

Supervised learning setup:

Training data, D comes in pair (x, y) where $x \in R^d$ is input instance and y is its label.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \subseteq R^d \times C$$

Where,

R^d = d dimensional feature space

x_i = input vector of i^{th} sample

y_i = output vector of i^{th} sample

C = label space

Different types of label Space

Binary classification $\Rightarrow C = \{0, 1\}$ or $\{-1, 1\}$

Multiclass classification $\Rightarrow C = \{0, 1, \dots, K\}$

Regression $\Rightarrow C = \mathbb{R}$, Real number set

Hypothesis class

Before we find the function ' h ' such that $h(x) \approx y$ we have to find to which class this fn belongs to. Different classes are, AI, Neural networks, Decision tree, Regression, etc.

tree, Regression, etc.

The set of all hypothesis functions is known as hypothesis class.

No free lunch theorem -

This theorem states that every successful ML algo. must make assumptions. This also implies that there is no single ML algorithm that can solve every problem.

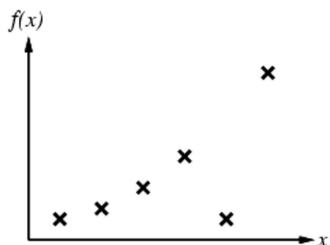
pure induction task is Supervised (Inductive) Learning

- Given a set of hypothesis function h find such h that approximate the function f .

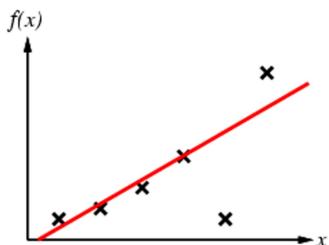
i.e if we are given Training data, D as

$$D = \{(x_i, f(x_i)) \mid f(x_i) = y_i\}$$

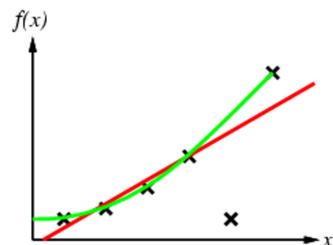
what we have to do is that find a h such that $h \approx f$ with high probability.



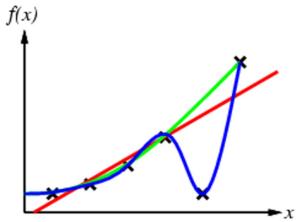
Given D



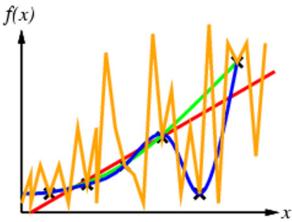
One possible h



Another possible h



One more h



another one

So, here we have different h 's but we have to take one which approximate f . i.e. $h \approx f$

Ockham's razor: prefer the simplest hypothesis consistent with data

Supervised Learning

Regression

Learning a continuous function f

Classification

Learning a discrete function f

linear regression

$$\begin{array}{l} \rightarrow \text{Simple } h_0(x) = \theta_0 + \theta_1 x \\ \rightarrow \text{Multiple } h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_k x_k \end{array}$$

Simple linear Regression :

Here, our hypothesis h is,

$$h(x) = \theta_0 + \theta_1 x \quad , \quad \Theta = (\theta_0, \theta_1)$$

Now, we have to adjust ' h ' means adjust Θ such

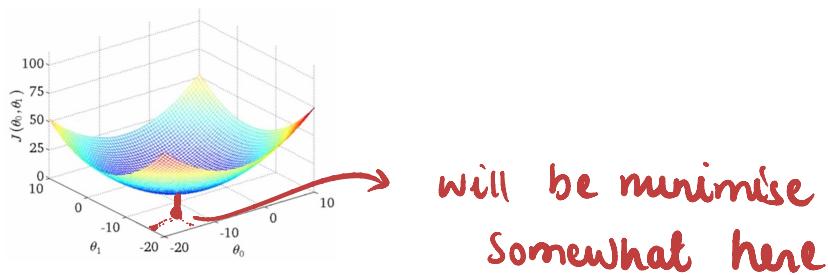
Now, we have to adjust 'h' means adjust θ such that $h \approx f$ with high probability.

