



Suppose that the data for analysis includes the attribute age. The age values for data tuples are (in increasing order)

$$\{13, 15, 16, 16, 19, 20, 20, 21, 22, 22, \underline{25}, \underline{25}, \underline{25}, \\ 30, 33, 33, \underline{35}, \underline{35}, \underline{35}, 36, 40, 46, 52, 70\}$$

find Mean, Median, Mode, Range, Quartiles (Q_1, Q_2, Q_3)

give five number Summary of the data
Show a box plot of the data

$$\sum_{i=1}^n \frac{x_i}{n} = \frac{809}{27} = 30$$

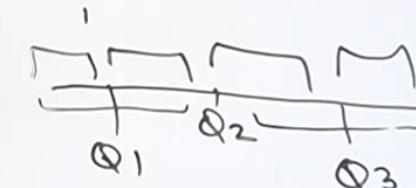
$$\begin{aligned}\text{Median} &= \left(\frac{N+1}{2}\right)^{\text{th}} \text{ term} \\ &= \left(\frac{27+1}{2}\right)^{\text{th}} \text{ term} \\ &= \left(\frac{28}{2}\right)^{\text{th}} \text{ term} \\ &= 14^{\text{th}} \text{ term} \\ &= \underline{\underline{25}}\end{aligned}$$

3. Mode: Since two numbers are repeating highest no. of times So the dataset is a bimodal dataset. Modes of the dataset are 25 & 35.

4. Midrange:

$$\frac{13 + 70}{2} = 41.5$$

5. Range: $70 - 13 = 57$



6. Quartile 2 (Q_2) = Median
 $= 25$

Quartile 1 (Q_1) = $\left(\frac{N+1}{4}\right)^{\text{th}}$ term
 $= 20$

Quartile 3 (Q_3) = $\left(\frac{(N+1)*3}{4}\right)^{\text{th}}$ term
 $=$



)

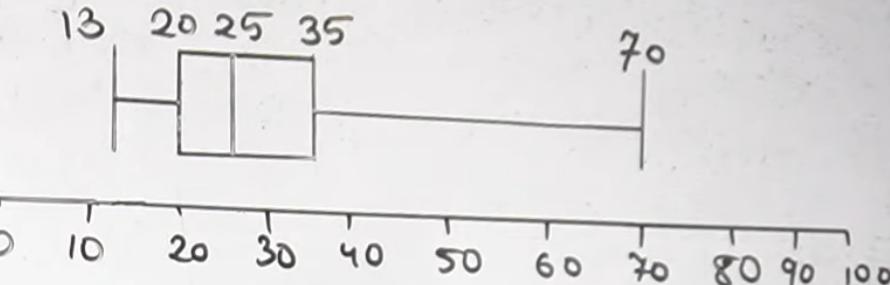
25, 25, 25, 25,
45, 46, 52, 70

range, Range,
five number
box Plot of

$$= \frac{809}{27} = 30$$

n

rm



4. Midrange: $\frac{13 + 70}{2}$
 $= 41.5$

$$\begin{aligned}Q_1 &= 20 \\Q_2 &= 25 \\Q_3 &= 35 \\Min &= 13 \\Max &= 70\end{aligned}$$

