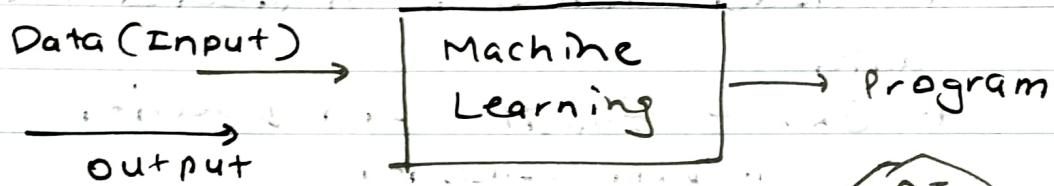
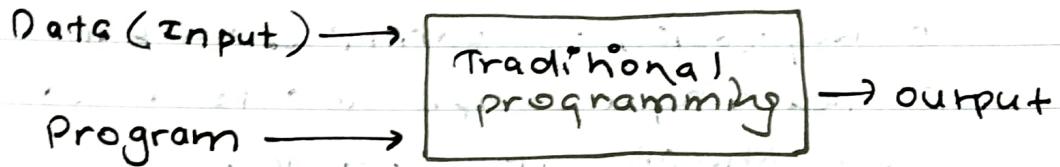
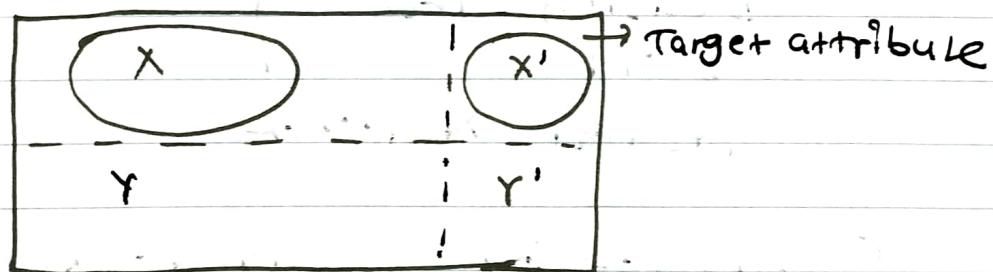


Machine Learning

It is the field of study that gives computers the capability to learn without being explicitly programmed.



Training & Testing



- + we take x & x' for training purpose. Then create one model.
- + we will test our model for y
- + verify y' values (predicted values) from model to actual y' values to check accuracy.
- + Handwriting Recognition System

ML approaches

① Supervised learning

+ classification - categorical output

+ Regression - continuous output

+ teaching a comp. how to do something by showing it examples.

+ model learns from labeled training data

+ Ex. - Email Spam Detection

1. Data Collection

2. Training

3. Testing

+ KNN

Decision Tree

SVM

Neural Networks

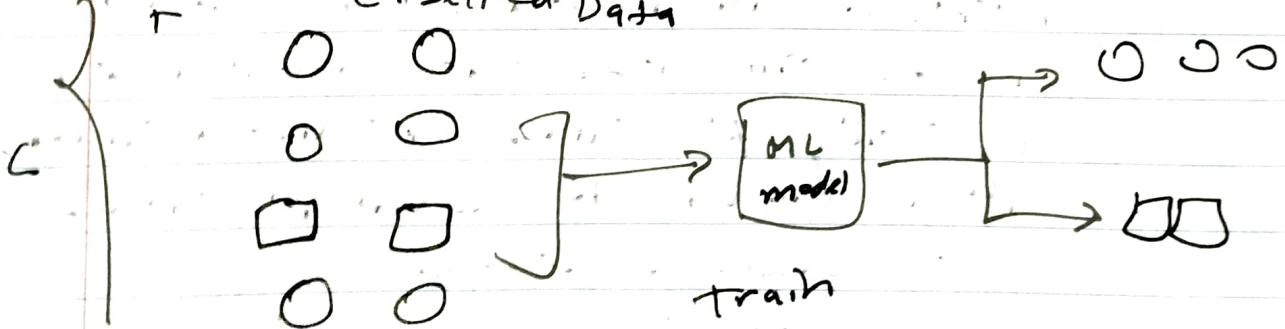
+ Regression

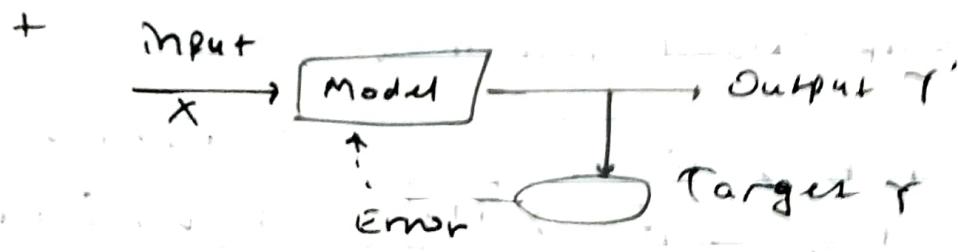
Linear Reg

Polynomial Reg.

Reg. Trees

Labeled Data





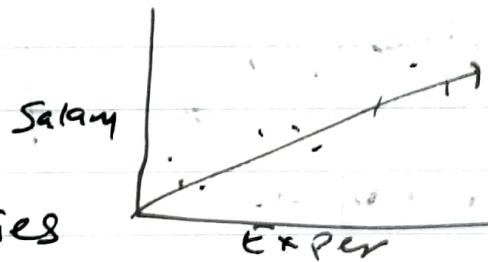
+ Ex.

1. Credit Card fraud Det
2. Sentiment analysis
3. Churn Prediction
4. Medical Diagnosis
5. Spam mail detection

class.

+ Regression

Ex.



1. Time series prediction.
 - a. Rainfall on certain region
 - b. Spend on voice call
2. Price prediction
3. Trend analysis
 - a. Linear or Exponential

② Unsupervised learning.

-

no labels

- + Clustering



- + Clustering applications

- ## 1. Customer data

- a. discover classes of art

- ## 2. Image pixels

- ### a. Discover regions

- ### 3. Word synonym

- ## 4. Document

- ## + Predicting co-occurrences

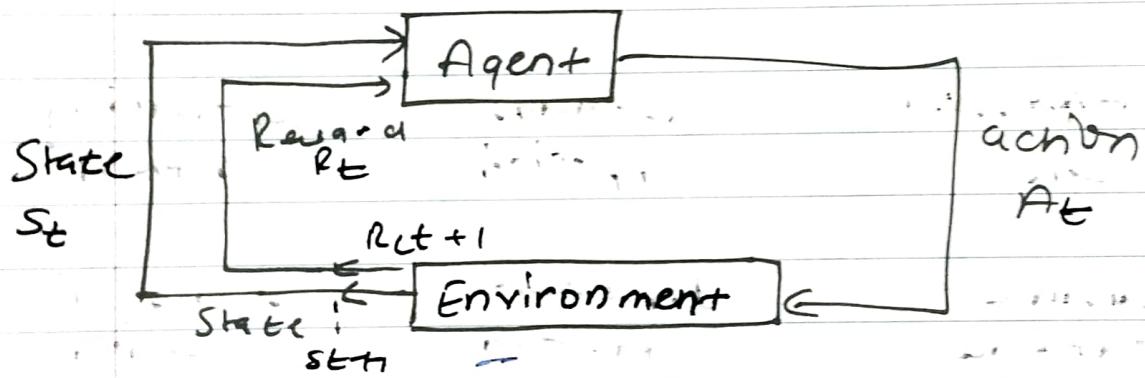
- ## + Market basket analysis

- ## + Time series analysis

- + k-means clustering

- + Association Rule mining
 - a. mining frequent pattern
 - b. association rules
 - c. conditional dependencies
 - d. find frequent pattern
 - e. derive association rule

③ Reinforcement Learning



+ Applications

- a. Business - Strategizing
- b. Gaming - Playstation
- c. Recommender System
- d. Intelligent tutorial system.