Parameters are θ_0, θ_1 .

$$\text{Cost function, } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

Squared error

Goal : Minimise $J(\theta_0, \theta_1)$



Gradient Descendent :

Let some function $J(\theta_0, \theta_1)$

- we have to find $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$
- Start with some $\theta_0 \in \theta_1$
- Change θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully reach minimum.

Algorithm:

while (not converged)

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad \text{for } j = 0, 1 \\ \downarrow \end{array} \right.$$

this shows simultaneous update

}

Simultaneous update :

— Right —

$$\left. \begin{array}{l} \text{temp}_0 = \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \\ \text{temp}_1 = \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \end{array} \right\} \text{first calculate new values}$$

$$\left. \begin{array}{l} \theta_0 = \text{temp}_0 \\ \theta_1 = \text{temp}_1 \end{array} \right\} \text{then update}$$

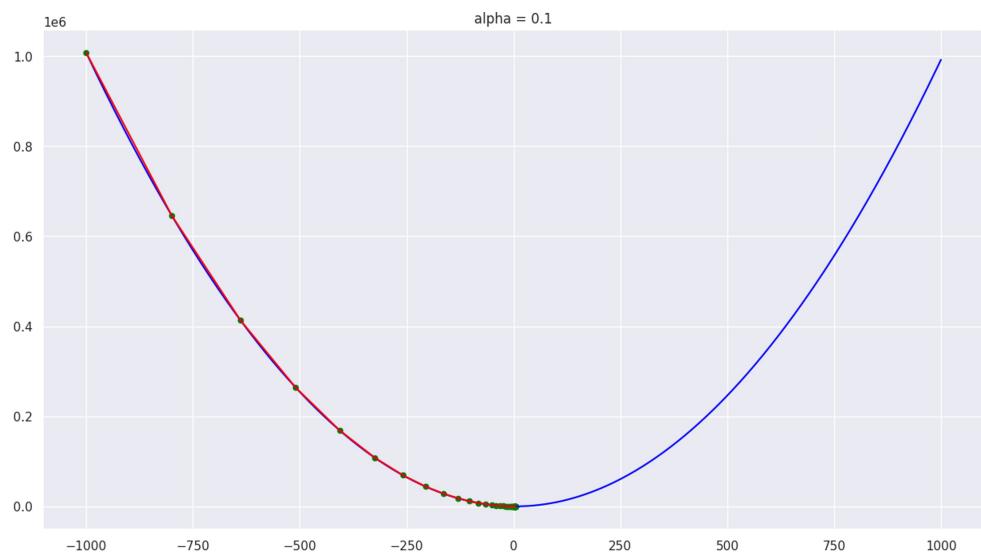
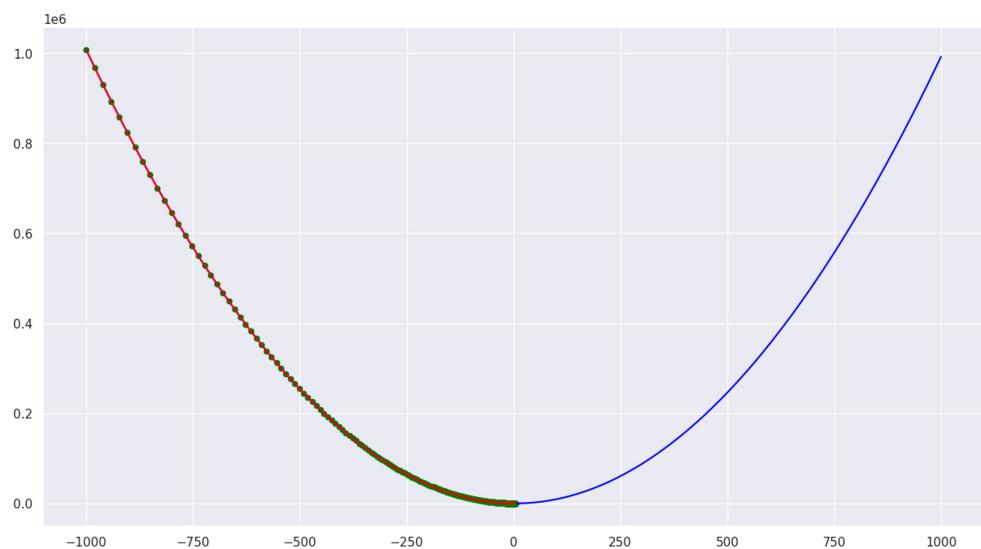
— Wrong —

