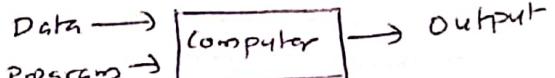
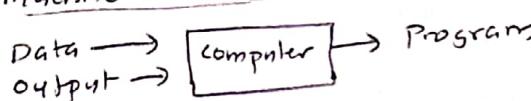


Module I

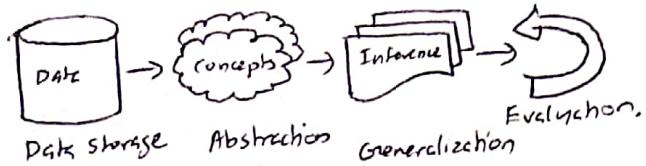
Introduction to machine Learning:

- Machine Learning: field of study interested in the development of computer algorithms to transform data → intelligent action
- It is about ~~teach~~ predicting the future based on the past.
- Algorithms enables computers to 'learn' from data
- Related to many areas such as
 - data mining
 - statistics etc
- Traditional programming,

Machine Learning

- Applications,
 - identification of unwanted spam messages in email
 - weather forecast & long-term climate change.
 - development of self-driving cars & auto-piloting drones
 - preventing fraudulent credit card transactions
 - Prediction of election outcomes
 - optimization of energy use in homes and offices
- Limitations,
 - recognize speech and handwriting accurately
 - requires massive data to train on.
 - High error-susceptibility.

How machines Learn:

- A formal definition proposed by computer scientist Tom M. Mitchell states "A machine learns whenever it is able to utilize an experience, such that its performance improves on similar experiences in the future."
- Human brains are naturally capable of learning from birth (computers \Rightarrow explicitly done)

→ The basic learning process is divided into 4 components,



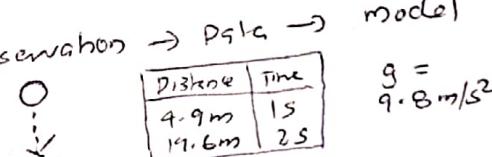
(i) Data Storage

[utilizes observation, memory and recall to provide a factual basis for further reasoning]

- all learning begins with data
- In human beings, brain uses electrochemical signals of biological cells to store & process observations for future call
- In computers, recall is done using HDD, flash memory and RAM in combination with a CPU

(ii) Abstraction

[involves translation of stored data into broader representations & concepts]

- process which assigns meaning to stored data (raw data comes to have a more abstract meaning)
- computer summarizes stored data using a model (explicit description of patterns within data)
- Training: process of fitting a model to a dataset.
- Different types of models,
 - 1) mathematical equations
 - 2) Relational diagrams (trees/graphs)
 - 3) Logical IF/else rules
 - 4) Clusters (grouping of data)
- Eg:
 Observation \rightarrow Pts \rightarrow model


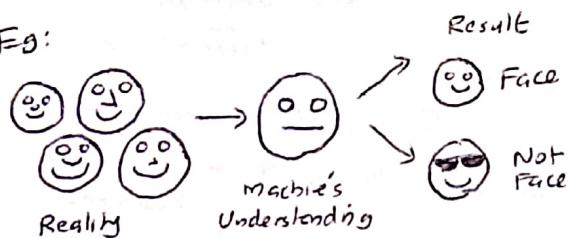
$$g = 9.8 \text{ m/s}^2$$

(iii) Generalization

[uses abstracted data to create knowledge & inferences — that drive action in new contexts]

- turning abstract knowledge into a form that can be utilized for future actions.
- It is a useless process (difficult to describe)
- involves reduction of hypothetical set in manageable number of important findings

- patterns are limited to only those that will be most relevant to future tasks
- it is not feasible to reduce no. of patterns by examining them one-by-one & ranking them
- so, algorithms employ shortcuts that reduce search space more quickly
- Towards end, algorithms will employ heuristics - (technique to solve problems more quickly)
- Eg:

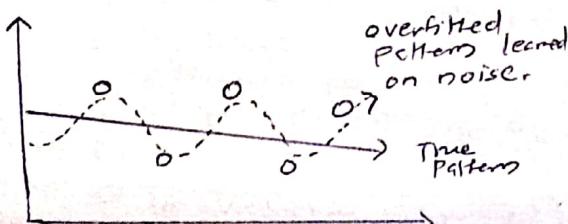


(iv) Evaluation

- [provider a feedback mechanism to measure utility of learned knowledge & inform potential improvements]
- evaluate / measure the learners success and use this info to inform additional trainings
 - models fails to perfectly generalise due to problem of noise
(noise \Rightarrow unexplained variations of data)
 - causes of noise
 - \rightarrow error due to sensors
 - \rightarrow issue with human subjects (Eg: Survey respondents)
 - \rightarrow data quality problems (missings, corrupted values etc)
 - \rightarrow complex phenomena (unsystematic impact to data)
 - Overfitting, used to reduce model noise

Solutions to the problems of overfitting \Rightarrow are specific to particular machine learning approaches

- Eg:



Machine Learning in practice (5 steps)

- 1) Data collection
- 2) Data exploration & preparation (checking quality of data)
- 3) Model training (selection of an appropriate algorithm to form a model)
- 4) Model Evaluation (evaluating the accuracy of a model using test data)
- 5) Model Improvement (using advanced strategies & switching model if required)

Types of Input Data

- \rightarrow Unit of analysis \Rightarrow smallest unit from which inference is made
 - \rightarrow Datasets (stores unit of observations and their properties)
 - collection of data includes,
 - (i) Examples [instances of unit of observation]
 - (ii) Features [recorded properties of examples]
 - \rightarrow Eg: To build an algorithm to identify spams e-mail;
 - Unit of observation \Rightarrow email messages
 - Examples \Rightarrow specific messages
 - Features \Rightarrow words used in emails
 - \rightarrow Examples and features are commonly gathered in Matrix Format Automobile
- | year | model | price | miles | color |
|------|-------|-------|-------|-------|
| | | | | |
| | | | | |
| | | | | |
- } examples (row)
Features (columns)

- \rightarrow Various forms of feature,
 - 1) Numerical (characteristic measured in numbers)
 - Eg: no. of black pixels
 - 2) Categorical (attribute consists of a set of categories)
 - Eg: Gender.
 - 3) Ordinal (follows an ordered list)
 - special case of categorical
 - Eg: cloth size (S, M, L, XL)
- Eg 2 - measurement of customer satisfaction
[not at all happy \rightarrow very happy]

Types of Machine Learning Algorithms

- Machine Learning algorithms are divided into categories according to their purpose
- They are,
 - Predictive model**
 - involve the prediction of one value using other values in the dataset.
 - Attempts to discover & model the relationship b/w target feature (feature being predicted) and other features
 - Supervised Learning algo. used
 - Attempts to optimize a model - to find combo of feature values that result in target output
 - We have input variables (X), and an output variable (Y), and we use an algorithm to learn mapping function

$$Y = f(X)$$

- Classification : predicting the category an example belongs to
- Using classifier, we can predict, (whether an email is spam, whether a person has cancer etc)

2) Descriptive model

- no single feature is more important than the others.
- Here, there is no target to learn.
- Descriptive modeling tasks include,
 - Pattern Discovery** [identify useful associations within data]
 - Eg: Market-Basket Analysis
 - clustering** [dividing datasets to homogeneous groups]
 - Eg: ads can be tailored for particular audience
- Unsupervised Learning algo. used
 - Attempts to model the underlying structure in data, in order to learn more about data]
 - We have input data (X), and no corresponding output variable (Y)

Eg: Clustering (K-means)
Association (Apriori algorithm)

Note ⇒ supervised learning algo. ⇒

- Decision Tree - classification
- Linear Regression - Numeric prediction

Exploratory Data Analysis

Numeric variables

1) Measuring Central Tendency - Mean, Median and Mode

- A class of statistics used to identify a value that falls in the middle of a set of data

[Mean]:

- also known as average
- Def: $\frac{\text{sum of all values}}{\text{number of values}}$

- Eg: mean income of 3 people having salaries 10,000, 5000, 15000

$$\Rightarrow \frac{10000 + 5000 + 15000}{3} \\ \Rightarrow \underline{\underline{10000}}$$

- In R, we use mean() function
- not always appropriate.

[Median]:

- Def: value that occurs halfway through an ordered list of values

- Eg: median of the salaries in ascending order (5000, 10,000, 15,000)

$$\Rightarrow \underline{\underline{10000}}$$

- To find median, first arrange the data in ascending/descending order.

- If no. of observations (n) is odd, median value is at position $\left(\frac{n+1}{2}\right)^{\text{th}}$ term

- If (n) is even, median value $\Rightarrow \left(\frac{n}{2}, \frac{n+1}{2}\right)^{\text{th}}$

- In R, we use median() function

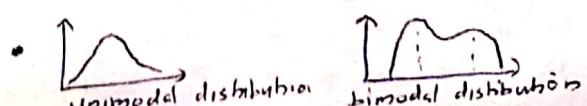
[Mode]:

- Def: number in a set that appears most frequently.

- A set can have more than one mode.

→ single mode \Rightarrow Unimodal
→ Two modes \Rightarrow Bimodal

- Eg: 12, 8, 7, 15, 7
mode = 7



2) measuring spread - quartiles and 5-number summary

- used to measure diversity
- denotes how tightly / loosely the values are spaced.
- specifies whether most values are like or unlike mean & median.

5 number summary

- A set of five statistics that roughly depict the spread of a feature's values
- All 5 are included in the output of `summary()` fn.

- 1) Minimum (Min.)
- 2) First quartile, or Q_1 (1st Qu.)
- 3) Median, or Q_2 (median)
- 4) Third quartile, or Q_3 (3rd Qu.)
- 5) Maximum (max)

- R provides `min()` & `max()` fn to calculate minimum and maximum
⇒ `range()` ⇒ returns both min & max.
⇒ `quantile()` ⇒ returns 5-number summary.

