

Intro. ML, Types, Linear Regression, classification & logistic regression, DT & RF, Naive Bayes & SVM, Applications of ML

Regression Analysis (RA)

A predictive modelling technique that investigate relationship b/w dependent (target) & independent variable (predictor).

Used for forecasting & finding relation b/w variables.

E.g. relationship b/w cash driving & no. of road accidents by driver

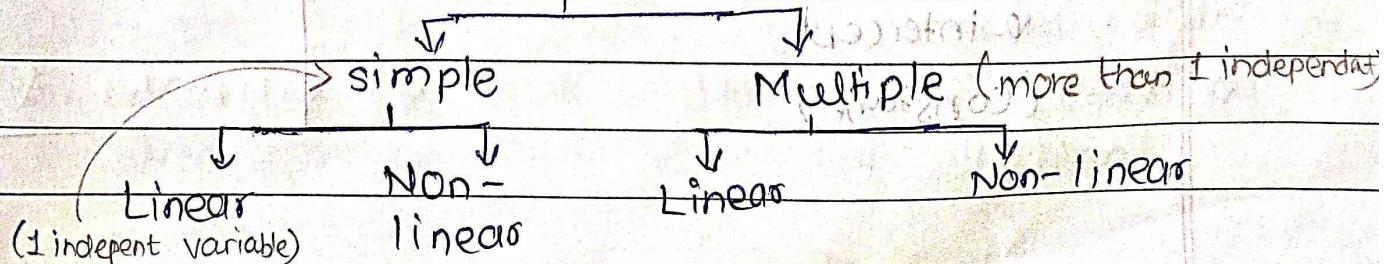
In RA, curve/line is fitted to data pts., in such manner that differences b/w distances of data pts. from curve or line is minimized.

Benefits

Indicates significant relationship b/w dependent & independent variables.

Indicates strength of impact of multiple independent variables on dependent variables.

Types of RA



$$\hat{y} = b_0 + b_1 x$$

$$b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

Random error term

Linear Regression - Continuous

Most widely known modelling technique.

A linear approach to modelling of the relationship b/w a dependent variable & one or more independent variables using best fit straight line (known regression line).

In this technique dependent variable is continuous & independent variable is discrete.

Simple Linear Regression

Objective: Fit a straight line through the data points.

To find best-line called regression line through the points which is nearest to most of pts.

Example

Predict height from age

Independent var x

House price predictor

Sales in upcoming months predictor

Y-intercept (slope) of a line

$y = b_0 + b_1 x + \epsilon$ random error

\downarrow error b_0 slope b_1

\downarrow Y-intercept

(constant)

\downarrow error

Note: LSE method that builds model & RMSE is metric that evaluate model's performance

classmate

Date _____

Page _____

$$R^2 = \frac{\sum (y - \bar{y})^2}{\sum (y_i - \hat{y}_i)^2}$$

If R^2 close to 1, \rightarrow good fit
If R^2 close, 0 \rightarrow no relationship

RMSE: tells measure of dispersion of predicted values from actual values.

$$RMSE_{\text{Error}} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n}}$$

$y_i - \hat{y}_i$ = Residual or difference b/w actual & predicted values.

y_i = observed value

\hat{y}_i = predicted \rightarrow

Use cases of LR

- (1) Sales of product, pricing, performance parameter.
- (2) Generating insights on consumer behavior, profitability & other business factors.
- (3) Evaluation of trends, making estimates & forecasts.

- 4) Determining marketing effectiveness, pricing & promotions on sales of product.
- 5) Assessment of risk in financial services & insurance domain.
- 6) Studying engineer performance from test data in automobiles.
- 7) Calculating causal relationship b/w parameters in biological systems & environment.
- 8) Conducting market research studies & customer survey results analysis.
- 9) Astronomical data analysis.
- 10) Predicting house price based on sizes of houses.

