

→ Data mining - Process of extracting extra knowledge & hidden patterns from large amounts of data

Steps in data mining

- ① Data Selection - extracting relevant data from various sources.
- ② Data Preprocessing - Converting data to consistent form
Removal of unnecessary info.
- ③ Data Transformation suitable format
- ④ Data mining - finding hidden patterns
- ⑤ Pattern evaluation
- ⑥ Presenting

* Data mining applications

- (i) Customer segmentation
- (ii) Risk management
- (iii) Fraud detection

Data Mining

→ Types of attributes

① Categorical (Qualitative)

(i) Nominal :- They are just some kind of symbols used to distinguish one object from other.
eg - zip codes, employee ID.

(ii) Ordinal :- They define an ordering of objects apart from just distinguishing.
eg - grades.

② Quantitative attributes (having properties of numbers)

(i) Numeric - ~~discrete~~

(a) Interval : Unit of measurement exist. (+, - can be applied. not \div & \times)
eg - temp. in °C & F

(b) Ratio : \times & \div also can be applied.
eg - length, mass, age

(ii) Discrete :- It has a finite or countably infinite set of values. Binary attributes are special types of discrete attributes.
eg - Profession, zip code

(iii) Continuous :- It has an infinite no. of states.
eg - Infinite no. of float values b/w 2 & 3.

* For asymmetric attributes, it is a special type of binary attribute where both the values are not equally important.

→ Precision, Bias & accuracy

(i) Precision - Closeness of repeated measurement to one another. It is measured by S.D of set of values.

(ii) Bias - Variation of measurements from quantity being measured. Measured by Diff b/w mean of set of values & known value of quantity being measured.

→ Outliers

* The data points that have characteristics or value that are diff from most of the other data objects in data set.

* Unlike noise, they can be sometimes be of interest. (in fraud detection).

→ Missing values

Ways to handle it:

(i) Eliminate data objects or attributes:-
Reliable if missing values is less. Very risky to eliminate attribute as it may be imp. for analysis.

- (ii) Estimate missing values : we can use method of KNN. (K nearest neighbour)
- (iii) Ignore the missing value during analysis

→ Issues related to application

- 1) Timeliness : If data provides a picture of ongoing phenomena, such as purchasing behaviour, this is relevant for only a limited time
- 2) Relevance : Common problem of sampling bias: When a sample does not contain diff types of objects in proportion to actual occurrence in population
eg - Survey responses are not by whole population

→ Data pre processing

Topics -

- (i) Aggregation
- (ii) Sampling
- (iii) Dimensionality reduction
- (iv) Feature subset selection
- (v) Feature creation
- (vi) Discretization
- (vii) Variable transformation

① Aggregation

- Combining of two or more objects into a single object
- How ~~as~~ the aggregation is done is something to look out for
- Aggregation leads to smaller data sets which require less memory & time
- ~~Adv~~ → A group of objects generally have less variation ~~than~~ than individual objects.

⇒ Data Quality

* Examples of data quality problems

(i) Noise & outliers

(Regression, binning, clustering) (box plot & SD method)

(ii) missing values

(KNN method, remove)

(iii) duplicate data

(iv) Inconsistent values

(i) Noise :- Refers to modification of original values. eg - distortion of person's voice.

② Sampling

- * Used because whole dataset is too expensive & time consuming to look at so we look at a sample
- * Sample size should be large enough
- * Using a sample should work as using the entire data set
- * Sample is representative of dataset if it has a mean close to that of original data.
- * Sampling approach -
 - ① Simple random sampling
 - Randomly selecting the objects
 - equal prob. of selecting any object
 - ② Sampling without replacement
 - as each item is selected, it is removed from whole pop"
 - ③ Sampling with replacement
 - Objects are not removed from pop as they are selected
 - More simpler to analyze as prob. of selecting any obj remains constant throughout
 - * These approaches do not adequately represent all object types.

④ Stratified Sampling : Splitting the data into several partitions, draw random samples from each partition.

→ Equal no. of obj. are drawn from each partition ensuring nothing is left out.

