

Popular Attribute Selection Measures are:-

1. Information Gain
2. Gini Index.
3. Gain Ratio.

Information Gain :-

$$P_i = |C_{i,D}| / |D|$$

Expected Information :-

$$\text{Info}(D) = - \sum_{i=1}^m P_i \log_2(P_i)$$

Information Needed :-

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times I(D_j)$$

Info. Gain :-

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D).$$

Step 1:- Entropy of Entire Dataset ($4+9 \rightarrow 9$)
 $(10+5 \rightarrow 15)$

$$S_{4+9} = -\frac{4}{14} \log_2 \frac{4}{14} - \frac{9}{14} \log_2 \frac{9}{14} = 0.94$$

Step 2:- Entropy of All Attributes (Weather)

$$\text{Entropy of sunny } S_{2,1,-3} = -\frac{2}{5} \log_2 \left[\frac{2}{5} \right] - \frac{1}{5} \log_2 \left[\frac{1}{5} \right]$$

$$= 0.97$$

$$\text{Entropy of Rainy Cloudy } S_{1,4,-0} = -\frac{1}{5} \log_2 \left[\frac{1}{5} \right] - \frac{4}{5} \log_2 \left[\frac{4}{5} \right]$$

$$\text{Entropy of Rainy } S_{3,1,-2} = -\frac{3}{5} \log_2 \left[\frac{3}{5} \right] - \frac{2}{5} \log_2 \left[\frac{2}{5} \right]$$

$$\therefore \text{Info Gain} = \text{Entropy}(\text{whole data}) - \frac{5}{14} \text{Entropy}(S) - \frac{1}{14} E(C) - \frac{3}{14} E(R)$$
$$= 0.246.$$

Calculate IG of Temp.

\rightarrow Entropy at entire Dataset

$$\Sigma \xi + q, -S_g = 0.94$$

$$S_{21} = \text{Entropy of Hot } \varphi + 2, -S_g = -2/2 \log_2 2 + -2/4 \log_2 4 = 0.1$$

$$\text{Entropy of Mild } \varphi + 1, -S_g = -4/6 \log_2 1/6 = -2/3 \log_2 1/3 = 0.91$$

$$\text{Entropy of Cold } \varphi + 3, -S_g = -3/4 \log_2 3/4 = -1/3 \log_2 1/3 = 0.81$$

$$\therefore \text{Info Gain} = E(\text{whole}) - 4/14 E(H) - 4/14 E(M) - 4/14 E(C)$$

$$= 0.029$$

calculate IG of humidity

$$S_{21} = E(H) \xi + 3, -S_g = 0.98$$

$$E(N) \varphi + 6, -S_g = 0.59$$

$$IG = 0.94 - 7/14 E(H) - 7/14 E(N)$$

$$= 0.15$$

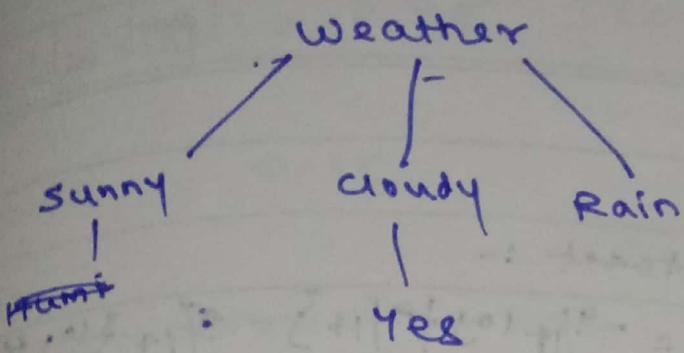
IG of wind

$$S_{21} = E(S) = 1.0$$

$$E(N) = 0.81$$

$$\therefore 0.0478$$

\therefore Greater Entropy is Weather



Now, Again for sunny for entire table.

$$S1 := \text{Entropy of sunny} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.97$$

IG of Temp :-

$$E(Hot) = 0$$

$$E(Mild) = 1.0$$

$$E(Cold) = 0$$

$$\therefore \text{Info Gain} = \text{Entropy(sunny)} - \frac{2}{5} Ent(H) - \frac{2}{5} Ent(M) - \frac{1}{5} Ent(C)$$