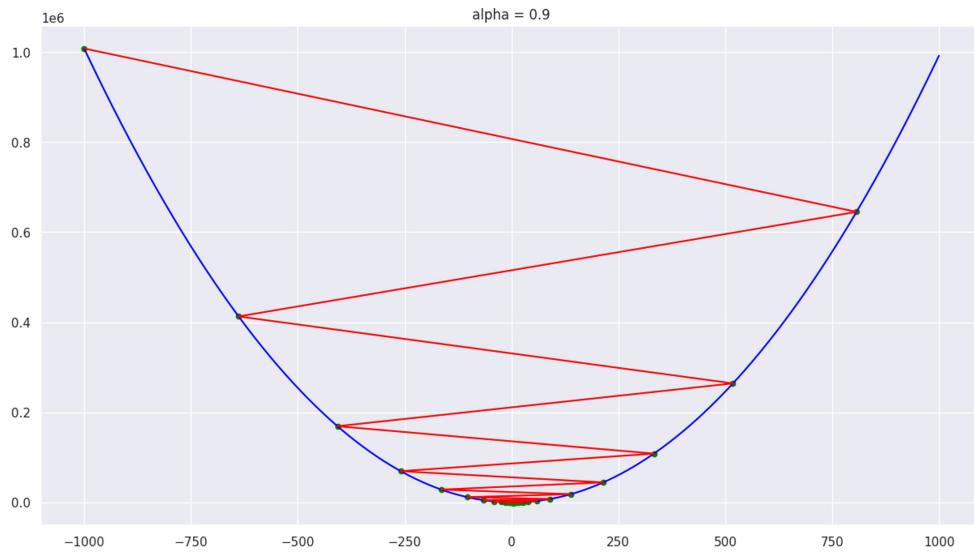
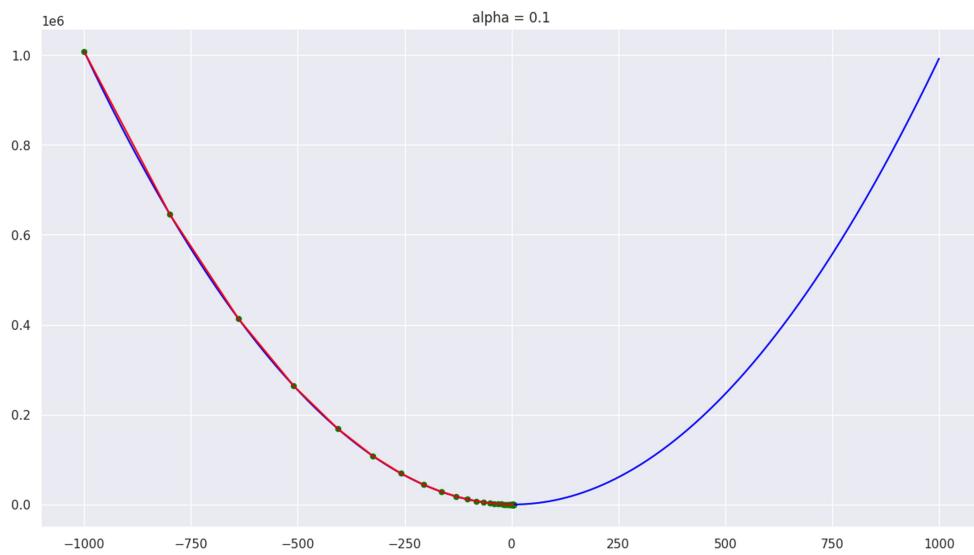
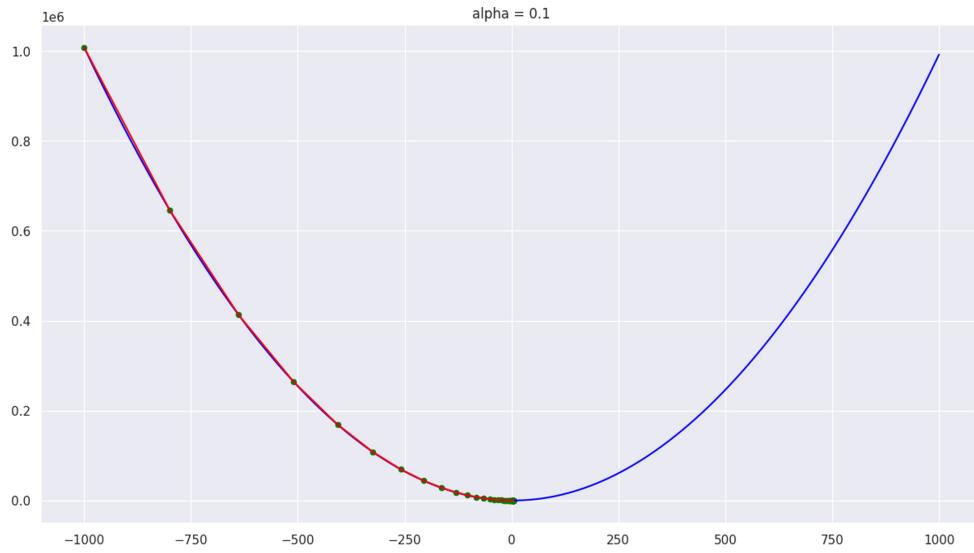
$$\theta_0 = \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

