

ML

↳ Supervised learning (both input and output) → regression
 data given and finding relationship in between them → classification

↳ Unsupervised learning (creating groups just from input data) → clustering

Semi-supervised ML → Reg + classification + clustering

{Combo of supervised and unsupervised}

↳ Linear Regression ($y = mx + c$)

↳ When to apply linear regression
 Error term must showcase constant variance (homoscedasticity)
 No multicollinearity

→ y vs $x \rightarrow$ linear tendency
 mean of the residual should be zero.
 error terms are not supposed to be correlated
~~error terms are supposed to be uncorrelated~~
 x and residuals must be uncorrelated
 error terms are supposed to be normal dist

$$\text{↳ residual} = r = y - \hat{y} = y - (mx + c)$$

$$\text{↳ } r^2 = (y - \hat{y})^2$$

$$\text{Sum of squares} = \sum_{i=1}^m r_i^2 = \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$$\frac{dC(\sum_{i=1}^m r_i^2)}{dc} = 0 = \frac{d(\sum_{i=1}^m r_i^2)}{dc}$$

{ \hat{y} : predicted; y : observed}

must be made minimum for best fit

$$0 = m_{\text{new}} = m_{\text{old}} - \eta \frac{\partial}{\partial m} \left(\sum_{i=1}^m (y_i - \hat{y}_i)^2 \right)$$

$$0 = C_{\text{new}} = C_{\text{old}} - \eta \frac{\partial}{\partial c} \left[\sum_{i=1}^m (y_i - \hat{y}_i) \right]$$

slope = $m_{\text{old}} - \eta \nabla E_m$
 not grad w.r.t. c

gradient descent method

η controls the rate of convergence

{exogeneity: x and residuals must be uncorrelated}
 learning rate

{Multi-collinearity: no collinearity expected from two or more independent variables (x_1 and x_2 not correlated)}

↳ Accuracy of the model $R^2 = \frac{1 - RSS}{TSS}$

(\rightarrow RSS: residual summation of squares)

$$\sum_{i=1}^m (y_i - \hat{y}_i)^2$$

(\rightarrow TSS: total summation of squares)

$$\sum_{i=1}^m (y_i - \bar{y})^2$$

because it will otherwise start learning of relations among feature-feature and will affect feature-label correlations

homoscedasticity

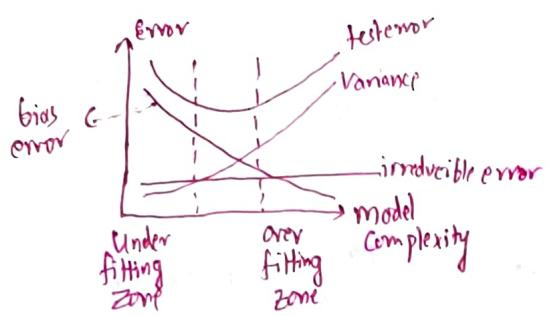
all random variables have the same finite variances

y x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9 x_{10}

homoscedastic heteroscedastic

Bias-Variance trade-off

↳ test error (error in test data) dependent on



bias error (error from wrong model assumption)

Variance (error from sensitivity to small fluctuations in train data)

irreducible error (inherent noise in problem itself)

(2)

↳ Comes from oversimplification of the model (taking less parameters)

↳ Comes from overcomplexity of the model (taking more parameters)

↳ Adjusted R-square → Corrected gof(model accuracy) measure for linear models.

% of variance in the data (dependent variable) explained by independent variables(s) in regression models

adjusted for the no. of attributes taken (minimizes the apparent effect of overfitting)

$$R_{adj}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

m → no. of samples
p → no. of features(attributes)

↳ Selecting features for doing ML (linear regression) → $\{y = m_1x_1 + m_2x_2 + \dots + c\}$

↳ Check for multicollinearity: if two features multicollinear \Rightarrow take only one of them (check correlation values) Otherwise overfitting

↳ if to take a particular feature for regression

↳ check p-values $\Rightarrow H_0: \text{feature } x_i \text{ does not affect } y \Rightarrow y = mx \text{ with } m=0$

TS: test statistics
ts: observed test statistics from the data

↳ if pvalue ($P(TS \leq ts | H_0)$, $P(TS \geq ts)$, $P(TS \geq |ts|)$) $< \alpha$ ↓ significance value

↳ lower tail test upper tail test 2tail test

↳ dependent on x_i Cdf(ts) $1 - Cdf(|ts|)$

↳ different for different test, e.g. z-statistics for z-test

↳ calculated using the various statistics tables (e.g. z-table)

* Note → Student's t-distribution (t -test)
used in statsmodels formula, api (Python)

↳ Variance Inflation Factor (VIF) to check multicollinearity

$VIF = 1 / (\text{not correlated})$

= 1-5 (moderately correlated)

= >5 (high correlation)

$$VIF = \frac{1}{1 - R_i^2}$$

obtained by regressing one feature (x_i) with another (x_j)

Dummy variables trap

→ Dummy variables → Creating additional attributes with values (0 or 1) for each class of a categorical value

③
one-hot encoding

→ Dummy variable trap → high correlation among variables

→ need to take only independent dummy variables!!
drop one of the dummy variable

→ One dummy variable can be predicted from other dummy variables (due to one-hot encoding) and if all the dummy variables included ⇒ dummy variable trap

→ Regularization → technique to discourage learning a more complex or flexible model, reducing the risk of overfitting → trying to incorporate noise hence less predictability for test dataset.