- + sup & , unsup x
- + react to envir. on their own
- + rapid grow , adapt

Supervised

labelled data

problems Regression

~~Classif~~

extra supervision

KNN, SVM

LR

calculate outcomes

risk evalution
forecast

Unsupervised

unlabelled data, no guidance

clustering
associat

no super

k-means
A prriori

discover pattern

Recommendations
anomaly detect

Reinforcement

environment

Exploraton
Exploitation

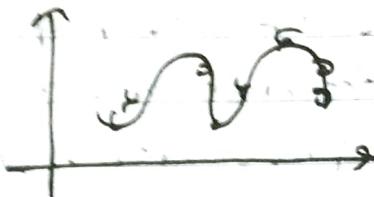
no super.

q-learning
SARSA

series of action

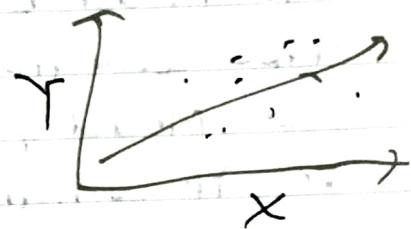
selfdriving cars
gaming
healthcare

Overfitting



+ Extra features

Underfitting



+ less features

Underfitting -

- + ML algo. said. to be underfitting when a model is too simple to capture data complexities
- + inability of model to learn the training data effectively
- + It result in poor performance both on training & testing.
- + inaccurate when applied to unseen ex.

Reas.

- + simplified assumptions
- + high Bias & low variance
- + size of training dataset is not enough
- + features are not scaled.
- + model complex
- + ↑ feature
- + remove noise
- + start of training.

Overfitting

1. Model is said ... when it does not make accurate predictions of testing data.
2. Model gets trained from so much data so it contains (noise & inaccurate entries).
3. Non-parametric & non-linear
4. Linear Algor. (sol.) - DTree

Reasons

1. Low Bias & High Variance
2. model is too complex
3. size of training data.

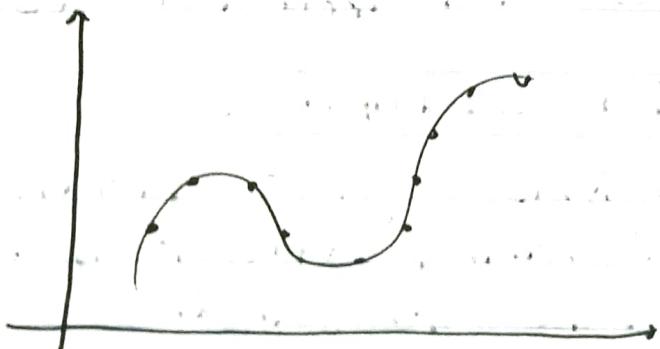
Tech.

1. Reduce model complexity
2. Ridge regularization & lasso regularization
3. Add accurate depth to the parameters

Ex.

Ball

Diag.



Feature Selection Techniques

~~1) Filter methods~~

- + Information gain
- + chi-square test
- + correlate coeff.

~~2) Wrapper methods~~

① Variance threshold

- + simple approach
- + removes features whose variance doesn't meet some threshold.
- + remove all 0 variance features
- + features that have same value in all samples.

from sklearn.feature-selection import VarianceFilter

$$n = [[0, \dots, 3], [0, \dots, 3], \dots]$$

selector = VarianceThreshold(?)

selector.fit_transform(n)

array ([[0, \dots, 3], [0, \dots, 3], \dots]])

② k-best

- + uses statistical tests like chi-square, ANOVA F-test to rank features based relationship → output.
- + select k Features (highest score)
- + Parameters - score_fn & k
- + Score_fn evaluate → feature imp.
- + f-regression - LRP, F-value b/w T & Fcs
- + mutual_info_regression - 2 random variables
- + f-classif - classification, ANOVA fvalue F & T
- + mutual_info_classif - 2 random
- + chi2 - class, F & T chi: sta.
- + Select_Percentile - hi% score band.

```
f = SL.ds import load_digits
```

```
f = SL.Feature_selection import SelectKBest
```

```
X,y = load_digits (return_X_y = True) (chi^2)
```

x.shape

(1532, 64)

```
X-new = SelectKBest(chi2, k=20).
```

fit_transform (x, y)

③

Select Percentile

- + select feature according to percentile of us.

Principal Component analysis

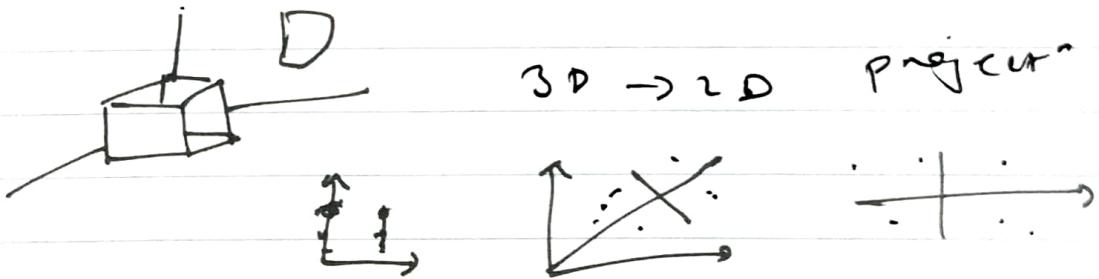
- + Dimensionality reduction tech
- 2. Higher Dim - Overfitting - not generalize
- 3. to project High D \rightarrow Low Dim & also retain reltn, Info, pattern
- 4. set of features \rightarrow principal components

Advantages -

- + remove correlated features
- + improve mlalgo performance
- + reduce overfitting

Disadvantages -

- + Info loss
- + difficult to select features ^{PCA} own



Steps

1) Standardization

- + ~~small~~ similar range - all variables equal contribn

2) + Covariance matrix

- + understand variance

- + find relat low Fctns

3) Eigen values & eigen vectors

- + CM sc compute \rightarrow Principal components

Properties -

- + linear combination of variables of original ds.
- + PCs are orthogonal to each other
- + O. correlation b/w pair of features
- + PC1 $\perp \perp$ PCN, Imp \downarrow

MT unit 2.

+ Linear Regression

+ statistical method used to model relationship between dependent variable y and one or more independent variable x .

+ It seeks to find best fitting straight line

$$y = B_0 + B_1 x + \epsilon$$

B_0 - intercept B_1 = slope.

ϵ = error term - difference between error term and predicted values.

