

Machine Learning

PAGE NO.:

★ Pandas data structure

- ⇒ Pandas has 2 types of data structures:
 - a) Series: It's a one dimension array with indexes, it stores a single column or row of data in a DataFrame.
 - b) DataFrame: It's tabular spreadsheet like structure representing row of each of which contain one or multiple columns.
- ⇒ A 1D array (labeled) capable of holding any type of data-Series
- ⇒ A 2D data (labeled) structure with columns of potentially different types of data - Dataframes

| index | Column | | |
|----------------|----------------|----------------|----------------|
| | c ₁ | c ₂ | c ₃ |
| x ₁ | 31 | 52 | 6 |
| x ₂ | gr | 3 | u |
| x ₃ | | | |

whole DataFrame

dataframe is form by combination of series

C₁ series
C₂ series ..

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

PAGE NO.:

* Precision, Recall, confusion matrix & F1 score

| | | Predicted Negative | | Predicted Positive | |
|-----------------|--|--------------------|----------------|--------------------|----------------|
| | | T ₁ | T ₀ | F ₁ | F ₀ |
| Actual Negative | | 7 | TN | 8 | FP |
| Actual Positive | | 1 | 5 | True Negative | Type 2 |
| | | 2 | 2 | FN | 2 |
| | | | | TYPE I | 2 |

Precision: what percent of positive predictions made are correct

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

TPR

Recall: what percent of actual positive value were (sensitivity) correctly classified by classifier

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

Here we chose TP but miss some of them which are actual positive but we marked -ve so we recall them

F1-score: harmonic mean of precision and recall

$$\text{F1-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{1}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \quad \frac{1}{2}$$

* Linear regression

Linear regression cost function $J = \frac{1}{n} \sum_{i=0}^n (\text{Pred}_i - y_i)^2$

J is fun of m, c Here $\text{Pred}_i = mx_i + c$

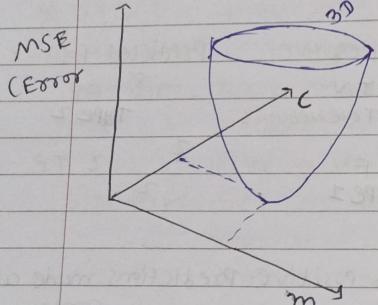
$$J(m, c) = \frac{1}{n} \sum_{i=0}^n (mx_i + c - y_i)^2 \quad y_i \text{ pred}$$

⇒ Our goal is to minimize cost to get best fit line
we are not going to try all permutation and combinations of m & c but use gradient descent algo to best fit line

⇒ gradient decent: algorithm to find best fit for a given training dataset in smaller number of iterations

For linear regression feature must be dependent of each other
 because $\beta = (x^T x)^{-1} x^T y$ has lots of computation
 $(x^T x)^{-1}$ is very expensive computationally

Plot m & c against MSE



For some combination of m & c we get min error this combination of m & c gives best fit line

Then we reduce value of m & c by some amount (learning we will notice decrease in MSE - we will continue doing the same until loss function is very small value or ideally 0 (which means 100% accuracy))

Step by step algorithm

1) Let $m = 0, c = 0 \text{ & } LR = 0.01$

LR gives the speed where the gradient moves during the gradient decent. Setting LR to high make our path unstable to low makes converge slow zero means our model is not learning anything from gradients

2) calculate partial derivative of cost function w.r.t m.
 (with little change of m how cost function changes)

$$\frac{\partial J}{\partial m} = \frac{1}{n} \frac{\partial}{\partial m} \left(\sum_{i=0}^n (y_i^2 - m^2 x_i^2 + c^2 + 2mx_i c - 2y_i m x_i) \right)$$

$$\frac{\partial J}{\partial m} = -2 \sum_{i=0}^n x_i (y_i - (m x_i + c))$$

$$\frac{\partial J}{\partial m} = -2 \sum_{i=0}^n x_i (y_i - y_{\text{pred}})$$



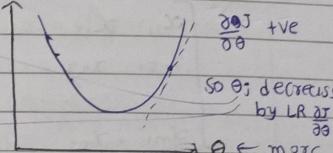
Similarly for c

$$\frac{\partial J}{\partial c} = -2 \sum_{i=0}^n (y_i - y_{\text{pred}})$$

3) repeat update value of m & c in loop eqn

$$m = m - LR \frac{\partial J}{\partial m}$$

$$c = c - LR \frac{\partial J}{\partial c}$$



4) repeat this process until cost function J is very small (ideally 0)

* Advantages of gradient descent

Flexibility: GD is used for various cost function and It can handle non-linear regression problems

Scalability: GD is scalable for large dataset since it updates parameters for each training example one at a time.

Convergence: GD can converge to the global minimum of cost function provided LR is set appropriately

* Disadvantages of gradient descent

Slow convergence: GD may require more iterations to converge to min. since it updates the parameters for each training example at a time

Local minima: GD gets stuck to local minima if cost function has multiple local minima

Noisy updates: Update in GD are noisy and have high variance which can make optimization process less stable and leads to oscillations around the minima

ML → model-based learning (generalized)
 Instance-based learning: give obj instance
 (memorized) Ex: K-nearest, no need to train
 RBF, Kernel fun

online
offline
out of core learning: Training is done like online (batch learning) but first we take model offline
 * For multilinear regression (two or more features)

x = feature matrix of $n \times p$ where x_{ij} is j th feature for i th observation

$$x = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & & & \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \text{response vector}$$

$y_i = \text{response for } i\text{th observation}$

$$y_{\text{pred}} = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi}, \quad \beta \text{ is regression coeffi}$$

$$y_i = y_{\text{pred}} + \epsilon_i \quad \leftarrow \text{residual error}$$

let we represent our feature matrix x as

$$x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \quad \text{where } \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

now we calculate β for which error is minimized

(using least squares method)

As already explained the Least squares method tends to determine $\hat{\beta}$ for which total residual error is so we present result directly [Error $E = \sum(y - \hat{y})^2$ minimized]