* Sample size should be large enough for accuracy but it eliminates purpose of sampling.

Small sample size leads to erroneous conclusions.

⑤ Progressive Sampling

* → Proper sample size can be difficult to determine, hence we start with small sample size & progresses to large one which is a sample of sufficient size.

→ When we need to look when we need to stop increasing sample size Note

③ Dimensionality reduction

- * Data sets may have large no. of features & attributes which are not necessary.
- * Data mining algos work better on data set with less attributes.

* Data can be more easily visualized.

* Curse of Dimensionality :-

As the dimension of data increases, it becomes more harder to analyze as chances of noisy & irrelevant data increases.

* Techniques :-

(i) Principal Component analysis (PCA)

Linear algebra technique to find new attributes that are linear combination of original attributes & are orthogonal to each other.

(ii) Singular value decomposition (SVD)

④ Feature Subset Selection

* Creating new attributes that are a combination or a subset of old attributes.

* Techniques :-

(i) Embedded approach - Feature selection occurs naturally as part of data mining algo. The algo itself decides the attributes.

(ii) Filter approach - Features are selected before data mining algo. is run through independent task.

(iii) Wrapper approach - The data mining algo is used as black box to find best subset of attributes.

⑤ Feature creation

* To create new set of attributes from original attributes

Techniques

(i) Feature extraction - The creation of a new set of features from original raw data

⑥ Discretization

* To replace raw values of numeric attribute by interval levels.

* Numerous continuous attributes are replaced by small intervals

Techniques

(i) Equal width discretization : It divides m data into equal intervals

(ii) Equal freq. discretization : Same no. of objects into each interval.

(iii) K means :- Clustering method, generally forms clusters of similar objects

⑦ Binarization

* If there are n. categorical values, $n = \log_2 m$ digits are reqd to represent it in binary form

⑧ ~~No~~ Variable Transformation

→ Techniques:

(i) Normalization or standardization: Attribute values are normalized by scaling their values so that they fall in specified range.

e.g. If we have our data which is scattered, we can use Z-transformation so that our values lie b/w 0 & 1.

$$\boxed{Z \text{ value} = \frac{x - \bar{x}}{s}} \quad \begin{array}{l} \bar{x} \rightarrow \text{mean} \\ s \rightarrow \text{SD} \end{array}$$

But if we have outliers in our data, we use median instead of mean in our calculation of SD.

→ Examples of Proximity decision

~~Question~~ → Simplest w/ n & y be two objects
that consist of n binary attributes

foo → no. of attributes where x is 0 & y is 0

f01 → " " " " where x is 0 & y is 1

f10 → " " " " where x is 1 & y is 0

f11 → " " " " where x is 1 & y is 1

(i) Simple matching coeff

$$\text{SMC} = \frac{\text{no. of matching attribute}}{\text{no. of attributes}} = \frac{f_{11} + f_{01}}{f_{01} + f_{10} + f_{11} + f_{00}}$$

(ii) Jaccard coeff

$$J = \frac{\text{no. of matching presence}}{\text{no. of attributes except 00}} = \frac{f_{11}}{f_{01} + f_{10} + f_{00}}$$

(iii) Cosine similarity

$$\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

$$\text{eg- } x = (3, 2, 0, 5, 0, 0, 0, 2, 0, 0)$$

$$y = (1, 0, 0, 0, 0, 0, 0, 1, 0, 2).$$

$$x \cdot y = 3 \times 1 + 2 \times 0 + 0 \times 0 + 5 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 2 \times 1 + 0 \times 0 + 0 \times 2 \\ = 5$$

$$\|x\| = \sqrt{3^2 + 2^2 + 5^2 + 0^2 + \dots} = 6.48$$

$$\|y\| = \sqrt{1^2 + 0^2 + \dots + 2^2} = 2.24$$

$$\cos(x, y) = 0.31$$

Classification :

→ Definition

- * Task of assigning objects to one of the predefined categories
- * It is a task of mapping an input attribute set x into its class label y .
- * It is usually done by a technique called decision tree induction.

e.g- classifying the species ~~are~~ by looking at various attributes like Body temp, skin cover, etc into classes like mammal, reptile, bird, fish, etc.

e.g- categorizing news as finance, weather, entertainment, sports

→ General approach

Training set

1
2
3
4
5
...
10

Induction

Learn model

Model

Test set

1
2
3
4
5

Deduction

Apply model

Model

→ We need to predict the classes of these five unknown records using my training set.