Classification: y is discrete variable.
Classification algo. is supervised ML technique used to identify category of new observations on training data basis.
Program learns from given dataset of observed & classifies new observation into no. of classes, such as "Yes" or "No" or 1, 0, spam or not spam, cat or dog. It uses variables as input.
Output variable is category, not value such as 'G' or 'B'. A discrete output function (y) is mapped to input variable (x). $y = f(x)$

	classmate		classmate	
	Date _____	Page _____	Date _____	Page _____
1) Algo. which implements classification logic dataset is called classifier.			Types of classification problems:	
2) 2 types of classifier are binary & multi-class			Linear models vs Non-linear models	
1) Binary classifier			Logistic regression vs K-NN	
If classification problem has only 2 possible outcomes			SVM vs Support Vector Machine	
2) Multi class classifier			Kernel SVM vs	
More than 2 outcomes are handled by			Random Forest vs Naive Bayes	
Classification of types of learners			Decision Tree vs Random Forest	
Learners in classification problem			Evaluating model	
Lazy learners: Firstly stores training dataset & wait until receives test dataset.			Log loss or cross-entropy	
The classification is done on basis of most related data stored in training dataset.			Confusion matrix	
Take less time in training but more time in prediction.			AUC - ROC curve	
Ex:- K-NN, naive bayes classifier			Log loss / cross entropy	
2) Eager Learners: They develop classification model based on training dataset by receiving test data. More time learning, less time prediction.			for evaluating performance of classifier, whose output is probability value b/w 0 & 1.	
Ex:- DT, NB, ANN, etc.			After good binary classification, value should be near to 0.	
			Other issues: Log loss represents higher accuracy or model is overfitted.	
			at it is better to use AUC - ROC curve.	

2) Confusion Matrix / error matrix: This chapter gives a table / matrix of O/P & describes the performance of predictions in summarized form, which has total no. of correct predictions & incorrect predictions.

	Actual +ve	Actual -ve
Predicted +ve	TP	FP
Predicted -ve	FN	TN

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

Recall (R) = $\frac{TP}{TP+FN}$ → measures model's ability to detect true samples

Sensitivity / precision = $\frac{TP}{TP+FP}$ → just positive prediction values that reflect how reliable model is in classifying

$$\text{F-1 score} = \frac{2 \times (P \times R)}{(P+R)}$$

Use 'cover of classification' Algo. to find out

Classification techniques to increase prior probability

Email spam detection, speech recognition

Speech recognition using hidden markov

Identification of cancer, tumor cells, etc

Drugs classification

Biometric identification, fingerprinting,虹膜识别

Face detection, handwritten digit recognition

AUC-ROC (Area Under ROC Curve)

This graph that shows performance of the classification model at diff thresholds.

To visualize the performance of multi-class classification, use AUC-ROC visualization

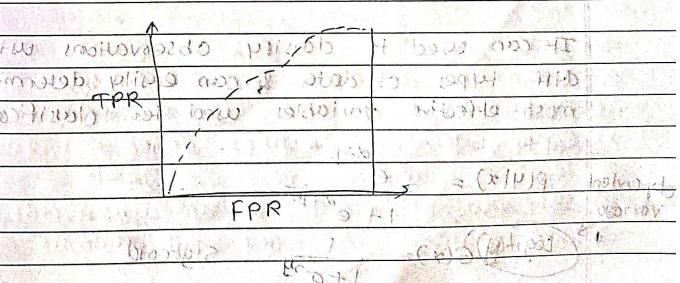
ROC curves plotted with TPR & FPR

$$\text{Recall} \rightarrow \text{TPR} \rightarrow \text{Y-axis} \rightarrow \frac{\text{TP}}{\text{TP}+\text{FN}}$$

$$\text{FPR} \rightarrow \text{X-axis} \rightarrow \frac{\text{FP}}{\text{FP}+\text{TN}}$$

ROC → probability curve

AUC → degree of separability



In logistic, logit transformation is applied on Odds - Prob. of success / prob. failure. (log odds)

LR estimates probability of event occurring based on data.