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

$$\hat{y} = x \hat{\beta}$$

Estimated response vector

$$\frac{dE}{d\beta} = 0 = \frac{d}{d\beta} [y^T y - 2 y^T x \beta + \beta^T x^T x \beta]$$

$$= 0 - 2 y^T x + \frac{d}{d\beta} (\beta^T x^T x \beta) \quad \text{For this there is multivariate} \quad \frac{dy}{dA} = 2 x A^T \quad \text{If } y = A^T x \quad \therefore \text{3D matrix}$$

$$= -2 y^T x + 2 x^T x \beta = 0 \Rightarrow \boxed{\beta = (x^T x)^{-1} x^T y} \quad (1) \quad A^T = A \text{ so } (x^T x)^{-1}$$

Odds: The proportion of an event's chances of happening to its chances of not happening
 ⇒ Regularization is also solved by direct form formulation & gradient
 ⇒ by regularizing high w & make them compare to low weight

* Regularization

Regularization is process where we ~~eliminate~~ some feature which is not useful

⇒ Regularization is a ~~tech~~ technique used to reduce the errors by fitting the function appropriately on the given training set and avoid overfitting

⇒ commonly used regularization technique-

L1 regularization = LASSO (Least Absolute Shrinkage and Selection Operator)

L2 " = Ridge regression

Dropout "

⇒ Lasso Regression adds "absolute value of magnitude" of coefficient as penalty term in loss function

$$\|w\|_1 = |w_1| + |w_2| + \dots + |w_n|$$

⇒ Ridge adds "square magnitude" of coefficient as penalty term to the loss function

$$\|w\|_2 = (\|w_1\|^2 + \|w_2\|^2 + \dots + \|w_n\|^2)^{1/2}$$

Note: During the regularization the output function y_{pred} does not change. The change is only in the loss function

The output function

$$y_{\text{pred}} = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Loss function before regularization

$$\text{Loss} = \text{Error}(y, y_{\text{pred}})$$

After Regularization

$$\text{For Lasso Loss} = \text{Error}(y, y_{\text{pred}}) + \lambda \sum_{i=1}^N |w_i| \quad \lambda \text{ is hyperparameter}$$

$$\text{For Ridge Loss} = \text{Error}(y, y_{\text{pred}}) + \lambda \sum_{i=1}^N w_i^2 \quad \text{regularization constant } \lambda > 0$$

$$\text{For Ridge Loss} = (xw - y)^T (xw - y) + \lambda \|w\|^2$$

$$L = (xw - y)^T (xw - y) + \lambda w^T w$$

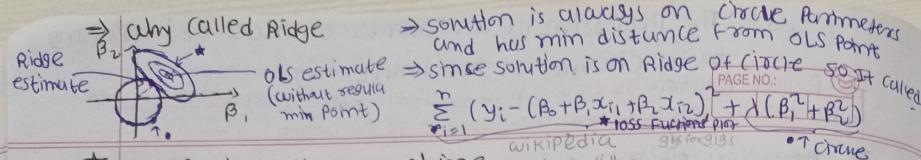
$$\frac{\partial L}{\partial w} = 2x^T(xw - y) + 2\lambda w = 0$$

$$\Rightarrow w = (x^T x + \lambda I)^{-1} x^T y$$

For gradient descent

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$\frac{\partial L}{\partial w} = x^T x w - x^T y + \lambda w$$

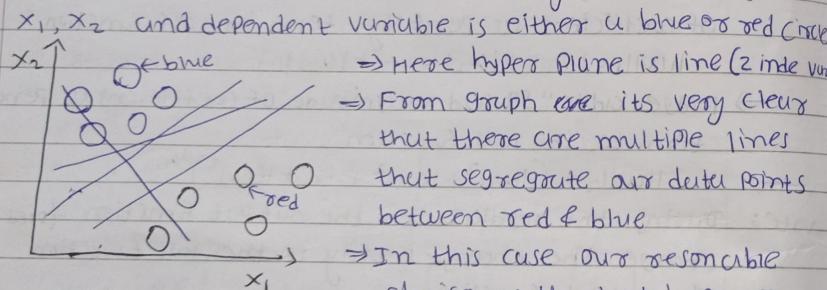


* Support vector machine

- SVM is a supervised ML algo used both for classification and regression
- objective is of SVM to find hyperplane in N-dimensional space that distinctly classifies the data points
- the dimension of hyperplane depends on no. of input features. If input features are 2 the hyperplane is just a line, If input features are 3 then the hyperplane become 2D plane.

It is difficult to imagine when the number exceed 3.

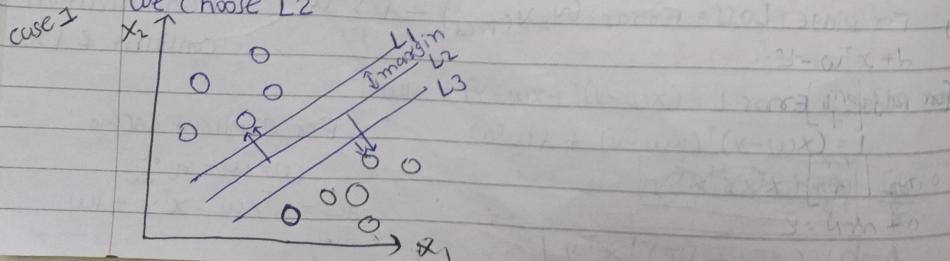
- lets consider input variable are only 2 (features)



Linearly separable data point is the one that represent largest separation or margin between the two classes

- so we choose the hyperplane whose distance from it to the nearest data point on each side is maximized

- If such a hyperplane exists it known as maximum-margin hyperplane or hard margin so we choose L2



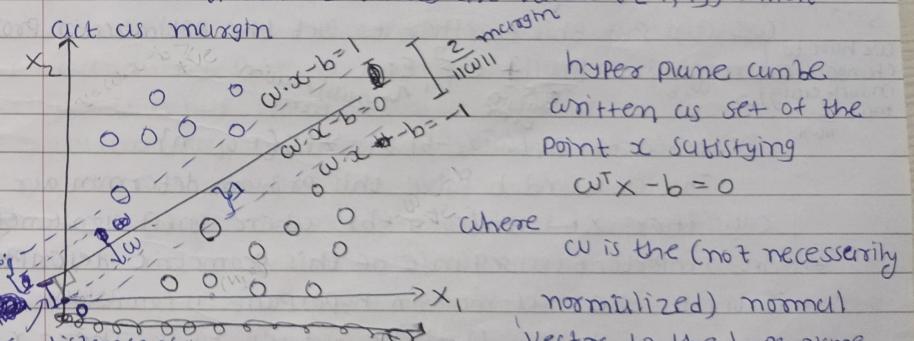
Hyperplane: In geometry, a hyperplane is a subspace whose dimension is one less than that of its ambient space.