0%	25%	50%	75%	100%
x	y	z	a	b

 ⇒ IQR() ⇒ returns Inter-Quartile Range ($Q_3 - Q_1$)

Quartiles

- ⇒ special case of quantiles (a type of statistics)
- ⇒ quantiles: numbers that divide data into equally sized quantiles
- ⇒ commonly used quantiles,
 - a) Tertiles (3-parts)
 - b) Quintiles (5-parts)
 - c) Deciles (10-parts)
 - d) Percentiles (100-parts)
- ⇒ quartiles divide dataset into 4 portions

⇒ Q_1 = middle number b/w (min) and median.

Q_2 = median

Q_3 = middle value b/w median and (max).

Eg 5, 7, 4, 4, 6, 2, 8

Ascending order → 2, 4, 4, 5, 6, 7, 8

Cut list into quarters

2, 4, 4, 5, 6, 7, 8

↑ ↑ ↑

Q_1 Q_2 Q_3

Eg 2

3, 5, 7, 8, 12, 13, 14, 18, 21

lower half

median

upper half

$$Q_1 = \frac{5+7}{2} = \frac{12}{2} = 6 //$$

$$Q_3 = \frac{14+18}{2} = \frac{32}{2} = 16 //$$

$$Q_2 = 12 //$$

Variance and std. Deviation

- Variance [average of squared differences between each value and the mean value]

$$\text{Var}(X) = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

- Standard deviation [square root of variance denoted by sigma]

$$\text{std Dev}(X) = \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2}$$

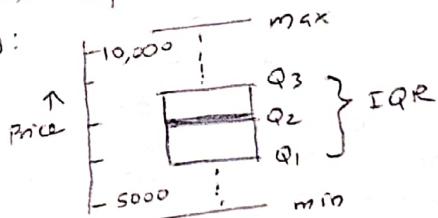
- var() and sd() functions in R can be used.

Visualising Numeric variables

1) Boxplot:

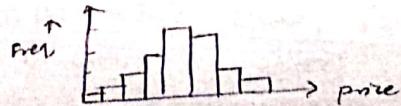
- `boxplot()` used in R
- depicts 5 ~~is~~ number summary using horizontal lines & dots.
- horizontal lines form the box represents Q_1 , Q_2 , & Q_3

• Eg:



2) Histograms

- `hist()` used in R
- divides values into predefined number of portions / bins
- composed of series of bars.
 - Height ⇒ freq. of values
 - vertical ⇒ start & end points of range of values.



Understanding numeric data

Uniform Distributions

- all values are equally likely to occur
- Eg: fair 6-sided die rolled on a floor.
- easy to detect with a histogram. (bars have same height)
- Histogram



Normal Distributions

- some values are more likely to occur than others.
- Eg: Distribution of height/weight in general populations.
- values are less likely to occur on both sides of central area.



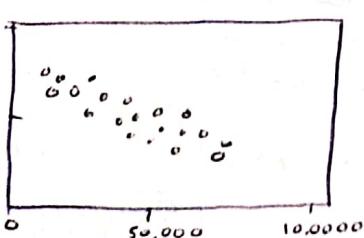
- Skewness \Rightarrow
 - a) Right skew
 - b) No skew
 - c) Left skew

Visualizations Relationships

Scatter plot

- visualizes a bivariate relationship (relationship b/w 2 variables)
- dots are drawn on a coordinate plane (values of one feature \Rightarrow x, values of other feature = y)
- plot() fn. used in R.
- Negative association \Rightarrow pattern of dots sloping downwards.
- Positive association \Rightarrow pattern of dots sloping upward.
- plot(x = ..., y = ..., main = "Headings", xlab = "...", ylab = "...")

Eg:



Categorical / Nominal Variables

- allows to assign categories
- Eg: used car dataset, 3 categorical variables [model, color, transmission]
- examined using table (one way table)
- table() in R.

Dimensionality Reduction

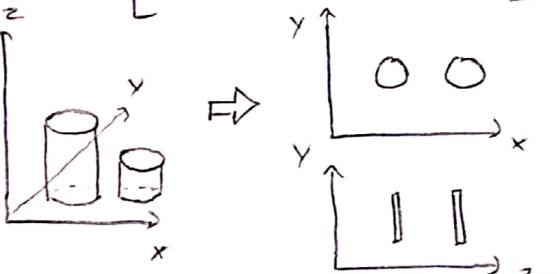
\rightarrow process of reducing the number of random variables, by obtaining a set of principal variables.

\rightarrow Divided into,

- Feature selection
 - find the subset of the original set of variables
 - involves 3 ways,
 - Filter
 - Wrapper
 - Embedded

- Feature extraction
 - reduces data in a high dimension space \rightarrow low dimension space

Eg:



Advantages

- helps in data compression (ie, reduced storage space)
- reduces computation time
- removes redundant features.

Disadvantages

- some data loss
- linear correlation is undesirable
- mean & covariance not enough to define datasets.

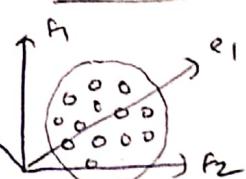
Method (Principal Component Analysis)

\rightarrow introduced by Karl Pearson

\rightarrow here variance of data in lower dimensional space \Rightarrow maximum

\rightarrow steps,

- 1) construct the covariance matrix
- 2) compute the eigenvalues of the matrix.
- 3) Largest eigenvalues are used to reconstruct a large fraction of variance of original data.



Module II

Lazy Learning

- Lazy Learning methods simply store the data.
- Generalizing beyond these data is postponed until an explicit request is made.
- Eg: Nearest Neighbor classifiers defined by their characteristic of classifying unlabelled examples by assigning them the class of similar labelled examples plus
- Nearest Neighbor methods are simple & extremely powerful
- Applications,
 - Predictions whether a person will enjoy a music recommendation
 - Facial recognition
 - Identifying patterns in genetic data.

The K-NN Algorithms

- Stands for K - Nearest Neighbors algorithms (K-NN)
- Simplest & widely used machine learning algorithms.
- Strengths & Weaknesses,

Strengths	Weaknesses
• simple & effective	• doesn't produce a model.
• makes no assumptions regarding data distribution	• selection of an appropriate 'K'
• Fast training phase	• slow classification phase.
	• Nominal features & missing data need additional processing

- K-NN algorithm uses information about an example's 'K' nearest neighbours to classify unlabelled examples
- K → any number of nearest number neighbors can be used.
- After choosing K, algorithm requires a training dataset, that have been classified into several categories (labeled by a nominal variable)

- Then for each unlabeled record, K-NN identifies K records in the training data that are the "nearest" in similarity.
- The unlabeled test instance is assigned the class of majority of K neighbors
- 2 features of each ingredient is only taken (simplifying on 2D graph)
- K-NN algorithm treats the features as coordinates in a multidimensional feature space.

Example of K-NN Algorithm

Ques: Is tomato a fruit or vegetable?



Measuring similarity with distance

- Locating tomato's nearest neighbors requires a distance function / formula
- Many ways to calculate distance
- K-NN uses Euclidean distance. (distance measured by a ruler to connect two points)
- Euclidean distance, P & Q ⇒ examples to be compared, having n features
 $P_1 \Rightarrow$ value of first feature of P,
 $q_1 \Rightarrow$ value of first feature of q

$$\text{dist}(P, Q) = \sqrt{(P_1 - q_1)^2 + (P_2 - q_2)^2 + \dots + (P_n - q_n)^2}$$

- We have the table of ingredient and features (sweetness & crunchiness)

Ingredient	Sweetness	Crunchiness
grape	8	5
green bean	3	7
nuts	3	6
orange	7	3

To calculate distance to tomato (6,4) from,

- a) grape $\Rightarrow \sqrt{(6-8)^2 + (4-5)^2} = 2\cdot2$
- b) green bean $\Rightarrow \sqrt{(6-3)^2 + (4-7)^2} = 4\cdot2$
- c) nuts $\Rightarrow \sqrt{(6-3)^2 + (4-6)^2} = 3\cdot6$
- d) orange $\Rightarrow \sqrt{(6-7)^2 + (4-3)^2} = 1\cdot4$

1 - NN classification

- To classify the tomato as a vegetable, protein or fruit, we'll begin by assigning the tomato \rightarrow to a single nearest neighbor
- Here orange is the nearest neighbor (distance = 1.4)
- As orange = a fruit, 1-NN would classify tomato also as a fruit

3-NN classification

- $K=3$ and performs vote against 3 nearest neighbours (orange, grape and nuts)
- majority class \Rightarrow Fruit ($\frac{2}{3}$ out of 3 votes)
- so tomato is again classified as a fruit

Choosing an appropriate K

- The decision of how many neighbors to use for K-NN \Rightarrow determines how well model will generalize future data.
- Bias variance Tradeoff: balance b/w overfitting & underfitting of training data
- Overfitting \Rightarrow occurs when algo. captures noise of data
- Underfitting \Rightarrow occurs when algo. neither model training data nor generalize to new data

- Very large K value \Rightarrow model always predict majority class
 - single nearest neighbor \Rightarrow noisy data / outliers influence classification
-
- Best value of K \Rightarrow somewhere b/w two extremes

Common practice, $K = \sqrt{\text{No of training examples}}$

Eg: For 15 example ingredients
 $K=4$, [since $\sqrt{15} = 3.87$]

Preparing data for use with K-NN

- one solution \Rightarrow rescale features by shrinking / expanding their range such that one contributes relatively equally to distance formula.