(Prediction algo)

* Evaluation of performance of classifier model is based on number of test records. correctly & incorrectly predicted by model. These numbers are represented in the confusion matrix.

Predicted class

	Class = 1	Class = 0
Actual Class = 1	f_{11}	f_{10}
Actual Class = 0	f_{01}	f_{00}

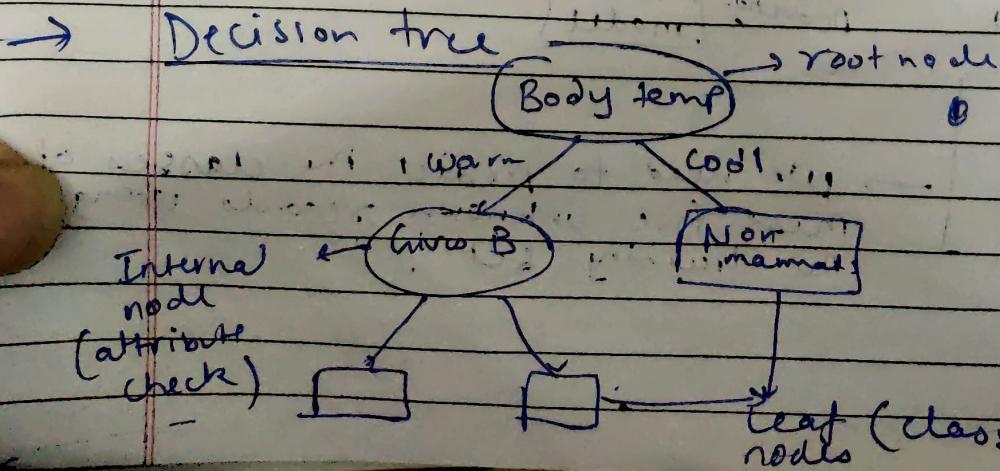
No. of correct predictions = $f_{11} + f_{00}$

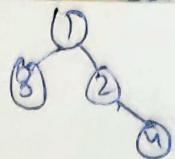
"incorrectly predicted as class 1" = $f_{10} + f_{01}$

where f_{01} = no. of records from class 0 incorrectly predicted as class 1

* Accuracy = $\frac{\text{No. of correct pred}}{\text{Total no. of pred}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$

* Error rate = $\frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$



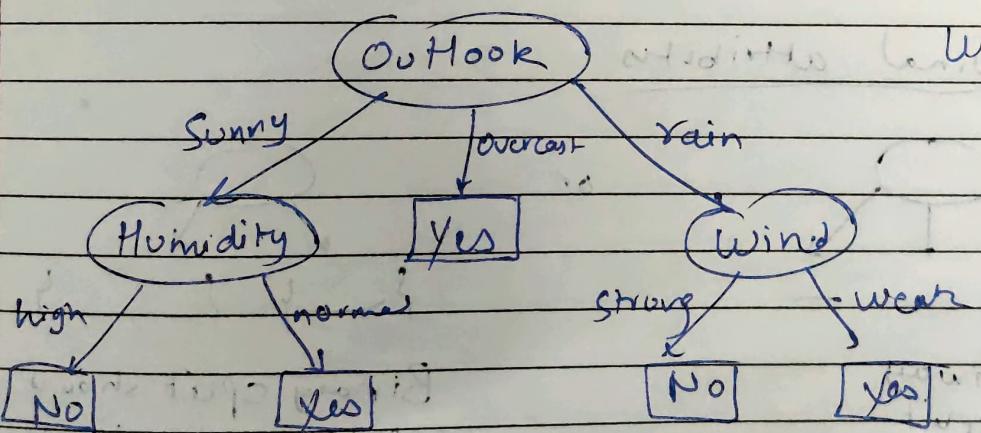
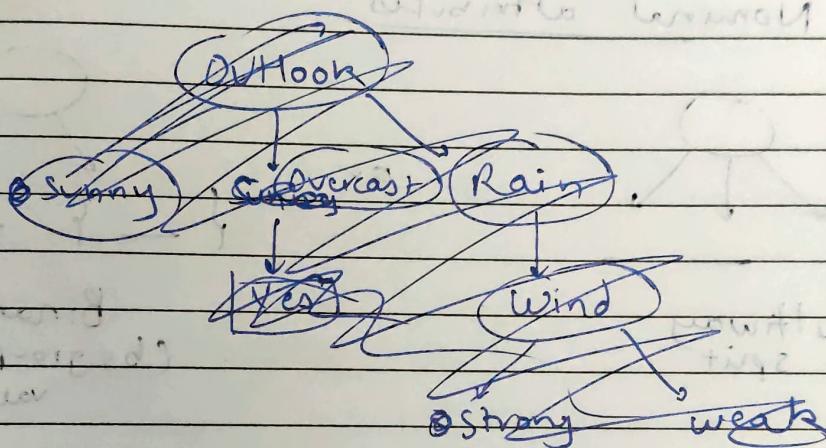