Logistic Regression - categorical

for predicting categorical dependent variable using given set of independent variables, O/P categorical or discrete (Y or N, 0 or 1). But instead of giving exact values of 0 or 1, it gives probabilistic values b/w 0 & 1.

In this regression, instead of fitting a regression line, we fit an "S" shaped logistic function which predicts 2 max. values (0 & 1) O/P.

The curve from logistic function indicates the likelihood of something e.g. whether cells are cancerous or not, whether a mouse is obese or not based on its weight etc.

It is mLE algo. bcoz it has ability to provide probabilities & classify new data using continuous & discrete datasets.

It can used to classify observations using diff. types of data & can easily determine most effective variables used for classification

$$P(Y|x) = \frac{e^{a+bx}}{1+e^{a+bx}}$$

$$\text{Logit}(y|x) = \ln \frac{y}{1-y} \quad \text{Sigmoid}$$

LR doesn't require linear relationship b/w dependent & independent variables. It can handle various types of relationships bcoz it applies log transformation to predicted.

It requires large sample size bcoz mle likelihood estimators are less powerful at low sample size than ordinary least square.

Logistic Function (Sigmoid)

Mathematical function used to map the range [0, 1] predicted values onto probabilities. It maps any real prob. value into another real value within a range of 0 to 1. Instead of 0 & 1, which cannot go beyond this limit. We use concept of threshold value.

Assumptions of LR: Logistic regression is obtained from LR eqn. Mathematical steps to get eqn are given.

Eqn of straight line is $y = b_0 + b_1 x + b_2 z_2 + \dots$ b/n only terms. In logistic, y can be b/w 0 & 1 only, so for this lets divide above eqn by $(1-y)$.

$$\frac{y}{1-y}; \quad 0 \text{ for } y=0 \\ 0 \text{ for } y=1$$

But need range b/w $-\infty$ to $+\infty$, take log.

$$\ln \left[\frac{y}{1-y} \right] = b_0 + b_1 x_1 + b_2 x_2 + \dots$$

MLE finds diff. values of b_i through multiple iterations to optimize for best fit.

Types

Binomial → only 2 possible types of dependent var. such as 0 or 1, pass or fail.

Multinomial → there can be 3 or more possible unordered types of dependent var., such as 'cat', 'dog' or 'sheep'.

Ordinal → 3 or more possible ordered types of dependent variables such as 'low', 'medium' or 'high'. on a scale 1 to 5.

Decision Tree → Non-linear classification & both

These are non-parametric supervised learning method used for classification & regression task. Mostly used in classification problems & works for both continuous and categorical input & output variables.

- Decision tree classifier is in form of tree structure with 2 types of nodes:
- Decision mode → specifies choice or test of some attribute with 1 branch for each outcome.
- Leaf node → indicates decision or classification of example.
- It is tree structured classifier, where internal nodes represent features of dataset, branches represent decision rules & each leaf node represent outcome.

The decisions or tests priori is performed on basis of features of dataset.