Eg: if sweetness & crunchiness are both measured 1-10 scale, we would like to measure spiciness on (1-10) scale.

several methods for scaling include,

1) Min-Max Normalization

- traditional method for rescaling
- transforms a feature such that all values fall b/w (0-1)
- Formula,

$$x_{\text{new}} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

2) Z-score standardization

- $x_{\text{new}} = \frac{x - \mu}{\sigma}$
- $\mu = \text{mean}(x)$
- $\sigma = \text{stdDev}(x)$
- resulting value is called z-score (falls on an unbound range of negative & positive numbers)
- no predefined maximum/minimum

3) Dummy Coding

- Euclidean distance not defined for nominal data.
- To calculate distance b/w nominal features \Rightarrow convert them to numeric form
- Eg: Dummy coding for gender,

$$\text{male} = \begin{cases} 1 & \text{if } x = \text{male} \\ 0 & \text{otherwise} \end{cases}$$

- Dummy coding of two category \Rightarrow results in a single new feature
- No need to construct a separate feature for female. (since two are exclusive & mutually)
- for n category nominal feature \Rightarrow dummy code = $(n-1)$ levels of feature.

For temperature,

$$\text{hot} = \begin{cases} 1 & \text{if } x = \text{hot} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{medium} = \begin{cases} 1 & \text{if } x = \text{medium} \\ 0 & \text{otherwise} \end{cases}$$

cold if both hot & medium = 0

- distance b/w dummy feature \Rightarrow $[0-1]$
- (no additional transformation required) (values fall on same scale)

Why k-NN algo. is lazy

- here no abstraction occurs.
- abstractions and generalizations are skipped altogether \Rightarrow undermines def. of learning
- process of making predictions tends to be relatively slow
- Instance based / Rule learning.

Probabilistic Learning

- Eg: A meteorologist says there is "70%" chance of rain in a weather forecast report ($7/10$)
- Use of past ~~data~~ events' data to predict future events

Probabilistic classifier

- A classifier that is able to predict, given a set of classes in probability distribution, a sample input

Naive Bayes

- uses probabilities
- mathematician Thomas Bayes - developed principles to describe probability of events.
- These principles formed the foundation for Bayesian methods
- Probability (b/w 0 and 1, i.e., 0% to 100%)
lower probability \Rightarrow event is less likely to occur
 $P(0) \Rightarrow$ definitely not occur
 $P(1) \Rightarrow$ event will occur with 100% certainty.
- Classifiers based on Bayesian methods \Rightarrow calculate observed probability from training data
- Later classifier applied to unlabelled data \Rightarrow to predict most likely class

Bayesian classifiers

- estimated likelihood / or a potential outcome should be based on evidence across multiple trials
- Eg: Trial \Rightarrow coin flip Event \Rightarrow Result = Head

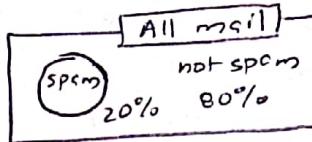
Understanding Probability

- Probability = $\left(\frac{\text{No. of trials in which event occurred}}{\text{Total no. of events}} \right)$
- Eg: If it rained 3 out of 10 days with similar conditions as today, Probability of rain = $\underline{\underline{3/10}} = \underline{\underline{30\%}}$
- $P(A)$: Probability of event (A)
- $P(\text{Rain}) = \underline{\underline{0.30}}$
- Sum of all possible outcomes is equal to 1

complement of $A = A^c/A'$

$P(\neg A) \Rightarrow$ probability of A not occurring

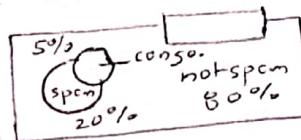
$$P(\text{spam}) = 0.20 \\ P(\neg \text{spam}) = 0.80$$



\uparrow event and their compliments

Joint Probability

- If certain events occur w.r.t. event of interest, we will be able to make predictions.
- consider a second event \Rightarrow [an email - message contains word "congratulations"]

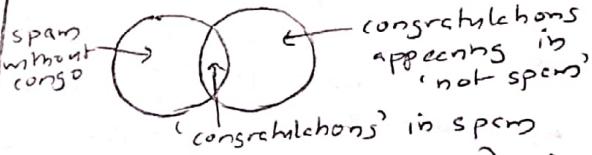


Not all spam messages contain word 'congratulations' & not every email with word 'congratulations' is spam.

- Probability of event X occurring at the same time event Y occurs

$$P(X \cap Y) \text{ or } P(X, Y)$$

Venn Diagram (By John Venn)



$P(\text{spam} \cap \text{congratulations}) \Rightarrow$ now probability of one event is related to probability of other.

- If independent events,

$$P(A \cap B) = P(A) \times P(B)$$

$$P(A) = 0.2 ; P(B) = 0.05 \\ \therefore P(A \cap B) = 0.05 \times 20 = \underline{\underline{0.01}}$$

- If dependent events,

- using predictive modelling
- $P(\text{spam})$ and $P(\text{congratulations})$ are likely to be highly dependent \Rightarrow calculations above may be incorrect

- more careful formulation of relationship b/w two events required

\Rightarrow use of Advanced Bayesian methods

Conditional Probability - with Bayes' Theorem

→ The relationships b/w dependent events can be described using Baye's theorem.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

→ Conditional probability of event A is the probability that event will occur, given B is already occurred.

[$P(A)$ is dependent on what happened after event B]

→ Baye's theorem tells that estimate of $P(A|B)$ should be based on, $P(A \cap B)$ & $P(B)$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A \cap B) = P(A|B) \cdot P(B)$$

$$\text{Also } P(A \cap B) = P(B \cap A)$$

Substitutions

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

→ Prior probability \Rightarrow probability of an even computed before collection of new data

$$\text{Ex: } P(\text{spam}) = 20\%$$

Likelihood $\Rightarrow P(\text{congo}|\text{spam})$

Marginal Likelihood $\Rightarrow P(\text{congratulations})$

Posterior probability $\Rightarrow P(\text{spam}|\text{congo})$

$$P(\text{spam}|\text{congo}) = \frac{P(\text{congo}|\text{spam})P(\text{spam})}{P(\text{congo})}$$

$$\text{If } P(\text{congo}|\text{spam}) = 4/20, \\ \text{Then } P(\text{spam} \cap \text{congo}) = (4/20) \cdot (20/100) \\ = 0.04$$

To compute posterior probability,

$$P(\text{spam}|\text{congo}) \\ \Rightarrow (4/20) \cdot \frac{(20/100)}{(5/100)} = 0.80$$

\Rightarrow Probability is 80% that a message is spam, given that it contains word congratulations

The Naive Bayes algorithm

- simple method to apply Bayes theorem to classification problems
- Strengths & weaknesses,

Strengths	Weakness
<ul style="list-style-type: none"> • simple, fast & effective • works well with noisy & missing data • Requires only few training examples • easy to obtain estimated probability. 	<ul style="list-style-type: none"> • sometimes rely on faulty assumptions of features • not ideal for datasets with many numeric features • Estimated probabilities are less reliable than predicted classes

→ it makes some "naive" assumptions about data. (assumes all features are equally important & independent)

→ assumptions rarely true in most real world applications.

→ To identify spam by monitoring email messages

- sender maybe an important indicator of spam (many messages)
- words are not independent etc.
- other words like 'earn' & 'money' (subscribers) will be here.

Classification with naive Bayes

• Suppose a message contains — congo & money, but no earn/subscribe

• Using bayes', difficult to solve

• Easier in Naive Bayes' since it assumes independence among events

[cross conditional \Rightarrow conditioned on independence \Rightarrow be same class value]

$$\Rightarrow P(A \cap B) = P(A) \cdot P(B)$$

$$\Rightarrow P(\text{spam} | w_1 \cap w_2 \cap w_3 \cap w_4) \\ = P(w_1 | \text{spam}) \cdot P(\neg w_2 | \text{spam}) \\ \cdot P(\neg w_3 | \text{spam}) \cdot P(\neg w_4 | \text{spam}) \\ P(\text{spam})$$

[5 terms]

Formula

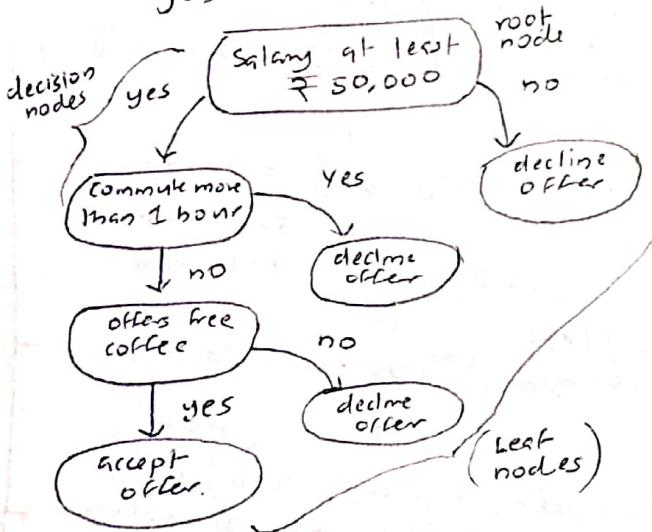
$$P(C_L | F_1 \dots F_n) = \frac{1}{Z} \prod_{i=1}^n P(F_i | C_L)$$

- Laplace estimator : if a message with all 4 terms occur, $P(\text{spam})$ becomes 0. \rightarrow so add a small no. to each counts in frequency table to ensure non-zero probability (typically set to 1)
- Note : To use Naive Bayes, features must be categorical, not numeric. Solution \rightarrow bins

Module III

Decision Trees

- machine learning model that can make complex decisions from sets of simple choices.
- useful for business strategy.
- Decision Tree Learners are powerful classifiers, which utilizes a tree structure to model relationships among features and potential outcomes.
- It uses a structure of branching decisions which guides to a final predicted class value.
- Types of Nodes;
 - (i) Root Node [tree starts from here]
 - (ii) Decision Node [choices to be made based on the attributes]
 - (iii) Branches [potential outcomes of a decision (Yes/No)]
 - (iv) Leaf Nodes [final decisions (also called terminal nodes)]
- Eg: Should I accept a new job offer?



