|-----------------|
| | <u>Yes</u> | <u>No</u> | <u>Yes</u> | <u>No</u> | |
| Hot | 2/9 | 2/5 | Y | N | |
| Cold | 4/9 | 2/5 | High | 3/9 | 4/5 |
| Cool | 3/9 | 1/5 | Non | 6/9 | 1/5 |

wind

| | Y | N |
|--------|-----|-----|
| Strong | 3/9 | 3/5 |
| weak | 6/9 | 2/5 |

* $P(\text{yes}|x) = P(\text{yes}) \cdot P(\text{sunny}|\text{yes}) \cdot P(\text{cool}|\text{yes})$

$$P(\text{high}|\text{yes}) \cdot P(\text{strong}|\text{yes})$$

$$= \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} = \frac{3}{9}$$

$$= 0.053$$

$$P(\text{no}|n) = P(\text{no}) \cdot P(\text{sunny}|\text{no}) \cdot P(\text{cool}|\text{no}) \cdot P(\text{high}|\text{no})$$

$$P(\text{strong}|\text{no})$$

$$= 0.206$$

here

$$P(\text{yes}|n) < P(\text{no}|n)$$

So No Doesnt Play.

| <u>Age</u> | <u>income</u> | <u>credit</u> | <u>Class</u> |
|------------|---------------|---------------|--------------|
| H | J. | L | Y |
| L | J. | H | N |
| L | E. | H | N |
| L | E. | H | N |
| H | J. | H | Y |

Prudit for $A = \text{Low}$, $\text{Income} = \text{fair}$, $Cn = \text{Low}$ wstrg Native

Bayesin may use Laplacian correction to avoid computing probability value of zero.

$$P(\text{Yes}) = \frac{2}{5} \quad P(\text{No}) = \frac{3}{5}$$

| <u>Age</u> | | <u>income</u> | | <u>credit</u> | |
|------------|---------------|---------------|--|---------------|---------------|
| | | | | | |
| H | $\frac{2}{2}$ | $\frac{0}{3}$ | | $\frac{Y}{2}$ | $\frac{N}{3}$ |
| L | $\frac{0}{2}$ | $\frac{3}{3}$ | | $\frac{J}{2}$ | $\frac{1}{3}$ |
| | | | | $\frac{E}{2}$ | $\frac{2}{3}$ |

Before correction.

yes

$$P(\text{Age} = L | \text{yes}) = \frac{0}{2}$$

$$P(\text{Cm} = \text{fa} | \text{yes}) = \frac{2}{2}$$

$$P(Cn = \text{low} | \text{yes}) = \frac{1}{2}$$

No

$$P(\text{Age} = L | \text{No}) = \frac{3}{3}$$

$$P(\text{Cm} = \text{fa} | \text{No}) = \frac{1}{3}$$

$$P(Cn = \text{low} | \text{No}) = \frac{2}{3}$$

Avoid Computing O used.

Laplacian Correction

- Given 5 to 2² 4% decision So k=2

and add ① Count (Geo⁰)

After

for yes

$$P(\text{Age} = \text{Low} | \text{geo}) = \frac{0+1}{2+2} = \frac{1}{4}$$

$$P(\text{line} = \text{fa} | \text{yes}) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$P(\text{cn} = \text{Low} | \text{yes}) = \frac{1+1}{2+2} = \frac{2}{4}$$

for no

$$P(\text{Age} = \text{Low} | \text{No}) = \frac{3+1}{3+2} = \frac{4}{5}$$

$$P(\text{line} = \text{fa} | \text{No}) = \frac{1+1}{3+2} = \frac{2}{5}$$

$$P(\text{cn} = \text{Low} | \text{No}) = \frac{0+1}{3+2} = \frac{1}{5}$$

So Avoid O successful

Now Posterior probabilities

yes

$$\begin{aligned}
 P(\text{yes}|x) &= P(\text{yes}) \cdot P(A=\text{low}) \cdot P(\text{in}=\text{f}) \cdot P(C=\text{low}) \\
 &= 0.2 \times \frac{1}{4} \times \frac{3}{4} \times \frac{2}{4} \\
 &= 0.0375
 \end{aligned}$$

for No

$$\begin{aligned}
 P(\text{No}|x) &= P(\text{No}) \cdot P(A=\text{low}) \cdot P(\text{in}=\text{f}) \cdot P(C=\text{low}) \\
 &= 0.6 \times 0.8 \times 0.2 \times 0.6 \\
 &= 0.384
 \end{aligned}$$

hence

$$P(\text{No}|x) > P(\text{yes}|x)$$

Ano 

Gaussian Naive Bayes

| Person | Height | weight | foot size |
|--------|--------|--------|-----------|
| M | 6 | 186 | 12 |
| M | 5.92 | 190 | 11 |
| M | 5.58 | 170 | 12 |
| M | 5.92 | 168 | 10 |
| f | 5.00 | 106 | 6 |
| f | 5.50 | 150 | 8 |
| f | 5.42 | 130 | 7 |
| f | 5.75 | 150 | 9. |