+ Assumptions of linear regression

+ linearity

+ Independence - observations are independent of each other

+ Homoscedasticity - variance of residuals - constant difference between observed and predicted value - across all independent variable.

+ Normality - residuals are normally distributed

+ No multicollinearity - Independent variables are not highly correlated with each other.

+ Bidimensional example.

+ Two variables involved (dependent and independent)

+ e.g. study hours and marks gained.



+ Linear regression

$$y = \alpha + \beta x_i$$

$$y_i = B_0 + B_1 x_i + \epsilon_i \quad i = 1, 2, \dots, n.$$

+ Assumptions.

+ y_i depends only on x_i and variation is y_i is random.

+ Variance of ϵ or y does not depend on x_i .

+ $E(\epsilon_i) = 0$ for all i

$$\text{Var}(\epsilon_i) = \sigma^2 \quad \text{for all } i$$

$$\text{corr. Var}(y) = \sigma^2$$

+ $\text{cov}(\epsilon_i, \epsilon_j) = 0$, $\text{cov}(y_i, y_j) = 0$.

+ ϵ variables or y variables are unrelated to each other.

$$+ \quad b_1 = \frac{\sum_{i=1}^n (x_i y_i - \bar{x} \bar{y})}{\sum_{i=1}^n (x_i^2 - \bar{x}^2)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

Hours studied (x)	Test scores (y)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$
2	65	-2	-10	4
3	70	-1	-5	1
4	75	0	0	0
5	80	+1	5	1
6	85	+2	10	4

$$\bar{x} = \frac{2+3+4+5+6}{5} = 4$$

$$\bar{y} = \frac{3+5}{5} = 75$$

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{-2 \times -10 + -1 \times -5 + 0 \times 0 + 1 \times 5 + 2 \times 10}{4+1+1+4}.$$

$$= \frac{20+5+5+20}{10}$$

$$= \frac{50}{10} = 5$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_0 = 75 - 5 \times 4 = 55$$

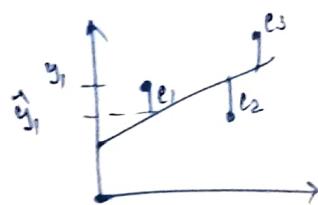
$$\therefore y = 55 + 5x$$

+ Predicted and residuals.

Predicted - find using linear regression

residual - deviation of observed from predicted value.

$$e_i = y_i - \hat{y}_i$$



+ sum of square error (SSE)

$$SSE = \sum (e_i)^2 = \sum (y_i - \hat{y}_i)^2$$

predicted
observed

low SSE \rightarrow better fit model.

explains unexplained variation in y.

+ SST (Total sum of squares)

$$SST = \sum (y_i - \bar{y})^2$$

total variability in y.

+ R^2 = coefficient of determination.

$$R^2 = 1 - \frac{SSE}{SST}$$

$$+ (y_i - \bar{y}_{\text{mean}})^2 = (y_i - \hat{y}_{\text{predi}})^2 + (\hat{y}_{\text{predi}} - \bar{y}_{\text{mean}})^2$$

$$+ TSS = RSS + ESS$$

$$\text{but } R^2 = 1 - \frac{RSS}{TSS}$$

$$R^2 = \frac{ESS}{TSS}$$

+ $R^2 \rightarrow$ determines proportion of variance in dependent variable that can be explained by independent variable.

$$R^2 = \frac{\text{Variance explained by model}}{\text{Total variance}} = \frac{SSR}{SST}$$

$R^2 \rightarrow$ coefficient of determination.

+ Adjusted $R^2 \rightarrow$ modified $R^2 \rightarrow$ adjusts to number of independent variable.

+ only those independent variable which helps in explaining proportion of variation in dependent variable

$$R^2_{\text{adjusted}} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

R^2 = sample R square.

P = number of predictors. (no of independent variables).

N = Total sample size.

can be negative.

\rightarrow worse model.

+ Regularization.

+ prevent overfitting

+ penalty term

+ improve generalization performance of model.

+ Ridge regression

+ adds penalty to squared magnitude of coefficient.

+ large coefficient \rightarrow heavy penalty.

+ shrinking the coefficient.

$$y = 0.9 + 15x \rightarrow y = 0.9 + 2.09x.$$

+ reduces variance of predictions $\rightarrow \sum (y_i - \hat{y})^2$.

$$\text{Ridge F} = \text{Loss} + \alpha \|w\|^2 = \text{RSS} + \alpha \|w\|^2.$$

$$\|w\|^2 = w_1^2 + w_2^2 + \dots + w_n^2.$$

$\alpha \uparrow \rightarrow$ more shrinking.



+ predictor variables \rightarrow highly correlated
 \rightarrow multicollinearity

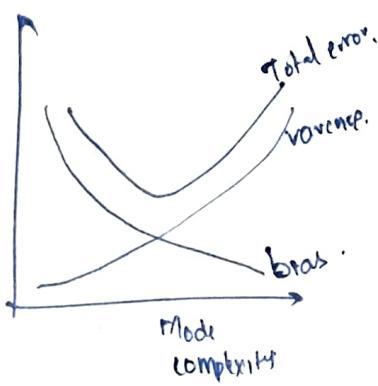


takes care \leftarrow high variance
of that predictor.

variance \rightarrow change in F if we train on different set.

bias \rightarrow error introduced

Inevitable error \rightarrow cannot be reduced.
that



+ Lasso Regression :

+ adds penalty to absolute value .

+ shrinks coefficient to 0

+ feature selection and removing irrelevant features

$$\text{Lasso R} = \text{Loss} + \alpha \|\mathbf{w}\| \text{ penalty}$$

$$\|\mathbf{w}\| = |w_1| + |w_2| \dots |w_n|$$

+ Elastic Net Regression :

$$\text{+ Elastic Net Regression} = \text{loss} + \alpha_1 \|\mathbf{w}\|_1^2 + \alpha_2 \|\mathbf{w}\|_2$$

+ for feature selection

+ for multicollinearity .

+ emerged because \rightarrow lasso is too dependent on data .

$$f_{\text{net}}(\mathbf{B}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i \mathbf{B})^2 + \lambda \left(\frac{1-\alpha}{2} \sum_{j=1}^m B_j^2 + \alpha \sum_{j=1}^m |B_j| \right).$$

$$\alpha = \begin{array}{l} \text{mixing parameter} \\ \text{ridge } (\alpha=0) \\ \text{lasso } (\alpha=1) \end{array}$$

+ Robust Regression with Random Sample Consensus (RANSAC)

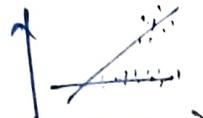
+ Outlier containing dataset .

+ outliers \rightarrow biased result .

+ RANSAC \rightarrow identifies outliers \rightarrow ignores them

+ RANSAC \rightarrow perform iterations \rightarrow until inliers remain

+ noise ignored



+ threshold \rightarrow defined \rightarrow for outlier identification .

f. Gradient descent.

+ descending slope to reach lowest point

+ objective \rightarrow find x such that $\rightarrow y = \text{minimum}$

+ iterative update parameters \rightarrow minimize cost function.

+ minimize sum of square residuals \rightarrow errors.

+ starts with random point \rightarrow descend the slope \rightarrow reach required position.