→ Loss function for models → $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$

$$+ \text{RSS} = \sum_{i=1}^n \left[Y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right]^2$$

→ Ridge regression (L₂ norm) → Loss function = RSS + $\lambda \sum_{j=1}^p \beta_j^2$ → minimize this func

→ λ : tuning parameter to decide how much to penalize the abrupt change in β_j

→ Lasso regression (L₁ norm) → Loss function = RSS + $\lambda \sum |\beta_j|$



→ Ellipses are the equivalent RSS values
→ in L₂ RSS and $\sum \beta_j^2$'s will never touch at axes
→ no variable can ever be reduced to zero
→ no chance of dimensionality reduction

→ But in L₁ RSS and $\sum \beta_j^2$'s can touch at axes ⇒ dimensionality reduction

→ Elastic net regression : linear combination of both L₁ and L₂

→ Standardization (using `sklearn.preprocessing.StandardScaler`)

→ to bring all the attributes to the same scale

→ This helps the model to have more robust relations.

$$(Z = \frac{X - \mu}{\sigma})$$

Used as a method to find the best possible tuning parameter in regularization (or elsewhere) (4)

Cross-validation

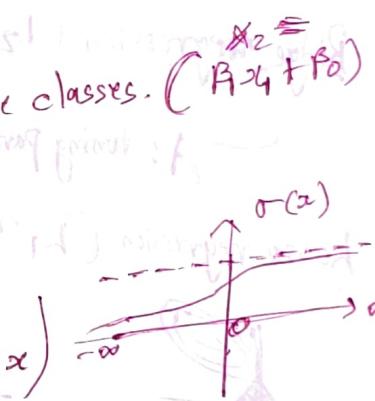
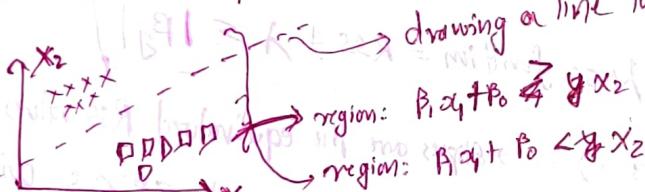
a model validation technique using resampling

Used for tuning the hyperparameters everywhere (like λ)

- $CV=5 \Rightarrow$ divides the data-set into 5 chunks and performs iteratively testing and training by taking different portions of the data \rightarrow one-round of CV
- multiple rounds of ~~itera~~ CV done taking into account the number of iterations given.

Logistic Regression

- is often used for classification and predictive analysis.
- estimates the probability of an event occurring, such as voted/didn't vote
- logistic regression not a classification algorithm on its own. Only classification in combination with a decision rule. It is regression because it estimates the probability of class membership
- estimates the probability of class membership drawing a line to divide the classes. ($\beta_1x_1 + \beta_0 = 0$)



Probability model \Rightarrow Sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$

Binary Classification

$$P(x) = \frac{1}{1+exp[-(\beta_1x_1 + \beta_0)]}$$

Can be rewritten as

$$P(x) = \frac{1}{1+exp[-(\beta_1x_1 + \beta_0)]} \quad \text{category-1}$$

$$P(x) = \frac{e^{\beta_1x_1 + \beta_0}}{1 + e^{\beta_1x_1 + \beta_0}} \quad (1)$$

$$P(\bar{x}) = 1 - P(x) = \frac{1}{1+exp[\beta_1\bar{x} + \beta_0]}$$

here learnable parameter is β_1, β_0

logit function (logistic function)

$$\log\left(\frac{P(x)}{P(\bar{x})}\right) = \beta_1x + \beta_0$$

Used to get in which region the parameter is present.

need to find β_1, β_0 based on $P(x)$ to get $P(x)$

$$y = \begin{cases} 0, & P(x) \leq 0.5 \\ 1, & P(x) > 0.5 \end{cases}$$

To get β_1, β_0 we use MLE

(S)

↳ Cost function for logistic regression (Binary)

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

no. of data points

$$\text{Cost} = -\frac{1}{m} \sum_{i=1}^m [y_i \log(P(x_i)) + (1-y_i) \log(1-P(x_i))]$$

↳ log-loss equation.

(NT) The eqn $(m_1x_1 + \dots + m_nx_n + c = 0)$ defines the dividing line in the parameter space where at either side the two classes exist. at that line $P(x) = \frac{1}{2}$

↳ Minimize Cost function $\rightarrow C = -\frac{1}{m} [y \log(\sigma(m\theta)) + (1-y) \log(1-\sigma(m\theta))]$

Gradient descent \rightarrow

$$\theta_j^{new} = \theta_j^{old} - \eta \nabla E_m$$

$$m_j^{new} = m_j^{old} - \eta \nabla E_m$$

$$\nabla E_m = \frac{1}{m} \sum_{i=1}^m [\sigma(\theta) - y] \theta_j$$

$$= \frac{\partial C}{\partial \theta_j}$$

→ Evaluation metrics for logistic regression

↳ Confusion matrices

		Actual values	
		Predicted (1) True Positive (TP)	(0) False Positive (FP)
Predicted (1)	True Positive (TP)		
	False Negative (FN)		True negative (TN)