5. Range: $70 - 13 = 57$

6. Quartile 2 (Q2) = Median
 $= 25$

Quartile 1 (Q1) = $\left(\frac{N+1}{4}\right)^{\text{th}} \text{ term}$
 $= 20$



Confusion Matrix

1% logo KO
TB hai

99% logo KO
TB nahi hai

Actual

True	False
True +ve	False +ve
-ve	True -ve
TP	FP
N	TN

Actual

Have TB	Don't have TB
True +ve	False +ve
-ve	True -ve
TP	FP
N	TN

TP = True +ve = correctly identified +ve

TN = True -ve = correctly identified -ve

FN = False -ve = Incorrectly identified +ve

FP = False +ve = Incorrectly identified -ve

1. Accuracy = $\frac{\text{Total correct Predictions}}{\text{Total Predictions}} = \frac{TP+TN}{TP+FP+FN+TN}$

2. Recall = $\frac{\text{Correctly identified +ve}}{\text{Total actual +ve}} = \frac{TP}{TP+FN}$

3. Specificity = $\frac{\text{Correctly identified Negative}}{\text{Total actual -ve}} = \frac{TN}{TN+FP}$

Details of 100 Patients → 99% logo K0 TB nahi

Predicted

	True	False
True	$\frac{TP}{TP+FP}$	$\frac{FN}{TP+FP}$
False	$\frac{FP}{TP+FP}$	$\frac{TN}{TP+FP}$

$$TP = T$$

$$TN$$

$$F$$

Score
Recall

$$\text{Precision} = \frac{\text{Correctly Identified tues}}{\text{Total tues Predicted}} = \frac{TP}{TP+FP} = 1$$

$$\text{Recall} = \frac{\text{Correctly Identified tues}}{\text{Total actual tues}} = \frac{TP}{TP+FN} = 1$$

F1 Score

$$F_1 \text{ Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

When Should we use F1 Score?

Apriori Algorithm

Given the following data, apply the Apriori algorithm. Given Support threshold = 50%, Confidence = 60%

Transaction	List of Items	Count
T ₁	I ₁ , I ₂ , I ₃	3
T ₂	I ₁ , I ₃ , I ₄	3
T ₃	I ₄ , I ₅	2
T ₄	I ₁ , I ₂ , I ₄	3
T ₅	I ₁ , I ₂ , I ₃ , I ₅	4
T ₆	I ₁ , I ₂ , I ₄	3

Support count

Association rules, minimum confidence

$$\{I_1\} \rightarrow \{I_2, I_3\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1\}}$$

$$= \frac{3}{4} \times 100 = 75\%$$

$$\{I_2\} \rightarrow \{I_1, I_3\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_2\}}$$

$$= \frac{3}{4} \times 100 = 75\%$$

$$\{I_3\} \rightarrow \{I_1, I_2\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_3\}}$$

$$= \frac{3}{4} \times 100 = 75\%$$

does not meet min-sup after discarding

Step 5: Join Step

Itemset	Count
I ₁ , I ₂ , I ₃	3
I ₁ , I ₂ , I ₄	2
I ₁ , I ₃ , I ₄	1
I ₂ , I ₃ , I ₄	2

Step 6: {I₁, I₂, I₄}, {I₁, I₃, I₄} & {I₂, I₃, I₄} does not meet min-sup = 3 So after discarding them we get :

Itemset	Count
I ₁ , I ₂ , I ₃	3

Thus {I₁, I₂, I₃} is frequent.

Step 7: Generate Association Rules

$$\{I_1, I_2\} \rightarrow \{I_3\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1, I_2\}} \\ = \frac{3}{4} \times 100 = 75\%$$

$$\{I_1, I_3\} \rightarrow \{I_2\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1, I_3\}} \\ = \frac{3}{3} \times 100 = 100\%$$

$$\{I_2, I_3\} \rightarrow \{I_1\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_2, I_3\}} \\ = \frac{3}{4} \times 100 = 75\%$$

Apriori Algorithm

Given the following data, apply the Apriori algorithm. Given Support threshold = 50%, Confidence = 60%

Transaction	List of items
T ₁	I ₁ , I ₂ , I ₃
T ₂	I ₂ , I ₃ , I ₄
T ₃	I ₄ , I ₅
T ₄	I ₁ , I ₂ , I ₄
T ₅	I ₁ , I ₂ , I ₃ , I ₅
T ₆	I ₁ , I ₂ , I ₃ , I ₄

$$\circ \{I_1\} \rightarrow \{I_2, I_3\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_1\}} \\ = \frac{3}{4} \times 100 = 75\%$$

$$\circ \{I_2\} \rightarrow \{I_1, I_3\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_2\}} \\ = \frac{3}{5} \times 100 = 60\%$$

$$\circ \{I_3\} \rightarrow \{I_1, I_2\}$$

$$\text{Confidence} = \frac{\text{Support}\{I_1, I_2, I_3\}}{\text{Support}\{I_3\}} \\ = \frac{3}{4} \times 100 = 75\%$$

Step 2: I₅ itemset does not meet min-sup = 3 so after discarding it we get

Itemset	Count
I ₁	4
I ₂	5
I ₃	4
I ₄	4

Step 3: Join Step

Itemset	Count
I ₁ , I ₂	9

Ans: All the above association rules are strong. Minimum confidence is 60%.