+ Algorithm

+ find slope \rightarrow with respect to each function

+ pick a random initial value.

+ differentiate y with respect to x .

+ update the gradient function.

+ calculate step size for each feature

$$\text{Step size} = \text{gradient} * \text{learning rate}.$$

+ New parameter = old params - step size.

+ repeat it until gradient is almost 0.

+ stochastic gradient descent

+ for random small part of observation
rather than all values.

+ convergence \rightarrow reached reaching a threshold of iterations

\rightarrow thus reaching threshold of
finding magnitude of gradient.

+ Hyper Parameters

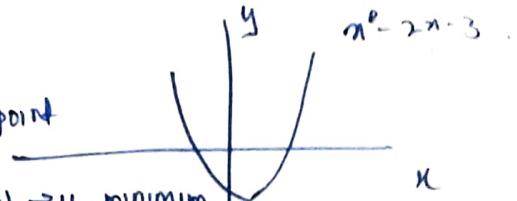
+ cannot be determined from data

+ not directly learned from training process

+ set prior to training

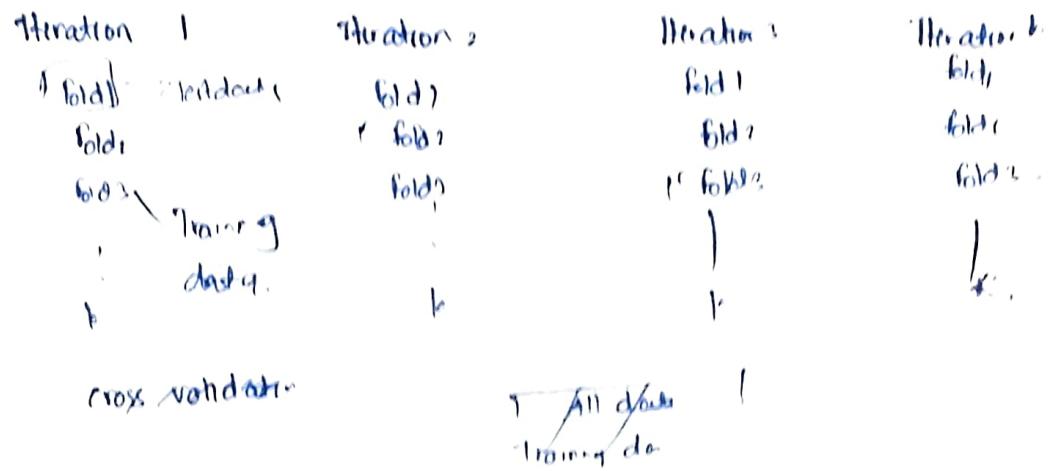
+ influences \rightarrow complexity, optimization algorithm used etc.

+ e.g. learning rate in gradient descent, number of hidden layers,
Regularization parameter.



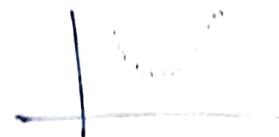
+ Grid Search.

- + specify grid of hyperparameters
- + evaluate each
- + finds best performance.
- + Repeated K-fold → → Training data into K folds.
 - * for each hyper param., train on $K-1$ fold.
 - + select best parameters.
- + K models fit on $\frac{K-1}{K}$ data (training split).
- + evaluated on $\frac{1}{K}$ of data (called test split).



+ Polynomial regression.

+ linear model in non-linearly dataset.



+ $y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p$

(\Rightarrow polynomial functions)

$$\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_p$$

+ Covariance matrix.

relationship between variables

and their variability.

$$\begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) \end{bmatrix} = \begin{bmatrix} \text{Var}(x_1) + \frac{\sum (x_i - \bar{x})^2}{n-1} & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{Var}(x_2) + \frac{\sum (x_i - \bar{x})^2}{n-1} \end{bmatrix} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Properties:

- m-symmetric
- m n-square

+ Examples:

Ridge \rightarrow Host price \rightarrow Squar feet/1000, logarithm function.

Tanso \rightarrow DR \rightarrow Rent price \rightarrow Per price

Ridge \rightarrow Farm marks

Gradient descent \rightarrow min prediction

Ore rock \rightarrow Mineral \rightarrow crop yield.

Q.	Height	Weight	Σxy	Σx^2
	x	y		
	50	15	750	2500
	20	10	200	400
	40	15	600	1600
	30	12	360	900
	50	15	750	2500
	190	67	2660	7900

$$y = b_n x + a.$$

$$b = \frac{n(\Sigma xy) - \bar{x}\bar{y}}{n\Sigma x^2 - (\bar{x})^2}$$

$$= \frac{5(2660) - 190 \times 67}{5 \times 7900 - (190)^2}.$$

$$= \frac{570}{3400} = 0.167.$$

$$a = \bar{y} - b \bar{x}$$

$$a = 13.4 - 0.167 \times 38.$$

$$a = 7.054.$$

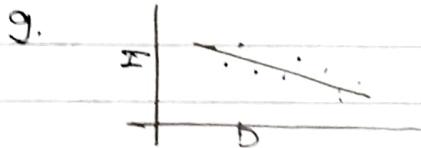
$$y = 7.054 + 0.167x.$$

$$\begin{array}{r}
 y \quad y_{\text{hat}} \cdot (y_{\text{hat}} - y)^2 \quad (y - \bar{y})^2 \\
 \hline
 5 \quad 6 \quad +1 \quad 1 \cdot 4 \cdot 4 \\
 2 \quad 1 \quad +1 \quad 3 \cdot 2 \cdot 4 \\
 4 \quad 5 \quad 1 \quad 0 \cdot 04 \\
 3 \quad 4 \quad 1 \quad 0 \cdot 64 \\
 5 \quad 7 \quad \frac{4}{8} \quad \frac{1 \cdot 44}{6 \cdot 8} \\
 \hline
 \frac{19}{5} \cdot 1 - \frac{\text{SSE}}{\text{SST}} = 1 - \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}
 \end{array}$$

$$R^2 = 1 - \frac{8}{6 \cdot 8} \\
 = 0 \cdot 176$$

Linear Regression

1. A linear approach that models the relationship b/w a DV & one or more IVs.
2. solved regression problems
3. estimates the DV when there is a change in the IV.
4. Output (DV) continuous in nature
5. Linear relationship
6. Uses a straight line
7. 1-D & 1-I \rightarrow Simple Linear Reg.
1-D & >1 -I \rightarrow Multiple Linear Reg.
8. $y = \beta_0 + \beta_1 x + \epsilon$
 $\downarrow \quad \downarrow \quad \rightarrow$ error term
 Intercept slope

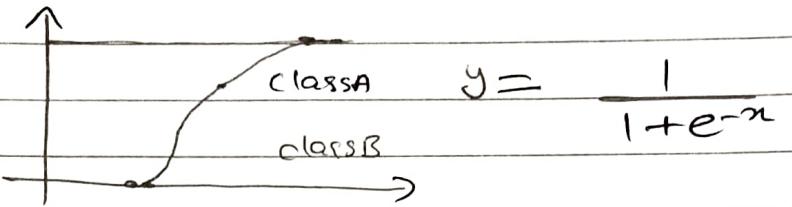


10. Ex. a. Predicting the GDP of a country.
- b. Predicting product price
- c. Score prediction.

Logistic Regression

1. A statistical model that predicts the probability of an outcome that can only have two values
2. solve classification problems (binary class.)
3. calculates the possibility of an event occurring.
4. Output (DV) is discrete
5. Goal is to find best fitting model for IV & DV relationship.