Accuracy

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

not always reliable

↳ Recall (sensitivity) $= \left(\frac{TP}{TP+FN} \right) \rightarrow P(\text{True} | \text{it is True})$
how much probable you are to detect a positive case

$$(10/40) \rightarrow \text{Accuracy} = \frac{10}{130} = \frac{10}{130}$$

$$\text{but } \frac{TP}{FP} = \frac{10}{50}$$

↳ Precision $= \left(\frac{TP}{TP+FP} \right) \rightarrow P(\text{it is True} | \text{predict true})$
given the true predictions how many are actually positive

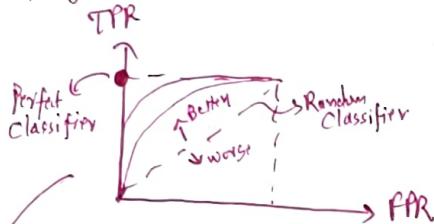
↳ F1 score $F1 \text{ score} = \left(\frac{2 \cdot P \cdot R}{P+R} \right) \rightarrow$ takes both precision and recall into account

↳ ROC (Receiver Operating Characteristic) \rightarrow

 A graph showing the relationship between True positive rate (Y-axis) and False positive rate (X-axis). The Y-axis ranges from 0 to 1, and the X-axis ranges from 0 to 1. A diagonal line from (0,0) to (1,1) represents a random classifier. A curve above and to the left of this line represents a classifier with better performance.

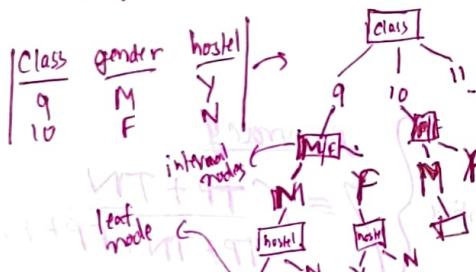
graphical plot that illustrates the diagnostic capability of a binary classifier system as its discrimination threshold is varied

(6)

↳ TPR(true positive rate) \rightarrow sensitivity (recall)↳ FPR(false positive rate) \rightarrow $1 - \underbrace{\text{specificity}}_{\text{TN} + \text{FP}}$ ↳ for different thresholds draw ROC curve \rightarrow take the point (threshold) where TPR high and FPR low→ AUC (Area Under the Curve)for different models (different (m_i, c) values) we have different ROC curves.higher AUC \Rightarrow better model

Decision Tree

flowchart-like tree structure



Gini Impurity Index

$$G(t) = 1 - \sum_{i=1}^j [P(i|t)]^2$$

j : no. of classes of the label (prediction)
 i : one category of the given attribute (whose impurity needs to be calculated)

ranges from 0 \rightarrow 0.5
 (Pure) \quad (most impure)

Gini for an attribute \Rightarrow

$$G(\text{attribute}) = \sum_t \frac{n_t}{n} G(t)$$

t : runs over all the categories of the attribute

n_t : no. of instances having that attribute
 $\sum n_t = n$ = total no. of instances

Root node \Rightarrow attribute with least $G(\text{attribute})$

helps in uncovering non-linear and non-clustered relationships among the attributes.
 with each internal nodes: a test on an attribute
 each branch: an outcome of the test (or record)
 each leaf node: Class label (terminal node)
 class of the predictor

How to build a Decision tree?

using Gini
 to increase the certainty in finding the final leaf node we must always try to zero down on it by starting from least cluttered decisions. This will help to get high purity nodes.
 root nodes should be of attributes with least disorder

Entropy and Information Gain

$$E(t) = \sum_{i=1}^j P(i|t) \log_2 [P(i|t)]$$

higher entropy \Rightarrow higher impurity.pure mode $\Rightarrow E(t) = 0$ Entropy for an attribute \Rightarrow

$$E(\text{attribute}) = \sum_t \frac{n_t}{n} E(t)$$

Information gain: $IG = E_{\text{initial}} - E_{\text{split}}$ E_{initial} = entropy without the split (ie $\sum P(i) \log_2 (P(i))$) E_{split} = entropy after the split using the attributeRoot node \Rightarrow attribute with highest information gain should be preferred as root node.

(7)

CART algorithm (Classification and regression tree)

for attributes having categorical values, the gini index is used to build the decision tree (seen earlier)

Pruning → to reduce the likelihood of overfitting of decision tree in CART

Post-pruning (or just Pruning)

after building tree

by minimum error

pruned back to a point where cross-validation error is minimum.

by smallest tree

tree pruned back slightly further, further than the minimum error, is still increase in error, at the expense of more intelligible error.

Ensemble Learning

↳ Bagging (Similar models)

↳ Stacking (different models)

↳ Boosting (similar models)