⇒ Support vectors are data point are closer to hyperplane and influence the position and orientation of hyperplane. Using these support vectors, we maximize the margin of the classifier. Deleting the support vector will change position of hyperplane. These point can help us to build our SVM.

* Large margin intuition

⇒ In logistic regression we take the linear function's output and squash the value within range of $[0, 1]$ and if the o.p > threshold value (0.5) then label it as 1 otherwise 0

⇒ In SVM we take the o.p of linear function and it is greater than 1 we identify it with one class and if output is less than 1 we identify it with another class since the threshold values are values ($[-1, 1]$) which act as margin



where w is the (not necessarily normalized) normal vector to the hyperplane

Hard margin margin along w vector (normal to hyperplane)

⇒ If training data is linearly separable then we can select two parallel hyperplane that separates two classes of data. The region bound by those is called "margin"

⇒ The maximum margin hyperplane is the ~~margin~~ hyperplane that lies half way between them.

⇒ with normalized or standardized data set those hyperplanes can be described by the equations

if $w^T x - b = 1$ (anything on or above this boundary is one class)

if $w^T x - b = -1$ (anything on or below this boundary is other class)

Now let point a is a point in plane from where normal to the plane passes through origin so $\vec{c} \cdot \vec{n} = |c| |n| \cos \theta = |c| \cdot 1 \cdot 1$.
 PAGE NO.: Angle between \vec{n} and \vec{c} is zero
 cmy point on plane so $d = c \cdot n$
 Point a from origin

distance between those hyperplane is $\frac{2}{\|w\|}$ so the maximum distance between them we need to minimize $\|w\|$. This distance is computed using the distance from a point to plane eqn.

we also have to prevent data points from falling into the margin, we add the following constraint for each i either

$$w^T x_i - b \geq 1, \text{ if } y_i = 1$$

or

$$w^T x_i - b \leq -1 \text{ if } y_i = -1$$

This constrains state that each data point must lie on correct side of margin. This can be written as

$$y_i(w^T x_i - b) \geq 1, \text{ for all } 1 \leq i \leq n$$

actual value predicted value

soft margin

We can put this together to get the optimization problem we have to change w, b minimize $\|w\|^2 + C \sum_{i=1}^n \epsilon_i$
 max $\|w\|$ margin error classification error

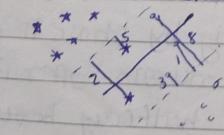
$$\text{Subjected to } y_i(w^T x_i - b) \geq 1 \quad \forall i \in \{1, 2, \dots, n\}$$

The w , b and ϵ_i solve this problem determine our classifier $x \mapsto \text{sgn}(w^T x - b)$ where $\text{sgn}(\cdot)$ sign function

An important consequence of this geometric description is that the max-margin hyperplane is completely determined by all that lie nearest to it. These x_i are called support vectors.

Here $C =$ How many errors model can consider

soft margins (written) $\epsilon_i =$ Value of the errors



$\epsilon = 2 + 5 + 3 + 1 + 8$
 error is sum of distance between point inside and their corresponding line

$(\vec{c} - \vec{a}) \cdot \vec{n} = 0 \quad \vec{c} \cdot \vec{n} = 0 \quad (\because \text{divide both sides by } |n|)$

$$\vec{c} \cdot \vec{n} - \vec{a} \cdot \vec{n} = 0$$

$$\vec{c} \cdot \vec{n} - d = 0$$

$$\vec{c} \cdot \vec{n} = d \quad \text{distance from origin}$$

when normal not necessarily unit vector

here d is distance of plane from the origin

but d_{plume} is distance of plume from origin

dist of point x to plane or line

$$d = w^T x + b \leftarrow \text{eqn.}$$

$$\|w\| \leftarrow \text{stop}$$

PAGE NO.: Features changing for each class

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

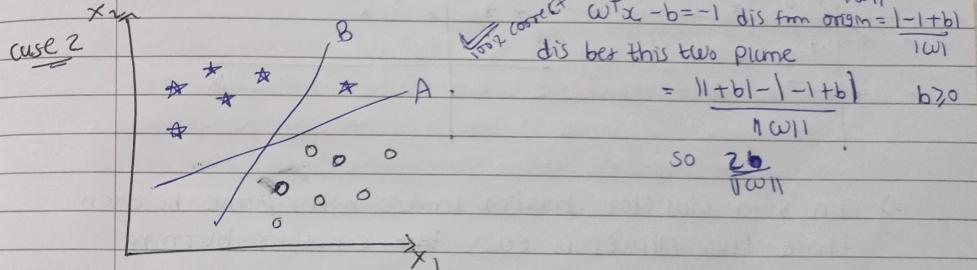
$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$

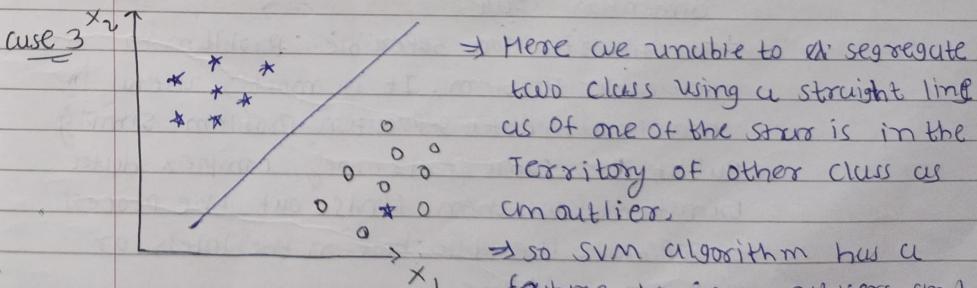
$w^T x - b = 0$ \downarrow $w^T x - b = 0$

