

History →

i) first machine learning system was developed in the 1950's.

ii) 1952, Arthur Samuel was at IBM developed a program for playing checkers.

According to Samuel → It is a field of study that gives computers the ability without being explicitly programmed.

1960's

→ Neural N/D → Rosenblatt's perceptron.

→ pattern recognition.

→ Minsky & Papert proved limitation of perceptron.

1970's

→ Symbolic Concept induction

→ Expert systems and knowledge acquisition bottleneck

→ Quinlan's ID3

→ Natural language processing.

1980's

→ Advanced decision tree and rule learning

→ learning and planning and problem solving

→ Resurgence of Neural N/D

→ Focus on Experimental Methodology

90's ML & Statistics →

→ support vector Machines

→ Data Mining

→ Adaptive agent & web Applications.

→ Text learning.

→ Reinforcement learning

→ Ensembles

→ Bayes Net learning

1994 → self driving car road test

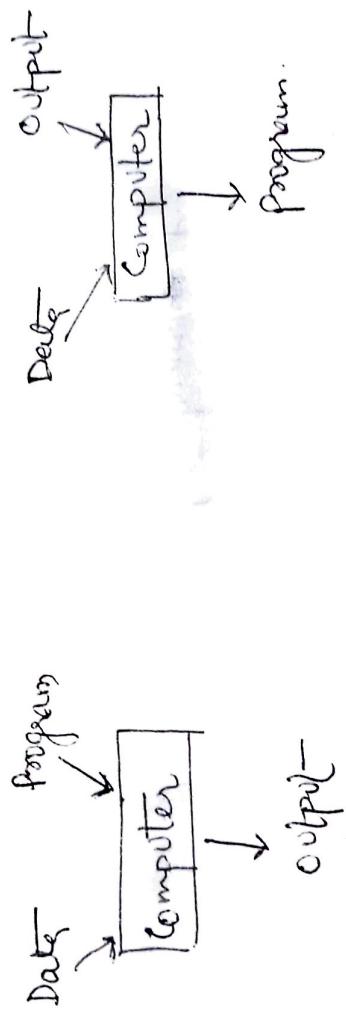
1997 → Deep blue beats Gary Kasparov

Popularity of this field in recent time and the reasons behind that -

- New Software / Algorithm
 - Neural Network
 - Deep learning
- New Hardware
 - GPU's
- Cloud Enabled
- Availability of big data

How Machine Learning solutions differs from "Programmatic Sol'n" →

Algorithm



Algorithm → The program takes data as input and the program produces output.

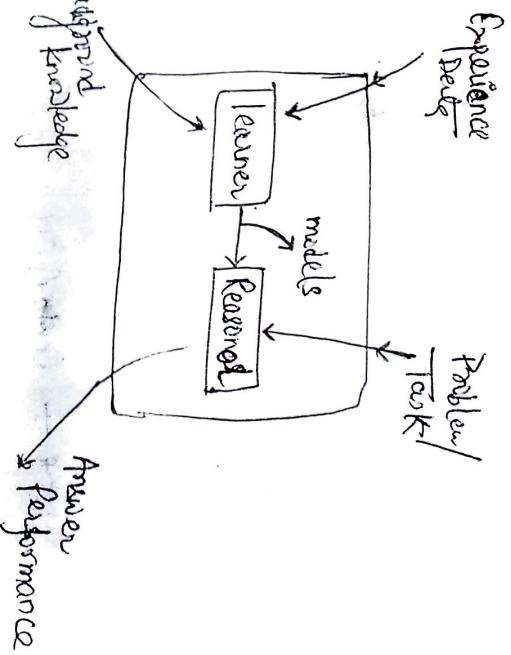
ML → you are feeding the data as input as well as output data and you are getting a program or a model with which you can solve subsequent tasks.

ML → learning is the ability to improve one's behavior with experience. So, it is about building computer system, that automatically improves with experience and we have to discuss what are the fundamental laws that govern the learning processes.

Machine learning explores algorithm that learn from data, build models from data and this model can be used for different tasks.

Exm - prediction, decision Making or solving tasks.

Tom Mitchell → def → A computer program is said to learn from experience E with respect to some task T, if its performance at tasks in T, as measured by P, improves with experience E.



Domain & Applications →

Medicine:-

- Diagnose a disease →
- Input - symptoms, lab measurements, test results, DNA tests, ...
- output - one of set of possible diseases or 'none of the above'
- Data-mining historic medical records to learn which future patients will respond best to which treatment.