α is known as the learning rate that actually decide how slow, fast we are reaching towards minima

Given below figures show a graph of function $f(x) = (x-3)(x-5)$ and Gradient descendent with different values of alpha .





Multivariate linear regression:

Multivariate linear regression:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\ = \theta^T x$$

where,

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, \quad X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Cost function :

$$J(\theta) = \frac{1}{2m} \sum \epsilon_i^2$$

where ϵ_i = error in i^{th} instance of prediction

$$\epsilon_i^2 = (h_{\theta}(x_i) - y_i)^2$$

$$h_{\theta}(x_i) = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \dots + \theta_n x_{in}$$

differentiating ϵ_i^2 wrt θ_j

$$\frac{\partial}{\partial \theta_j} \epsilon_i^2 = 2(h_{\theta}(x_i) - y_i)^2 \cdot x_{ij}$$

$$S_0, \quad \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_{ij}$$

Feature Scaling:

- Make sure that features are on the similar scale.

$x_1 \rightarrow$ Number of bedrooms : 1-5

$x_2 \rightarrow$ Size of room : (0 - 2000)

$$x_1 = \frac{x_1}{5} \quad \Delta \quad x_2 = \frac{x_2}{2000}$$

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1$$

Logistic regression: classification:

linear regression $\Rightarrow h_{\theta}(x) = \theta^T x$

Logistic regression $\Rightarrow h_{\theta}(x) = g_{\theta}(\cdot \theta^T x)$

$$\text{where , } g(x) = \frac{1}{1+e^{-x}} \quad \left\{ \begin{array}{l} \text{Sigmoid or} \\ \text{logistic fn.} \end{array} \right.$$

Due to this $h_{\theta}(x) \in [0,1]$

if $h_{\theta}(x) < 0.5 \rightarrow$ predict 0

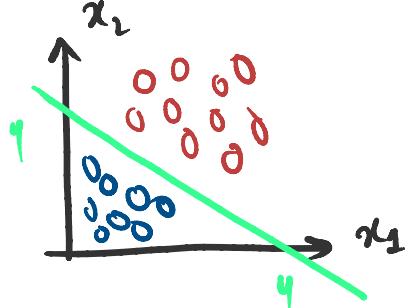
else predict 1.

Here, $h_{\theta}(x)$ represent the probability of any feature being salient.

$$h_{\theta}(x) = P(y=1 | x; \theta)$$

probability that $y=1$ given x & parametrised by θ .