Step 5: Join Step

Itemset	Count
I ₁ , I ₂ , I ₃	5

Itemset	Count
I ₁ , I ₂ , I ₄	4

Itemset	Count
I ₁ , I ₃ , I ₄	3

Itemset	Count
I ₂ , I ₃ , I ₄	3

Step 6: {I₁, I₂, I₄}

does not meet min-sup = 3 so after discarding them we get :

Itemset	Count
I ₁ , I ₂ , I ₃	5

Itemset	Count
I ₁ , I ₃ , I ₄	3

Itemset	Count
I ₂ , I ₃ , I ₄	3

Thus {I₁, I₂, I₃} is frequent itemset.

Step 7: Generate Association Rules

o {I₁, I₂, I₃} → {I₁}

Confidence = Support{I₁, I₂, I₃} / Support{I₁}

= (3/5) × 100 = 60%

o {I₁, I₂, I₃} → {I₂}

Confidence = Support{I₁, I₂, I₃} / Support{I₂}

= (3/4) × 100 = 75%

o {I₁, I₂, I₃} → {I₃}

Confidence = Support{I₁, I₂, I₃} / Support{I₃}

= (3/4) × 100 = 75%



FP-Growth Algorithm

⇒ A database has five transactions.
 Let min-sup = 60% of min support.
 Find all the frequent items.

FP-Growth

TID	Items
T ₁₀₀	{M, O, K}
T ₂₀₀	{D, O, N, K}
T ₃₀₀	{M, A, K}
T ₄₀₀	{M, U, C, K}
T ₅₀₀	{L, D, Y}

⇒ min-sup = 60%.
 ∴ Sup-count to be calculated.

Step 1: Count

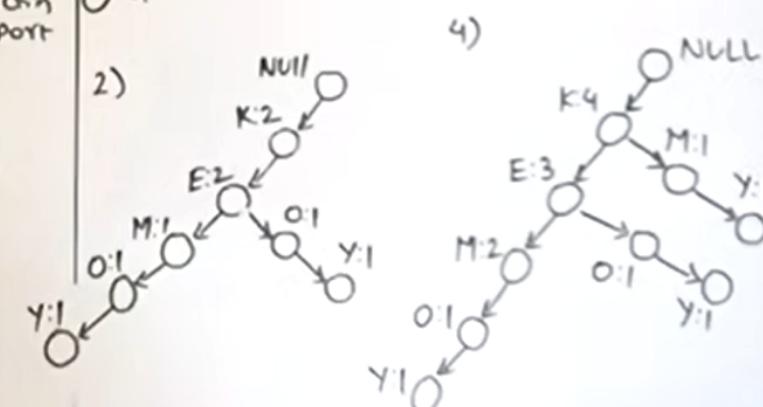
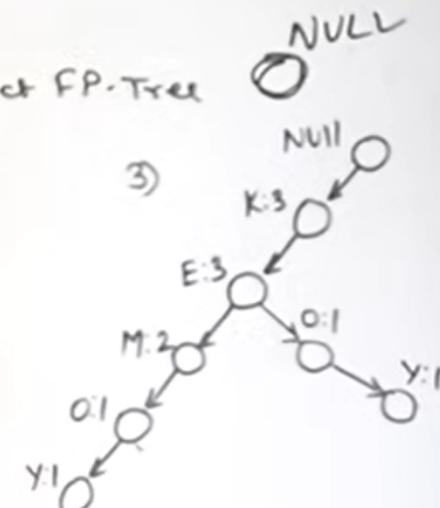
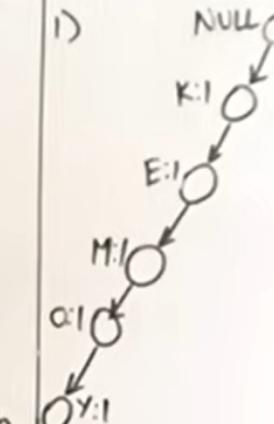
FP-Growth
 Step 1: Record A, C, D, I, N, U as frequent items & calculate sup-count & sort remaining items in descending order of sup-count.