Vision :-

Robot control

NLP

Speech recognition

Machine Translation

Financial

- ↑ features
1. choose the training experience.
 2. choose the target function (that is to be learned) → hypothesis language
 3. choose how to represent the target function class of f^n on features.
 4. choose a learning algorithm to infer the target f^n .

Types of Learning →

- Supervised learning → I/P (labeled training data)
 - O/P (de-classified training examples)
 - Given an observation x , what is the best label for y ?

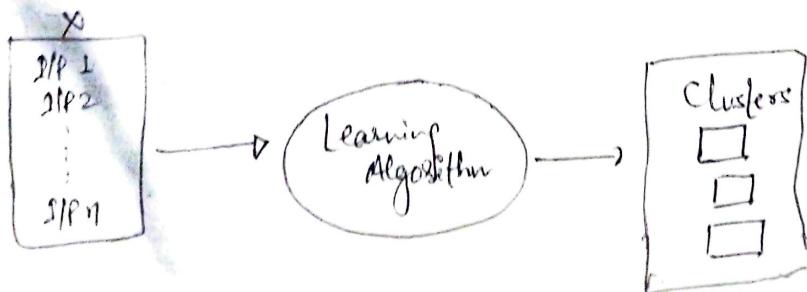
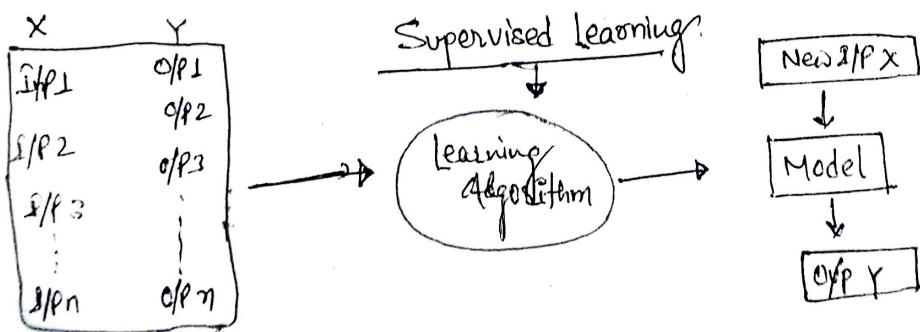
Unsupervised learning → (unlabeled training data).

- x
- Given a set of x 's, cluster or summarise them

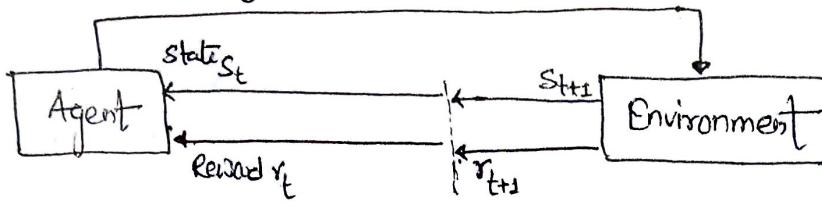
Reinforcement learning —

- Determine what to do based on rewards & punishments.

Semi-Supervised Learning — Combo of supervised & unsupervised learning.



reinforcement learning →



Supervised learning →

(input feature)

$$x_1, x_2, \dots, x_n$$

discrete value feature (classification)

continuous value feature (regression)

y (target feature)

I_1	a_1, a_2, \dots, a_n	y_1
I_2	b_1, b_2, \dots, b_n	y_2
I_3	c_1, c_2, \dots, c_n	y_3
-	-	-
-	-	-

Test instance.

$$z_1, z_2, \dots, z_n$$

- A set of input features x_1, x_2, \dots, x_n .
- A target feature y .
- A set of training exam where the values for the input features and the target features are given for each exam.
- predict the values for the target features for the new exam.
 - Classification when y is discrete.
 - Regression when y is continuous.

Classification →

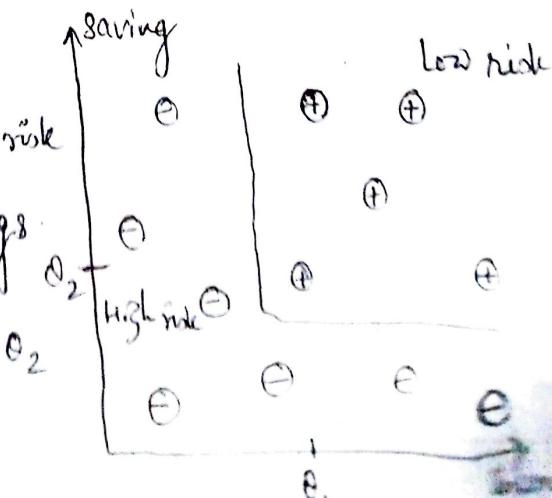
Ex → Credit Scoring

Differentiating between low risk & high risk customers from their income & savings.