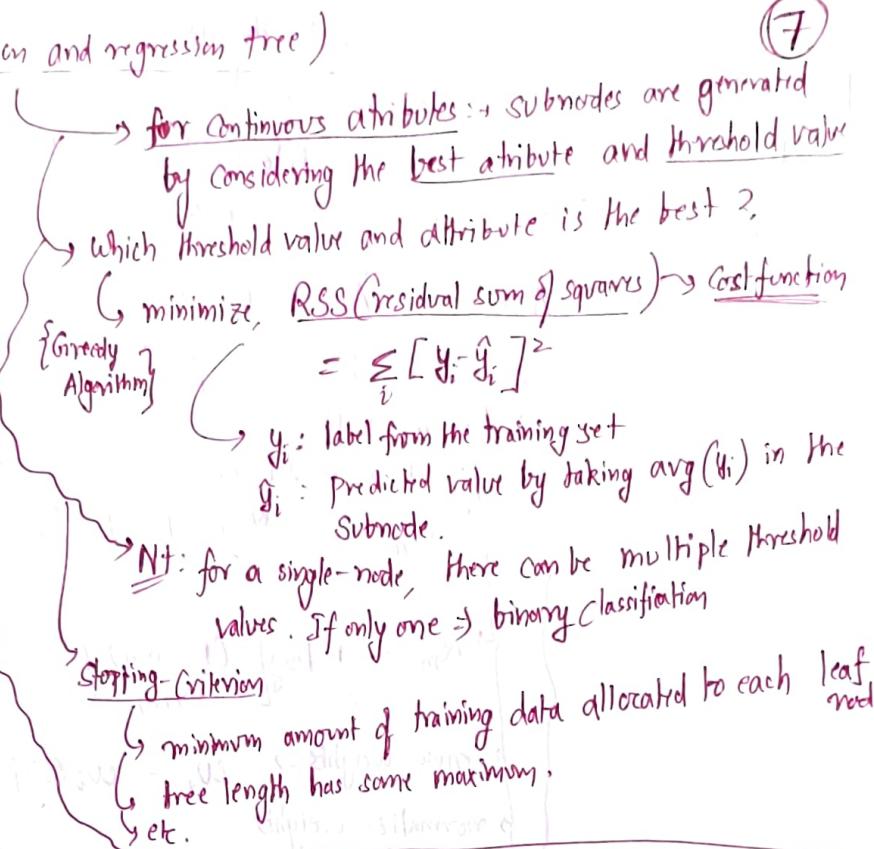
Random Forest

one kind of bagging of decision trees (DT) in which the correlation between the individual DTs is reduced by restricting the no. of attributes taken into consideration for node splitting

if m total attributes present

$p = \sqrt{m}$ (classification) default no. of random attributes that

$p = m/3$ (regression) can be taken for splitting



Early stopping (pre-pruning)

at each stage of tree-splitting, cross-validation error checked if error does not decrease significantly enough → we stop

if may underfit as it may not allow to have subsequent splits where the error reduce significantly

meta approach to ML which seeks better predictive performance by combining the predictions from multiple models.

↳ Bagging (Bootstrap AGgregATING) → many samples generated from the given data by bootstrapping (drawing samples with replacement)

↳ simple voting or averaging of the outcomes (in case of decision trees)

Smaller bootstrap samples help in case of large datasets.

difficult to visualize individual trees.

leads to lower correlations

improvement in accuracy

↳ DT → handling of missing values is easy

↳ high comp. time.

B) Boosting

A method in which the training dataset for each subsequent model is increasingly focussed on instances misclassified by previous model.

- ① AdaBoost
- ② Gradient Boosted Trees
- ③ XGBoost

Use of very simple DTs (called weak learners) in case of DTs as models
Combine prediction using a weighted average of models

Ada Boost (Adaptive boost)

- initial weights $\rightarrow w_i = \frac{1}{N}$ to all samples $\{\bar{x}_i, y_i\}$
- for $m = 1 - M$
 - fit a model $G_m(x)$ to the training data
 - $err_m = \sum_i^N w_i I(y_i \neq G_m(\bar{x}_i))$ } $I(\text{True}) = 1$
 $I(\text{False}) = 0$
 - gives weights of the wrong predictions
 - $\alpha_m = \frac{1}{2} \log \left[\frac{(1 - err_m)}{err_m} \right]$
 - new weights $\rightarrow w_i = w_i \exp[\pm \alpha]$ + \rightarrow incorrect classified
 \rightarrow correct classified
 - normalize weights
 - Draw new sample by weighted random draw and ascribe equal initial weights.
 - Combine the predictions using a weighted average of models.

Gradient boosting → updating the decision trees based on minimization of the loss function $L(y, F(x))$

$$\text{Steepest descent approach: updated model: } F_{m+1}(x_i) = f_m(x_i) - \gamma \left(\frac{\partial L}{\partial f_m(x_i)} \right) \\ = f_m(x_i) + \gamma h_m(x_i)$$

different loss functions for different models

$$\text{Least square regressor: } h_m(x_i) \propto \frac{\partial L}{\partial f_m(x_i)} \left[\frac{1}{n} \sum (y_i - f(x_i))^2 \right] \propto (y_i - F_m(x_i))$$

$$\text{residual} \rightarrow h_m(x_i) = y_i - F_m(x_i)$$

Algorithms:

(1) get the base model, $f_0(x_i) = \text{avg}(y_i)$

(2) fit decision tree to the residuals (ie find a model for residuals)

$$h_0(x_i) = y_i - \text{avg}(y_i) \rightarrow DT-1$$

$$(3) \text{ get the new value: } f_1(x_i) = f_0(x_i) + \gamma h_0(x_i)$$

(4) get the next updated model by finding a new model for new residuals

$$h_1(x_i) = y_i - f_1(x_i) \rightarrow DT-2$$

$$F_2(x_i) = f_0(x_i) + \gamma [h_0(x_i) + h_1(x_i)]$$

$$(5) \text{ till } M-\text{step} \rightarrow h_{M+1}(x_i) = y_i - F_M(x_i) \quad F_M(x_i) = f_0(x_i) + \gamma \sum_{m=1}^{M-1} h_m(x_i)$$