$$= 0.57$$

IG of Humidity

$$E(H) = 0$$

$$E(N) = 0$$

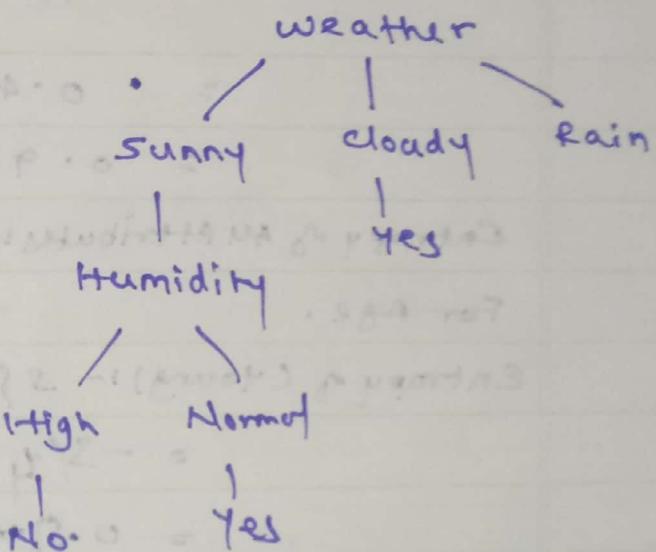
$$\therefore \text{Info Gain} = 0.97$$

IG of Wind :-

$$\text{Entropy}(S) = 1$$

$$E(W) = 0.918$$

$$IG = 0.019$$

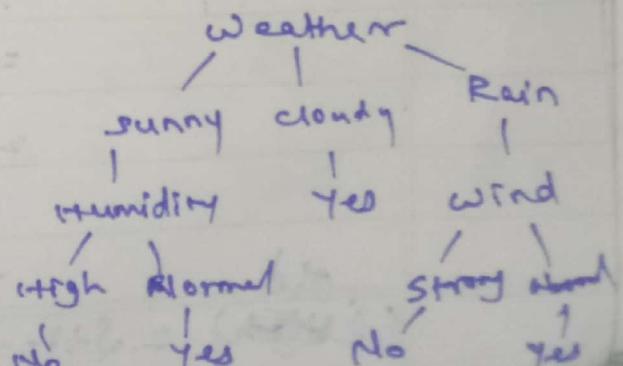


Now calculate Again for Rain :-

$$IG \text{ of temp} = 0.019$$

$$IG \text{ of Humidity} = 0.019$$

$$IG \text{ of Wind} = 0.97$$



Gain Ratio.

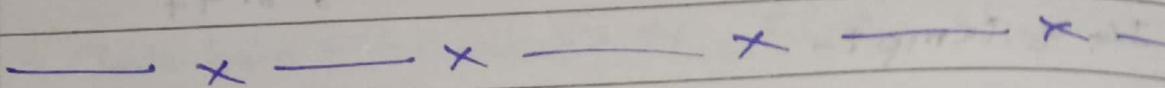
Gain for Entire Dataset :-

$$S_{\{1, -5\}} = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$= 0.94$$

Info Gain for Income :-

$$0.94 - 1.557$$



Class → Owns a House \rightarrow Yes
 \rightarrow Rented

Gain for Entire Dataset

$$S_{\{1, -5\}} = -\frac{7}{12} \log_2 \left(\frac{7}{12} \right) - \frac{5}{12} \log_2 \left(\frac{5}{12} \right)$$

$$= 0.45 + 0.52$$

$$= 0.9798$$

Entropy of All Attributes:-

For Age.

Entropy of (Young) :- $S_{\{1, 0\}}$

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

$$= 0.81 + 0.25 = 0.25$$

Entropy of (Medium) = $S_{\{1, -2\}}$

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

$$= 0.25 + 0.302 = 0.552$$

Old ≈ 1.115

~~$E(\text{Old}) = S_{\{1, -3\}} = -\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right)$~~

$$\text{Info Gain (Age)} = 0.97 - \frac{4}{12}(0.52) - \frac{5}{12}(1.115) - \frac{3}{12}(0.600)$$