- more than 2 outcomes is also possible
- leaf nodes \Rightarrow actions to be taken as the result of series of decisions.
- Benefits \Rightarrow Model is created in a human readable format.
- provides insight into - why the model works or doesn't work well for a particular task (for future business practices)

Uses of decision trees

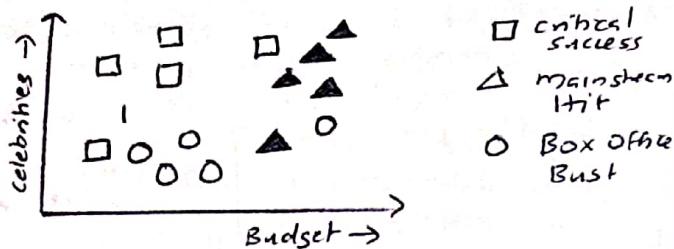
- Applicant rejection [credit scoring models should mention criteria in a doc free of bias]
- Customer Behaviour [whether satisfied/not - send to marketing agencies]
- Medical Diagnosis [based on symptoms, laboratory measurements etc]

Divide and Conquer

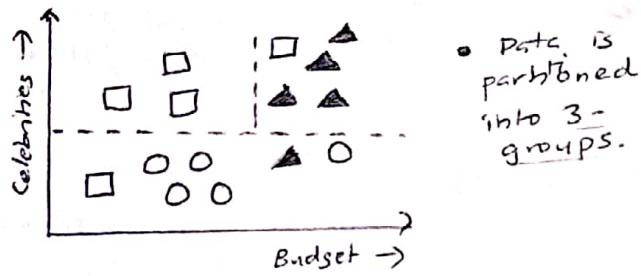
- also called recursive partitioning
- Decision trees are built using this heuristic
- splits data into subsets \Rightarrow and are then split repeatedly into even smaller subsets and so on.
- Algorithms stopped when subsets are sufficiently homogeneous / stopping criterion has met.
- In a decision tree,
 - at root node \Rightarrow entire dataset, no splitting done.
 - now choose a feature to split (most predictive feature is used)
 - examples are partitioned into groups according to distinct values
- Stopping Criteria;
 - (i) All the examples at the node have same class
 - (ii) No remaining feature to distinguish among the examples
 - (iii) Tree has grown to a predefined size limit.

- Eg: Decision Tree algorithms to predict whether a movie would fall into one of the three categories
 - critical success
 - mainstream hit
 - Box office Bust (flop)
- company's last 15 most recent releases are analyzed (success/failure)
- Relation (budget & no. of major celebrities)

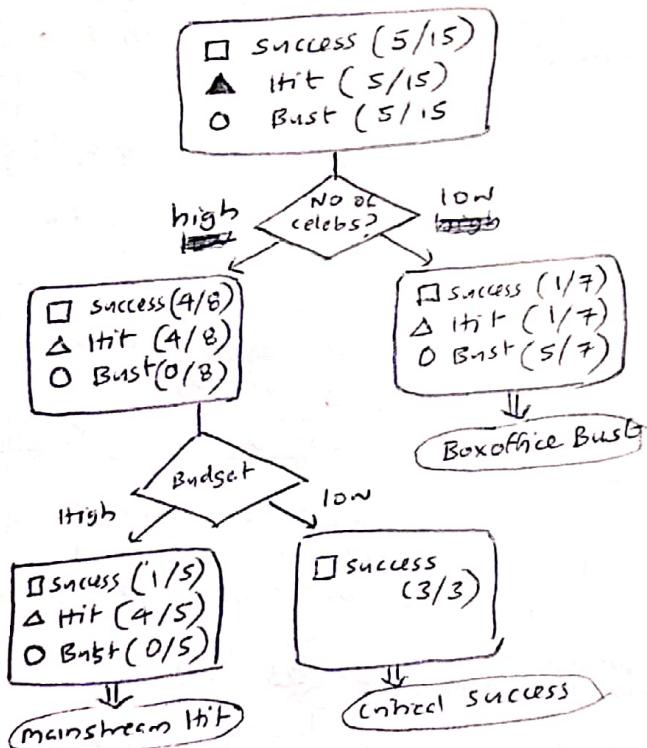
→ Scatter Plot



- Firstly split the dataset using the feature (no. of celebrities) \Rightarrow create root node.
- second split based on the feature (budget of movie - high/low)



- Decision Tree: predicting the future success of movies



The C5.0 decision tree algorithm

- popular algorithms to build decision tree models automatically.
- well known \Rightarrow C5.0 Algorithm (developed by Ross Quinlan)
 - improved version of prior algo. C4.5, which itself is an improvement over ID3 algorithm
 - easier to understand & deploy

Strengths

- all purpose classifier (works well on most problems)
- highly automatic learning process (handles missing, numerical / nominal)
- Excludes unimportant features
- used on small / large datasets
- results can be interpreted without mathematical background
- more efficient

Weakness

- biased towards splits (having a large no. of levels)
- easy to overfit / underfit model
- Trouble modeling relationships (reliance on axis-parallel splits)
- small changes in training data \Rightarrow large changes in decision logic
- Large trees can be difficult to interpret

Choosing the best split:

(to identify which feature to split upon)

\rightarrow Purity \Rightarrow degree to which a subset contains only a single class

\rightarrow Pure subset \Rightarrow subset composed of only a single class.

\rightarrow Entropy \Rightarrow C5.0 uses entropy as a measurement of purity.

- quantifies randomness / disorder in a set of values.
- sets with high \Rightarrow very diverse

- decision trees hope to find splits that reduce entropy.

measured in bits (only 2 classes possible)

$$\text{Range (Entropy)} = [0 - \log_2(n)]$$

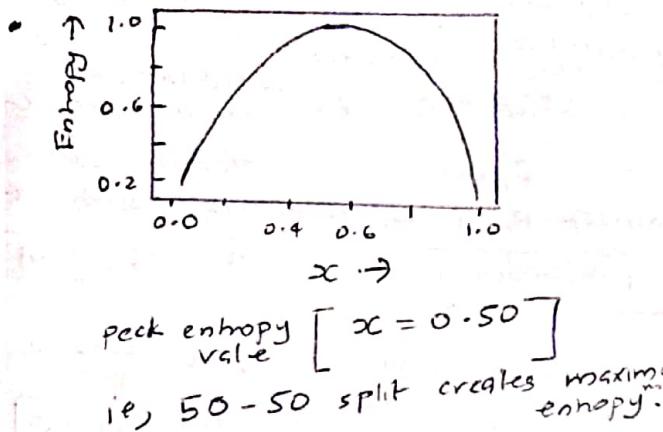
- Entropy, $\frac{\text{min value}}{\text{max value}} \Rightarrow$ completely homogeneous
- $\frac{\text{max value}}{\text{min value}} \Rightarrow$ diverse as possible

$$\text{Entropy}(S) = \sum_{i=1}^c - p_i \log_2(p_i)$$

S - segment of data
c - no. of class levels
 p_i - proportion of values falling into class level i.

- Ex: A partition of data into two classes
 - red (60%) \rightarrow white (40%)
 - \rightarrow Entropy = $-0.60 \times \log_2(0.60) - 0.40 \times \log_2(0.40)$
 - $= 0.9709$

- If proportion of examples in one class is \underline{x} , then proportion (other class) $\Rightarrow \underline{(1-x)}$



- If one class increasingly dominates the other, entropy reduces to zero

Information Gain

(For a feature F is calculated as the difference b/w the entropy in the segment before split (S_1) and partitions resulting from split (S_2))

$$\text{InfoGain}(F) = \text{Entropy}(S_1) - \text{Entropy}(S_2)$$

- After a split, data is divided into more than one partition
- So $E(S_2)$ must consider Total Entropy (weighting each partitions entropy)
- ie, Total entropy,

$$\text{Entropy}(S) = \sum_{i=1}^n w_i \text{Entropy}(P_i)$$

- ie, total entropy resulting from a split is sum of the entropy of each of n partitions weighted by proportions

• Higher Information Gain [a feature is better at creating homogeneous groups]

• Information Gain = 0 [no reduction in entropy for splitting on a particular feature]

Pruning the decision tree

→ involves reducing the size of decision tree

→ 2 types

- a) pre-pruning (early stopping)
 - stop growing tree further, before it classifies training set
 - condition
 - i) reaching a certain no. of decisions
 - ii) decision node contains only small no. of examples.

Post-pruning

- allows tree to perfectly classify training set and then post-prune the tree.
- pruning leaf nodes to reduce the size of the tree to a more appropriate level
- more effective approach (difficult to determine the optimal depth without growing) it first

Regression methods

Regression

→ concerned with specifying the relationship b/w a single dependent variable and one/more numeric independent variables.