Discriminant :- if $\text{income} > 0$, AND $\text{savings} > 0_2$

then low-risk

else High risk



Regression :-

Price of used car

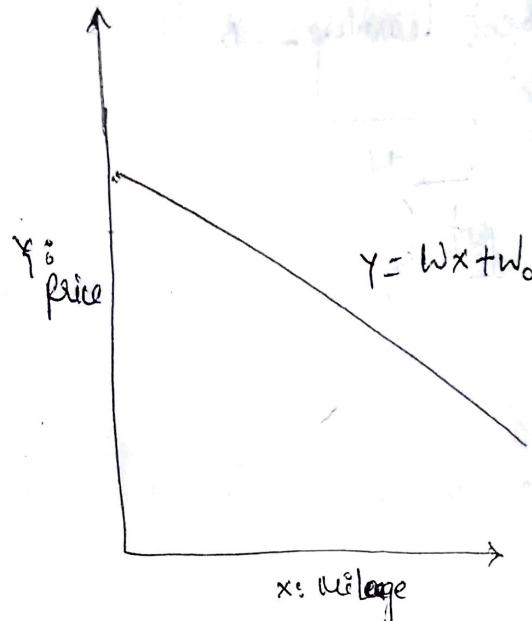
x : Car attributes

y : Price

$$Y = g(x, \theta)$$

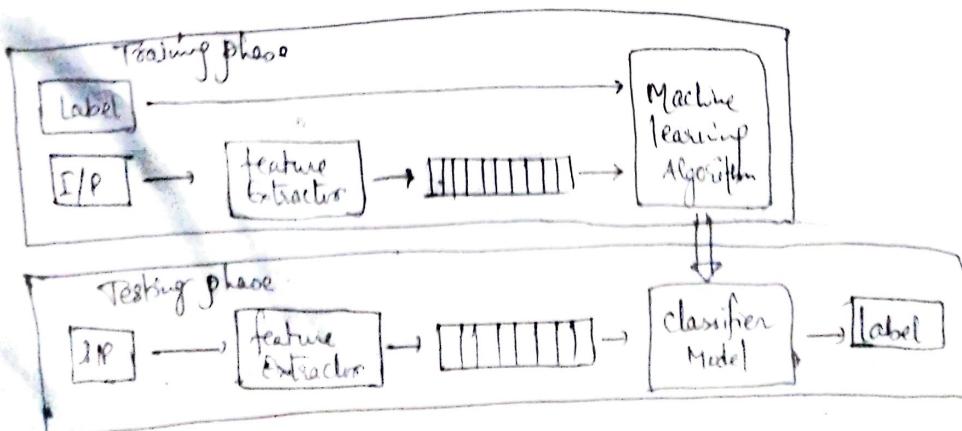
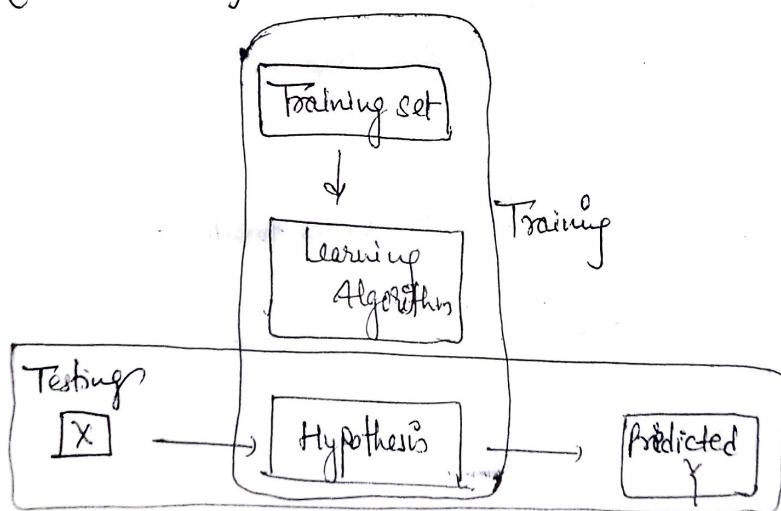
gU Model

θ Parameters



The individual observations are analyzed into a set of quantifiable properties which are called features.