↳ XGBoosting (Extreme Gradient Boosting) ↳ very powerful model
 ↳ high speed and high performance ↳ developed by Tianqi Chen

↳ steepest descent approach to minimize the

objective function: $\text{Obj}(\theta) = \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^T w(f_i)$ ↳ at t^{th} iteration

$$\hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i)$$

(e.g. RMSE, log logistic error)

$$w(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

Controls no. of leaves Controls the score

here, $T \rightarrow$ no. of leaves

↳ building the $f_{k=(t+1)}$ tree:

↳ not built based on traditional methods like gini index or entropy

↳ taylor expansion of loss function (upto 2nd order): $\sum_{i=1}^n \left[\mathcal{L}(y_i, \hat{y}_i^{(t-1)}) + \frac{1}{2} h_i f_t^2(x_i) \right]$

$$\text{gradient} \leftarrow g_i = \frac{\partial}{\partial \hat{y}_i^{(t-1)}} \mathcal{L}(y_i, \hat{y}_i^{(t-1)})$$

$$\text{hessian} \leftarrow h_i = \frac{\partial^2}{\partial \hat{y}_i^{(t-1)} \partial \hat{y}_i^{(t-1)}} \mathcal{L}(y_i, \hat{y}_i^{(t-1)})$$

$$\text{Similarity at one leaf: } \frac{G_j}{H_j + \lambda}$$

↳ def: divide the θ function leaf wise

so, for j^{th} leaf:

$$G_j = \sum_{i \in I_j} g_i \quad I_j = \text{indices of the data points in the leaf } j$$

$$H_j = \sum_{i \in I_j} h_i$$

↳ Algo 2: (i) Initialize the tree with one leaf

and (ii) leaf split based on gain (maximize the gain)

$$\text{gain} = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$$

score of new left leaf score of new right leaf score of original leaf regularization term

↳ Pruning the tree: if gain < $\gamma \Rightarrow$ do not split the leaf

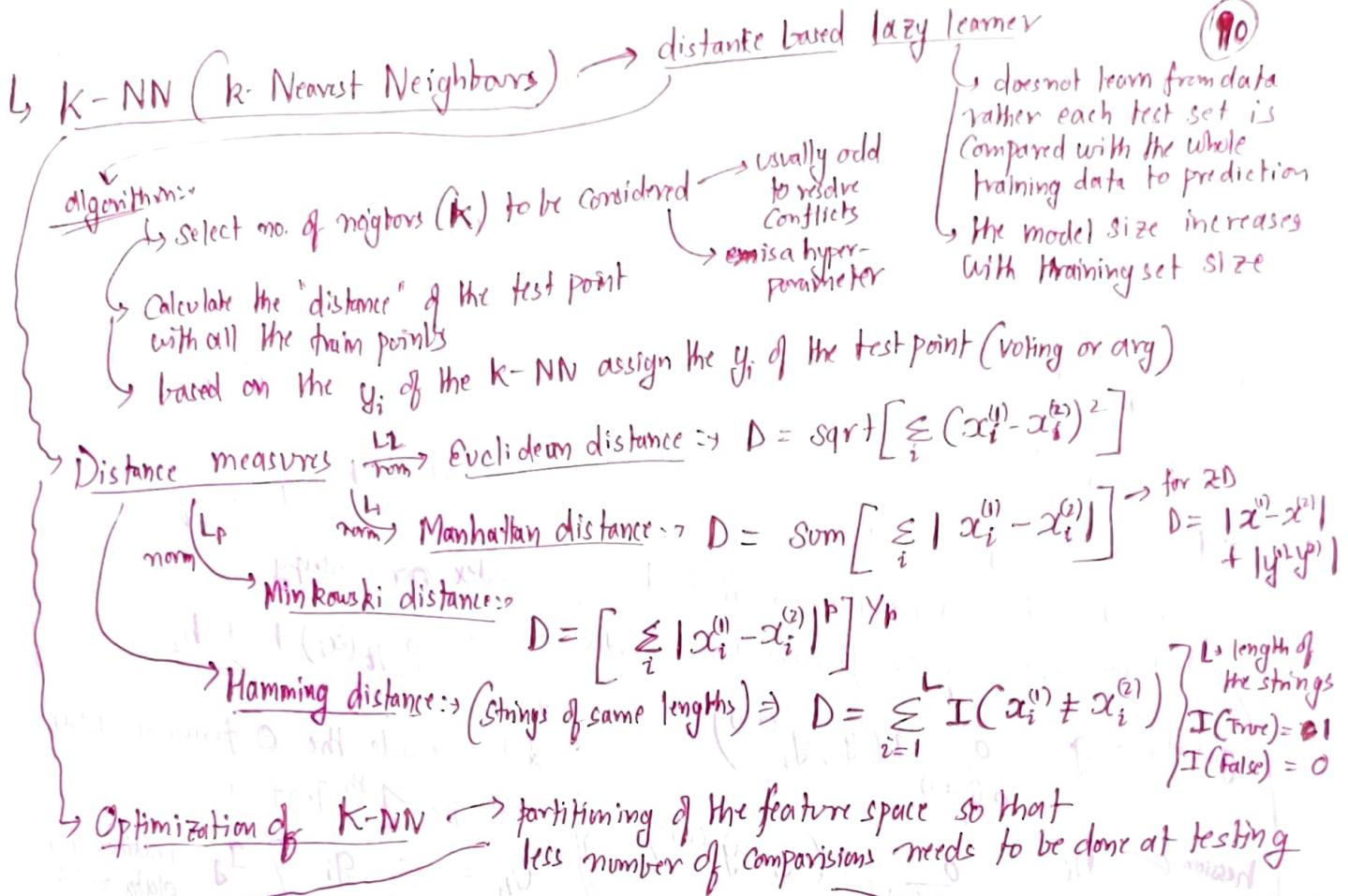
(iii) stopping criterion

↳ (M1) minimum leaf size
 ↳ (M2) minimum leaf frequency
 ↳ (M3) minimum leaf entropy