$w^T x - b = 0$ \downarrow $w^T x - b = 0$



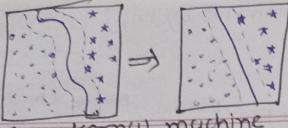
\Rightarrow SVM select the hyperplane which classifies the class accurately prior to maximizing margin

\Rightarrow Plane B has classification error and plane A classified all correctly so plane A is correct hyperplane



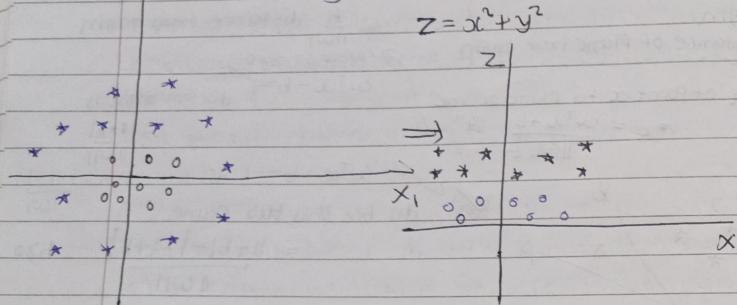
\Rightarrow Here we unable to segregate two class using a straight line as one of the stars is in the territory of other class as outlier.

\Rightarrow so SVM algorithm has a feature to ignore outliers and find hyperplane that has max margin



\Rightarrow SVM can solve this problem ~~more easily~~ easily by introducing additional feature

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



\Rightarrow In SVM classifier, having linear hyperplane between those two classes is easy but another burning question arises that if we should ~~be~~ need to add this feature manually to have a hyperplane. No. SVM algorithm has a technique called the Kernel trick
 \Rightarrow Kernel trick: is a function that take low dimension vector space and transfer it to higher dimension space

\Rightarrow It converts non-separable problem to separable problem. It is mostly useful in non-linear data separation problem. Simply put it does some extremely complex data transformations, then finds out the process to separate the data base on the labels or output you define.

* Advantage of SVM

- \Rightarrow effective in the higher dimension cases
- \Rightarrow It's memory efficient as it uses a subset of training point in the decision function called support vectors
- \Rightarrow different Kernel function can be specified for the decision function and it's possible to specify custom kernels

* Disadvantage of SVM

- \Rightarrow difficult to choose Kernel function & C & gamma
- \Rightarrow It's does not perform well on large dataset because the required training time is large
- \Rightarrow It's does not perform well when dataset has more noise i.e., target set classes are overlapping
- \Rightarrow SVM doesn't directly provide probability estimates; these are calculated using an expensive five-fold cross-validation. It is included in the related SVC method of the Python Scikit-learn library
- \Rightarrow SVM is work best on small dataset which contain the number of dimension is greater than number of samples.

Soft margin

To extend the SVM to cases in which data are not linearly separable the hinge loss function is helpful

$$\epsilon_i = \text{Loss} = \max[0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)] \quad y_i \text{ is } i\text{th target } (\pm 1)$$

$\mathbf{w}^T \mathbf{x}_i - b$ is the i th output

⇒ This Loss function is zero when Condition

$$y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \text{ for all } 1 \leq i \leq n$$

satisfied - (means x_i lies on correct side of margin)

⇒ For data on the wrong side of the margin, the function's value is proportional to the distance from x_i the margin

⇒ The goal of the optimization then is to minimize

$$\lambda \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \epsilon_i$$

where, $\lambda > 0$ determines the trade-off between increasing margin size and ensuring that x_i lie on the correct side of the margin, by deconstructing the hinge loss

⇒ This optimization problem can be massaged into the following:

$$\text{so we have calculate minimum value of } \left(\frac{1}{n} \sum_{i=1}^n \epsilon_i + \lambda \|\mathbf{w}\|^2 \right)$$

$$\lambda \|\mathbf{w}\|^2 + \sum_{i=1}^n \epsilon_i \quad (C \propto \frac{1}{\lambda})$$

with respect to

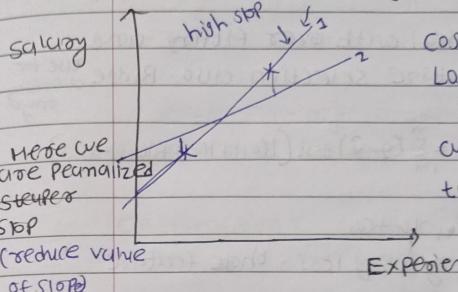
$$\text{(subject to)} \quad y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0 \quad \forall i \in \{0, 1, 2, \dots, n\}$$

For the large value of C it behaves similar to the hard margin SVM.

feature which is not highly correlated is minimized but not zero

Regularization

★ Ridge Regression & Lasso Regression to avoid overfitting



cost function for Ridge regression

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\text{slope})^2$$

$$\lambda > 0$$

we have to reduce this cost function to deal with overfitting

$$\text{If there are more than one feature then } \text{Loss} = n + \lambda (w_1^2 + w_2^2 + \dots + w_n^2)$$

$$\text{Loss} = n + \lambda \text{ (slope)}^2$$

$$\Rightarrow w_i \text{ is slope with respect to } x_i \text{ If increase } x_i$$

by one unit then how much effect in y is calculated from w_i

$$\text{in case of overfitting}$$

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\text{slope})^2$$

$$\text{Loss} = \lambda (\text{slope})^2$$

$$\text{let assume line 1 has slope } 1.3 \text{ & } \lambda = 1$$

$$\text{so Loss}_1 = 1.69$$

but we can find another line 2 with less slope then

$$\text{Line 1 let this line has slope } 1 \quad \lambda = 1$$

$$\text{so Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\text{slope})^2$$

value of $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ has small value let 0.39

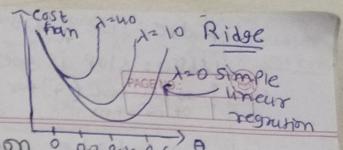
$$\text{so Loss} = 0.39 + 1(1)$$

$$\text{Loss}_2 = 1.39$$

$$\text{Loss}_1 > \text{Loss}_2$$

⇒ Here we reduce overfitting by reducing slope so there will be more cost for training data then overfitting but for test data its give generalized solution

⇒ and we reduce slope with has higher value so in units increasing in λ is less change in y



Ridge: slope goes to zero ~~not~~ but it is shrink
 Lasso: slope goes to zero
 Why? \Rightarrow because $(\partial \text{L})/\partial w_i$ for slope of Ridge & Lasso can explicitly kill this

\Rightarrow If λ is higher value of slopes are goes toward zero

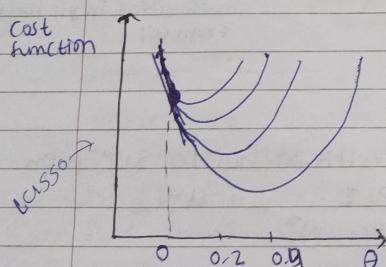
* Lasso regression: use to deal with over fitting and use for feature selection also Ridge use for overfitting only

$$\text{Loss} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |\text{slope}| = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (|w_1| + |w_2| + |w_3| + |w_4|)$$

$$y = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$$

\Rightarrow where ever slope value is very very less those feature will be removed.