- Categorical ("A", "B", "O", blood group).
- ordinal ("large", "small")
- Integer-valued (Age)
- Real valued (height)



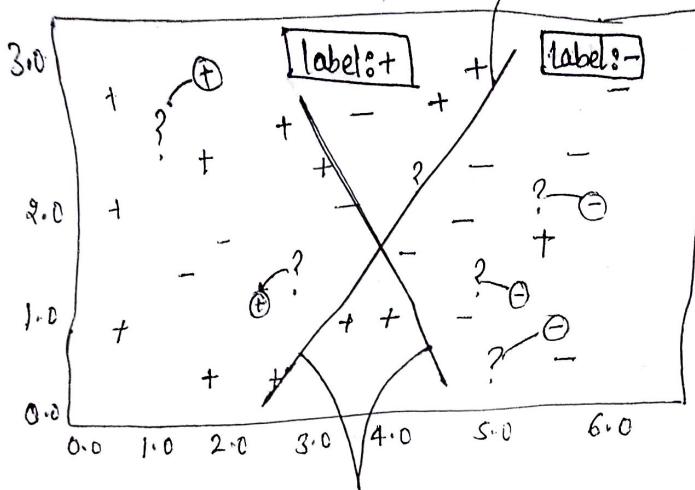
Given exm $(\hat{x}, y), (\hat{x}, f(x))$

Classification $f(x) \rightarrow$ discrete

Regression $f(x) \rightarrow$ Continuous.

Probability of estimation $f(\hat{x}) = \text{probability of } \hat{x}$.

functions.



Hypothesis Space:
Set of legal Hypothesis

Hypothesis Space:-

- The space of all hypothesis that can, in principle, be output by a learning algorithm.
- We can think about a supervised learning machine as a device that explores a "Hypothesis Space".
 - Each setting of the parameters in the machine is a different hypothesis about the f that maps input vectors to output vectors.

exm (x, y)

Training data \Rightarrow set of examples,
feature space $\rightarrow X$

Concept $\rightarrow C, C \subseteq X$

Target function $f, [f : X \rightarrow Y]$

Concern about choosing H.S.

- size.
- randomness
- parameters.

Inductive learning / prediction :-
Need to make assumptions.
experience alone doesn't allow us to make Conclusion about unseen data instances.

Two types of Bias -

- Restriction - limit the Hypothesis space. (Polynomials - specifying the form of the f)
- Preference - Impose ordering on hypothesis space.

Inductive learning :- Inducing a general function from training examples.

- Construct hypothesis h to agree with C on the training set.
- A Hypothesis is Consistent if it agrees with all training examples.
- A Hypothesis said to be recognized generalize well if it correctly predicts the value of y for novel example.

Inductive learning Hypothesis :- Any Hypothesis 'f' found to approximate the target function 'c' well over a sufficiently large set of training examples 'D' will also approximate the target function well over other unobserved examples.

Types of Inductive Biases -

- Occam's Razor :- the simplest consistent Hypothesis about the target fn. is actually the best.
- Minimum description length :-
- Maximum Margin :- When drawing the boundary b/w two classes, attempt to maximize the width of the boundary (SVM).

Generalization :-

Components of Generalization Error -

- Bias - how much the average model over all training sets differ from the true model.
 - Error due to inaccurate assumption / simplifications made by the model.
- Variance - how much models estimated from different training sets differ from each other.

Underfitting and Overfitting:-

Underfitting → model is too simple to represent all the relevant class characteristics.

- High bias and low variance. (low)
- High training error and High test error.

Overfitting → Model is too 'complex' and fits irrelevant characteristics of data.

- low bias and high variance.
- low training error and High test error.

Evaluation:-

H - Hypothesis Space

S - Training data

h - Learning algorithm.

$h \in H$

Experimental Evaluation → Matrix based

Error Matrix

Accuracy

Precision / Recall.

Evaluating Predictions:-

Suppose we want to make a prediction of a value for a target feature on example

- y is the observed value of target feature on example x .

- \hat{y} is the predicted value of target feature on example x .

• How is the error measured?

Absolute Error :- $|h(x) - y|$ for n training examples - $\frac{1}{n} \sum_{i=1}^n |h(x) - y|$

Sum of Square :- $\frac{1}{n} \sum_{i=1}^n (h(x) - y)^2$
→ linear regression.