New instance, H = 130, f = 8

Avg

$$\text{here } (M) = \frac{4}{8} = 0.5 \quad | \quad \text{female} = \frac{4}{8} = 0.5$$

Male

$$\text{mean (height)}_M = \frac{6.00 + 5.92 + 5.58 + 5.92}{4} = 5.855$$

$$\text{variance } \sigma^2(M) = \frac{(6-5.85)^2 + (5.92-5.85)^2 + (5.58-5.85)^2 + (5.92-5.85)^2}{4-1}$$

(6)

is man

$$\left[\frac{\sum n_i - \bar{n}}{n-2} \right]$$

$$= 0.0350$$

Same way (to find out mean, variance find out)

~~female weight~~

$$\text{mean} = \frac{180 + 190 + 170 + 165}{4} = 176.25$$

$$\text{variance} = \frac{(180 - 176.25)^2 + (190 - 176.25)^2 + (170 - 176.25)^2 + (165 - 176.25)^2}{4-1}$$
$$= 122.92$$

male foot size

$$\text{mean (FS)} = \frac{12 + 11 + 12 + 10}{4} = 11.25$$

$$\text{variance} = \frac{(12 - 11.25)^2 + (11 - 11.25)^2 + (12 - 11.25)^2 + (10 - 11.25)^2}{4-1}$$
$$= 0.91667$$

Female

$$\text{Mean (H)} = 5.417$$

$$\text{variance (H)} = 0.097$$

$$\text{mean (W)} = 132.5$$

$$\text{variance (W)} = 558.333$$

$$\text{mean (FS)} = 7.5$$

$$\text{variance (FS)} = 1.6667$$

Gaussian distribution

$$g(n, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(n-\mu)^2}{2\sigma^2}}$$

$$P(\text{Male} | H_{1g}) = \frac{1}{\sqrt{2 \cdot 3.14 \times 0.035}} \times e^{-\frac{(6 - 5.85)^2}{2 \times 0.035}}$$
$$= 1.578$$

Same way

$$P(\text{Male} | w) = 5.98 \times e^{-6}$$

$$P(\text{Female} | fs) = 1.311 \times 10^{-3}$$

$$\underline{\text{Female}}$$

$$P(F|H) = 2.23e^{-1} \quad | \quad P(F|w) = 1.67e^{-2} \quad | \quad P(F|fs) = 2.86e^{-1}$$

$$\text{Posterior male} = \underbrace{P(M) * P(H|M) * P(w|M) * P(fs|M)}_{\text{evidence}} = 6.1984 e^{-9}$$

$$= 6.1984 e^{-9}$$

$$\text{value} = 6.1984 e^{-6} \cdot e^{-3} = e^{-9}$$

$$\text{Posteriori female} = \frac{P(S) \cdot P(f|S) P(w|f) P(fs|s)}{\text{Evidence}} \quad \text{w} \} \text{ unknown} \\ \approx 5.37 e^{-4}$$

here $P(f) > P(M)$

So test sample is classified as female.

Naive Bayesian classification known as class condition.

Independence

Became

Independence feature: easy decision \approx 3 or 2 vs independent \approx 20

Simpler math: assumption and simpler math calculation

Fast Prediction: allows quick and easy probability.

Efficiency: one feature doesn't change another feature make it more efficient.

* easy implementation

[Adv within equal Gains (from 20)]

Model Evaluation & Selection

Evaluation → How can we measure accuracy.

Methods for estimating accuracy:

— holdout method.

— cross validation

— Bootstrap.