\Rightarrow when ever we use $|\text{slope}|$ magnitude in loss function (cost fun) slope will be moving toward zero so some of slope values of feature is reduce to zero.



Features with non zero.

\Rightarrow Features which not highly correlated is become zero

\Rightarrow Feature Selection

\Rightarrow Reduce features which are not important

$$\text{For direct formula } \frac{\partial L}{\partial m} = 0$$

$$\text{Ridge: } m = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) / \sum_{i=1}^n (x_i - \bar{x})^2$$

$m = w$ = weights
 $=$ slopes

$$L = w^T x^T x w - 2w^T x^T y + y^T y + \lambda w^T w$$

$$\frac{\partial L}{\partial w} = x^T x w - x^T y + \lambda w = 0$$

$$w = (x^T x + \lambda I)^{-1} x^T y$$

$$\text{Lasso: } m = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) / \sum_{i=1}^n (x_i - \bar{x})^2$$

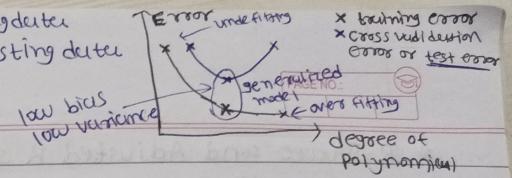
$$\text{also called } \alpha$$

$$0 < \lambda < \infty$$

- formula
 + formula
 $(\because \frac{\partial L}{\partial m} \text{ is } \neq 0)$

So that term will be -ve

Bias: accuracy in training data
 Variance: accuracy in testing data



* Elasticnet Regression (Ridge + Lasso)

$$\text{Cost function: } \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{i=1}^n (\text{slope})^2 + \lambda_2 \sum_{i=1}^n |\text{slope}|$$

MSE

Ridge

Lasso

\Rightarrow Reduce overfitting

\Rightarrow Feature selection

\Rightarrow use for big data set and we haven't knowledge about which regression is better

\Rightarrow use when data has multi collinearity

$$\text{L1 ratio} = \frac{\lambda_1}{\lambda_1 + \lambda_2}$$

$$\text{L1} = \frac{\lambda_1}{\lambda}$$

$$\lambda_1 = \lambda \times L_1$$

$$\lambda_2 = \lambda - \lambda_1$$

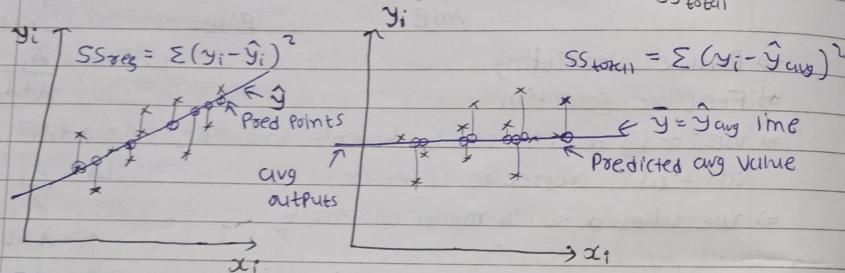
★ R-Square and Adjusted R-Square

- ⇒ In classification we check accuracy using confusion matrix, accuracy score, TP value, Recall value, precision, F1-score
- ⇒ But in regression problem there are different methods called R-square & ~~R-squared~~ technique.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}} \quad \begin{aligned} SS_{\text{res}} &= \text{Sum of Residual or Error} \\ &\text{goodness of best fit line} \\ &\text{Real Points} \\ &SS_{\text{total}} = \text{Sum of Avg total} \end{aligned}$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad \begin{aligned} \text{Predicted Points} \\ \text{Predicted avg Value} \end{aligned}$$

$$R^2 = \frac{SS_{\text{res}}}{SS_{\text{total}}}$$



⇒ If R^2 is less than zero then predicted model is very bad (\because It is worse than avg)

⇒ $R^2 \rightarrow 1$ model is best

$$y = \beta_0 + \beta_1 x_1 \quad (\text{simple})$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (\text{multiple})$$

$$+ \beta_3 x_3$$

as we adding new and new independent features $\Rightarrow R^2$ is usually increases.

⇒ let's add ~~β_3~~ β_3 feature so our linear regression model assign β_3 some value such that $SS_{\text{res}} \downarrow$ so $R^2 \uparrow$

⇒ so $R^2 \uparrow$ as we added new independent features

⇒ so adding independent feature R^2 always increases it will never decrease, (\because even by adding feature if weights is not change but $y_i - \hat{y}$ always $\geq y_i - \bar{y}$ so...)

★ ⇒ even if new added independent feature has no correlation with target output even though $R^2 \uparrow$ so it is not penalizing the new added features. To overcome this problem we introduce adjusted R^2

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

where R^2 = sample R-square

P = Number of predictors (Independent features)

N = Total sample size

If independent whenever Attributes increases R^2_{adj} ↓

⇒ If independent feature is correlated with output then $R^2 \rightarrow 1$ and R^2_{adj} also high

⇒ If variables are not correlated R^2_{adj} slightly decreases compare to original R^2 value.