$$= 0.15$$

Info gain for Income

very high :- $\$5(+210)$

$$= -\frac{2}{3} \log_2(2/3) - \frac{1}{3} \log_2(1/3)$$

$$= 0$$

High :- $\$54103$

$$= -\frac{1}{3} \log_2(4/3) - \frac{2}{3} \log_2(2/3)$$

$$= 0$$

medium :- $\$1123$

$$= -\frac{1}{3} \log_2(1/3) - \frac{2}{3} \log_2(2/3)$$

$$= 0.52 + 0.389$$

$$= 0.909$$

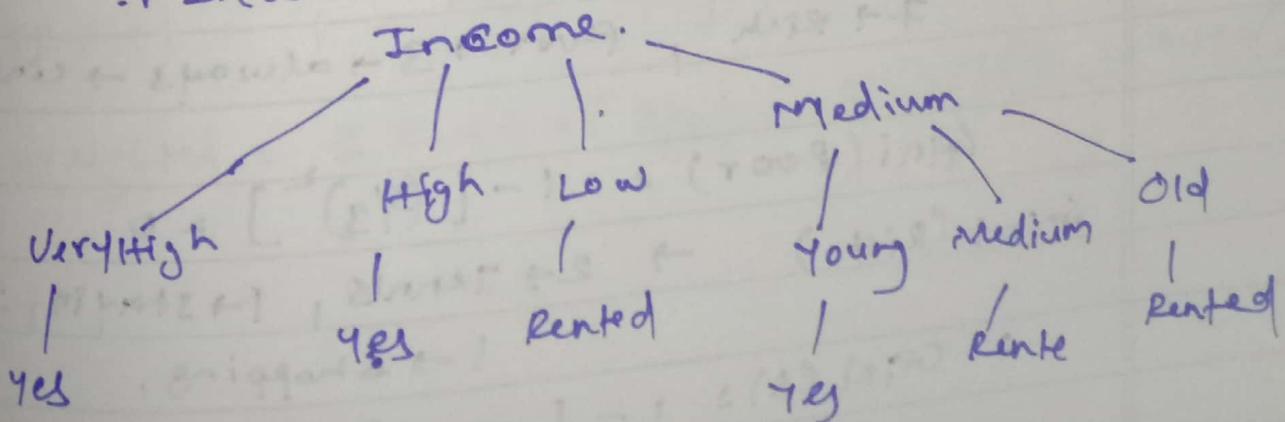
Low :- $\$24913$

$$= 0.$$

$$\text{Info (Income)} = 0.97 - \frac{3}{12}(0.909)$$

$$= 0.78$$

\therefore income



Gini Index

weekend	weather	parents	Money	Decision
w ₁	sunny	yes	Rich	Cinema
w ₂	sunny	No	Rich	Tennis
w ₃	windy	yes	Rich	Cinema
w ₄	Rainy	Yes	Poor	Cinema
w ₅	Rainy	No	Rich	Stay in
w ₆	Rainy	yes	Poor	Cinema
w ₇	windy	No	Poor	Cinema
w ₈	windy	No	Rich	Shopping
w ₉	windy	Yes	Rich	Cinema
w ₁₀	sunny	No	Rich	Tennis

Compute Gini Index for overall collection.

Cinema → 6, Tennis → 2, Stay in → 1, Shopping → 1

$$Gini(S) = 1 - \left[\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right]$$

$$\rightarrow = 0.58$$

Gini for Money Attribute.

7 → Rich, Poor → 3 → always → Cinema

$$Gini(Poor) = 1 - \left[\left(\frac{3}{3}\right)^2 \right] = 0$$

Gini(Rich) → 2 → Tennis, 1 → stay in, 3 → Cinema
1 → Shopping,

$$\therefore Gini(S) = 1 - \left[\left(\frac{2}{7}\right)^2 + \left(\frac{3}{7}\right)^2 + \left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 \right] \\ = 0.594$$

$$\text{Weighted Avg} = 0 * \left(\frac{3}{10}\right) + 0.594 * \left(\frac{7}{10}\right) = 0.486$$

Gini for parents

$$Yes \rightarrow 5, No \rightarrow 5$$

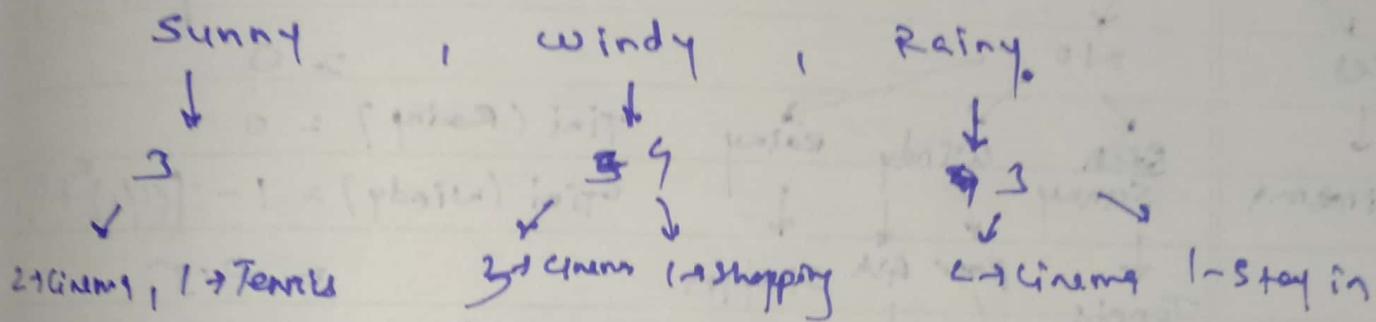
for Yes $5 \rightarrow \text{Cinema}$ for No, $\rightarrow \text{Tennis}$, $1 \rightarrow \text{Stayin}$
Gini (Parents = Yes) $= 1 - \left[\left(\frac{5}{5}\right)^2 + \left(\frac{0}{5}\right)^2 \right]$ $1 \rightarrow \text{Shopping}$
 $= 0$ $1 \rightarrow \text{Unrest}$

Gini (Parents = No) $= 1 - \left[\left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right]$
 $= 0.72$

$$\text{Weighted} = 0.2 \times \left(\frac{5}{10}\right) + 0.72 \times \left(\frac{5}{10}\right)$$

0.36 ← since, minimum max info.