Classification → number of Misclassifications + $\frac{1}{n} \sum_{i=1}^n S(h(x), y)$

		True class		Hyp. Class
		Pos	Neg	
Pos	TP	FP		
	FN	TN		
	P			N

$$\text{True Positive Rate} = \frac{\text{TP}}{\text{TP} + \text{TN}}$$

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

$$\text{Recall} = \frac{\text{TP}}{\text{P}}$$

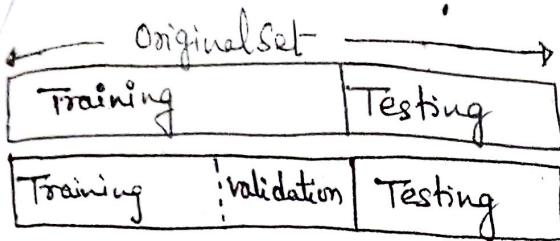
False Positive Rate

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

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

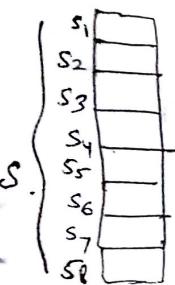
$$\text{Recall} = \frac{\text{TP}}{\text{P}}$$

Validation Set

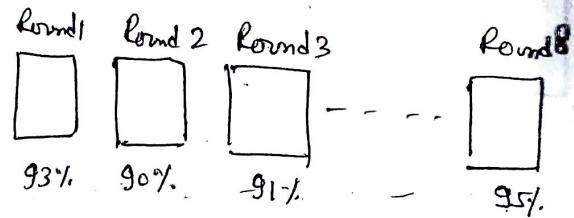


validation set is used during tuning training data to tune the model parameters.
check for the accuracy of your hypothesis on the test set. So, when you do validation if your splitting data into these parts validation fails to all the above data. to overcome this we use cross validation.

k-fold Cross Validation:-



Round i : Use s_i for testing
 $s - s_i$ for training
Validation Accuracy \rightarrow

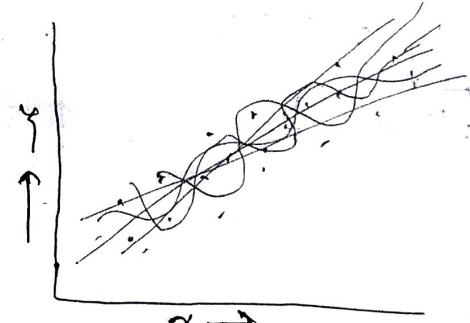


$$\text{Final Accuracy} = \text{Average}(\text{Round 1}, \text{Round 2}, \dots)$$

Linear Regression:-

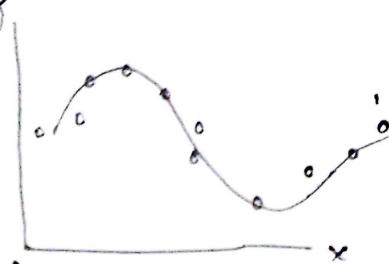
Instance x, y
Given x predicts y $x \rightarrow \text{Continuous}$.

Simple Regressions \rightarrow single feature x .

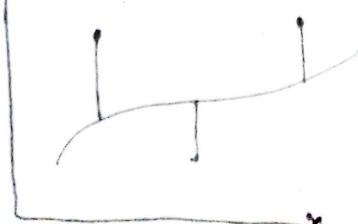


x & y both are Continuous

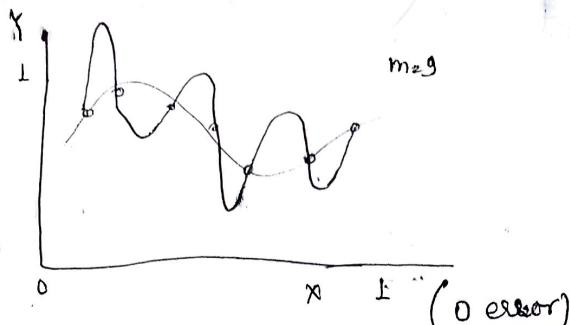
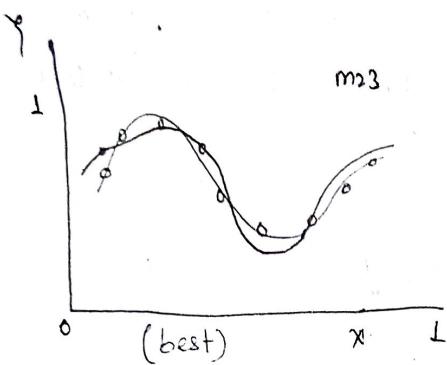
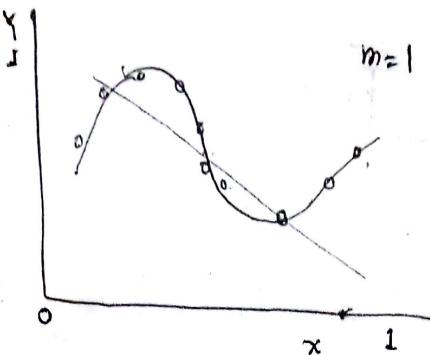
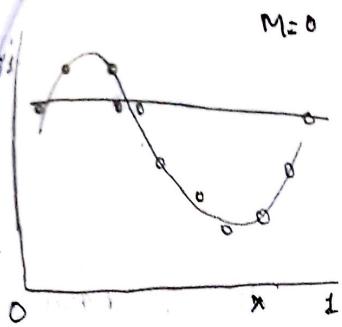
Example \rightarrow