⇒ If it is completely correlated there will be slight increase in R^2_{adj}

⇒ Every time you add a independent variable to model the $R^2 \uparrow$ even if the independent variable is insignificant. It never declines. whereas R^2_{adj} ↑ only when ind variable is significant and affects dependent variable

$$R^2_{\text{adj}} \leq R^2$$

Hypothesis Testing

PAGE NO.:

→ Evaluate z (can also more than z) mutually exclusive statement (either this true or that) on population data using sample data

Null hypothesis (let innocent)

1) make Initial Assumption (H_0)

H_0 = Alternative

2) Collect Data (evidences)

hypothesis

3) gather evidence to Reject OR non Reject null hypothesis

(guilty)

Actual $\rightarrow H_0 \quad H_1$

Truth

| H_0 Do not Reject ↑ predicted | OK | TP | FP | |
|---------------------------------------|-----------------|----|----|----|
| H_1 Reject (so H_1) | TYPE I error | FN | OK | TN |

⇒ Type 1 error is more harmful or type 2 error is more harmful is depends on Problem Statement

★ P value : Probability for null hypothesis to be true

space key

p=0.8

If p=0.01

that means If we repeat
this experiment 100 times
number of touching this
area 0 1 time

Null hypothesis: Treats everything
same or equal

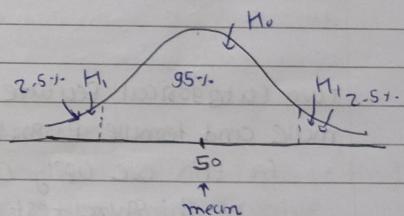
Toss experiment:

H_0 = Null hypothesis = fair coin

Toss 100 times

alternative hypothesis

H_1 = coin is not fair



$$\text{Chi-square} = \sqrt{\frac{(\text{Actual} - \text{Expected})^2}{\text{Expected}}}$$

PAGE NO.:

★ T-test, chi square test, Anova test

| Gender | Age group | Weight (kg) | Height (m) |
|--------|-----------|-------------|------------|
| M | Elderly | 70 | 1.4 |
| F | Adult | 65 | 1.2 |
| M | Adult | 65 | 1.4 |
| M | Child | 20 | 1.1 |
| F | Adult | 75 | 1.3 |
| M | Elderly | 80 | 1.3 |

H_0 = There is no diff
 H_1 = There is diff
 one categorical
 ratio of
 M & F is
 diff or not M F

$P \leq 0.05$ ← significance value

one sample → (P value)
 proportion Test If $P \leq 0.05$ then Reject H_0 and
 There is diff in male & female
 ratio

Two categorical feature : Is there any difference in
 male and female Proportion base on Age group
 for this we apply Chi square test

If Chi square test gives value < 0.05 then reject H_0 .

one numerical variable
 Continuous variable (let consider Height)

now what is statement that comes wrt continuous variable
 height? ~~say~~

avg height = 1.3 is there any diff wrt mean height

Here we apply → Test: T-test

base of T test If values > 0.05 then reject H_0 .

two numerical variable
 Continuous variable

Test: Correlation (-1 to +1) / T test

If correlation value $\rightarrow 0$ so no correlation

H_0 = There is relation betw two variable

H_1 = " " no relation " " "

base on this we apply T test If $P \leq 0.05$ reject H_0 .

more then 2 category Test: Anova $P \leq 0.05$ reject H_0 .

(Ex: Income (elderly, adult, child)) ⇒ correlation between categories

★ Matrix In classification

Classification Problem

default Probability
 Class Lables is 0.5

Probabilities

In binary classification
 there are two classes A, B

In this method threshold
 value is taken 0.5

If \geq greater than 0.5

then Its one class otherwise
 another class

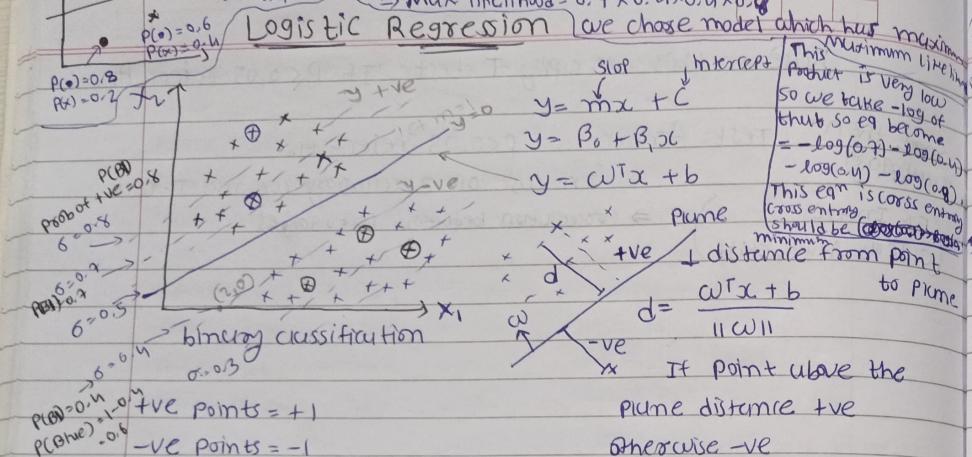
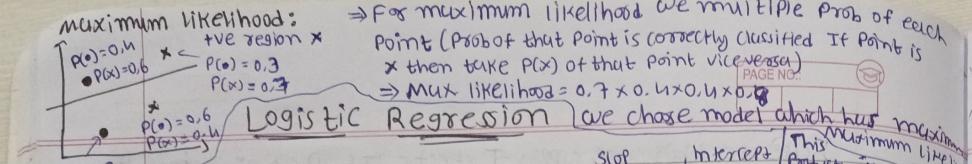
$$F_{\text{beta}} = \frac{(1+\beta)^2 \text{ Precision} \times \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}}$$

$$\text{If } \beta=1 \quad F_{\text{beta}} = \frac{2 \text{ Pre} \times \text{Re}}{\text{Pre} + \text{Re}}$$

when FP & FN both close equally important then we
 select $\beta=1$

\Rightarrow If FP having more impact then FN then we reduce β
 (Type 2 errors)