Now, on the basis of decision tree, we can classify the new objects that come.

* How to build a decision tree Hunt's algo :-

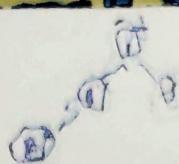
We have to build an optimal decision tree starting from an attribute which gives the most pure results as compared to other attributes.

eg -



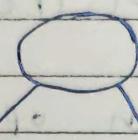
* Design issues of decision tree

- How should training records be split?
- How should the splitting procedure stop?

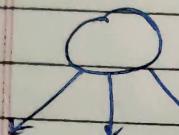


* Methods for expressing Attribute test conditions

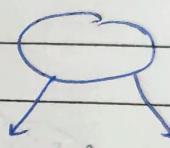
- (i) Binary attributes Two outcomes



- (ii) Nominal attributes

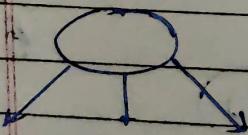


Multway split

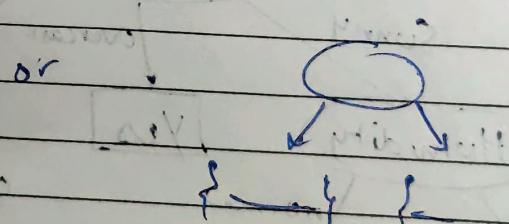


Binary split
(by grouping attribute values)

- (iii) Ordinal attributes



Multway split



Binary split shows
ensure order is preserved
eg - small cannot be combined
with large

- (iv) Continuous attributes

No Binary / Multway split.

or
Through discretization, we can make intervals

III II I

classmate

Date _____

Page _____

* Decision Tree using Gini index

Example

Weekend	Weather	Parents	Money	Decision	1	2	3	4	5
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					
w ₁₁									
w ₁₂									
w ₁₃									
w ₁₄									

To decide root node.

① Gini index for overall sample → 3 possibilities

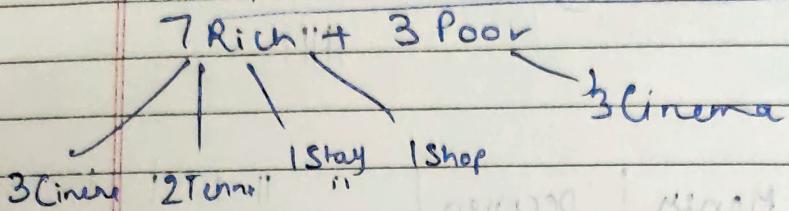
→ Gini index of overall sample

$$6C + 2T + 1ST + 1SH$$

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

$$= 0.58$$

→ Gini index of money attribute

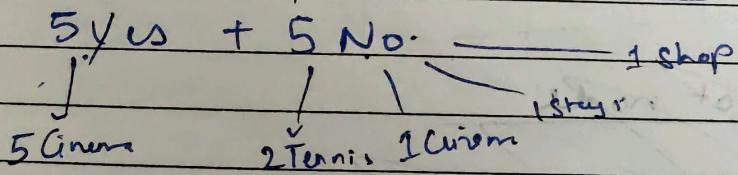


$$\text{Money} = \frac{\text{Rich}}{\text{Rich}} = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{2}{7} \right)^2 + \left(\frac{1}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] = 0.694$$