Set	Sup-count
A	5
C	4
D	3
I	3
N	3
U	3

Step 2: Sort items in each transaction according to descending support count.

TID	List of items
T ₁₀₀	{K, E, M, O, Y}
T ₂₀₀	{K, E, O, Y}
T ₃₀₀	{K, E, M}
T ₄₀₀	{K, M, Y}
T ₅₀₀	{K, E, O}

Step 4: Construct FP-Tree





FP-Growth Algorithm

=> A database has

Let min_sup =

Find all the frequent itemsets

FP-Growth

TID

T₁₀₀

T₂₀₀

T₃₀₀

T₄₀₀

T₅₀₀

=> min-sup = 60%
 i.e. sup-count to be

Step 1: Count

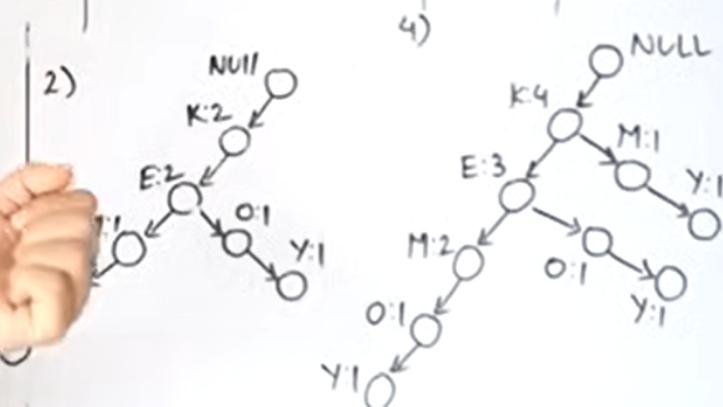
Step 2: Discard A, C, D, I, N, U
 Itemset does not satisfy sup-count &
 arranging remaining items in
 descending order of sup-count

Itemset	Sup-count
{K}	5
{E}	4
{M}	3
{O}	3
{Y}	3

Conditional pattern base
 Conditional FPTree

Frequent Patterns Generated

{Y}	{K, E, M, O: 1}, {K, E, O: 1}	(K: 3)	{K, Y: 3}
	{K, M: 1}		
{O}	{K, E, M: 1}, {K, E: 2}	(K: 3, E: 3)	{K, O: 3}, {E, O: 3}, {K, E, O: 3}
{M}	{K, E: 2}, {K: 1}	(K: 3)	{K, M: 3}
{E}	{K: 4}	(K: 4)	{K, E: 4}
{K}	-	-	-



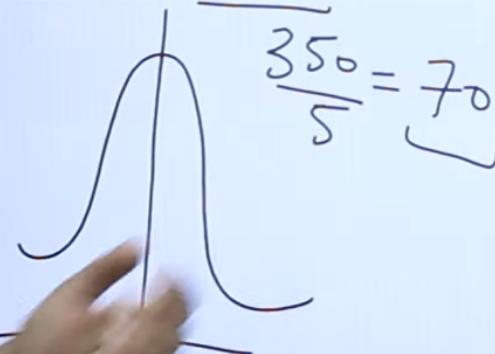
Min Max Technique [0-1]