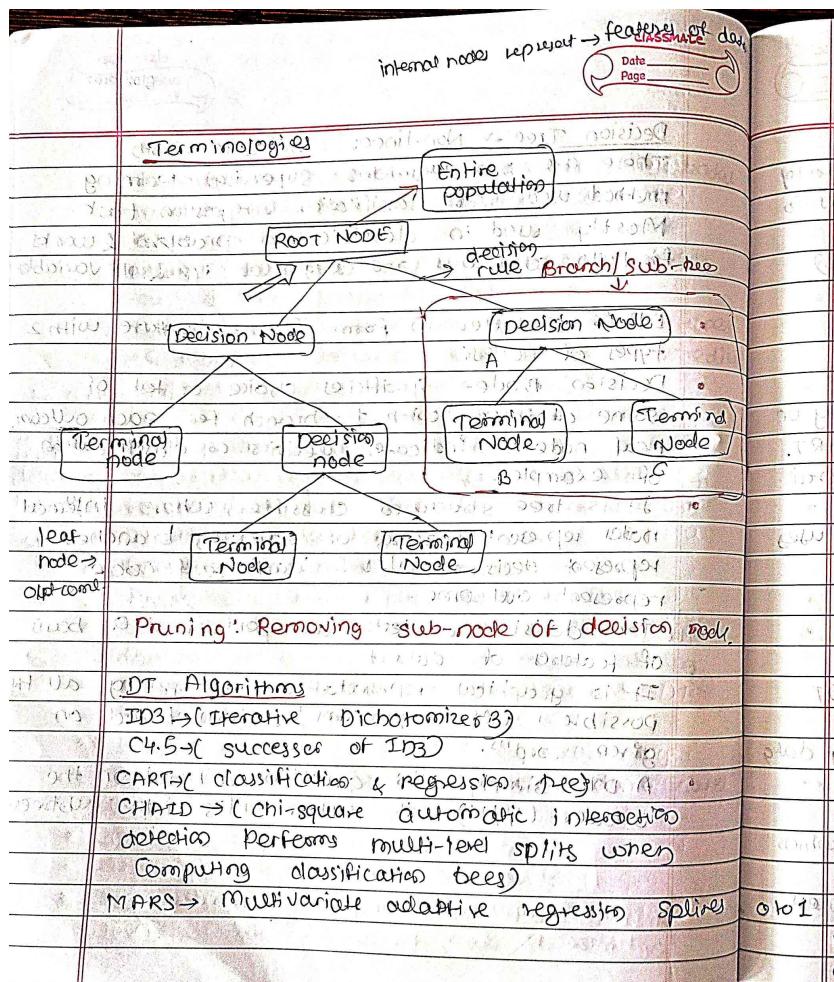
It is graphical representation of getting all the possible splits to problems / decision based on given condns.

- A DT is simply yes/no question & based on the answer (Y/N) it further splits into subtrees.

- 1) Trying to figure out who will win election
- 2) Whether student will pass or fail exam
- 3) Whether customer will come back
- 4) Whether email is spam
- 5) A credit card company want to know whether transaction amount & credit score impact a prob. of given transaction being fraudulent

Advantages

- 1) Makes no assumption about distribution of data in feature space.
- 2) Easily extended to multiple classes
- 3) Quick to train
- 4) Very fast at classifying unknown records
- 5) Good accuracy for many simple datasets



Attribute Selection Measures

While implementing DT, main issue arises that how to select best attribute for root node & for sub-nodes.

So, to solve such problems there is a technique which is called Attribute selection measure or ASM.

By this measurement, we can easily select the best attribute for nodes of tree.

- Techniques used for testing of code in function of
1) Info. gathering
2) Client analysis

Info. Gain is measure of changes in entropy after segmentation of dataset based on attribute.

It calculates how much info. feature provides us about class.

According to value of info. gain, we split node & build DT.

Decision algorithm tries to maximize Node or info gain & node/attribute having highest info gain is split first.

Infor Gain = Entropy (S) - [(weighted avg)* Entropy
into: theory metric that measures impurity / uncertainty in observations. It determines how DT chooses to split data

Entropy: It is metric to measure impurity in given attribute. It specifies randomness in dataset, so it's the value of uncertainty.

$$\text{Entropy}(S) = -P(\text{yes}) \log_2 P(\text{yes}) - P(\text{no}) \log_2 P(\text{no})$$

It is used to measure the entropy of the dataset.

Total no. observations must be divided among samples and each remaining part must be

to construct DT which can be divided into

Gini Index: It is a measure of impurity or purity used while creating DT in CART.

An attribute with low Gini index is preferred compared to high Gini.

It only creates binary splits, & CART uses

Gini index to calculate the purity of the

Gini Index formula = $1 - \sum p_i^2$

Major factors considered to improve DT performance.

Choosing splitting attribute: By examining data in training set & getting input from domain experts.

Ordering splitting attributes: Order in which attributes are chosen is important.