Decision boundary



Let's take $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} = \theta_{\text{min}}$

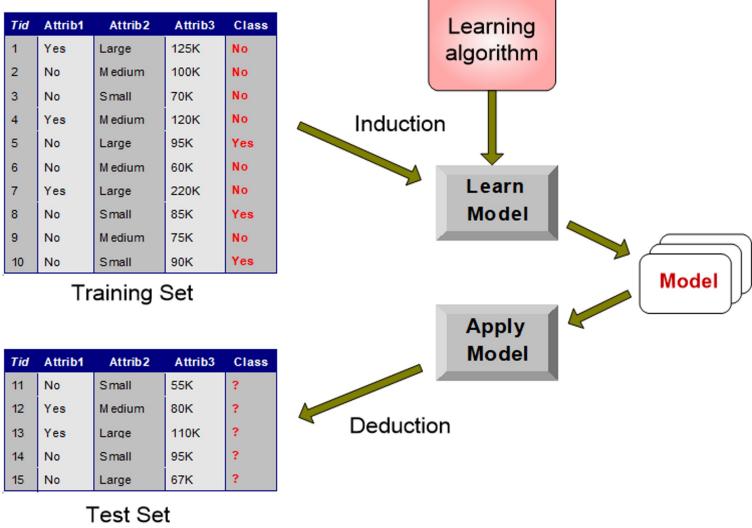
then

$$\theta^T x = -3 + x_1 + x_2$$

so, $g(\theta^T x) = 0.5$ if $\theta^T x = 0$

i.e. $x_1 + x_2 = 3$ → This is our decision boundary.

Decision boundary is a property of hypothesis function.



Loss function :

- Determine how good or bad our hypothesis fn is

a) zero-one loss function { generally used in classification }

$$L_{0/1}(h) = \frac{1}{m} \sum_{i=1}^m \delta_{h(x_i) \neq y_i}$$

where, $\delta_{h(x_i) \neq y_i} = \begin{cases} 1 & \text{if } h(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases}$

b) Squared loss function { generally used in regression settings }

$$L_{\text{squared}}(h) = \frac{1}{m} \sum (h(x_i) - y_i)^2$$

- 1) If prediction is far away from actual then loss suffered will grow quadratically and that leads to better prediction.

2) If prediction is so close to actual then its square will be more less and then less attention will be given to that example.

Absolute loss

$$L_{abs}(h) = \frac{1}{m} \sum |h(x_i) - y_i|$$

To optimise we should reach

$$h(x) = \text{Median}_{p(y|x)} [y]$$

optimal hypothesis fn:

$$h = \operatorname{argmin}_{h \in \mathcal{H}} L(h)$$

h will be the hypothesis for which the loss is minimum.

cost function for logistic regression:

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$$

We have a point 'x' that has to classified into some class
 Let $S_x \subseteq D$ and $|S_x| = K$ be a set of K Nearest
 neighbours of point X.

then our hypothesis function is -

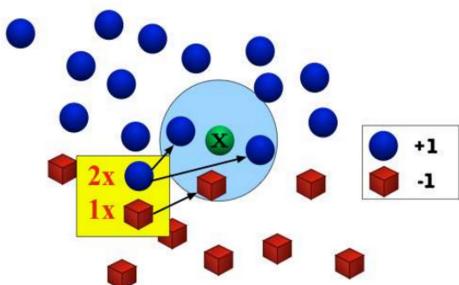
$$h(x) = \text{mode}(\{y' : (x', y') \in S_x\})$$

i.e select the class of high frequency.

Assumptions taken for this problem is -

Similar inputs have Similar Output.

Example: $K=3$



Distance function:

$$\text{dist}(x_1, x_2) = \left\{ \sum_{i=1}^d |x_i - x_{i2}|^p \right\}^{\frac{1}{p}}$$

This is Known as "Minkowski distance"

This is known as "Minkowski distance"

$p=1 \Rightarrow$ Manhattan Distance

$p=2 \Rightarrow$ Euclidean Distance

curse of dimensionality :

- As the dimension of feature space increases the closeness of points decreases.

Let's illustrate the above statement with an example:

Lets take a cubic of edge length = 1 as our dataset i.e it encloses all the points in our dataset.