House	Square Foot	Bedrooms	Price(in Lakh)
1	1200	3	50
2	1500	4	60
3	1000	2	40
4	1800	5	80

$$X' = \frac{X - \min}{\max - \min} = \frac{1200 - 1000}{1800 - 1000} = \frac{200}{800} = 0.25$$

Z-score Technique

Student	Height(in inches)
1	64 = -1.5
2	70
3	72
4	68
5	76



$$X' = \frac{X - \text{mean } \mu}{\text{Std. dev } \sigma}$$
$$\frac{(64-70)^2 + (70-70)^2 + (72-70)^2 + (68-70)^2 + (76-70)^2}{5}$$

$$36 + 0 + 4 + 4 + 36 = \frac{80}{5} = \sqrt{16} = 4$$

$$\frac{64-70}{4} = -\frac{6}{4} = -1.5$$

Data Smoothing A Data Reduction

Suppose a group of age records

has been sorted as follows:

$\{5, 12, 13, 14, 15, 18, 26, 28, 42\}$

Perform data smoothing &
data reduction.

=> Equal freq Bin:

Bin1: $\{5, 12, 13\}$

Bin2: $\{14, 15, 18\}$

Bin3: $\{26, 28, 42\}$

Equal width Bin:

Bin1: $\{5, 12, 13, 14, 15\}$

Bin2: $\{18, 26, 28\}$

Bin3: $\{42\}$

* Smoothing by Mean(1):

Bin1: $\{10, 10, 10\}$

Bin2: $\{16, 16, 16\}$

Bin3: $\{32, 32, 32\}$

* Smoothing by Boundary(1):

Bin1: $\{5, 13, 13\}$

Bin2: $\{14, 14, 18\}$

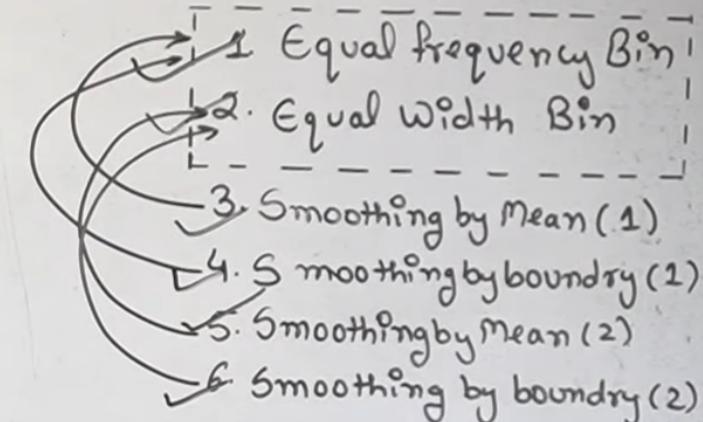
Bin3: $\{26, 26, 42\}$

* Smoothing by Mean(2):

Bin1: $\{12, 12, 12, 12, 12\}$

Bin2: $\{24, 24, 24\}$

Bin3: $\{42\}$



Smoothing by boundary(2)

Bin1: $\{5, 15, 15, 15, 15\}$

Bin2: $\{18, 28, 28\}$

Bin3: $\{42\}$

Chi-Square Test

Watch Actor ↓	Y	N	(R) Total ↓
	M	140	44
F	178	38	216
(C) Total →	318	82	400

✓ $\alpha = 0.05$ -

If $= (Rows-1)*(cols-1) = 1 \quad 3.841 <$

$$\chi_c^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

$$E = \frac{RT \times CT}{T} / = \frac{184 \times 318}{400} = 146$$

H₀: No connection

H₁: full connection

$$\checkmark = (184 \times 82) / 400 = 38$$

$$\checkmark = (216 \times 318) / 400 = 172$$

$$\checkmark = (216 \times 82) / 400 = 44$$

Actor Watch	O	E	$(O-E)^2$	$(O-E)^2/E$
M Y	140	146	36	0.246
M N	44	38	36	0.947
F Y	178	172	36	0.209
F N	38	44	36	0.818

Cal $\rightarrow \frac{2.22}{2.220}$