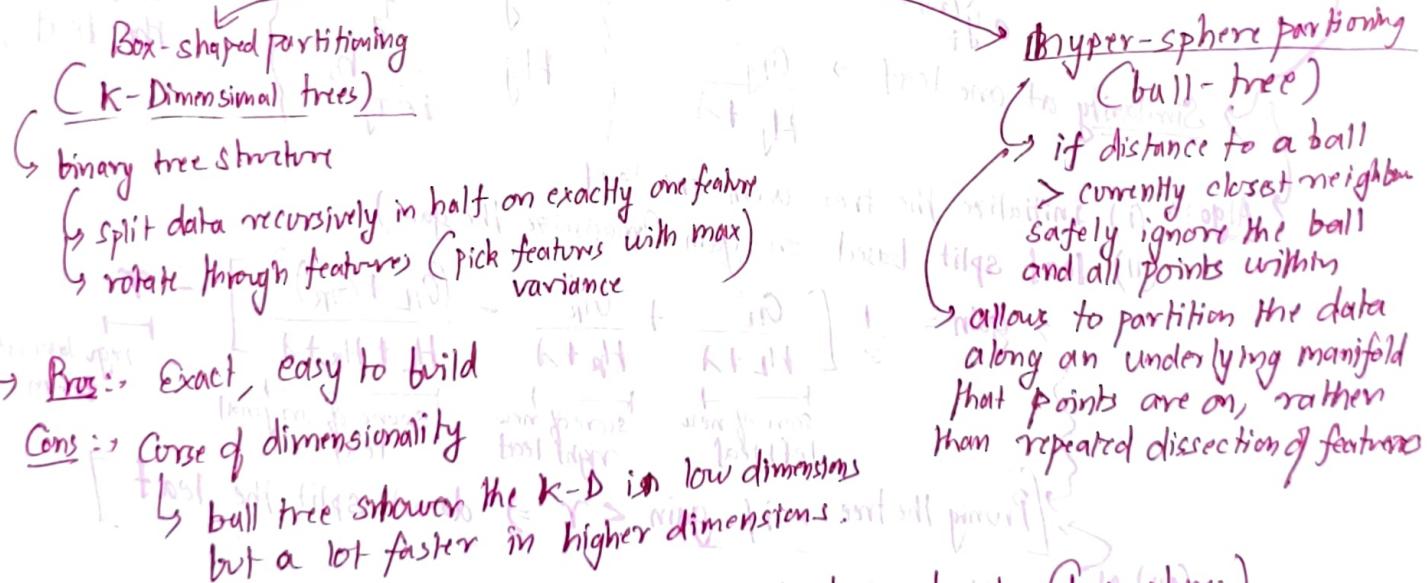
all leaf of size multiple to M3
 represent leaf of min size

min leaf size only off
 leaving off backward

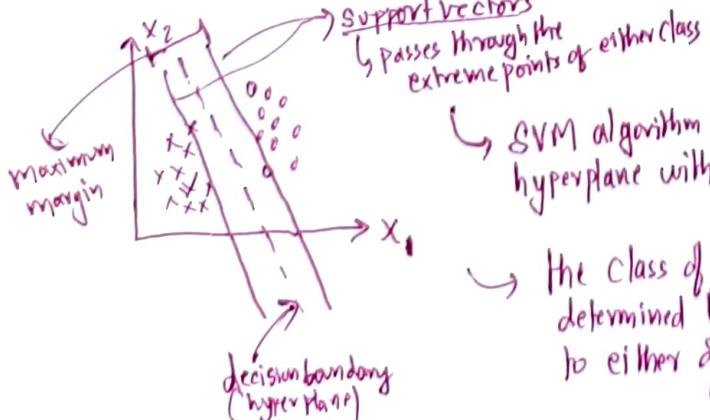
pruning leaf



Optimization of K-NN → partitioning of the feature space so that less number of comparisons needs to be done at testing



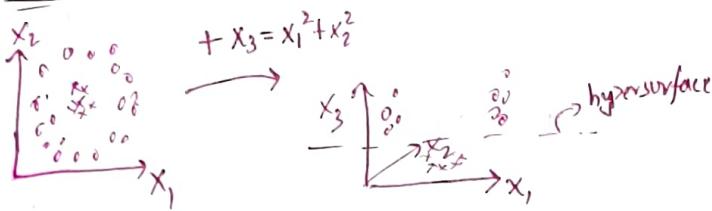
Support Vector Machine (SVM) → to create decision boundary (hyperplanes) separating classes.



SVM algorithm tries to find the hyperplane with highest margin

The class of the test data is determined by its proximity to either of the support vectors.

SVM Kernel trick (non-linear SVM)



when the data points are separated by non-linear hypersurface.