Tree structure: A balanced tree with fewest levels is desirable.

Work of decision tree is to classify it uniformly.

Stopping criteria: The creation of tree stops when training data are perfectly classified. In addition, stopping earlier may be performed to prevent overfitting.

Training data: If training dataset is too small, tree may not work properly.

With more data, it is possible to remove unwanted branches.

If tree is too large, tree may overfit.

Pruning: Once tree is constructed, pruning phase might remove redundant comparisons or remove subtrees to achieve better performance.

Cost of complexity in pruning.

Reduced error pruning.

Advantages: It is simple, understanding, simple, etc.

Simple to understand as it follows same process that human follows while making any decision in real life.

Useful for solving decision-related problems.

It helps to think about all possible outcomes for problem.

There is less requirement for data cleaning compared to other algo.

	<p><u>Disadvantages</u></p> <ul style="list-style-type: none"> ① The DT contains lots of layers, which makes it complex. ② It may have an overfitting issue, which can be resolved using RF algorithm. ③ For more class labels, computational complexity of decision tree may go up. ④ If dataset has many categorical variables, then it may lead to overfitting. <p><u>Uses</u></p> <ul style="list-style-type: none"> • Building knowledge management platforms for customer service that improve 1st call resolution, average handling time & customer satisfaction rates. • In finance for forecasting future outcomes & assigning probabilities to those outcomes. • Binomial option pricing predictions & real options analysis. • Customer's willingness to purchase a given product in a given setting, online & offline both. • Product planning, for example, Gerber products, Inc. used decision trees to decide whether to continue planning PNC fees manufacturing toys or not & to increase general business decision-making. • Loan approval.
	<p><u>Random Forest</u> - Non linear classifiers (categorical) used for both classification & regression.</p> <p>Based on concept of ensemble learning, in which is process of combining multiple classifiers to solve complex problems to improve performance.</p> <p>RF is classifier that contains no. of decision trees on various subsets of dataset & take average/majority to improves predictive accuracy.</p> <p>Instead of relying on 1 DT, RF takes predictions from each tree, & based on majority votes of predictions, it predicts final output.</p> <p>The greater no. of trees in forest leads to higher accuracy & prevents overfitting.</p> <p>but fails to generalize unseen testing data.</p> <p>RF model splits data into training & testing sets.</p> <pre> graph TD TD[Training data] --> TD1[Training Data] TD --> TD2[Testing Data] TD1 --> DT1[DT 1] TD1 --> DT2[DT 2] TD1 --> DTn[DT n] DT1 --> V[Voting] DT2 --> V DTn --> V V --> P[Prediction] </pre>

Assumptions

- Since RF combines multiple trees to predict classes of dataset, it is possible that some DT may predict correct OIP, while others may not. But together, all the trees predicts correct OIP.
- Below 2 assumptions:
 - There should be some actual value in features variable of dataset so that the classifier can predict accurate results rather than guessed result.
 - The predictions from each tree must have very low correlations.
- Need for RF more than other classifier algorithms due to its less training time, compared to others.
- It predicts OIP with high accuracy, even for large dataset, it runs efficiently.
- It can also maintain accuracy when large proportion of data is missing.

Advantages:

- Capable of performing both C & R.
- Handling large datasets with high dimensionality.
- It enhances accuracy of model & prevents overfitting issue.

Disadvantage

- Although RF uses both C & R, it is not more suitable for Regression type problems.
- User interface is not much user friendly.
- Banking lot uses this algo. for identification of borrower's loan risk, work with algorithm.
- Medicine: Disease trends & risks of disease can be identified.
- Land uses: We can identify areas of similar land use.
- Marketing: Marketing trends can be identified.
- Wood as fertilizer & fuel can be used.
- Marketing, bank, retail, stock market, etc.
- Alg. can be used for classification & regression.
- No programming skills required.