→ simplest \Rightarrow follows straight line form

→ model data using a slope-intercept format.

$$y = a + bx$$

y - dependent
x - independent
b - slope
a - intercept

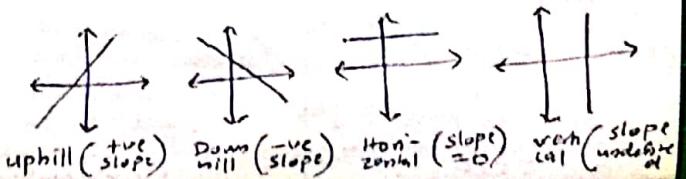
• slope (how much line rises for increase in x)

• intercept (value of y, when x=0)

• positive values \Rightarrow slope upward

• negative values \Rightarrow slope downward

$$\rightarrow \begin{cases} a=0 \\ b=1 \end{cases} \quad \begin{cases} a=0 \\ b=-1 \end{cases}$$



- Regression analysis - uses.
- ② modeling complex relationships among data elements.
- ③ estimating the impact of a treatment on an outcome.
- ④ Quantifying relationship b/w event & response (clinical drug tests, engineering safety tests)
- ⑤ Identifying patterns (used to forecast future behaviour given known criteria — natural disaster damages, election results)
- ⑥ Statistical hypothesis testing (whether a premise is likely to be true/false in light of the observed data)

Basic Linear Regression models

- Simple Linear Regression (only single independent variable)
- Multiple Linear Regression (two or more independent variables)

Other Regression Methods

- Logistic Regression (models a binary categorical outcome)
- Poisson Regression (models integer count data)
- Multinomial Logistic Regression (models a categorical outcome)

Simple Linear Regression

→ relationship b/w a dependent variable and a single independent predictor variable.

→ line equation,

$$y = \alpha + \beta x$$

intercept, α → line crosses the y axis

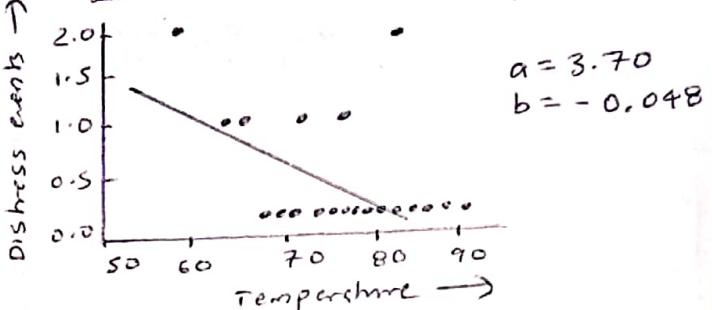
slope, β → change in y given to increase in x

→ Example; (Rocket Failure)

- regression model demonstrated a link b/w temperature & O-ring failure
- could forecast the change of failure given expected temperature at launch.

- component of distress (one of the two types of problem)
 - (i) Erosion [excessive heat burns up the O-ring]
 - (ii) Blowby [hot gases leak through poorly sealed O-rings]

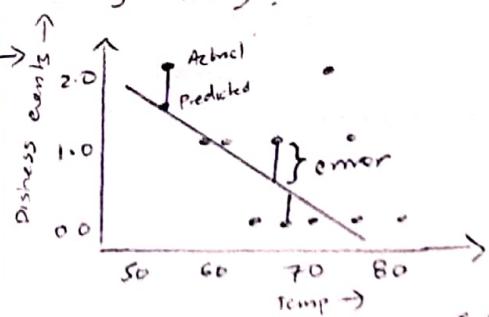
Scatterplot (23 launches)



- at 60°F , just one O-ring distress 70°F , expect around 0.3 failure
- Some interpretations;
 - Launches at $31^{\circ}\text{F} \rightarrow 3$ times more risky than at 60°F
 - $\Rightarrow 8$ times risky than at 70°F

Ordinary Least Squares (OLS) Estimation

- to determine optimal estimates of α and β
- slope & intercept chosen so that they minimize the sum of squared errors (vertical distance b/w predicted & actual y value).



→ Goal \Rightarrow minimizing $\sum (y_i - \hat{y}_i)^2$

→ solution for ' α ' depends on value of ' β '.

$$\alpha = \bar{y} - \beta \bar{x}$$

→ value of β

$$\beta = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

→ break this eqn to simplify it a bit.

→ denominator for $b \Rightarrow$ similar to variance of x

$$\text{Var}(x) = \frac{\sum (x_i - \bar{x})^2}{n}$$

average squared deviation from mean of x

→ numerator for $b \Rightarrow$ similar to covariance for x & y

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

→ 'n' gets cancelled as we divide covariance by variance of x .

$$\therefore b = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

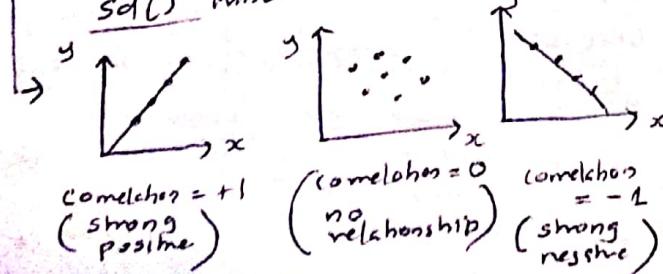
- easy to calculate in R language
- $x = \text{temp}$ & $y = \text{distress event}$
- use `cov()` & `var()` of R. to estimate ' b '
- use `mean()` to estimate ' a '

Correlations

- a number that indicates how closely their relationship follows in a straight line.
 - refers to Pearson's correlation coefficient (developed by Karl Pearson)
 - Ranges b/w $[-1 \text{ } \& \text{ } +1]$
 - extreme values \Rightarrow perfectly linear relationship.
 - close to zero values \Rightarrow absence of linear relationship.
 - Pearson's correlation - Formula
- $$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}}$$

→ Eg: To calculate correlation b/w launch temperature & no. of distress events (-0.51)

→ In R language, `cov()` and `sd()` functions are used



- Interpretation of correlation strengths
- weak $\Rightarrow [0.1 \text{ to } 0.3]$
 - moderate $\Rightarrow [0.3 \text{ to } 0.5]$
 - strong \Rightarrow above 0.5
for both positive & negative values.
- context-based
- 0.5 for human beings [extremely high]
 - 0.5 for mechanical processes [weak]

Multiple Linear Regression

- more than one independent variable.
- Extension of simple linear regression.
- similar \Rightarrow goal is to minimize prediction error.
- difference \Rightarrow additional independent variables involved.

→ General form:

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

ϵ (epsilon) \Rightarrow added since predictions are not perfect.

ϵ : residual term

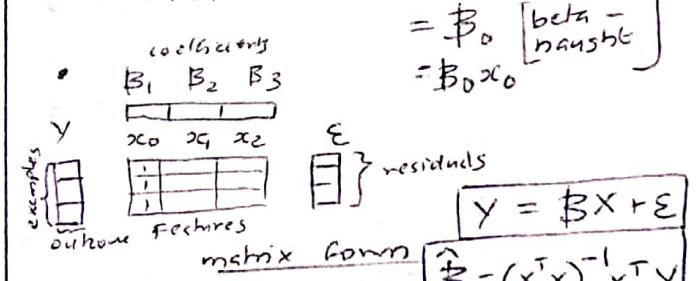
- y changes by amount β_i for each unit increase in x_i .

- when all independent variables = 0

$$y = \alpha \quad (\text{intercept})$$

$$= \beta_0 \quad (\text{beta - intercept})$$

$$= \beta_0 x_0 \quad (\text{beta - feature})$$



Strengths

- common approach for modelling numeric data
- can be adapted to any modelling task
- estimates strength & size of relationships among features & outcomes

Weakness

- makes strong assumptions on data.
- models' form \Rightarrow must be specified in advance
- doesn't handle missing data.
- only works with numeric features.
- requires some knowledge of statistics to understand model.

Module IV

Neural Networks

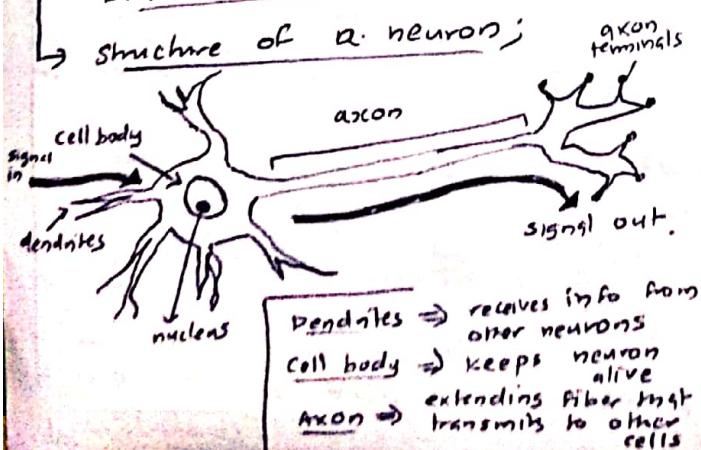
- mimics the structure of animal brains to model arbitrary functions.
- Artificial Neural Network (ANN)
- models the relationship b/w a set of input signals & an output signal.
- it uses a model derived from our understanding of how a biological brain responds to stimuli from sensory inputs.
- brain uses a network of interconnected cells (neurons)
- ANN uses a network of artificial neurons to solve learning problems
- human brain is made up of 85 billion neurons

Applications

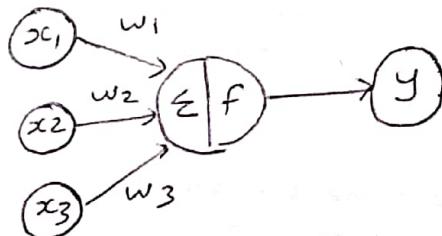
- speech and handwriting recognition programs (used in postal mail sorting machines)
- automation of smart devices (self driving cars & self-piloting drones)
- sophisticated models (weather and climate patterns)
- can be also applied to numeric prediction & unsupervised pattern recognition

Biological motivation

- ANNs were designed as conceptual models of human brain activity.