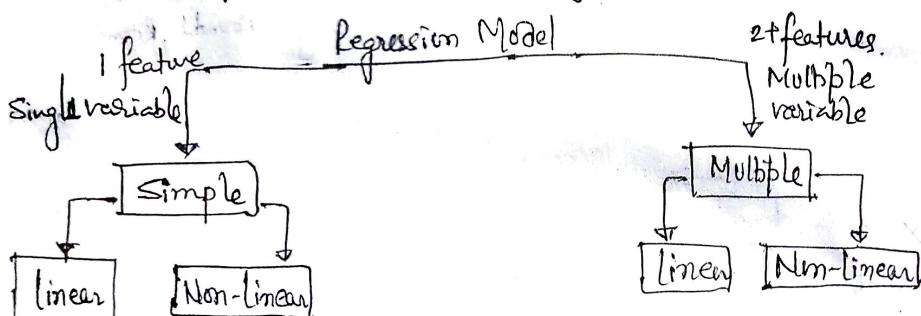
- The curve is true function (which is not polynomial)
- We may use a loss function that measures the square error in the prediction of $y(x)$ for x .



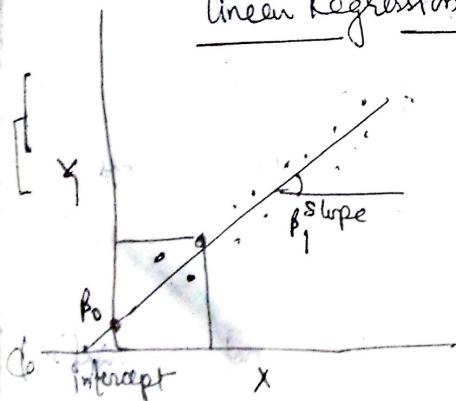
m (degree polynomial).



Sum of square error is decreasing.



Linear Regression:



$y = \beta_0 + \beta_1 x + \epsilon$

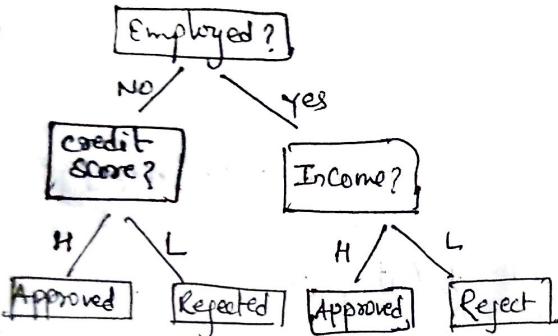
point where it meet to x -axis (y -intercept)
slope of line
error

Decision Tree \rightarrow It is a tree structured classifier, and types of nodes, decision nodes and leaf.

decision Node, they specify a choice or a test based on this you can decide which direction you can go.

leaf node indicates the classification of example or the value of example.

* Decision Trees can be used both for classification and regression. However, it is more popularly used for classification though they can be used for regression.



Prefer Smaller Trees
 — low depth
 — Small no. of nodes.

Top down Induction of Decision Trees ID3 \rightarrow

1. $A \leftarrow$ the 'best' decision attribute for next node.
2. Assign A as decision attribute for node.
3. for each value of A create new descendant
4. Sort training examples to leaf node according to the attribute value of ~~the~~
5. if all training examples are perfectly classified stop, else iterate over new leaf nodes.

Choices

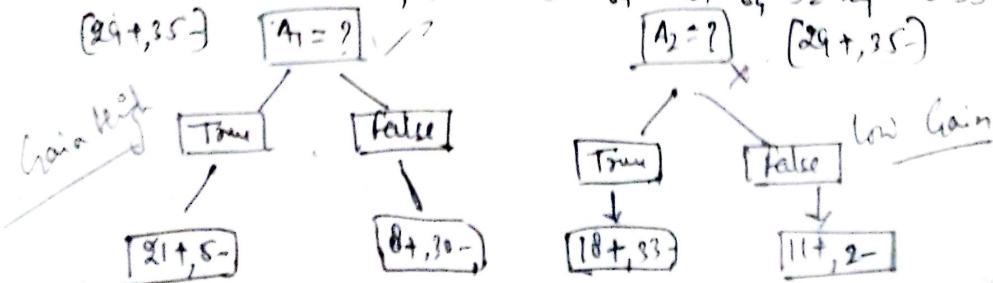
When to Stop

- No more input feature.
- all examples are classified the same
- too few samples to make an informative split.

which test to split on.