6. Uses S curve or sigmoid function
 7. also called Logit Regression
 8. EV can be continuous or binary \rightarrow (DV)
 9. Output \Rightarrow 1 (true, success)
 10.



11. Ex. a. predicting whether an email is spam or not
 b. predicting whether customer will take loan or not
 c. credit card transactions fraud or not.

12. $\log \left(\frac{y}{1-y} \right) = b_0 + b_1 x_1 + \dots + b_n x_n$

Confusion Matrix

		Predicted		
		NO	YES	
Actual	NO	50	10	60
	YES	[TN]	[FP]	
	YES	5	100	105
		[FN]	[TP]	
		65	110	

$$\text{Accuracy} = \frac{TP + TN}{\text{Total}} = \frac{100 + 50}{165} = 0.91$$

$$\text{Error rate} = 1 - \text{accuracy} = 1 - 0.91 = 0.09$$

$$\text{Precision} = \frac{TP}{\text{Predicted Yes}} = \frac{100}{110} = 0.91$$

$$\text{Recall} = \frac{\text{TP}}{\text{actual yes}} = \frac{100}{105} = 0.95$$

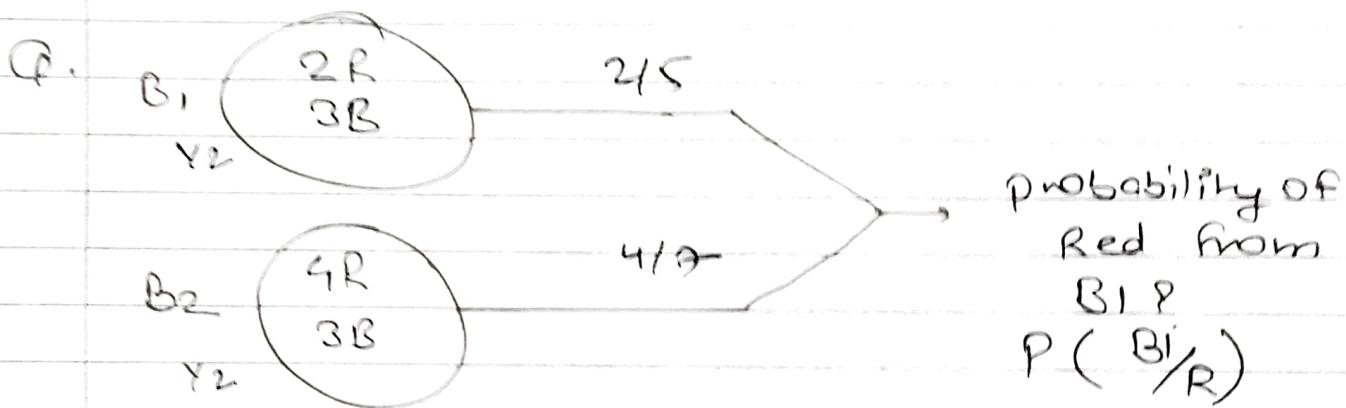
\rightarrow Independent

Naive Bayes Classifier

1. supervised learning

2. classification

3. Based on Bayes theorem



$$P(Y/x) = \frac{P(x/Y) * P(Y)}{P(x)}$$

$$Y = B1$$

$$x = R$$

$$P(R/B1) = 2/5$$

$$P(YB1) = Y_2$$

$$P(R) = Y_2 * 2/5 + Y_2 * 4/7$$

$$P(B1/R) = \frac{Y_2 * 2/5}{Y_2 * 2/5 + Y_2 * 4/7}$$

$$P(Y/X) = \frac{P(X/Y) * P(Y)}{P(X)}$$

$$P(Y/x_1, x_2, x_3 \dots x_n)$$

$$= \frac{P(x_1/Y) * P(x_2/Y) * \dots * P(x_n/Y) * P(Y)}{P(x_1) * P(x_2) * \dots * P(x_n)}$$

Q.

1] Prior Probability

$$P(\text{Fever} = \text{Yes}) = 7/10$$

$$P(\text{Fever} = \text{No}) = 3/10$$

2] Conditional Probability

Fever

	Yes	No
Covid	4/7	2/3
Flu	3/7	2/3

$$P(\text{Yes} / \text{Flu, Covid})$$

$$= \frac{P(\text{Flu}/\text{Yes}) * P(\text{Covid}/\text{Yes}) * P(\text{Yes})}{P(\text{Covid}) * P(\text{Flu})}$$

$$P(\text{No} / \text{Flu, Covid})$$

$$= \frac{P(\text{Flu}/\text{No}) * P(\text{Covid}/\text{No}) * P(\text{No})}{P(\text{Covid}) * P(\text{Flu})}$$

$$P(\text{Yes} / \text{Covid, Flu}) = \frac{3/7 * 4/7 * 7/10}{12/70}$$

$$P(\text{No} / \text{Covid, Flu}) = \frac{2/3 * 2/3 * 3/10}{4/30}$$

Q. $X = \text{Chills} = Y$, Runny Nose = N, Headache = Mild, Fever = Y, Flu = ?

3) Prior Probability (flu)

$P(\text{Flu} = \text{Yes})$	$= \frac{5}{8}$
$P(\text{Flu} = \text{No})$	$= \frac{3}{8}$

$$P(Y/\text{rest}) = P(\text{rest}/Y) \times P(Y)$$

2) Conditional Probability

Chills	Y	N
Flu (Y)	5/5	1/3
Flu (N)	1/5	1/3

	Flu	
	Y	N
✓ Chills (Y)	3/5	1/3
Chills (N)	2/5	2/3

	Y	N
RN (Y)	4/5	1/3
RN (N)	1/5	2/3

	Y	N
H (S)	2/5	1/3
✓ (M)	2/5	1/3
(N)	1/5	1/3

	Y	N
✓ Fever (Y)	3/5	0/3
(N)	0/5	3/3

$$\text{Yes} \Rightarrow (3/5 \times 1/5 \times 2/8 \times 1) \times 5/8 = \frac{6}{1200} = 0.005$$

$$\text{No} \Rightarrow 1/3 \times 0 \Rightarrow 0$$

$$\text{Normalized} \Rightarrow \frac{0.005}{0.005+0} = 1$$

① $P(H=\text{Yes}) = 3/8$
 $P(H=\text{No}) = 5/8$

	Y	Happy	N
Weather			

Good	$1/3$	$3/5$
Bad	$2/3$	$2/5$

② Study

Pass	$3/3$	$1/5$
Fail	$0/3$	$4/5$

③ Neighbor

Home	$2/3$	$2/5$
Out	$1/3$	$3/5$

a) $H = \text{Good}$, $S = \text{PASS}$, $N = \text{OUT}$

$$\text{Yes} \Rightarrow \frac{1}{3} \times \frac{3}{3} \times \frac{1}{3} \times \frac{3}{8} = \frac{1}{24} = 0.0416$$

$$\checkmark \text{No} \Rightarrow \frac{2}{3} \times \frac{1}{3} \times \frac{2}{5} \times \frac{5}{8} = \frac{1}{200} = 0.005$$

$$\text{Yes} \Rightarrow \frac{0.0416}{0.0416 + 0.005} = 0.898$$

$$\checkmark \text{No} \Rightarrow \frac{0.005}{0.0416 + 0.005} = 0.52$$

b) $\checkmark \frac{1}{3} \times \frac{3}{3} \times \frac{2}{5} \times \frac{3}{8} = \frac{1}{6} = 0.16$

$$\frac{2}{3} \times \frac{1}{3} \times \frac{2}{5} \times \frac{5}{8} = \frac{1}{200} = 0.005$$

$$\checkmark \text{Yes} \Rightarrow \frac{0.16}{0.16 + 0.005} = 0.89$$

$$\text{No} \Rightarrow 0.11$$

ENN (K-Nearest Neighbor)

NN classifiers are defined by their characteristic of classifying unlabeled ex. by assigning them the class of most similar labeled ex.