- incoming signals are received by dendrites (through bio-chemicals)
- impulses are weighted according to relative importance / freq.
- A cell fires only when it reaches a threshold (output signal transmitted via electro-chemical process)
- Electric signal is again processed as chemical signal at axon's terminals (to be passed to the neighboring neurons across a tiny gap [synapse])
- Artificial neurons \Rightarrow similar to biological model
- directed network diagram (defines relationship b/w input signals [x_1, x_2, x_3] & output signals [y variables])
- each dendrite's signal weighted (as in biological neuron)
- signal is passed according to activation function (f).



can be represented by the formula

$$y(x) = f \left(\sum_{i=1}^n w_i x_i \right)$$

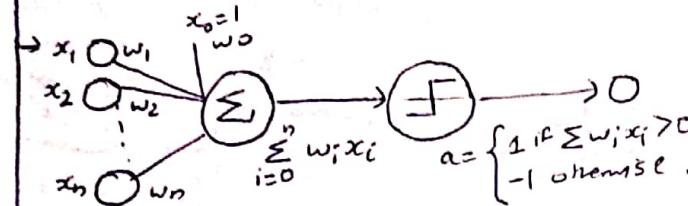
Characteristics of neural networks

- An activation function
- Network topology (architecture)
- Training algorithms.

- Activation Functions [transforms combined input signals into a single output signal to be broadcasted]
- Network topology [describes no. of neurons and no. of layers in which they are connected]
- Training algorithms [specifies how weights are set in order to inhibit/excite neurons in proportion to input signals]

Perceptrons

- one type of ANN systems is based on unit called Perceptron



- Takes a vector of real valued inputs, calculates a linear combination of these inputs.

$$O(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

- Another form of function;

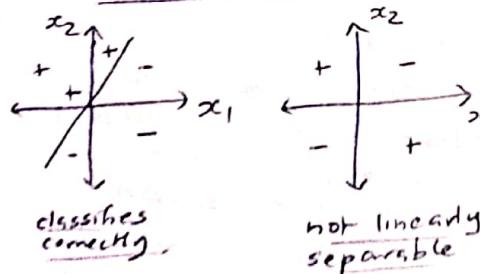
$$O(\vec{x}) = \text{sgn}(\vec{w} \cdot \vec{x})$$

$$\text{sgn}(y) = \begin{cases} 1 & \text{if } y > 0 \\ -1 & \text{otherwise} \end{cases}$$

- A single perceptron can be used to represent many boolean functions

- Assume value of 1 (True) & -1 (False)
- To implement AND function $\Rightarrow w_0 = -0.8$, $w_1 = w_2 = 0.5$
- OR function $\Rightarrow w_0 = -0.3$, $w_1 = w_2 = 0.5$
- ie can represent all primitive boolean functions (AND, OR, NOR, NAND)
- XOR can't be implemented using single perceptrons.

- Decision surface represented by two-input perceptron;



Perceptron Training Rule:

- Two well known algorithms are,
- (i) The perceptron rule &
- (ii) The delta rule
- provide the basis for learning networks of many units.

(i) Perceptron Training Rule

(revises weight w_i associated with input x_i)

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

t = target output
 o = output gen. by perceptron.
 η = +ve constant (learning rate)

- conditions \Rightarrow

a) training examples are linearly separable.

b) sufficiently small η used.

(ii) Delta Rule (for examples not linearly separable)

- key idea \Rightarrow gradient descent (to find weights that best fits examples) (basis of Backpropagation algorithm)

- Output for first stage of perceptron

$$O(\vec{x}) = \vec{w} \cdot \vec{x}$$

- measure of training error of a weight vector,

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

D = set of training examples
 t_d = target output of d
 o_d = output of linear unit d.

Activation Functions

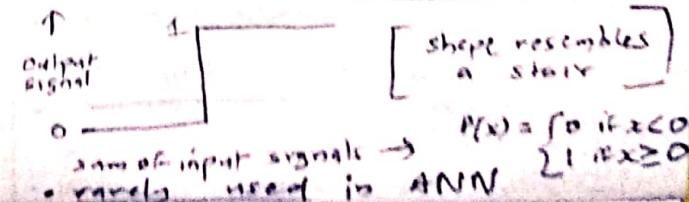
- mechanism by which the artificial neurons processes incoming information & passes it throughout the network

- Biological cases;
summing process of total input & determining whether it meets firing threshold

- In ANN;

Threshold activation function: results in an output signal only once a specified input threshold has been attained.

- Also called Unit step activation function.



Sigmoid Activation Functions

- most commonly used.

$$f(x) = \frac{1}{1 + e^{-x}}$$

$e \approx 2.72$
(base of natural log)

- step / "S" shape
- output signal is not binary
- values can fall anywhere in the range [0 to 1]
- differentiable (ie, calculate the derivative)
[crucial to create efficient ANN algorithms]
- used as default activation fn.

Linear Activation Functions

- results in neural network very similar to a linear regression model.

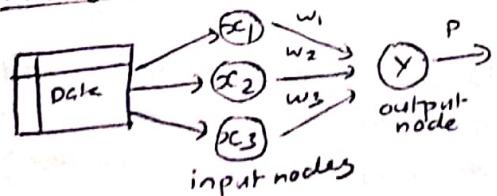
Gaussian Activation Functions

- results in a model called a Radial Basis Function (RBF) network



Network Topology

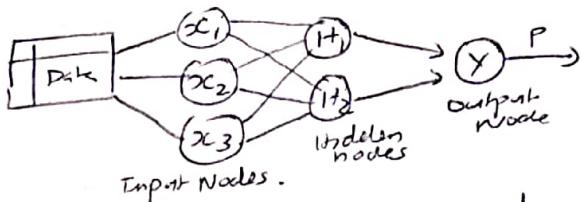
- The patterns and structures of interconnected neurons.
- The three key characteristics:
 - The number of layers
 - whether info is network is allowed to travel backward
 - The number of nodes within each layer of network
- complex networks are capable of identifying more subtle patterns.



(i) No of Layers

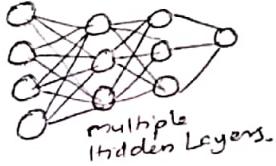
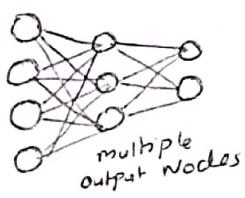
- single layer network
 - only one set of connections weights
 - process incoming data exactly as received
 - used for basic pattern classification

- multi layer network
 - adds one/more hidden layers
 - generally fully connected
 - hidden layers process input signals prior reaching output node.



(ii) direction of information travel

- feedforward networks
 - forms connections -to- connections until reaches to output layer
- multiple outcomes can be modeled
- multiple hidden layers can be applied.
- feedback network
 - allows signals to travel in both directions
 - use: weather forecasting



(iii) No. of nodes in each layer.

- No of input nodes \Rightarrow no. of features in the input data.
- No. of output nodes \Rightarrow no. of outcomes to be modeled.
- No. of hidden nodes \Rightarrow user decides.

Cost Function

- the closer the predicted values to corresponding real values the better the model.
- cost function used to measure how close predicted values are to their corresponding real values
- weights of model \Rightarrow responsible for predicting new values
- Eg: model $\Rightarrow y = 0.9 + x + 0.1$
Cost function measures \Rightarrow how close $(0.9 + x + 0.1)$ is to actual y
- we are responsible for finding better weight for our model (0.9 and 0.1) to come up with a lowest cost.

Back Propagation Algorithms

- ANN must be trained with experience.
- an algorithm which uses strategy of back propagation errors
- relatively slow (compared to other ML algos)
- common in the field of determining
- iterates through many cycles of two processes (each iteration \Rightarrow epoch)
- algorithm cycles until a stopping criterion is reached.
- 2 process cycle includes;
 - A forward phase.
 - A backward phase.

Forward phase

- neurons are activated in sequence from input layer to output layer.
- applying each neuron's weights and activation function along the way.
- upon reaching \Rightarrow output signal is produced.

Backward phase

- result of output signal from forward phase compared to the true target value.
- The difference \Rightarrow results in an error
- error is propagated backwards in the network \Rightarrow to modify the connection weights (to reduce future errors)
- uses derivative of each neuron's AF to identify gradient. (suggests how deeply error will be reduced / increased)
- learning rate \Rightarrow reduction in error by an amount

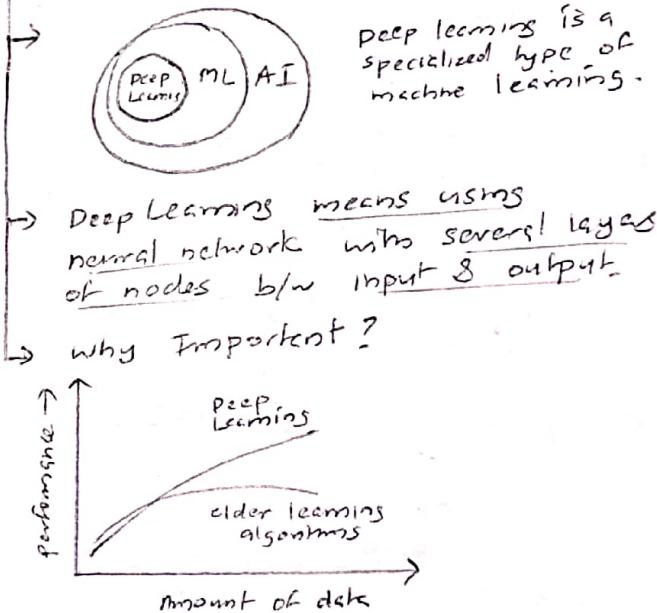
Backpropagation Algorithms

- ① Initialize Network with random weights.
- ② For all training cases (examples)
 - a) Give training inputs to network and calculate output.
 - b) For all layers (output \rightarrow input layer)
 - i) compare network output with correct output
 - ii) update weights in current layer.

Strengths	Weaknesses
<ul style="list-style-type: none"> • adapted to classification or numeric prediction problems • capable of modeling more complex patterns. • makes few assumptions about data's relationship 	<ul style="list-style-type: none"> • computationally intensive • slow to train • prone to overfitting data • difficult to interpret.