Gini for weather:-



$$\text{Gini}(\text{Sunny}) = 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 0.44$$

$$\text{Gini}(\text{Cinema}) = 1 - \left[\left(\frac{3}{10}\right)^2 + \left(\frac{4}{10}\right)^2 \right] = 0.375$$

$$\text{Gini}(\text{Rainy}) = 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 0.44$$

$$\therefore \text{Weighted} = 0.44 \times \left(\frac{3}{10}\right) + 0.375 \times \left(\frac{4}{10}\right) + 0.44 \times \left(\frac{3}{10}\right)$$
$$= 0.416$$

Parents
Yes No

↙
Cinema

w ₂	Sunny	No	Rich	Tennis
w ₅	Rainy	No	Rich	Stayin
w ₇	Windy	No	Poor	Grocery
w ₈	Windy	No	Rich	Shopping
w ₁₀	Sunny	No	Rich	Tennis

Now for weather

$$\text{sunny} \rightarrow 2, \text{windy} \rightarrow 2, \\ \text{rainy} \rightarrow 3$$

$$\text{Gini(Sunny)} = 1 - \left[\left(\frac{2}{2}\right)^2 \right] = 0$$

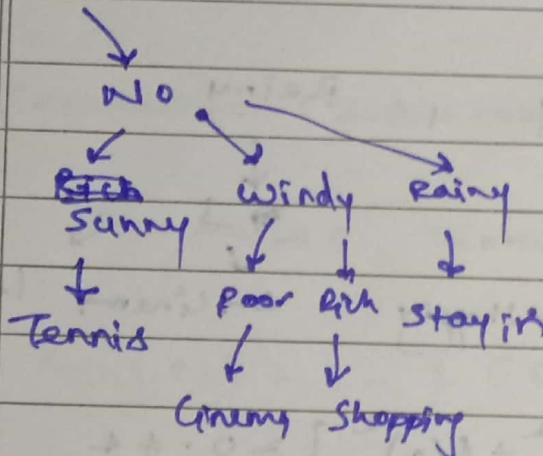
$$\text{Gini(Rainy)} = 0$$

$$\text{Gini(Windy)} = 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = 0.5$$

$$\text{Weighted} = 0 + 0 + 0.5 * \frac{2}{5} = 0.2$$

Parents

↙
Yes
↓
Cinema



For Money

$$\text{Rich} \rightarrow 4, \text{Poor} \rightarrow 1$$

$$\text{Gini(Rich)} = 1 - \left[\left(\frac{2}{2}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] = 0.625$$

$$\text{Gini Poor} = 1 - \left(\frac{1}{1}\right) = 0$$

$$\text{Weighted} = 0 * \frac{1}{4} + 0.625 * \frac{4}{5} = 0.5$$

Naive-Bayes Classification:-

Bayes theorem

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

Find $P(Y)$ if probability of X is given.

$$P(Y | X_1, X_2, \dots, X_n) = \frac{P(X_1|Y) * P(X_2|Y) * P(X_3|Y) \dots P(X_n|Y) * P(Y)}{P(X_1) * P(X_2) * P(X_3) \dots P(X_n)}$$

person	covid (Yes/No)	flu (Yes/No)	fever (Yes/No)
1	Yes	No	Yes
2	No	Yes	Yes
3	Yes	Yes	Yes
4	No	No	No
5	Yes	No	Yes
6	No	No	Yes
7	Yes	No	Yes
8	Yes	No	No
9	No	Yes	Yes
10	No	Yes	No

Given Person (flu, covid)

$$P(\text{Yes} | (\text{flu, covid})) = P(\text{flu/Yes}) * P(\text{covid/Yes}) * P(\text{Yes})$$

$$P(\text{No} | (\text{flu, covid})) = P(\text{flu/No}) * P(\text{covid/No}) * P(\text{No})$$

S1:- $P(\text{fever} = \text{Yes}) = 7/10$
 $P(\text{fever} = \text{No}) = 3/10$

		conditional prob.	
		Yes	No
	Covid	4/7	2/3
	Fever	3/2	2/3

Play-tennis

Outlook	Temp	Humidity	Windy	Class
sunny	Hot	High	False	N.
sunny	Hot	High	True	N
overcast	Hot	High	False	P.
rain	mild	High	False	P.
rain	cool	High Normal	False	P.
rain	cool	Normal	True	N
overcast	cool	Normal	True	P.
sunny	Mild	Normal High	False	N
sunny	Cool	High Normal	False	P.
rain	Mild	Normal	False	P.
sunny	Mild	Normal	True	P
overcast	Mild	Normal High	True	P
overcast	Hot	Normal	False	P
rain	Mild	High	True	N

S:- $P(C=P) = 9/14$
 $P(C=N) = 5/14$

outlook

	Yes	No
Hot	2/9	2/5
Mild	4/9	4/5
Cool	3/9	1/5

	P	N
Sunny	2/9	3/5
Overcast	4/9	0
Rain	3/9	2/5

	P	N
High	3/9	4/5
Normal	8/9	4/5

	P	N
False	6/9	2/5
True	3/9	3/5

$$P(X|Y) = \frac{2}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{2}{9} \times \frac{6}{9} \times \frac{4}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{1}{9} \times \frac{7}{9} \times \frac{1}{9} \times \frac{9}{14}$$

$$\frac{2}{9} \times \frac{4}{9} \times \frac{1}{9} \times \frac{2}{9} \times \frac{2}{9} \times \frac{4}{9} \times \frac{8}{9} \times \frac{1}{9} \times \frac{1}{14}$$

$$\frac{2}{9} \times \frac{4}{9} \times \frac{1}{9} \times \frac{4}{9} \times \frac{2}{9} \times \frac{4}{9} \times \frac{4}{9} \times \frac{1}{9} \times \frac{1}{14}$$