Add a new feature x_{n+1} with the equation of the non-linear hypersurface. Can then separate the data using hyperplane and support vectors.

One vs One (OvO) and One vs Rest (OvR) for multiclass classification

each class vs every other class

e.g. red vs blue, red vs green, blue vs green

each class vs rest of classes

e.g. red vs [blue, green]
blue vs [red, green]
green vs [red, blue]

binary classification algorithms can be used for multiclass classification.

Multiclass problem - splitted into many binary class problems.

Semi-supervised Learning

for supervised learning: dataset needs to be handlabelled
for unsupervised learning: application spectrum is limited

Algorithm is trained upon a combination of labelled and unlabelled data and then uses the existing labelled data to label the rest of the unlabelled data.

First use clustering to cluster similar data and then use the existing labelled data to label the rest of the unlabelled data.

Assumptions: Continuity \Rightarrow closer points \Rightarrow same output label
Cluster assumption \Rightarrow data can be divided into discrete clusters having same labels.
Manifold assumption \Rightarrow data lie approx on a manifold of much lower dimension than input space so that distances and densities can be defined.

K-means Clustering

K determines the number of clusters to be made
Centroid based algorithm, where each cluster is associated with a centroid \rightarrow (average of all the points)

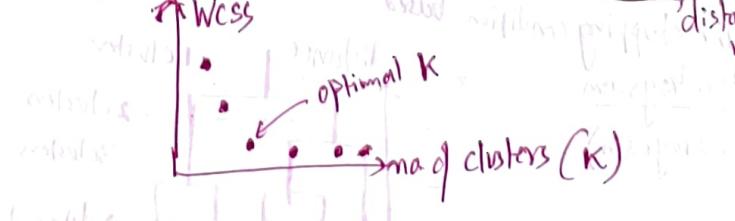
- ↓
Algorithm steps:
1. Select K random points in the input space as centroids.
- 2. Assign each data point to their nearest K -point to form K -clusters.
- 3. Recalculate the position of K -points as the centroid of the points in clusters.
- 4. End the loop when no more reassignment of points occur.

choosing the right $K \rightarrow$ Elbow method

for each value of K calculate the WCSS (Within Cluster Sum of Squares)

$$WCSS = \sum_{\text{Clusters}} \sum_{\text{in cluster}} \|P_i - C_{\text{cluster}}\|^2$$

Centroid of cluster c
distance metric.



(12)

K means algorithm

Smart initialization algorithm

Drawback of k-means
(susceptible to initialization of centroids)

- ↳ if far-off point initialized \Rightarrow might end up with
 - ↳ no points in cluster
 - ↳ and more than one cluster linked to same centroid
- ↳ if centroids too close
 - ↳ same cluster having > 1 centroids

↳ random first centroid from data point

for each data point compute its distance from the nearest chosen centroids (up till now) \leftarrow \rightarrow Stop: when K -centroids are initialized.

assign the next centroid based on a distance based probability from the nearest previously chosen centroid ($P(\text{next centroid}) \propto 1 / (\text{last centroid} - \text{new centroid})$)

→ Mini-batch k-means clustering \rightarrow high speed performance, good for large volume of data

↳ to use small random samples of fixed size data,
↳ each time a new random sample is taken to update the cluster until the convergence
↳ clusters are updated based on a learning rate which decreases as the no. of iterations increases.

→ Agglomerative vs Divisive Clustering (Top-down) \rightarrow all observations start in one cluster, and splits are performed recursively

(bottom-up) \rightarrow each observation starts in its own cluster and pairs of clusters are merged over time into a single large group.

e.g. Hierarchical clustering

don't need pre-determined number of clusters.

→ Hierarchical clustering analysis (HCA) \rightarrow linkages to measure the distance b/w clusters

Algorithm

take every point as a cluster

club the two nearest clusters together to reduce the no. of clusters by one

repeat the clubbing of clusters until the stopping condition based on dendrogram

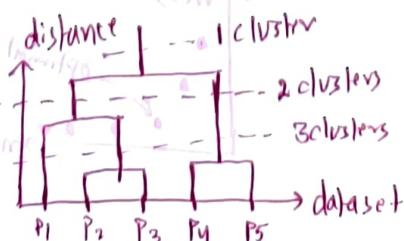
(dendrogram)

single linkage \rightarrow distance b/w closest points considered

complete linkage \rightarrow distance b/w farthest points considered

Average linkage \rightarrow distance b/w each pair of datasets added up and then divided by total number of datasets. one of most popular

Centroid linkage \rightarrow distance b/w centroids considered.



→ DBSCAN (density based spatial clustering of applications with noise)

↳ drawbacks of k-means and hierarchical clustering: (i) only good for convex clusters
 (ii) can't work with noise as they are highly susceptible to them

↳ DBSCAN can克服

↳ Two parameters:

↳ epsilon (ϵ or ep): defines the neighborhood around a data point, two points considered neighbors if their distance $< \epsilon$
 if ϵ too small \Rightarrow large no. of outliers // ϵ too large \Rightarrow clusters will merge

↳ Minimum points (MinPts): minimum no. of data points within ep radius.

↳ larger dataset \Rightarrow larger MinPts; general rule: $\text{MinPts} \geq D + 1$
 ↓ dimension of dataset

↳ Types of data points

↳ Core-point

if it has $\geq \text{MinPts}$
 within its ep radius

↳ border point

if it has $< \text{MinPts}$
 within its ep radius
 but is in a neighborhood of Core-point

↳ Noise/outlier

$\leq \text{MinPts}$ in its ep radius
 and not in any neighborhood
 of core points.