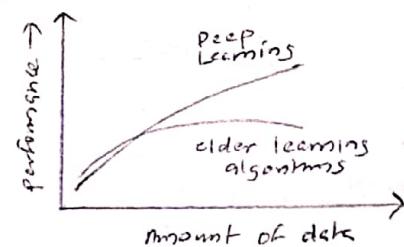
Deep Learning

- branch of machine learning
- set of algorithms that attempt to model high level abstractions in data.
- deep network \Rightarrow many layers b/w input & output.
- Def: Deep learning is a collection of statistical machine learning techniques used to learn feature hierarchies.



- Deep Learning means using neural network with several layers of nodes b/w input & output.

→ why important?



Advantages	Disadvantages
<ul style="list-style-type: none"> • reduces the need for feature engineering. • architecture can be adapted to new problems • best-in-class performance on problems (speech, language, vision etc.). 	<ul style="list-style-type: none"> • Requires large amount of data • computationally expensive to train • no strong theoretical foundations • no theory to guide topology etc. • what is learned is not easy to comprehend.

Module V

Support Vector Machines (SVM)

- A supervised machine learning algorithm
- used for both classification & regression
- can be imagineted as a surface that creates a boundary b/w points of data plotted & their feature values.
- Goal ⇒ To create a flat boundary (hyperplane) — to divide space to create fairly homogeneous partitions.
- SVM Learning combines,
 - instance-based nearest neighbor learners
 - linear regression modeling
- extremely powerful & can model complex relationships.

Applications:

- Gene Expression Data Classification [to identify cancer or other genetic diseases]
- Text Categorization [identification of language used in a doc]
- Detection of rare, but important events [Earthquakes, security breaches etc]

Hyper Planes

- boundary which partitions the data into groups of similar class values.
- SVMs use hyperplanes to partition data into groups of similar class values.
- Es: separate groups of circles & squares in 2D & 3D



Note ⇒ linearly separable (separated perfectly by a straight line/ flat surface)

- There can be more than one choice of dividing line b/w groups of circles and squares.
- Possibilities ⇒ a, b, c
- maximum margin
Hyperplane ⇒ chosen by algorithm from possibilities

Maximum Margin Hyperplane (MMH)

- creates greatest separation b/w two classes
- generalize best to future data.
- improve chance that in case of random noise, points will remain correct side of boundary.
- Es:

Support Vectors:

- points from each class that are closest to MMH.
- Each class must have at least one support vector
- Support vectors alone are possible to define MMH. (key feature of SVMs)

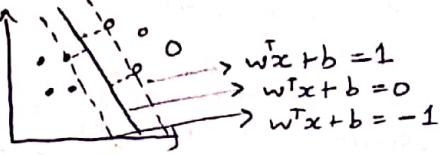
- z_1 z_2 z_3
- $h_1 \Rightarrow$ doesn't separate classes
- $h_2 \Rightarrow$ separates only by a small margin.
- $h_3 \Rightarrow$ separates with max margin

MMH - In case of linearly separable data

- in case of n-dimensional space
MMH is denoted by eqn.
- $\vec{w} \cdot \vec{x} + b = 0$ → implies vectors
 \vec{w} → vector of ' n ' weights $\{w_1, w_2, \dots\}$
- $b \Rightarrow$ a single no. (bias)
↳ intercept term.
- The goal of the process is to find set of weights that specifies two hyperplanes,

$$\begin{aligned}\vec{w} \cdot \vec{x} + b &\geq +1 \\ \vec{w} \cdot \vec{x} + b &\leq -1\end{aligned}$$

$w \rightarrow$ weight vector
 $x \rightarrow$ input vector
 $b \rightarrow$ bias.



All the points of one class fall above first hyperplane & all the points of other class fall below second hyperplane.

$$\text{Distance b/w two planes, } D = \frac{2}{\|\vec{w}\|}$$

$\|\vec{w}\| \Rightarrow$ Euclidean Norm
(distance from origin to vector w)

For max distance, we need to minimize $\|\vec{w}\|$

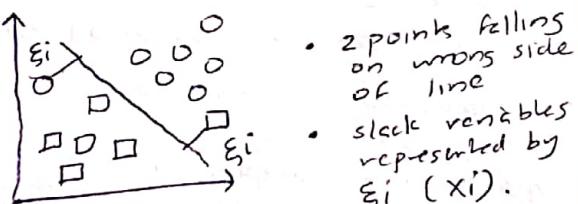
Task represented as a set of constraints

$$\begin{aligned} 1) \quad & \min \frac{1}{2} \|\vec{w}\|^2 \\ 2) \quad & \text{s.t. } y_i (\vec{w} \cdot \vec{x}_i - b) \geq 1, \forall \vec{x}_i \end{aligned}$$

- ① \Rightarrow to minimize Euclidean norm
- ② \Rightarrow subject to conditions that each y_i data points is correctly classified. $\forall \Rightarrow$ for all. $y \Rightarrow$ class value (+1/-1)

MMH - In case of non-linearly separable data.

use of slack variable (ξ_i)
[creates a soft margin that allows some points to fall on incorrect side of the margin]
i.e., allows misclassification of difficult/noisy examples.



- A cost value (C) is applied to all points that violate constraints
- Rather than finding maximum margin, algorithm attempts to minimize total cost.

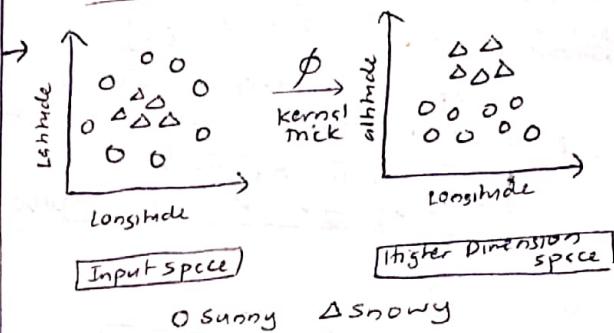
$$\begin{aligned} \min \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t. } y_i (\vec{w} \cdot \vec{x}_i - b) \geq 1 - \xi_i, \forall \vec{x}_i, \xi_i \geq 0 \end{aligned}$$

$C \Rightarrow$ cost parameter | Greater $C \Rightarrow$ harder the optimization
Lower $C \Rightarrow$ wider overall margin.

- balance b/w these 2 preferred to generalize future data.
- Real world applications \Rightarrow relationship b/w variables are non linear

Kernal Trick

Kernels used to map the problem into a higher dimension space. (ie, a nonlinear relationship may suddenly appear to be quite linear)



- Left scatter plot [non linear relationship b/w a weather class — 2 features — latitude & longitude]
- add additional dimensions in order to create separation

Kernal trick \Rightarrow process of constructing new features that express mathematical relationships between measured characteristics.

Altitude feature can be expressed mathematically as an interaction b/w latitude & longitude.

Strengths & Weaknesses (SVMs with nonlinear kernels)

Strengths	Weaknesses
<ul style="list-style-type: none"> used for numeric prediction problems not influenced by noisy data. not prone to overfitting. easier to use than neural networks high accuracy in data mining 	<ul style="list-style-type: none"> various combinations of kernels and model parameters needed testing. slow to train if dataset is large. Results in a complex black box model — difficult to interpret.

Kernal Functions - general form

- denoted by $\phi(\psi(x)) \Rightarrow$ mappings of data into another space.
- applies some transformations to the feature vectors x_i and x_j and combines using dot product
 \Rightarrow returns a single number.

$$K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

All SVM software packages will include these kernels.

Most commonly used,

(i) Linear Kernel

- simplest kernel function
- doesn't transform data at all
- simply expressed as the dot product of the features.
-

$$K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

(ii) Polynomial Kernel

- adds a simple non-linear transformation of the data.
- $K(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j + 1)^d$
 $d = \text{degree of polynomial}$

(iii) Sigmoid Kernel

- analogous to neural networks

$$K(\vec{x}_i, \vec{x}_j) = \tanh(K_0 \vec{x}_i \cdot \vec{x}_j - \delta)$$

$\tanh \rightarrow$ hyperbolic tangent function

(iv) Gaussian RBF Kernel

- general purpose kernel
- similar to RBF neural networks
- performs well on many types of data.

$$K(\vec{x}_i, \vec{x}_j) = e^{-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{2\sigma^2}}$$

$\sigma \Rightarrow$ adjustable parameter.

Note: No reliable rule to choose a kernel.

choice of kernel is arbitrary
 \Rightarrow performance may vary slightly.

Multiclass SVM

- classifications with more than two classes.
- depends on whether classes are mutually exclusive or not.
- 2 methods;

(i) One of classification

- classification not mutually exclusive
- The decision of one classifier has no influence on decision of other classes.
- Eg: Text classification
- Formal definition \Rightarrow

$$\gamma_j(d) \in \{c_j, \bar{c}_j\}$$

- Eg: A doc about '2008 Olympics' should be a member of 2 classes (china class & sports class)

classes	Regions		Industries		Subject	
	UK	China	Coffee	Tea	Electronics	Sports
training set	<input type="checkbox"/>					
	<input type="checkbox"/>					

• Steps:

- Build a classifier for each class
- Apply classifier separately.

(ii) One of classification

- classification mutually exclusive
- each document belongs to exactly one of the classes.
- called multinomial, multiclass functions
- Formal definition

$$\gamma(d) \in \{c_1, \dots, c_j\}$$

- Eg: KNN is a nonlinear one of classifier.
 [doc member of only one class]
 i.e., china.

• Steps:

- Build a classifier for each class
- Apply classifier separately.
- Assign doc to the class with
 - maximum score.
 - maximum confidence value.
 - maximum probability.