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

$$\text{Weighted index} = \frac{7}{10} \times 0.694 + 0 = 0.486$$

→ Gini index of Parents attribute



Parents =

Yes

$$1 - \left[\left(\frac{5}{5} \right)^2 \right] = 0$$

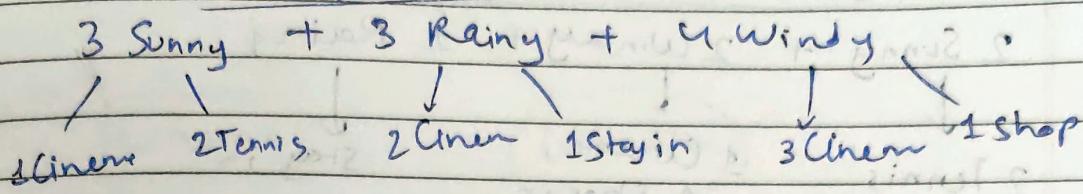
Parents =

No

$$1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \right] = 0.72$$

$$\text{Weighted index} = 0 + \frac{5}{10} \times 0.72 = 0.36$$

→ Cini index of weather attrin



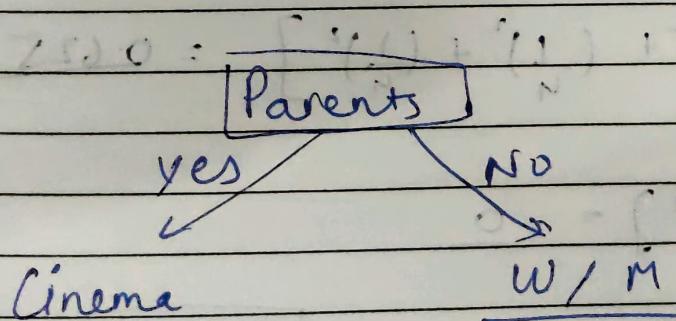
Sunny $= 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] = \underline{0.444}$

Rainy $= 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = \underline{0.444}$

Windy $= 1 - \left[\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] = \underline{0.375}$

Weighted $= 0.444 \times \frac{3}{10} + 0.444 \times \frac{3}{10} + 0.375 \times \frac{4}{10}$
 $= \boxed{0.416}$