\Rightarrow If FN having more impact (Type 1 errors)
 then FP then we increase β value then $\beta < 1$



$$y_i \sum_{i=1}^n w_i^T x_i > 0$$

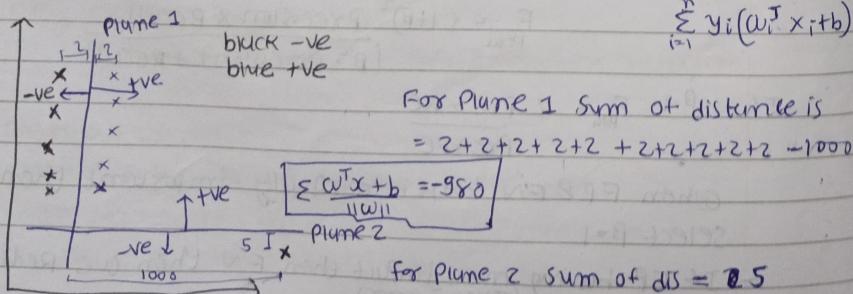
$$y_i (w^T x + b) > 0 \quad y = w^T x = [m] [x_1 \dots x_n]$$

If add b in matrix

$$\frac{y_i (w^T x + b)}{\|w\|} > 0 \quad y = [m] [x_1 \dots x_n]$$

$$\text{Cost function is } \max \sum_{i=1}^n y_i (w_i^T x_i + b)$$

for maximize $\sum y_i (w_i^T x_i + b)$ we need to update only w_i (given y_i , given x_i) until get max of



Impact of σ : $w_n = w_0 + n(y_i - \hat{y}_i)x_i$
 instead of using step fun it use or then $y_i - \hat{y}_i$ will be 0 and weight update every time either point is correctly classified or -ve point (+ve) move always

so here because of outliers we select plume 2 but plume 1 is best plume so overcome with this here we take best

s function to prevent this

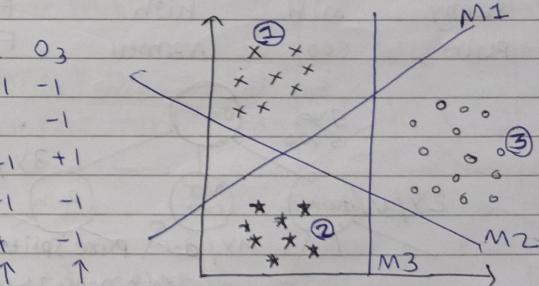
$$\max \sum_{i=1}^n \sigma(y_i + w^T x_i)$$

here $\sigma = \frac{1}{1 + e^{-z}}$, $z = 0 \Rightarrow \sigma = 0.5$

In pre calculation for plume 1 sum of distance is -980. This is big value so we transform this to between 0 to 1 by using σ function. by doing this we neglect large value due to outliers

* Logistic Regression (one vs Rest)

| f_1 | f_2 | f_3 | O/P | O_1 | O_2 | O_3 |
|----------|----------|----------|-------|-------|-------|-------|
| I_1 | I_2 | I_3 | O_1 | +1 | -1 | -1 |
| I_4 | I_5 | I_6 | O_2 | -1 | +1 | -1 |
| I_7 | I_8 | I_9 | O_3 | -1 | -1 | +1 |
| I_{10} | I_{11} | I_{12} | O_1 | +1 | -1 | -1 |
| I_{13} | I_{14} | I_{15} | O_2 | -1 | +1 | -1 |



m_1 will be able to find out that O_1 is +1 or not +1 for O_3

whenever we give new input feature we pass all input in M_1 model then pass all feature in M_2 or M_3

Respectively let m_1 give probability 0.2

m_2 " " 0.25

m_3 " " 0.55

Probability of m_3 highest so data belongs to O_3 category

~ continue in DL book

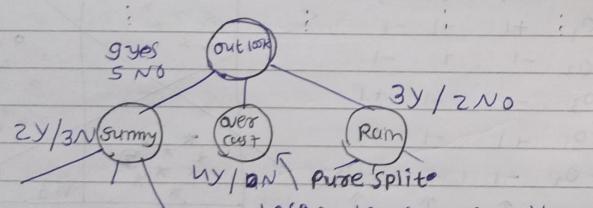
Decision Tree

- use for classification and regression both
- There are 2 technique in Decision Tree

ID3 \leftarrow CART

- Entropy and gini index \leftarrow Purity Split
- Information gain \leftarrow feature Decision Split

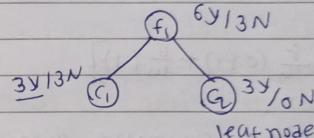
| outlook | Temp | Humidity | Wind | Play |
|----------|------|----------|-------|------|
| Sunny | Hot | High | False | No |
| Sunny | Hot | High | True | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | F | Yes |
| Rain | Cool | Normal | T | No |
| Overcast | Cool | Normal | Weak | Yes |
| Sunny | Mild | High | F | No |
| Rain | Cool | Normal | F | Yes |



leaf node If all are Y or N
 → If there is no pure split we take another feature and split again until leaf node

Purity Check 1) Entropy

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$



+ for Yes
- for No

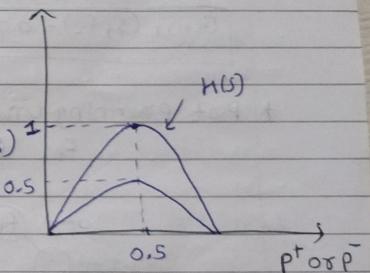
It can be more than 2 of category

$$\text{For } C_1, H(C_1) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$H(C) = 0 \rightarrow \text{Pure split}$$

$$\text{For } C_1, H(C_2) = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right)$$

$$(H(C_2) = 1) \rightarrow \text{Impure split}$$



2) gini Impurity

$$G.I. = 1 - \sum_{i=1}^n (P_i)^2$$

$$\text{Gini} = \sum_{i=1}^n (P_i)^2$$

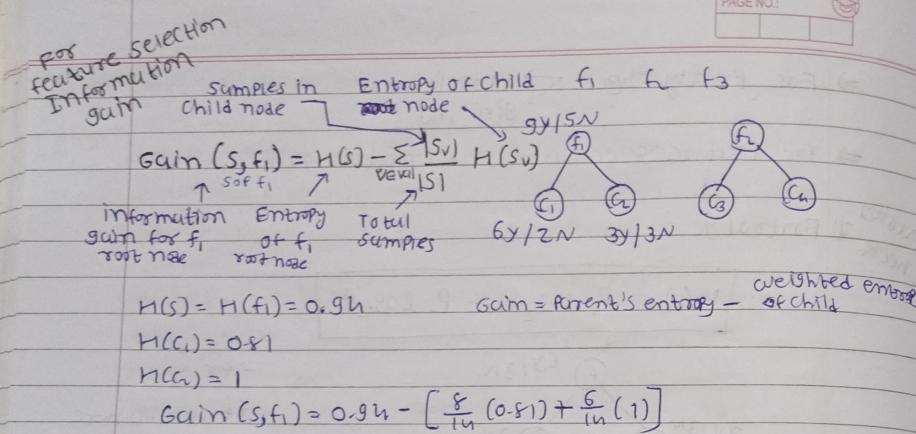
$$G.I. = 1 - \text{Gini}$$

$$G.I. \text{ for impure split} = 1 - \left((\bar{P}_+)^2 + (\bar{P}_-)^2 \right) = 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right]$$