+ can apply (facial)

+ movie enjoys

+ identify patterns in genetic data

Euclidean

$$\text{dist}(x, y) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

1) balance b/w overfitting & underfitting training data \rightarrow bias-variance

2) $k \uparrow \rightarrow$ noisy data impact \downarrow but small pattern \rightarrow imp & ignore



3) no. of records of difficulty

4) 3 to 10

5) $\sqrt{\text{no. of train. ex.}}$

6) 15 \rightarrow 4

163, 63

Q.	162	61	$\sqrt{5}$	S
	165	62	$\sqrt{5}$	S
	167	63	4	S
	170	66	$\sqrt{58} = 7.61$	
	165	63	$\sqrt{2}$	S
	167	64	$\sqrt{17}$	
	170	65	$\sqrt{53}$	
	173	68	x	
	165	23		
	163	63	2	S

II Metrics to evaluate the performance of the classification model.

- 1) Accuracy →
- 2) Precision
- 3) Recall
- 4) F1-score
- 5) AUC-ROC

Confusion matrix →

A confusion matrix is defined as the table that is often used to describe the performance of a classifier model on a set of the test data for which the true values are known.

	AC Values	
PV	P(1)	N(0)
P(1)	TP	FP
N(0)	FN	TN

TP → model predicted +ve & true

TN → ——————+—————→ & ——————+—————→

Type1 FP → ——————+—————→ +ve & False

Type2 FN → ——————+—————→ -ve & False

	A-V	
	P	N
PV	P	TP
N	FN	TN

① Accuracy

- 1) measures of how often classifier predicts correctly.
- 2) ratio of correct P / Total P.
- 3) DS \rightarrow Unbalanced \Rightarrow

$$A = \frac{TP + TN}{TP + TN + FP + FN}$$

② Precision

- 1) explains how many of the correctly predicted cases turn out to be positive.
- 2) case $FP > FN$
- 3) $TP / TP + FP$

$$P = \frac{TP}{TP + FP}$$

③ Recall (Sensitivity)

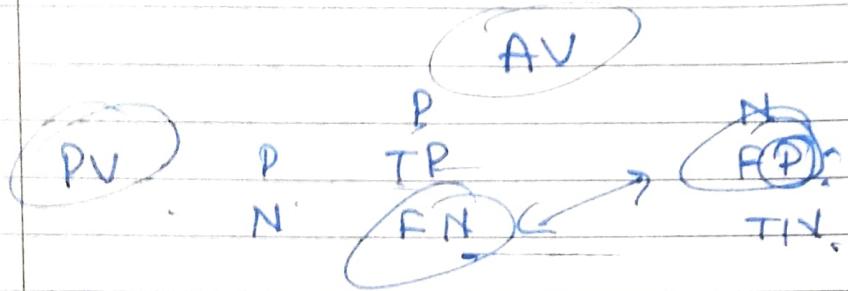
- 1) explains how many of the actual positive cases the model is able to predict correctly.
- 2) $FN > FP$
- 3) medical false X, true matter

$$R = \frac{TP}{TP + FN}$$

④ F1 Score

Combined idea, max $P=R$, HM of P & R
 $FP = FN$, $TN >$,

$$F1 = \frac{2 \times P \times R}{P + R}$$



B) $A = \frac{TN + TP}{All}$

measures how model predict correctly.

2) $P = \frac{TP}{TP + FP}$

measures how many OF correctly predicted are positive. $FP > FN$

3) $R = \frac{TP}{TP + FN}$ PPL FN

measures how many of the actual true cases model predicts correctly

4) F1 Score $2 \times \frac{P \times R}{P + R}$ $FP = FN$
 $TN > P$

5) AOC-ROC curve (receiver operating characteristic)

Showing the performance of a classi. model at all classi. threshold.

True +ve R $\Rightarrow \frac{TP}{TP + FN}$

FPR $\Rightarrow \frac{FP}{FP + TN}$

Q.

AV	PV		$TP = 3$	$FN = 8$	$TN = 87$
	P	N			

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

Actual Value

PV	P	N
P	3	2
N	8	87

AV	P	N	$TP = 3$	$FN = 8$	$TN = 87$
P	3	2			

$$\text{TPR} \Rightarrow \frac{TP}{TP + FN} = \frac{3}{3 + 8} = \frac{3}{11}$$

$$\text{FPR} \Rightarrow \frac{FP}{FP + TN} = \frac{2}{2 + 89} = \frac{2}{91}$$

Q.

PV	AV		$TP = 12$	$FN = 3$	$FP = 4$	$TN = 96$
	P	N				

$$\text{TPR} = \frac{TP}{TP + FN} = \frac{12}{15} = \frac{4}{5}$$

$$\text{FPR} = \frac{FP}{FP + TN} = \frac{4}{20} = \frac{1}{5}$$

Cross Validation

1. used to assess how well our ML models perform on unseen data.
2. resampling tech. DS → 2 parts

Test & Train

3. Train D → train model

Test D → unseen test & predict

4. TD → PV AT → ✓

Techniques

1. Hold out method
2. Leave one out cross-validation
3. k-fold cross validation
4. stratified k-fold cross validation

① Hold Out

1. simplest to evaluate a classifier
2. DS
Train ↑ Test
3. classifier performs function of assigning data items in a given collection to a target category
4. PS
[Tr.S | rs]