∴ So, we select Parent as root node



→ Gini index of Parents | Weather

2 Sunny + 2 Windy + 1 Rainy

↓ ↓ ↓
2 Tennis 1 Cinema,
 1 Shop in

$$\Rightarrow 1 - \left(\frac{2}{2}\right)^2 = 0$$

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

$$\Rightarrow 1 - \left[\left(\frac{1}{4}\right)^2\right] = 0$$

$$\text{Weights} = 0 + 0 + \frac{2}{5} \times 0.5 = \boxed{0.2}$$

→ Gini index of Parents | Money

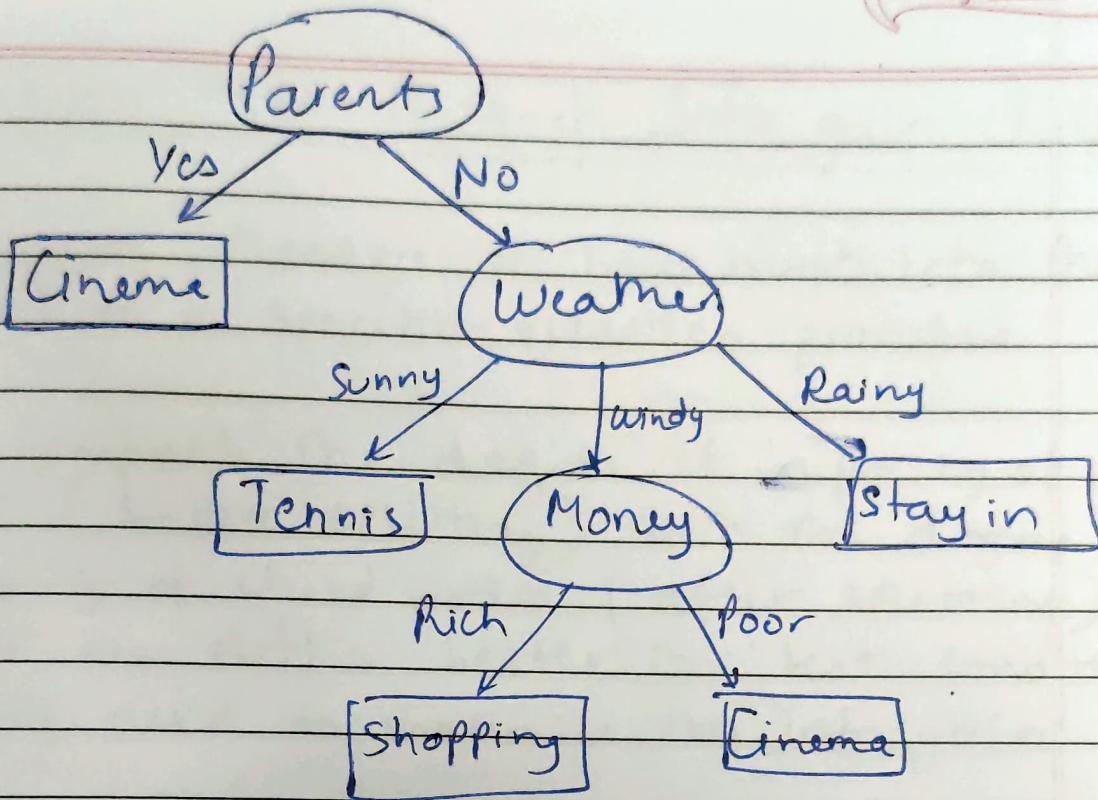
4 Rich + 1 Poor
↓ ↓
2 Tennis 1 Stay 1 Shop 1 Cinema

$$\Rightarrow 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{1}{5}\right)^2\right] = 0.625$$

$$\Rightarrow 1 - \left[\left(\frac{1}{4}\right)^2\right] = 0$$

$$\text{Weights} = 0 + \frac{4}{5} \times 0.625 = \boxed{0.5}$$

So, we select weather.



→ Decision tree through Info. gain & Entropy

- * Info' gain: Measure of how much info. the answer to a specific question provides.

We compare the degree of impurity of parent node (before splitting) with the degree of impurity of child nodes (after splitting). Larger the diff., better the test condition. So, we need to maximize the info.gain.

- * Entropy - measure of how much uncertainty is there in info.
- * Info gain $\uparrow \rightarrow$ Entropy \downarrow

* Formulas :-

$$\text{Info gain} : I(P_i, n_i) = -\frac{P_i \log_2 P_i}{S} - \frac{n_i \log_2 n_i}{S}$$

S: Total sample size

$$\text{Entropy} : E(A) = \sum_{i=1}^n P_i n_i I(P_i, n_i)$$

$$\text{Gain} : \text{Info gain} - \text{Entropy}$$

Example

Day	Outlook	Temp.	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Info gain:

$$\begin{aligned}
 &= - \left[\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right] \\
 &= - \left[\frac{9}{14} (3.16 - 3.80) + \frac{5}{14} (2.32 - 3.80) \right] \\
 &= - \left[(0.64)(-0.64) + (0.357)(-1.48) \right] \\
 &= \boxed{-0.4096 - 0.52832} \\
 &= \boxed{0.937} = 0.940
 \end{aligned}$$