$X = \{\text{rain, hot, high, false}\}$

$$P(\text{rain} | X) = \frac{2}{9} \times \frac{2}{9} \times \frac{3}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{1}{14}$$

$$\frac{2}{9} \times \frac{6}{9} \times \frac{1}{14}$$

$$\frac{12}{504} \rightarrow 0.023$$

$$\frac{2}{5} \times \frac{3}{5} \times \frac{4}{8} \times \frac{4}{5} \times \frac{5}{14}$$

$$\frac{3}{25} \rightarrow 0.036$$

Hence, X is classified with N.

Metrics for Evaluating Classifier Performance

+ve

TP (True+ve) :- refers to tuples that were correctly labelled by classifier

TN True -ve

~~FP~~ (False +ve) :- refers to ~~-ve~~ tuples that were correctly labelled by classifier.

FP (False +ve) :- These are -ve tuples that were incorrectly labelled

~~FN~~ (False -ve) :- These are +ve tuples that were incorrectly labelled.

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

$$\text{error rate} = \frac{FP + FN}{P + N}$$

$$\text{recall, } TPR = \frac{TP}{P}$$

$$TNR, \text{ specificity} = \frac{TN}{N}$$

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

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

		Actual		Total
		Yes	No	
Predicted	Yes	6954	46	7000
	No	412	258.8	3000
Total	7366	2634	10000	

Accuracy = $\frac{6954 + 258.8}{10000} = 0.9552\%$

Error = 0.0458

Sensitivity =

1. Nominal :- categories, states or "names of things"
 - The numbers are not intended to be used quantitatively
 - Also, values do not have any specific order
2. Binary Attributes → Two values
3. Ordinal Attributes → Values have meaningful order or a ranking among them but magnitude is not known.
- + Numeric Attributes:-

 - ↳ (4.1) Interval Scaled Attributes
(No true zero point)
 - ↳ (4.2) Ratio scaled Attributes
(Inherent zero point)

5. Discrete attributes :-
Has only a finite or a countably infinite set of values which may or may not be represented as integers.
6. Continuous Attributes -
Practically, real values can only be measured and represented using finite number of digits.

3 Areas of Statistical Descriptions :-

- Measuring central tendencies
 - ↳ Mean
 - ↳ Median $\rightarrow \frac{N+1}{2}$
 - ↳ Mode
 - ↳ Midrange
- Measuring dispersion,
- Graphic display of data.

SUPPOSE

$x = \{2, 4, 4, 5, 5, 7, 9\}$ mixtures

$\therefore \text{Mean} = 5$

$$\text{variance} = \frac{(2-5)^2 + (4-5)^2 + (4-5)^2 + (5-5)^2 + (5-5)^2 + (7-5)^2 + (9-5)^2}{8}$$

$$= \frac{32}{8} = 4$$

$\therefore SD = \sqrt{\text{var}}$

$= 2$

Proximity Measure for Nominal Attribute.

Object	test-1 (Nominal)
1	code A
2	code B
3	code C
4	code A A

$$d(i, j) = \frac{P - M}{P}$$

$P \rightarrow$ total attributes
describing objects

$M \rightarrow$ total no. of matched

$$\begin{bmatrix} 0 & d(2,1) & 0 \\ d(3,1) & 0 & d(3,2) \\ d(4,1) & d(4,2) & d(4,3) & 0 \end{bmatrix}$$

$$d(2,1) = \frac{7-0}{7} = 1$$

$$d(3,1) = \frac{7-0}{7} = 1$$

$$d(3,2) = \frac{7-0}{7} = 1$$

$$d(4,1) = \frac{7-0}{7} = 1$$

$$d(4,2) = \frac{7-0}{7} = 1$$

$$d(4,3) = \frac{7-0}{7} = 1$$

Proximity measures for binary attributes.

A contingency table for binary data

$$\text{for symmetric} \Rightarrow d(i,j) = \frac{r+s}{q+r+s+t}$$

Object j

	1	0	sum
1	$q+r$	$q+r$	$q+r$
	s		

$$\text{for asymmetric} = d(i,j) = \frac{r+s}{q+r+s}$$

	0	1	$s+t$	$s+t$
0	$q+s$	$t+t$	$q+s+t+t$	p
	sum	$q+s$	$t+t$	p

Name	Gender	Fever	Cough	Test1	Test2	Total	Total
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

* Only attribute
 * Gender is symmetric attribute
 rest all asymmetric

* Grade let Y and P be 1 and N = 0

$$M=1, F=0$$

	Gender	Fever	Cough	Test1	Test2	Total	Total
Jack	1	1	0	1	0	0	0
Mary	0	1	0	1	0	1	0
Jim	1	1	1	0	0	0	0

$$d(jack, mary) = q=1, r=0, s=0, t=1 = \frac{1}{3} = 0.33$$

$$d(jack, jim) = q=1, r=1, s=1, t=1 = \frac{2}{4} = 0.5$$

$$d(jim, mary) = q=1, r=0, s=1, t=1 = \frac{1}{4} = 0.25$$

Proximity Measured for Mixed Attributes.