$$CV(k) = \frac{1}{k} \sum_{i=1}^k MSE_i$$

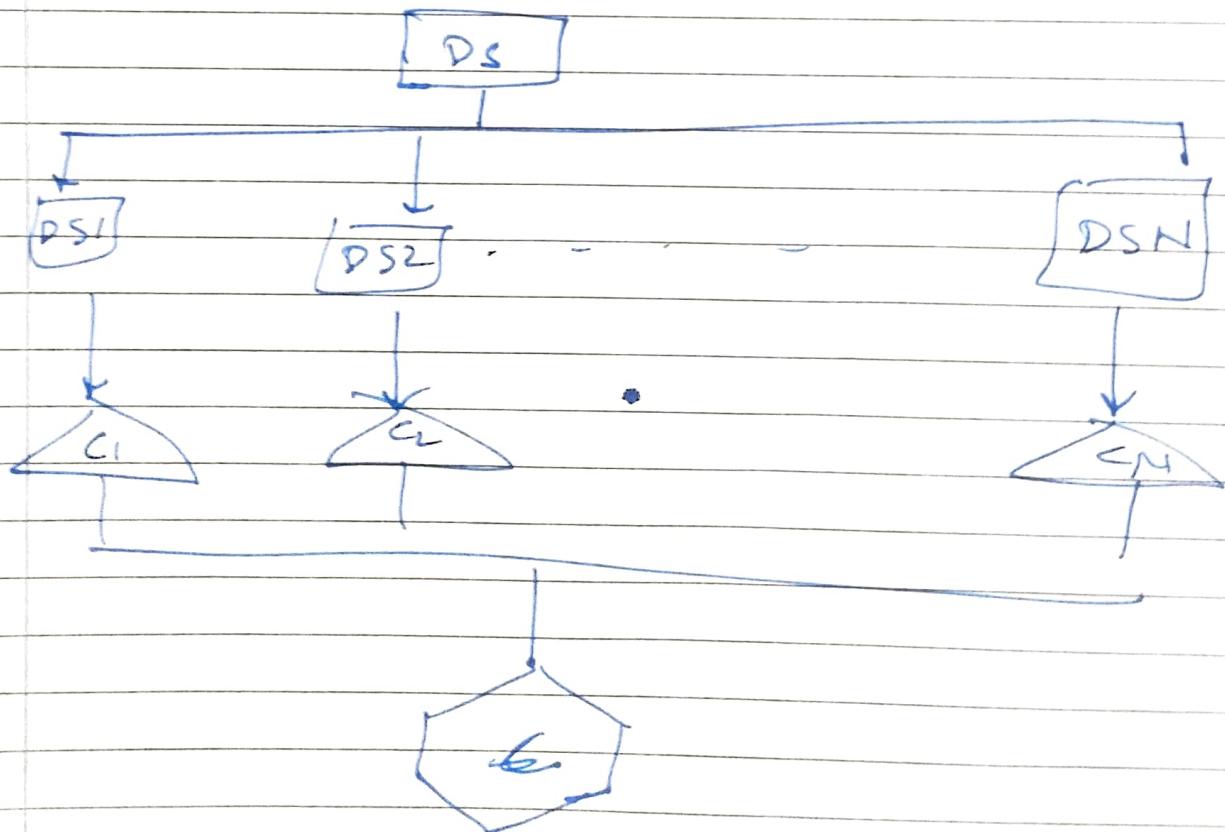
Weighted KNN

$$\hat{f}(n_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v, f(n_i))$$

$$0.45 d(m, n) + 0.15$$

Ensemble Learning (RF) (GB)

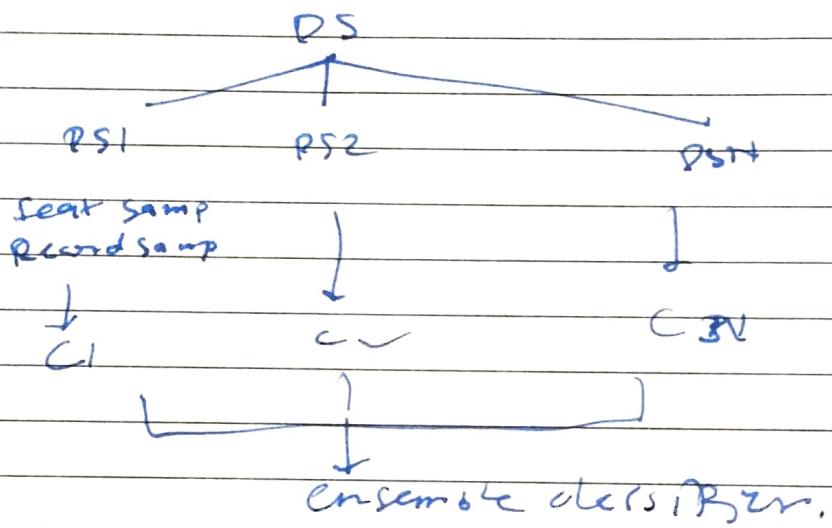
1. Improve ML results by combining several models.
2. better predictive performance compare to single model.
3. learn a set of classifiers (experts) and allow them to vote.



Types

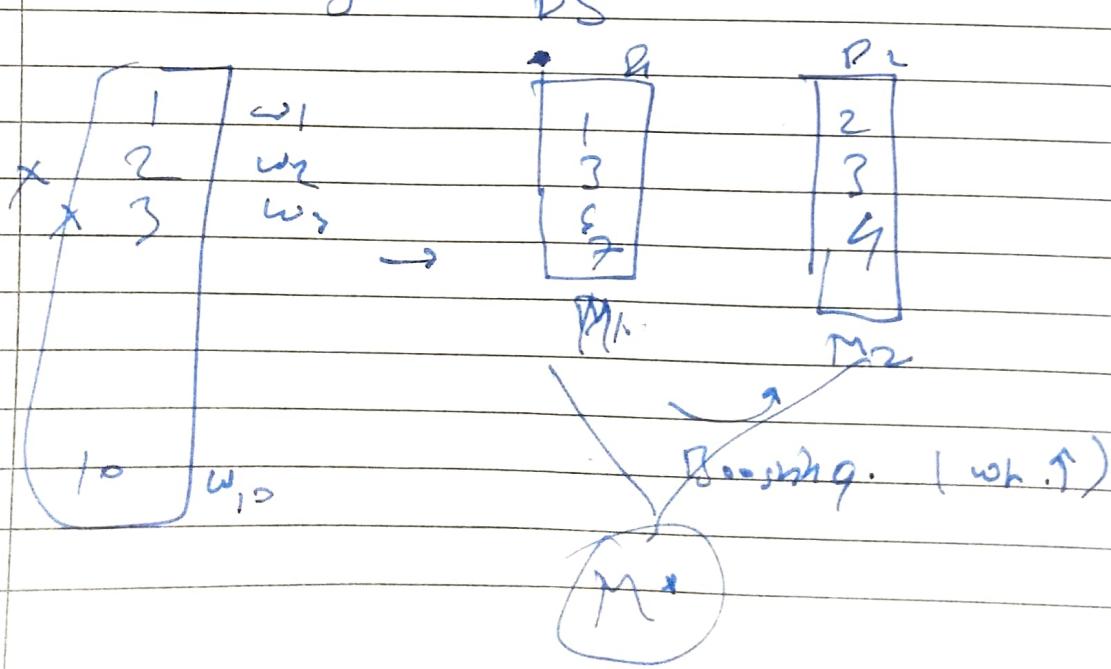
1. Bootstrap aggregating
2. Stacking
3. ~~the~~ Voting
4. Boosting