- split gives smaller error.
- with multi-valued features
 - ↳ split on all values or
 - ↳ split values into leaf

$$\text{Entropy } [29+35] = -\frac{29}{64} \log_2 \frac{29}{64} - \frac{35}{64} \log_2 \frac{35}{64} = 0.99.$$

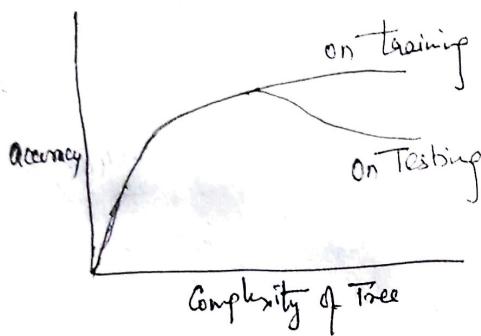


Gini Index

$$\text{Gini(Node)} = 1 - \sum_{\text{Node}} \left[P(C) \right]^2$$

$$\text{Gini}_{\text{split}}(A) = \sum_{\text{Value}(A)} \frac{|S_p|}{|S|} \text{Gini}(N_v)$$

Overfitting:- A Hypothesis 'h' is said to overfit the training data if there is another hypothesis 'h'' such that h' has more error than h on training data but h' has less error than h on test data.



Avoid overfitting \rightarrow

\rightarrow Pre-pruning - Stop growing when data split not statistically significant.

\rightarrow Post-pruning - Grow full tree then remove nodes.

K-nearest Neighbour:-

Training Method:- Save the training example.

At Prediction Time:- find the K training examples $(x_1, y_1), \dots, (x_k, y_k)$ that are closest to the test example x.

- \rightarrow Classification \rightarrow predict the most frequent class among those y_i 's.
- \rightarrow Regression \rightarrow predict the average of among the y_i 's.

Algorithm \rightarrow

Input:

D - Training data

k - Number of Neighbors

t - Input tuple to classify

Output:

c \rightarrow class to which t is assigned.

KNN Algo:-

$N=0$ 'N \rightarrow set of neighbors for'

for each $d \in D$ do

if $|N| < k$ then $N=N \cup d$;

else if $\exists u \in N$ such that $\text{sim}(t, u) > \text{sim}(t, N)$ then begin

$N=N \cup u$;

$N=N \setminus u$;

end.

Distance formulae :-

$$x_i = (x_{i1}, x_{i2}, \dots, x_{in})$$

$$x_j = (x_{j1}, x_{j2}, \dots, x_{jn})$$

dist - Euclidian (x_i, x_j)

$$= \sqrt{\sum_m^n (x_{im} - x_{jm})^2}$$

Manhattan

$$\sum_{i=1}^n |x_i - y_i|$$

Representing Hypothesis ->

most general hypothesis — every attribute is positive example $\langle ?, ?, ?, ? \rangle$

most specific hypothesis — no attribute is positive example $\langle \phi, \phi, \phi, \phi \rangle$

General to Specific Ordering of Hypothesis —

$$h_1 = \langle \text{sunny}, ?, ?, ?, \text{strong}, ? \rangle$$

$$h_2 = \langle \text{sunny}, ?, ?, ?, ?, ? \rangle \Rightarrow h_2 \text{ is more general than } h_1.$$

→ first, for any instance x in X and hypothesis h in H we say that x satisfies h iff $h(x) = 1$.

→ Given Hypothesis h_j and h_K .

→ h_j is more general than or equal to h_K iff any instance that satisfies h_K also satisfies h_j .

Find-S : finding A Maximally Specific Hypothesis —

1. Initialize h to the most specific hypothesis in H .

2. for each positive training instance x .

for each attribute a_i in h .

if the constraint a_i is satisfied by x
then, do nothing.

else replace a_i in h by the next more general constraint
that is satisfied by x .

3. Output hypothesis h .

Disadvantage ->

• Sensitivity to noise

• It's hard to determine good values

• Heuristic k = sqrt (No. of Examples)

Bayesian Learning

Bayesian probability — Partial belief.

Bayesian Estimation — Calculate validity of proposition.

↳ Prior estimate.

↳ New relevant evidence

Bayes Theorem →

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)}$$

• $P(h|D)$ is called posterior probability

• $P(h)$ is called prior probability

• $P(D)$ — probability of data, D.