No	Test I C Nominal	Test II Ordinal	Test III C Numeric
1	Code A	Excellent	98
2	Code B	Fair	22
3	Code C	Good	64
4	Code A	Excellent	28

For Mixed Attributes calculate values for each and do the average.

for Test I.

$$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

for Test II.

$$Z_{it} = \frac{(x_{it} - \bar{x}_i)}{\sigma_i}$$

$$2 \rightarrow 0.33$$

$$1 \rightarrow 0$$

$M_{it} \rightarrow$ Total no. of Attributes.

$$3 \rightarrow 1$$

$$\therefore 2/2 \rightarrow 1$$

$$2 \rightarrow 1$$

Code A (1)	0			
Code B (0)	1	0		
Code C (0.5)	0.5	0.5	0	
Code D (1)	0	1	0.5	0

For Test 3.

Code A (45)	0	17	17	17	17
Code B (22)	23	0	17	17	17
Code C (64)	19	42	0	17	17
Code D (18)	127	6	3.6	0	17

$$\text{Min} = 0$$

$$\text{Max} = 42$$

Divide by 42 for

Normalization.

Code A	0	0.41	0.41	0.41	0.41
Code B	0.55	0	0.41	0.41	0.41
Code C	0.45	1.0	0	0.41	0.41
Code D	0.40	0.14	0.86	0	0.41

why Data Preprocess in

- ii) incomplete
- iii) inaccurate or noisy
- iv) inconsistent

full accepted multidimensional view in

- ii) Accuracy
- iii) completeness
- iv) consistency
- v) Timeliness
- vi) Believability
- vii) Interpretability

ACC
TB Int

Handle NOisy
data (Binning
method)

- Smooth by bin mean
- Smooth by bin median
- Smooth by bin boundary

Major Tasks in Data Preprocessing in

1) Data Cleaning → fill in missing values
→ identify outliers and smooth out noisy data

2) Data Integration → correct inconsistent data
→ resolve redundancy caused by integration

3) Data Transformation

correlation analysis
can be used to handle
data redundancy.

4) Data Reduction

5) Data Discretization

Regression

• Linear regression.

• Helps in finding the best line
to fit b/w the two variables
so that one can be used to
predict other.

• Multiple reg:-

involves more than two
variables.

Dimensionality Redn:-

- Wavelet Transform
- Principal Component
- Attribute Selection

Dimensionality Redn:-

• Cube Aggregation

• Tree

• Clustering

• Discretization

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

χ^2 Test

There's also something
called as

covariance of Numerical Data

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

Explained in next page.

Stock price for AU Electronics and High Tech.

Time Point	AU Electronics	High Tech.
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5
t6		

$$E(\text{AUelect}) = \frac{6+5+4+3+2}{5} = \frac{20}{5} = \$4$$

$$E(\text{High Tech}) = \frac{20+10+14+5+5}{5} = \frac{54}{5} = \$10.80$$

$$\text{cov}(A \cdot B) = E(A \cdot B) - \bar{A} \cdot \bar{B}$$

$$= \frac{6 \times 20 + 10 \times 5 + 14 \times 4 + 5 \times 3 + 2 \times 5}{5} - 4 \times 10.80$$

$$= 50.2 - 43.2$$

$$= 7.$$

Since, covariance is +ve we can say that stock prices for both companies rise together.

Min-max Normalization

$$w_i' = \frac{v_i - \min_A}{\max_A - \min_A} (new_{-max}_A - new_{-min}_A) + new_{-min}_A$$

$$= \frac{200 - 200}{1000 - 200} (1 - 0) + 0$$

= 0.

$$0, 0.125, 0.25, 0.35, 1$$

For Data

$$200, 300, 400, 500, \\ 1000.$$

$$= \frac{300 - 200}{800} (1) + 0$$

$$= \frac{100}{800} (1) + 0$$

$$= \frac{400 - 200}{800}$$

$$= \frac{200}{800} \frac{1}{4}$$

$$\frac{600 - 200}{800}$$

$$= 1.$$

Naive Bayes Algorithm.