Entropy of Outlook

	Yes	No	Total	Prob
Sunny	3	2	5	5/14
Overcast	4	0	4	4/14
Rain	3	2	5	5/14

$$\begin{aligned} I(\text{Outlook, Sunny}) &= - \left[\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right] \\ &= - \left[(0.6) (1.58 - 2.32) + 0.4 (1 - 2.32) \right] \\ &= - [(0.6)(-0.74) + (0.4)(-1.32)] \end{aligned}$$

$$I(\text{Outlook, Overcast}) = 0$$

$$I(\text{Rain}) = 0.972$$

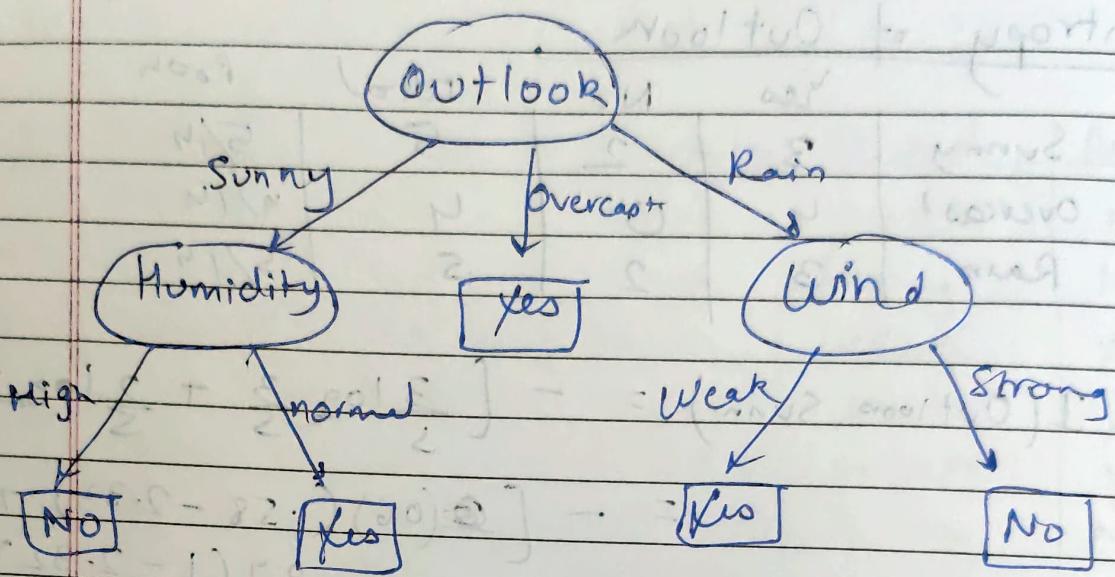
$$\begin{aligned} \text{Entropy} &= 0.972 \times \frac{5}{14} + 0 + 0.972 \times \frac{5}{14} \\ &= 0.70 \end{aligned}$$

Gain of Outlook = $0.940 - 0.70 = 0.24$

$$\text{Gain of Temp} = 0.029$$

$$\text{Gain of Humidity} = 0.151$$

$$\text{Gain of Wind} = 0.048$$



* Impurity measures like entropy & Gini index tend to favour attributes that have large no. of distinct values.

* Customer ID cannot be used as root node because it is unique for every customer & we will not be able to use decision tree

* Gain ratio = $\frac{\text{Info gain}}{\text{Split info}}$

Evaluating performance of classifier

① Hold out method

- * Original data is partitioned into two sets - training set & test set
- * Model is imposed on training set & tested on test set
- * Proportion of data for training & test set can be 50-50 or 75-25.
- * Accuracy of model on test set needs to be calculated
- * Limitations :-
 - (i) Fewer records are there for training set which affects accuracy of model
 - (ii) Smaller the training set, larger the variance of model ; larger the training set . We cannot judge accuracy from smaller test set

② Random subsampling

- * Hold out method repeated several times to improve estimation of classifier's performance is called random subsampling

③ Cross validation

- * Suppose we partition data into two equal subsets . First, we choose one of them for training ; other for testing . Then, we swap the roles . The total error is sum of error for both runs .
- * K-fold cross validation method generalises this by segmenting data into K equal partitions .