↳ Algorithm

↳ get all the neighbor points for each point and identify the core-points

↳ if a core-point not assigned to any cluster, assign it to a new one

↳ find all its densely connected points and assign them to the same cluster.

↳ two points a and b are densely connected if \exists a point c which is a core point
 and is in the neighborhood of both a and b

Stop: when visited all the points. Those points that don't belong to any cluster are noise.

* Check scikit-learn's page on Clustering for comparison on different clustering algorithms \rightarrow DBSCAN performs best both in terms of making clusters and time complexity.

↳ Performance Metrics for clustering

↳ Silhouette Score = $\frac{(n-i)}{\max(i,n)}$ || $n \rightarrow \text{avg}(\text{dist b/w each sample and its nearest cluster})$
 $i \rightarrow \text{avg}(\text{intra-cluster distance})$

Only applicable to convex clusters
 ↳ near +1 \Rightarrow cluster's samples further away from nearby clusters (good clustering)
 ↳ near 0 \Rightarrow sample very close to decision boundary (ambiguous clustering)
 ↳ < 0 \Rightarrow samples assigned to wrong cluster.

↳ When ground-truth clusters labels are present (gold-standards)

$P \rightarrow$ partition created by our clustering algorithm // $P' \rightarrow$ gold-standard partition of clusters

(true positive) SS: \rightarrow no. of pairs which are in same clusters in both P and P'

(false positive) SD: \rightarrow no. of pairs which are in same cluster in P and in different clusters in P'

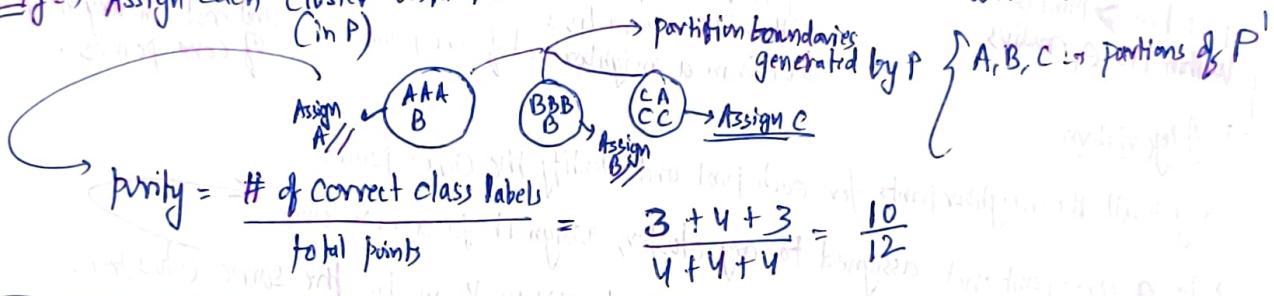
(false negative) DS: \rightarrow no. of pairs which are in different clusters in P and in same cluster in P'

(true negative) DD: \rightarrow no. of pairs which are in different clusters in both P and P'

$$\hookrightarrow \text{Rand index} \Rightarrow R = \frac{SS + DD}{SS + DD + DS + SD}$$

$$\hookrightarrow \text{Jaccard index} \Rightarrow J = \frac{SS}{SS + DS + SD}$$

↳ Purity: \rightarrow Assign each cluster with the class which comes the max time when compared with P'



→ Principal Component Analysis \rightarrow Converts observations of correlated features into

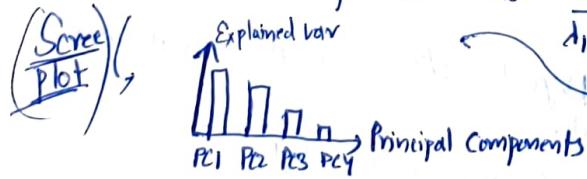
↳ used for dimensionality reduction \rightarrow a set of linearly uncorrelated features (Principal Components)

Step-1: Standardization: $\rightarrow X_i = \frac{\text{Value} - \text{mean}}{\text{stddev}}$ {to compare the variances of all attributes}

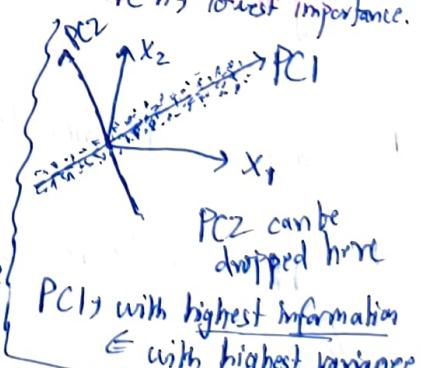
Step-2: Covariance matrix computation: $\begin{bmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) & \dots \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) & \dots \end{bmatrix}$

Step-3: Get eigenvalues and eigenvectors of the covariance matrix

↳ explained variance for each variance = $\frac{\lambda_i}{\lambda_1 + \lambda_2 + \dots + \lambda_n}$



The greater the explained variance, the more information the PC contains.



PC1 with highest information with highest variance

Step-4: get the new features \rightarrow multiply P^* matrix ($[PC_1, PC_2, \dots]$) to Z (original dataset) to get Z^*

↳ each column of Z^* is independent of one another

Step-5: remove the less important vectors