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

$$P(\text{Default} = \text{No}) =$$

Homeowner	Martial status	Job Experience	Refunded
Yes	Single	3.	No
No	Married	4.	No
No	Single	5.	No
Yes	Married	4.	No
No	Divorced	2.	Yes
No	Married	4.	No
Yes	Divorced	2.	No
No	Married	3.	Yes
No	Married	2.	No
Yes	Single	2.	Yes

$$P(D = \text{No}) = 7/10$$

$$P(D = \text{Yes}) = 3/10$$

Q

	Yes	No
Yes	1/3	3/7
No	2/3	4/7

	condition	
	Yes	No
Married	1/3	4/7
Single	1/3	2/7
Divorced	1/3	1/7

	Yes	No
1	2/3	4/7
2	1/3	2/7
3	1/3	2/7
4	3/7	0/7
5	0/7	1/2

$$\text{for } \langle \text{No, Married, } 3 \rangle = 0.022$$

$$2/3 \times 1/3 \times 1/3 \times 4/7 \times 1/10 \text{ (Yes)}$$

$$4/7 \times 1/7 \times 2/7 \times 1/7 \times 1/10 \text{ (No)} = 0.082$$

Gini Index

Gini Index for overall collection:-

$$\Rightarrow 1 - \left[\left(\frac{5}{14} \right)^2 + \left(\frac{9}{14} \right)^2 \right]$$

$$\Rightarrow 0.459$$

Gini Index for Age :-

Youth $\rightarrow 5$, middle-aged $\rightarrow 4$, senior $\rightarrow 5$

$$\text{Gini (Youth)} = 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right]$$

$$= 0.48$$

$$\text{Gini (middle-aged)} \approx 0$$

$$\text{Gini (senior)} = 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right]$$

$$= 0.48$$

$$\text{wt. avg} = 0.48 \times \frac{5}{14} + 0.48 \times \frac{5}{14}$$

$$= 0.342$$

Gini Index Income:-

high $\rightarrow 4$, medium $\rightarrow 6$, low $\rightarrow 4$

$$\text{Gini (high)} = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right]$$

$$= 0.5$$

$$\text{Gini (medium)} = 1 - \left[\left(\frac{4}{10} \right)^2 + \left(\frac{2}{10} \right)^2 \right]$$

$$= 0.44$$

$$\text{Gini (low)} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right]$$

$$= 0.375$$

$$\text{wt. avg} = 0.5 \times \frac{4}{14} + 0.44 \times \frac{6}{14} + 0.375 \times \frac{4}{14} = 0.438$$

student

$$\text{No} \rightarrow 7 \quad (\text{Avg} \rightarrow 7)$$

$$\text{Gini (No)} := 1 - [(4/7)^2 + (3/7)^2] \\ = 0.489$$

$$\text{Gini (Avg)} := 1 - [(6/7)^2 + (1/7)^2] \\ = 0.244$$

$$\text{wt. avg} = 0.489 \times 7/14 + 0.244 \times 7/14 \\ = \frac{7}{14} [0.733] \\ = 0.3665$$

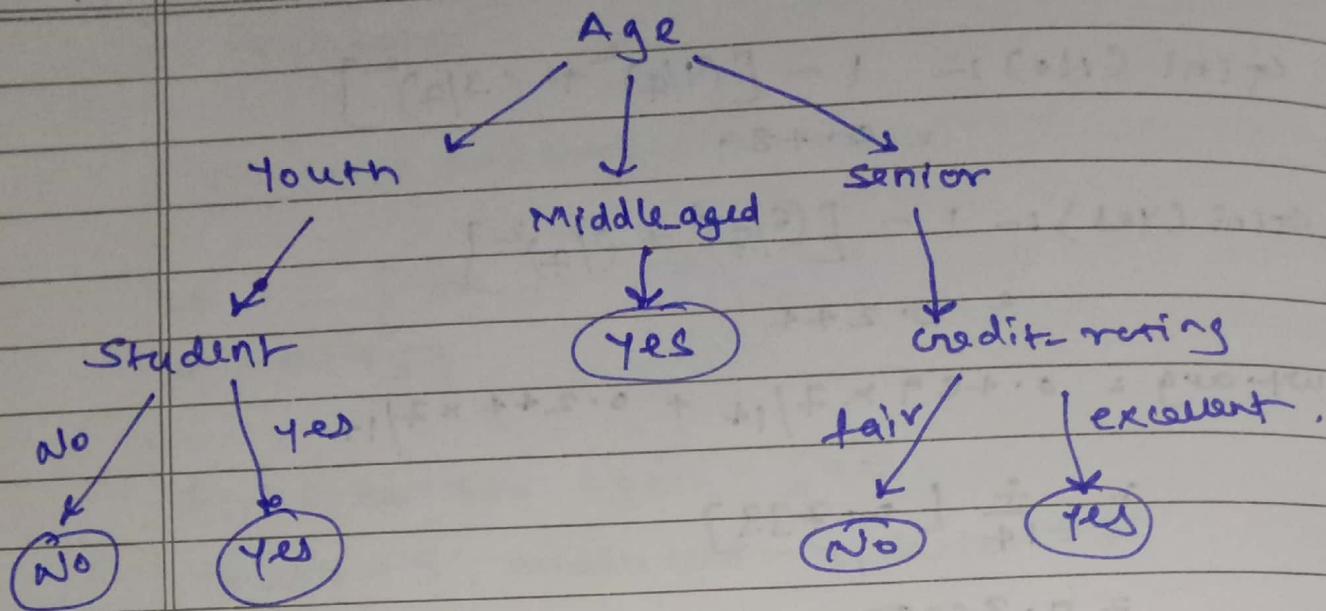
credit Rating :-

Fair :- 8 Excellent \rightarrow 6

$$\text{Gini(Fair)} := 1 - [(7/8)^2 + (6/8)^2] \\ = 0.375$$

$$\text{Gini(Excellent)} := 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right] \\ = 0.5$$

$$\therefore \text{wt. avg} = 0.375 \times \frac{8}{14} + 0.5 \times \frac{6}{14} \\ = 0.428$$