Confusion matrix \rightarrow represents accuracy of model

positive tuples (P)

main category yes \rightarrow (buy pe = yes)

buys pe = yes (does buy)

Negative tuples (N)

These are all other categories

buys pe = no (doesn't buy pe)

Y-harm

True positive (TP): positive tuples \rightarrow Identified as positive

Buy pe = yes

\downarrow P
He does buy pe

True Negative (TN): Negative tuples \rightarrow correctly negative

\oplus
buy pe = no

\circlearrowleft
he dont buy pe

True Positive (TP)
True Negative (TN)
False Positive (FP)
False Negative (FN)

False positive: Neg tuples get positive result \Rightarrow not buy PC
 buy PC = yes, actually don't buy.

False negative: Positive tuples get negative result \Rightarrow buy PC = NO → actual buy PC

| | | Example | | <u>Predicted</u> | | |
|---------------------|----|---------|----|------------------|----|----|
| <u>Actual class</u> | | +C | -C | +C | -C | +C |
| +C | TP | FN | AC | TN | FP | |
| -C | FP | TN | +C | FN | TP | |

Formal form

$$\text{accuracy} \rightarrow \frac{TP + TN}{P + N}$$

$$\text{error rate} \rightarrow \frac{FP + FN}{P + N} \quad (\text{1 - accuracy})$$

$$\text{Sensitivity, TP, recall} \rightarrow \frac{TP}{P}$$

$$\text{Specificity, TN, note} \rightarrow \frac{TN}{N}$$

$$\text{Precision} \rightarrow \frac{TP}{TP + FP}$$

$$\text{F1 score} \rightarrow \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

* Math

Predicted 0⁺

Predicted 1⁻

Actual 0⁺ 90 10 = 100

Actual 1⁻ 10 0 = 10

10 = 110

Find accuracy, error rate, precision, recall

Ans

hence

~~TP = 90, TN = 0, FP = 10, FN = 10~~

~~TP = 90, TN = 0, FP = 10, FN = 10~~

TN = 0, FP = 10, FN = 10, TP = 0

math solved by this.

$$\text{Accuracy} = \frac{TN + TP}{P+N} = \frac{0+0}{110} = 0.81 = 81\%$$

$$\text{error rate} = \frac{FP + FN}{P+N} = \frac{10+10}{110} = 0.18 = 18\%$$

(1 - accuracy)

$$\text{precision} = \frac{TP}{TP + FP} = \frac{0}{0+10} = 0$$

$$\text{recall} = \frac{TP}{TP + FN} = \frac{0}{0+10} = 0$$

Observations

④ Min precision, recall & f1m → model fails to correctly identify any positive instance.

⑤ accuracy high → model correctly classifies a good proportion of the overall data.

⑥ min precision high → model's strong ability.
to minimize false positive

⑦ Recall high → models capture all actual buyers
low → model missed many true positive cases.

⑧ F1 score high → Indicates Balanced performance.
low → poor balance precision & recall needs improvement

⑨ Accuracy high but recall low → miss leading → Medical Sector first time

Class Imbalance problem

when number of events in each class of dataset is not equal. Ex - fraud, spam detection

Dataset Imbalance:

when,

- unequal cls. distribution.

- if minority cls less than (10-20)% of the total data then dataset is imbalanced.

Measure the performance for ~~Imbalanced Algorithm~~

① confusion matrix: $TN, TP, FN, FP \rightarrow$ measured

Correct Neg
predict
 $N-N$

Correct Pos
predict
 $P-P$

Measured Neg
predict
 $P-N$

Pos
predicted
 $N-P$

② precision: Measure accuracy of positive predictions.

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

③ Recall: Measure the ability to identify actual positives

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

$$\boxed{\text{F1 Score}}, \boxed{\text{Specificity}} \rightarrow \frac{\text{TN}}{\text{TN} + \text{FP}}$$

ability of predict true negative.

ROC AUC: Plot the TP rate against the FP rate at the various threshold settings.

AUC : Measure overall performance higher is better

Evaluating Classifier Accuracy.

— Holdout vs Cross Validation Methods

Holdout 2 steps

Training sets → model construction

Test set → accuracy estimation.

Random Sampling

→ either the data split multiple times
→ final accuracy is the average all accuracies from this splits.

⑩ Cross validation: If $n=10$ then split into 10, 9 parts are training 1 are testing.

* Leave one out cross validation: each point test often is training

* Stratified cross → distribution similar to the original set. (balanced evaluation)

(*)

Roc Curve

It is a graph showing the performance of a classification model at all classification thresholds.

True Positive Rate False positive Rate
TPR \ominus FPR

AUC (Area under Curve)

measure accuracy

Closer to 1 \rightarrow Highly accurate model.

" " 0.5 \rightarrow poor model almost random prediction.