Naive Bayes:

- Based on Bayes theorem & used for solving classification problems (problems of identification)
- Mainly used in text classification, other includes high-dimensional training datasets.
- It is simple & most effective classifier algorithm, which helps in building fast (ML) models that can make quick predictions.
- It is probabilistic classifier, means: Predicting on the basis of prob. of object.

(Naive Bayes is called Naive bcoz it assumes that occurrence of certain feature is independent of occurrence of other features. Such as if fruit is identified on basis of color, shape & taste, then red, spherical & sweet fruit is recognized as apple. ∵ each feature individually contributes to identifying that it is apple w/o depending on each other.)

Bayes: called Bayes bcoz based on principle of Bayes theorem.
fundamental assumption is that each feature makes an independent, equal contribution to outcome.

at Bayes Theorem / Bayes rule / Bayes law
Used to determine prob. of hypothesis with prior knowledge. It depends on conditional probability in basic solution $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$. Prob. of evidence given that prob. of hypothesis is true.

prob. of hypothesis A prior prob. (prob. of hypothesis of observed event B on observing evidence)

Marginal prob. (prob. of evidence)

Types of Bayes Model: 1) Gaussian

1) Gaussian: - This model assumes that features follow a normal distribution. This means if predictions taken continuous values instead of discrete, then it assumes that these values are sampled from Gaussian distribution.

2) Multinomial: - This classifier used when data is multinomial distributed. It is primarily used for document classification problems, it means particular document belongs to which category such as sports, politics, education etc. The classifier uses freq. of words in input predictors.

<p>3) Bernoulli: This classifier works similar to Multi-nomial classifiers, but predictor variables are independent Boolean variables.</p>	<p>SVM - Linear classification - categorical.</p>
<p>Such as if particular word is present or not in document. This model is also famous for document classifications tasks.</p>	<p>It is one of most popular supervised learning algo., can be used for both C & R.</p>
<p>Advantages:</p> <ul style="list-style-type: none"> Fast & easy ML algo. to predict the class of datasets. Can be used for binary as well as multi-class predictions. It performs well in multi-class predictions compared to other algos. Most popular for text classification problems. 	<p>The goal of SVM is to create best line or decision boundary that can segregate n-dimensional space into classes so that can easily put new data pt. in correct category in future. The best decision boundary - Hyperplane.</p>
<p>Disadvantages:</p> <ul style="list-style-type: none"> Difficult with training dataset. It assumes that all features are independent & unrelated, so can't learn relationship b/w features. 	<p>SVM chooses extreme pts. vectors that help in creating hyperplane. These extreme cases called support vectors & algos. called SVM.</p>
<p>Uses:</p> <ol style="list-style-type: none"> For news sentiment classification For classifying articles in to different filters For credit scoring Medical data classification In real time predictions bcoz NB is aages learner In text classification such as spam filtering & sentiment analysis. 	<p>SVM Model: </p>

Hyperplane & SV in SVM algo.

Hyperplane: There can be multiple lines/decision boundaries to segregate classes in n-dimensional space, but we need to find out best decision boundary that helps to classify data pts.

The dimensions of hyperplane depend on features present in dataset, which means if there are 2 features, hyperplane will be straight line, 3 features - 1 dimensional plane.

We always create hyperplane that has max. margin, which means max. dist. b/w data pts.

Support Vectors: The data pts. or vectors that are closer to hyperplane & which affect the position of hyperplane termed Support vector. Since these vectors supports hyperplane called SV.

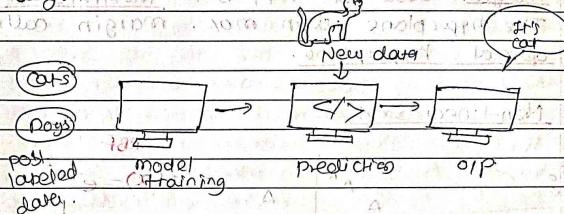
Example!

CAT & DOG.

Strange cat that has some features of dogs, so if we want model that accurately identifies whether it's dog or cat, so such model can be created by SVM algo.

We will start training our model with lots of images of cats & dogs, so it can learn about diff. features & then test with strange creature.