Thus, to enclose k Neighbours we take the edge length of min hypercube = l

Now,

$$l^d \approx \frac{k}{N}$$

Assuming distribution is uniform over volume.

$$l \approx \left(\frac{k}{N}\right)^{\frac{1}{d}}$$

So, as dimension increases value of ℓ will increases.

d	ℓ
2	0.1
10	0.63
100	0.955
1000	0.9954

k-NN summary

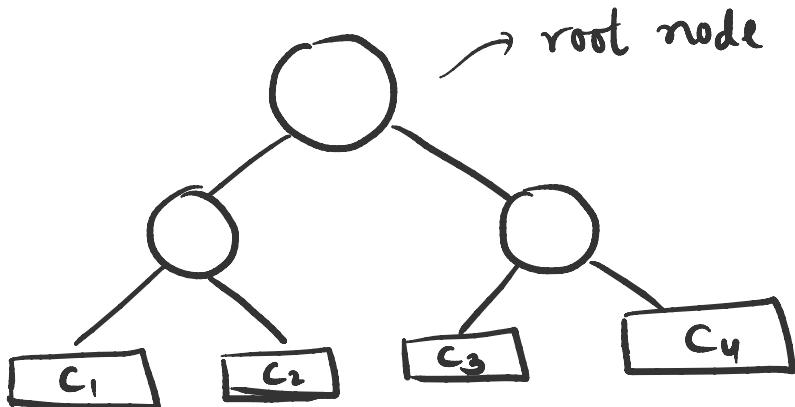
- k-NN is a simple and effective classifier if distances reliably reflect a semantically meaningful notion of the dissimilarity. (It becomes truly competitive through metric learning)
- As $n \rightarrow \infty$, k-NN becomes provably very accurate, but also very slow.
- As $d \gg 0$, points drawn from a probability distribution stop being similar to each other, and the kNN assumption breaks down.

Given a database $D = \{t_1, t_2, \dots, t_n\}$ where t_i is the tuple in database, which is defined by set of attribute $A = \{A_1, A_2, A_3, \dots, A_n\}$. Also given a set of classes $C = \{c_1, c_2, \dots, c_n\}$

Decision tree :

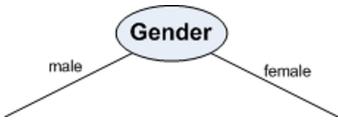
○ → internal nodes

□ → leaf nodes

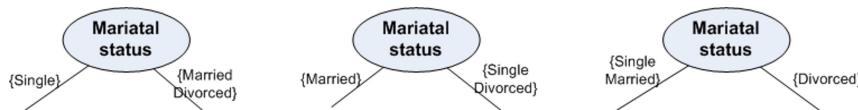


Node splitting in building binary tree:

a) Binary attribute

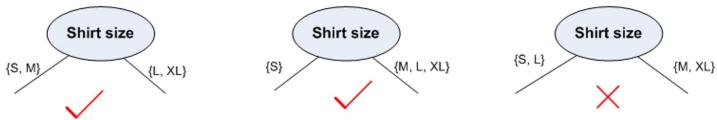


b) Nominal attribute



c) Ordinal attribute

- Order should be maintained properly:



d) Numerical attribute



Concept of Entropy:

- This is needed to chose the order of attribute selection
- First time it was used to measure information content in a message.

Entropy Calculation:

The minimum number of Yes/No question asked to reach to the answer.

Let suppose we have 'm' values

then its entropy is $H = \log_2 m$

Alternate definition of entropy:

The entropy of a set of m distinct values is the number of bits needed to encode all the values in the most efficient way.

Example:

00 : North	contains 4 values
01 : East	
10 : West	
11 : South	Number of bits - ?

10 : West
11 : South

→ Number of bits
needed = 2

If distribution of values is not uniform

Example: Let take four values A, B, C and D and with probability P_1, P_2, P_3 & P_4 respectively.

Then, for this we use huffman coding not simple Coding

$A \rightarrow b_1$: bits to represent A

$B \rightarrow b_2$

$C \rightarrow b_3$

$D \rightarrow b_4$

So, On an average number of bits required are

$$b_{avg} = \sum P_i b_i$$

So, if an object occurs with possibility P then Number of bits required to represent it is

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

$$\text{So, } b_{avg} = - \sum P_i \log_2 P_i = E_{\text{entropy}}$$

--,

varg

↳ entropy

Entropy of a dataset

$$E = - \sum p_i \log(p_i)$$

E is always non-negative

$m=2^k$

in case of
uniform
distribution

$$0 \leq E \leq \frac{1}{k} \log_2 K$$

→ when all data is
distributed
uniformly
over K classes

When all
data is distributed
in one class
only

- Different algorithms have been proposed to take a good control over

- Choosing the best attribute to be splitted, and
- Splitting criteria

ID₃ Algorithm:

→ Several terms

$$E(D) = - \sum p_i \log(p_i)$$

$$E_A(D) = \sum_{j=1}^m \frac{|D_j|}{|D|} E(D_j)$$

Number of sub dataset
'j' which is partition
on Attribute 'A'

Then information gain $\alpha(A, D)$ will be

Then information gain $\alpha(A, D)$ will be

$$\alpha(A, D) = E(D) - E_A(D)$$

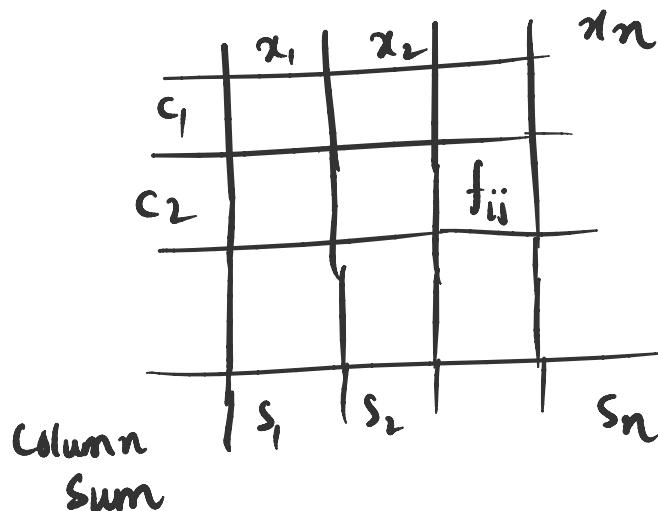
Frequency table : A shortcut method for calculating α

Data set <i>Table A</i>			X <i>Table X</i>				
X	Y	Class		1	2	3	4
1	1	A	A	1	1	1	1
1	2	B	B	1	1	1	1
2	1	A	C.Sum	2	2	2	2
2	2	B					
3	2	A					
3	1	B					
4	2	A					
4	1	B					

Table Y		
	1	2
A	2	2
B	2	2
C.Sum	4	4

Frequency table of X

Frequency table of Y



$$V = \sum_{i,j} f_{ij} \log f_{ij} \quad \text{for all } i \neq j$$

$$S_i = \sum_j f_{ij} \log f_{ij} \quad \text{for all } i$$

then, $E_x(D) = (-V+S)/N$

$$\alpha(x, D) = E(D) - E_x(D)$$

CART Algorithm:

In spite of using entropy of dataset it uses $G(D)$ i.e. the measure of impurity of dataset.

$$G(D) = \sum (1 - p_i^2)$$

$$0 \leq G(D) \leq 1 - \frac{1}{K}$$

\downarrow
when there
is only
single class

\downarrow
when data is
distributed
uniformly

Now, similar to ID₃ we divide D into D₁ and D₂ on attribute A then.

$$\underline{\gamma(A, D)} = G(D) - G_A(D)$$

Gini Index /

Impurity reduction

Calculation using frequency table

$$G(A, D) = 1 - \frac{1}{|D|} \sum_{j=1}^m \frac{1}{|D_j|} \sum_{i=1}^K f_{ij}^2$$

$$G(A, D) = 1 - \sum_{j=1}^m \frac{|D_j|}{|D|} \underbrace{\log}_{\substack{\text{Column} \\ \text{Sum}}} \sum_{i=1}^{|D_j|} \dots$$

Algorithm C4.5 -

$$\text{Split information} = - \sum_{j=1}^m \frac{|D_j|}{|D|} \log \frac{|D_j|}{|D|} = E_A^*(D)$$

then our golden ratio is , ↳ Here 'm' Shows
number of values
of Attribute 'A'

$$\beta = \frac{\alpha(A, D)}{E_A^*(D)}$$

$$0 \leq E_A^*(D) \leq \log_2 m$$



↓
Single

attribute value .

↳ m shows the number
of distinct values
of A .

calculation of $E^*(A, D)$ using frequency table

Split info.

$$\sum_{i=1}^m s_i \ln s_i$$

$$= - \sum_{j=1}^m \frac{s_j}{|D|} \log_2 \frac{s_j}{|D|}$$