Module VI

Evaluating Model Performance

- process of evaluating ML algorithms
- Algorithms have varying strengths and weaknesses.
- Classifiers are evaluated \Rightarrow reflects type of data.
 - Actual class values. [key to evaluation]
 - Predicted class values.
 - Estimated probability of the prediction.
- Goal \Rightarrow maintain two vectors (Actual & Predicted class values)
 - must have same no. of values.

Confusion Matrix (imp. measure of model performance)

- A table that categorizes predictions \Rightarrow whether they match actual value.
- One dimension \Rightarrow possible categories
other dimension \Rightarrow actual values.
- matrix can be created for models that predict any number of class values.

Eg:

	predicted.			predicted
	O	X		O X X
	X	O		X O X
	O	X		X X O
actual	X	O		X

Two classes model Three classes model

- Correct predictions [fall on the diagonals in confusion matrix (denoted by O)]
- Incorrect predictions [off diagonal matrix cells (denoted by X)]
- Performance [based on counts of predictions falling on/off the diagonal]

- Positive class** [class of interest]
- Negative class** [other except class of interest]

The relationship b/w positive & negative class predictions can be depicted as a 2x2 confusion matrix.

4 categories where prediction fall,

- True Positive (TP)
- True Negative (TN)
- False Positive (FP)
- False Negative (FN)

- TP [correctly classified as the class of interest]
- TN [correctly classified as not class of interest]
- FP [incorrectly classified as class of interest]
- FN [incorrectly classified as not class of interest]

Eg: spam classifier.

	no	yes
no	(TN)	(FP)
yes	(FN)	(TP)

Actually spam.
Predicted to be spam.

measuring performance [using 2x2 confusion matrix]

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Error Rate} = \frac{FP + FN}{TP + TN + FP + FN}$$

$$\text{or } \text{Error rate} = 1 - \text{accuracy}$$

i.e., if a model is 95% correct, it is 5% incorrect.

Precision and Recall

Precision [proportion of positive examples that are truly positive]

- trustworthy
- precise model will only predict positive classes

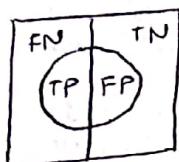
$$\text{Precision} = \frac{TP}{TP + FP}$$

→ Recall [measure of how complete the results are]

- high recall ⇒ large portions of positive examples.

• Eg: An SMS spam filter has a high recall ⇒ if majority of spam msgs are identified correctly.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$



$$\text{Precision} = \frac{\text{D}}{\text{D} + \text{B}}$$

$$\text{Recall} = \frac{\text{D}}{\text{D} + \text{C}}$$

Holdout method

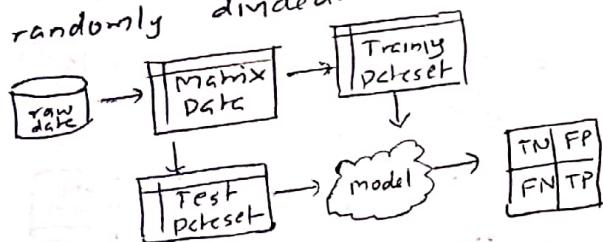
→ procedure of partitioning data into training & test datasets.

→ Training dataset (2/3) ⇒ used to gen. model

→ applied to, test dataset (1/3) ⇒ used to gen. predictions

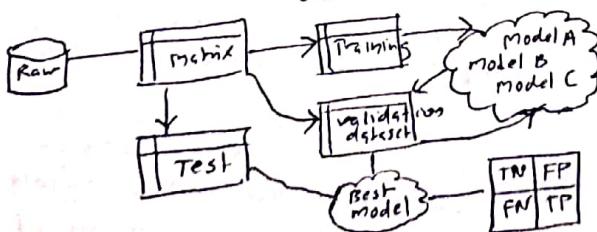
(proportions can vary)

randomly divided.



→ Validation dataset; used for iterating & refining model.

(Typical split ⇒ Training (50%) Test (25%) validation (25%)



• Problem: chance of a class being omitted from training dataset.

• Solution: [Stratified Random Sampling]

guarantees random partitions have nearly same proportion of each class (as the full dataset)

• `createDataPartition()` in R studio.

• Problem: substantial amount of data cannot be used for training.

Cross Validation

→ based on repeated holdout

→ known as k-fold cross validation

→ industry standard for estimating model performance.

→ randomly divides the data into 'k' random partitions called folds. (common = 10 fold CV)

→ For 10 folds, each fold comprises of 10% of data

• folds matching 10% sample used for model evaluation.

• process of training / evaluating done 10 times (10 diff. combinations) and average performance is reported

→ same class balance is maintained

→ createFolds() in R studio

Bootstrap Sampling

→ alternative to k-fold CV

→ less freq. used

→ method ⇒ creation of several randomly selected training & test datasets.

• Results of various datasets are then averaged to obtain final estimate

Difference

• In CV, each example can only appear once.

• In Bootstrap, examples can be selected multiple times process [sampling with replacement]

→ $P(\text{given instance included in training dataset}) = 63.2\%$
 i.e., test dataset = 36.8%

• Bootstrap sample ⇒ less representative of full dataset.

→ Error Rate: $0.632 \times \text{error}_{\text{test}} + 0.368 \times \text{error}_{\text{train}}$

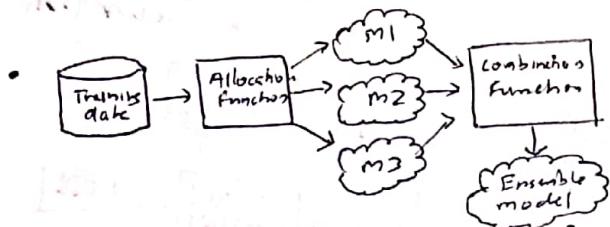
Advantage:

- works better on small datasets
- BS has applications beyond performance measurement
- used to improve model performance.

Improving Model Performance

- Meta Learning: Technique of combining & managing predictions from multiple models to form powerful model
- Ensemble methods: (uses meta-learning approach)

- uses multiple learning algorithms to obtain better performance
- helps improving machine learning results by combining several models
- Idea \Rightarrow combining multiple weaker learners \Rightarrow a stronger learner is created.
- Ideal ensemble \Rightarrow includes diverse set of models.



Input data \Rightarrow to build a number of models.

Allocation Function \Rightarrow dictates how much of training data each model receives.

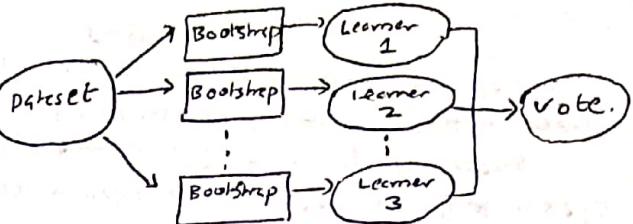
Combination Function \Rightarrow how disagreements are reconciled

• Advantages (over single models)

- (i) Better generalizability of future problems.
- (ii) Improved performance on massive or minuscule datasets.
- (iii) Ability to synthesize data from distinct domains
- (iv) More understanding of difficult learning tasks.

Bagging

- one of the first ensemble method.
- widespread acceptance.
- technique \Rightarrow bootstrap aggregating
- generates a no. of training datasets
- datasets used to generate models.
- models predictions are combined using voting (for classification) averaging (numeric prediction)



→ Bagging is often used with decision trees (tendency to vary)

→ ie, used with unstable learners. (one changes substantially when input changes slightly)

→ bagging() \Rightarrow to train model in R. ipred package \Rightarrow for bagged decision trees
n bag parameter \Rightarrow used to control no. of decision trees
 no. of votings (default = 25)

Boosting

- another ensemble-based method
- boosts performance of weak learners

→ similar to bagging
 uses models trained on resampled data & vote to determine final prediction

→ key differences to bagging

- resampled datasets are constructed specifically to generate complementary learners
- each learner's role is given a weight (based on past performance)

→ possible to increase performance \Rightarrow by adding additional classifiers

Adaboost / Adaptive boosting:

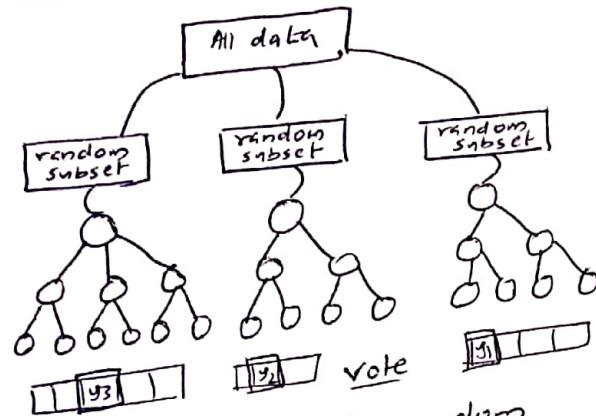
- A boosting algorithm
- generates weak learners that learn a large portion by paying attention to misclassified examples

Random Forests

- another ensemble-based method.
- focuses only on ensembles of decision trees.
- designed by Leo Breiman and Alejo Cutler.

→ combines (base principles of basins) with (random feature selection)
 ⇒ to add additional diversity to decisions tree models.

→ After ensemble of trees (forest) is generated, model uses a vote to combine the tree's predictions



→ uses only a small, random portion of full feature set

→ can handle extremely large datasets

→ error rates similar to any other method.

→ easier to use, less prone to overfitting
 powerful & versatile.

Strengths	Weaknesses
<ul style="list-style-type: none"> all purpose-model that performs well on most models can handle noisy/missing data as well as categorical/continuous features. selects only the most important features. can be used on data with large no. of features 	<ul style="list-style-type: none"> not easily interpretable require time to work to tune the model

→ randomForest used in R studio.
 package

→ randomForest() creates 500 trees by default
 → considers $\sqrt{(\text{no of features})}$ at each split.

→ Out of Bag (OOB) error rate \Rightarrow Total error rate of predictions.