	Income	Age	Own House
1.	Very-High	Young	Yes
2.	High	Medium	No
3.	Low	Young	Rented
4.	High	Medium	No
5.	Very-High	Medium	Yes
6.	Medium	Young	No
7.	High	Old	Yes
8.	Medium	Medium	Rented
9.	Low	Medium	Rented
10.	Low	Old	Rented
11.	High	Young	Yes
12.	Medium	Old	Rented

Gini Index :- Yes :- 7 Rented :- 5

$$1 - \left[\left(\frac{7}{12}\right)^2 + \left(\frac{5}{12}\right)^2 \right] = 0.486$$

For Income.

Very High \rightarrow 2, High \rightarrow 4, Medium \rightarrow 3, Low \rightarrow 3

$$\text{Gini}(\text{Very High}) \rightarrow 0 \quad \{1 - 1\}$$

$$\text{Gini}(\text{High}) \rightarrow 0$$

$$\begin{aligned} \text{Gini}(\text{Medium}) &\rightarrow 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] \\ &= 0.44 \end{aligned}$$

$$\text{Gini}(\text{Low}) \rightarrow 0$$

$$w \cdot \text{avg} = 0.44 \times \frac{8}{12} = 0.11$$

For Age.

Young \rightarrow 4, Medium \rightarrow 5, Old \rightarrow 3

$$\begin{aligned} \text{Gini}(\text{Young}) &\rightarrow \{1 - \left[\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] \\ &= 0.375 \end{aligned}$$

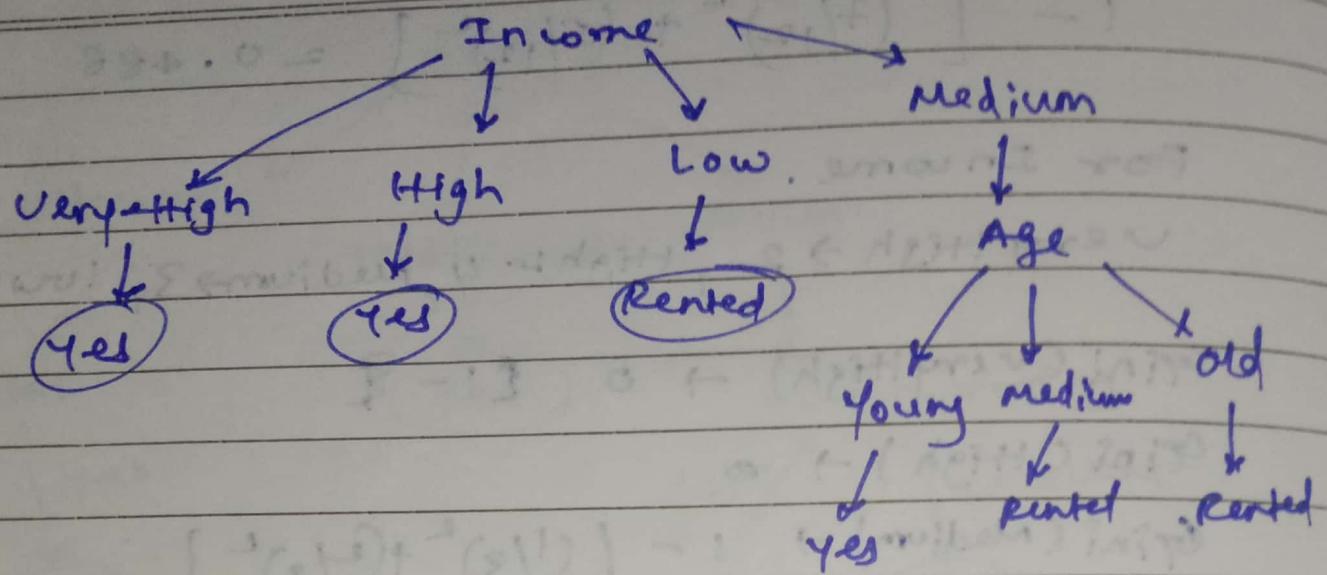
$$\text{Gini}(\text{Medium}) \rightarrow 1 - \left[\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right] \quad 0.4 \times 5$$

$$\text{Gini}(\text{Old}) \rightarrow 0.48$$

$$\begin{aligned} \text{Gini}(\cancel{\text{Old}}) &\rightarrow 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] \\ &\rightarrow 0.444. \end{aligned}$$

$$0.375 \times \frac{4}{12} + 0.48 \times \frac{5}{12} + 0.44 \times \frac{3}{12} =$$

$$0.375 \times \frac{4}{12} + 2.0 + 0.11 = 2.235.$$



Normalization

with z-score.