So as SV creates decision boundary b/w these 2 data (cat & dog) & chooses extreme cases (SV), it will see extreme case of cat & dog.



Types of SVM

Linear

Used for linearly

separable data.

Dataset can be classified into 2 classes by straight line.

Non-linear

Non-linearly separable data.

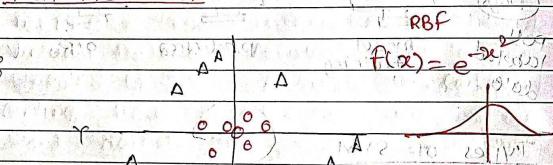
Dataset can't be classified using straight line.



Linear SVM

SVM algo. helps to find best line or decision boundary; this is best boundary / region called hyperplane. SVM finds closet pt. of lines from both classes. These pts. are SVs. The distance b/w vectors & hyperplanes called margin. Goal of SVM is to maximize margin. The hyperplane with max. margin called optimal hyperplane.

Non-linear SVM



So to separate these pts. we need to add 1 more dimension. For linear data, we add 1 dimension. For non-linear data, add 3rd dimension.

It can be calculated as

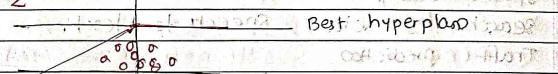
$$z = x^2 + y^2$$

By adding 3rd dimension, sample space will become as shown.

Decision rule
for any $y = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} + b \geq 0 \\ -1 & \text{if } \vec{w} \cdot \vec{x} + b < 0 \end{cases}$

$\vec{w} \cdot \vec{x} \geq c$ $\vec{w} \cdot \vec{x} - c \geq 0$, $\vec{w} \cdot \vec{x} + b \geq 0$
if $\vec{w} \cdot \vec{x} + b \geq 0$
if $\vec{w} \cdot \vec{x} + b < 0$

so, now SVM will divide data into classes in following way



∴ we are in 3-d. space, hence it is looking like parallel to Z-axis. If we convert it in 2d space with $Z=1$, it will become as follows

∴ we get circumference or radius, R , in case of non-linear data.

Hyperplane

Use $\vec{w} \cdot \vec{x} + b = 0$, $x = |\vec{w}| / \|\vec{w}\| \cos \theta = 0$

- 1) Face Detection
- 2) Image Classification
- 3) Text Categorization

Optimization function

$$\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 + C \sum_i \max(0, 1 - y_i(\vec{w} \cdot \vec{x}_i + b))$$

$$\vec{w} \cdot \vec{x} + b \geq 1 \quad \text{if } y_i(\vec{w} \cdot \vec{x}_i + b) \leq 1$$

$$\vec{w} \cdot \vec{x} + b \leq -1 \quad \text{if } y_i(\vec{w} \cdot \vec{x}_i + b) > 1$$

$$y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1 \quad \text{is constraint}$$

$$y_i(\vec{w} \cdot \vec{x}_i + b) \leq 1 \quad \text{is feasible region}$$

Applications of ML

1. Image recognition.

Face detection, automatic friend tagging suggestion

2. Speech recognition

Search by voice, speech to text

3. Traffic prediction

Real time location of vehicle from map.

4. Product recommendations.

Amazon, Netflix, etc.

5. Self-driving cars.

6. Email spam & malware filtering

MLP, DT, NB classifier used.

7. Virtual personal Assistant

Google Assistant, Alexa, Siri, Cortana

8. Online fraud detection.

9. Stock market trading

LSTM

10. Medical diagnosis

11. Automatic language translation.

(1) Suffer from overfitting.

if DT's allowed to grow,

with any control.

(2) Single DT is faster in computation.

(3) When dataset with feature is

taken, say if P by DT it will

(1) RF are created from subsets of data & final O/P based on

(2) Slower compare to DT.

(3) It randomly selects observations

builds DT & avg. result taken

No use of any set of rules

formulate some set of rules for prediction.)