1) Bagging / Bootstrap Aggregating



classification \rightarrow voting

2) Boosting



Ensemble Learning

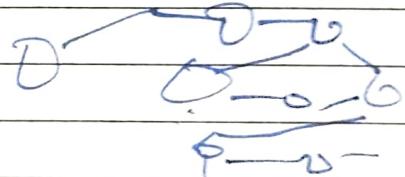
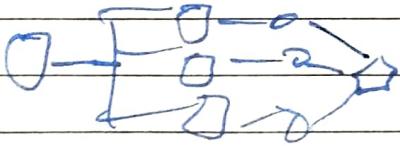
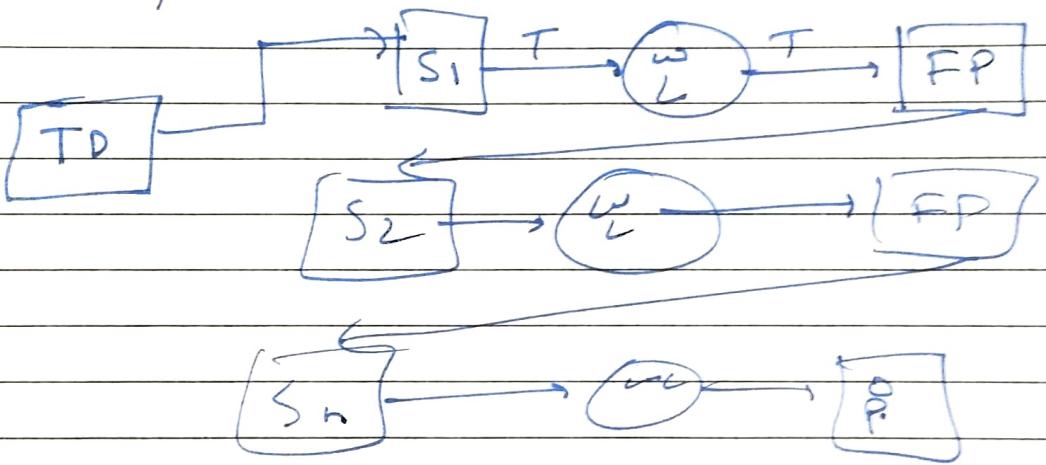
1. Type of ML uses multiple ML models to predict.
2. ex. RF, Gradient Boosting
3. Improve ML results by combining several models
4. Better result compare 1 m.
5. Individual model \rightarrow weak learner coz \uparrow variance or bias.
6. HB \rightarrow not learning from data
7. LV \rightarrow learning too well & varies for each data
8. Generalize x,
9. bias-variance trade off
Underfit \rightarrow HBLV
Overfit \rightarrow LBHV
ensemble balance parts has BVTQ
10. final strong model

① Bagging

1. Reduce the variance from predictions by generating additional data for training from DS
2. BA
3. Bootstrapping
resampling subsets of data with
4. Aggregating
5. Max voting

② Boosting

1. Sequentially training weak learners
2. Combining weak learners with HB.
3. Sample of data taken
4. Sample $\rightarrow 1m \rightarrow \text{Pre}$
- 5: CP ? IP
6. IP \rightarrow used again 2M
7. at each step error improve
8. weighted averaging
9. $M_1 - w_1, M_2 - w_2$ by predictive power.



Decision Tree (Sup)

1. Predict output of target variable
2. graphical presentation
3. all prob. solutions based on decisions on certain conditions.

Root Node

Decision Nodes

Leaf Node / Terminal

Subtree

Pruning.

Q. A (outlook)

values(outlook) = S, D, R

Entropy

1. Entropy is information theory metric that measures Impurity or uncertainty in a gr. of observations.

2. DT to help to split data.

3. N class

$$S = - \sum_{i=1}^N p_i \log_2 p_i$$

4. p_i probability of randomly selecting an example in class

5. categorical data.

6. $S \rightarrow 0$ if homogeneity

$$S = [S+, S-]$$

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$S_{\text{sunny}} \Rightarrow [2+, 3-] = 0.571$$

$$S_{\text{overcast}} \Rightarrow [4+, 0-] = 0$$

$$S_{\text{rain}} [3+, 2-] = 0.971$$

Information Gain

1. measures how much info. a feature provides about a class
2. order of attribute in the nodes of DT
3. $\text{Gain} = E_{\text{parent}} - E_{\text{children}}$

$$\text{Gain}(S, \text{outlook}) = \text{Entropy}(S)$$

$$= S \sum_{v \in \{S, O, R\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{outlook}) =$$

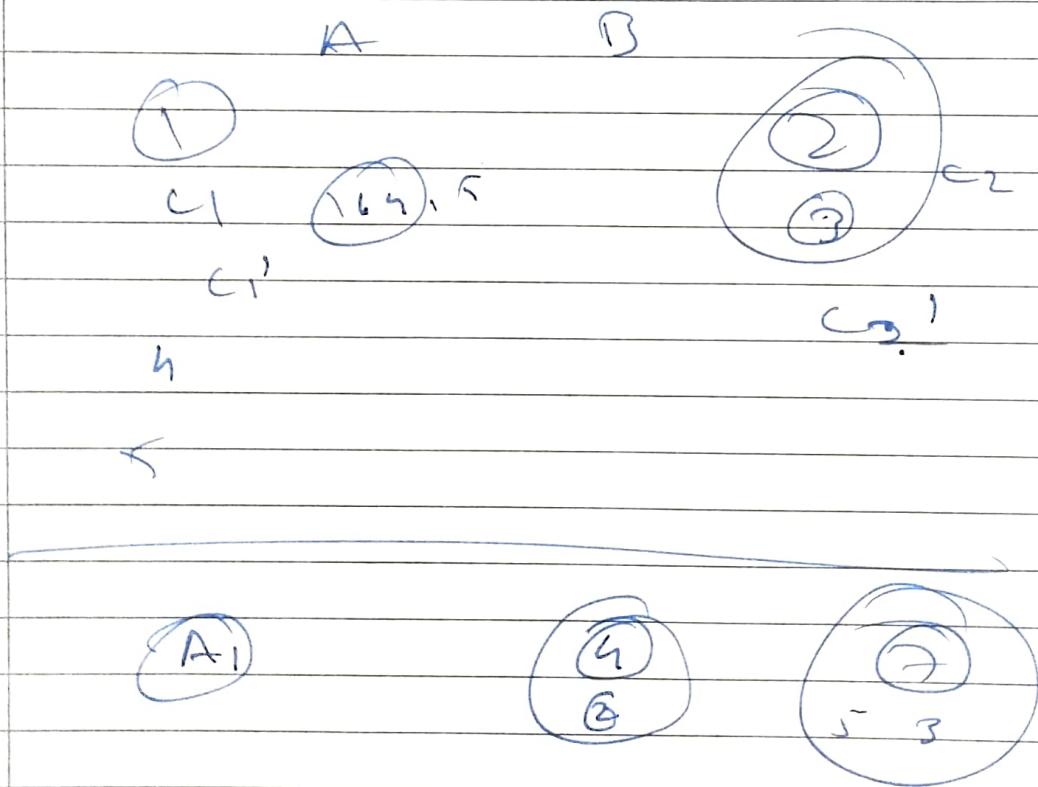
$$E(S) = \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.571$$

$$= 0.2464$$

RF

1. class & Reg
2. T do
3. AT
4. MIS ✓
5. H-Regu.
6. size ↑
7. B. Box

K-means (JS)



from SKL-metrics import accuracy_score

$$y_{true} = [0, 1, 1]$$

$$y_{pred} = [0, 0, 1]$$

$$acc = acc_score(y_{true}, y_{pred})$$