• $P(D|h)$ is Conditional probability

Q →

$$P(\text{Cancer}) = 0.008, P(-\text{cancer}) = 0.992$$

$$P(+|\text{Cancer}) = 0.98, P(-|\text{Cancer}) = 0.02$$

$$P(-|\neg\text{Cancer}) = 0.97, P(+|\neg\text{Cancer}) = 0.03$$

$$P(\text{Cancer}|+) = \frac{P(+|\text{Cancer})P(\text{Cancer})}{P(+)} \quad , \quad P(-\text{Cancer}|+) = \frac{P(+|\neg\text{Cancer})P(\neg\text{Cancer})}{P(+)}$$

$$= \frac{0.98 \times 0.008}{P(+)}$$

$$= \frac{0.03 \times 0.992}{P(+)}$$

Maximum A posteriori (MAP) Hypothesis —

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

$$h_{\text{map}} = \arg \max_{h \in H} P(h|D)$$

$$= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)}$$

$$= \arg \max_{h \in H} P(D|h)P(h)$$

Q →

$$\text{mach 1} = 30 \text{ wrenches/hr}$$

$$\text{mach 2} = 20 \text{ wrenches/hr}$$

Out of all produced parts — 1% are defective

out of all defective parts

5% came from mach 1.

50% came from machine 2.

What is prob that the part produced by
mach 2 is defective?

$$\begin{aligned} p(\text{mach 1}) &= 0.6 \\ p(\text{mach 2}) &= 0.4 \\ p(\text{defect}) &= 1\% \end{aligned}$$

$$\begin{aligned} p(\text{mach 1} | \text{Defect}) &= 50\% \\ p(\text{mach 2} | \text{Defect}) &= 50\% \end{aligned}$$

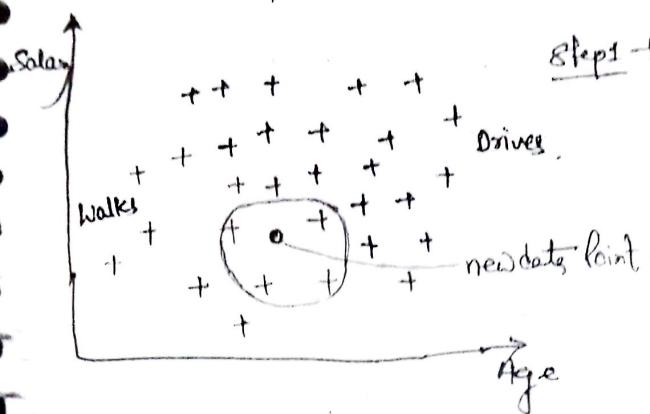
Now we have to find out.

$$P(\text{Defect} | \text{mach 2}) = ?$$

$$P(\text{Defect} | \text{mach 2}) = \frac{P(\text{mach 2} | \text{Defect}) P(\text{Defect})}{P(\text{mach 2})}$$

$$= \frac{0.5 * 0.01}{0.4} = 0.0125 = 1.25\%$$

Naive Bayes \rightarrow



$$P(\text{Walks}) = \frac{\text{Number of Walkers}}{\text{Total No. of Observations}}$$

$$P(\text{Walks}) = \frac{10}{30}$$

$$\begin{aligned} P(0) &= \frac{\text{No. of similar observations}}{\text{Total No. of observations}} \\ &= \frac{4}{30}. \end{aligned}$$

$$P(0 | \text{Walks}) = \frac{\text{No. of similar observations among those who walks}}{\text{Total. No. of workers}}$$

$$2 \frac{3}{10}$$

$$P(\text{Walks}|0) = \frac{10/30 * \frac{3}{10}}{\frac{4}{30}} = 0.75$$

$$\begin{aligned} \text{Step 1} \rightarrow P(\text{Walks}|0) &= \frac{P(0 | \text{Walks}) * P(\text{Walks})}{P(0)} \\ &\stackrel{\text{① Prior Probability}}{\uparrow} \quad \stackrel{\text{② Likelihood}}{\downarrow} \\ &\stackrel{\text{③ Marginal Likelihood}}{\leftarrow} \end{aligned}$$

$$\text{Step 2} \rightarrow P(\text{Drives}|0) = \frac{P(0 | \text{Drives}) * P(\text{Drives})}{P(0)}$$

$$\text{Step 3: } P(\text{Walks}|x) \text{ vs } P(\text{Drives}|x)$$

$$P(\text{Drive}) = \frac{20}{30}$$

$$P(0) = \frac{4}{30}$$

$$P(0 | \text{Drives}) = \frac{1}{20}$$

$$P(\text{Drives}|0) = \frac{\frac{1}{20} * \frac{20}{30}}{\frac{4}{30}} = 0.25$$

$$P(\text{Walks}|0) > P(\text{Drives}|0).$$

New data point should be walks to work.