$$G.I. = 0.5 \text{ for Impure}$$

$$G.I. = 0 \text{ for Pure}$$

→ we can use both Entropy and gini Impurity but if we have big dataset then for easy calculation GI is best.
 (∴ Entropy contain log)



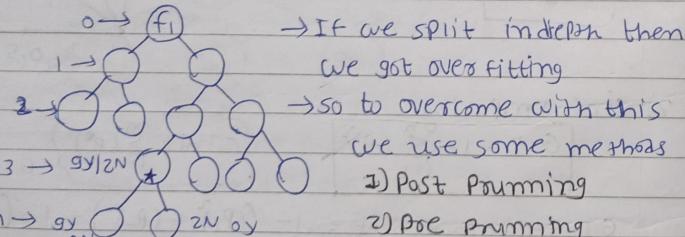
$$Gain(S, f_1) = 0.049$$

Here info. gain in f_2 is greater so we select f_2

$$Gain(S, f_2) = 0.052$$

* Post Prunning and Pre-Prunning Decision trees

$f_1 \quad f_2 \quad f_3 \quad O/P$



1) Post Prunning:

In Post Prunning we make decision tree and then Prunning (cut) it. using Parameters.

- These are many parameters let we take max-depth here for a node we see that probability of y is high so there is no need to further split. by doing this we can overcome with overfitting
- max-depth set us 3 level so there is no splitting after 3 level and there is no overfitting

- 2) Pre Prunning : we are applying hyper parameters tuning while making decision tree
- set max-depth in some range of values [.....]
 - and define max-features [.....]
 - now decision tree select and after applying all parameters find where we can get max accuracy

⇒ Post Prunning ← Small dataset
Pre Prunning ← Bigger dataset

* Decision Tree split for numerical

$f_1 \quad O/P$

$x_i \rightarrow 2-3 \quad Yes$

$3-6 \quad Yes$

$4 \quad No$

$5-2 \quad No$

$6-7 \quad Yes$

$8-9 \quad No$

$10-5 \quad Yes$

$14-2 \quad No$

disadvantage: time complexity
but in tot time it's low

1) sorting all values in ascending order

$x_i \leq 2-3$

$\leq 2-3 \quad y/y \quad n/n$

$\geq 2-3 \quad y/n \quad y/n$

1 yes

2) next threshold value is 3-6 and again split will happens

$\leq 3-6 \quad y/y \quad n/n$

$\geq 3-6 \quad y/n \quad y/n$

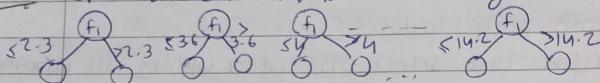
1 yes | 1 no

- 4) now for both tree we calculate Entropy and information gain and we select tree who has better entropy and information gain

- 5) this kind of division will happen for each x_i value (every value of f_1)

Step 1 Sort the data on basis of numerical column

Step 2 split the data on basis of every numerical values



Step 3 calculate entropy or information gain of each tree and select best tree on basis of entropy and IG.

Do

Step 4 Do it recursively for left and right nodes and finally we get decision tree.

\Rightarrow Splitter = best in random for best numerical column we choose those split which has high info gain

\Rightarrow Min Sample Split SPIT only when there is min n number of sample (rows)

\Rightarrow Min Sample Leaf If any node has min n sample then only split in every child

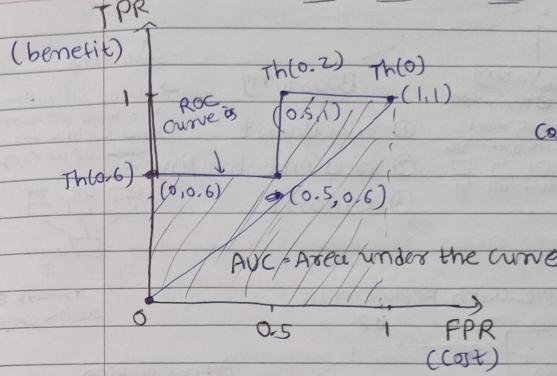
\Rightarrow Max Features = n every step choose n random features and split on it not split on every feature

\Rightarrow Max Leaf Nodes

\Rightarrow Min Impurity Decrease Parante entropy - weighted child entropy is $>$ min impurity then only split

(gini index)
gini impurity

ROC And AUC Curve



| Actual | 0 | 1 | Total |
|--------|----|----|----------|
| 0 | TN | FP | positive |
| 1 | FN | TP | actual |

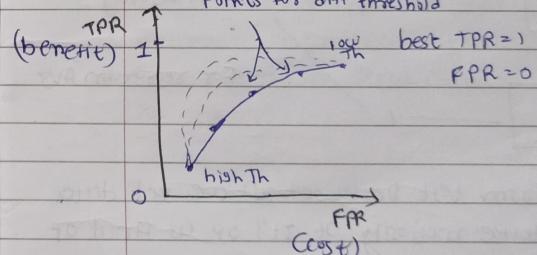
$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

\Rightarrow on base of what we actually need (high TPR or low or high FPR or low) we can select threshold value for given problem

\Rightarrow good model has always greater than Δ this area if area of model less than Δ that then our model is dumb model It guess randomly

Points for diff threshold



Ensemble Techniques

Combining multiple models

PAGE NO.:

BAGGING
make same but give diff data while training
(Bootstrap aggregation)
Ex random forest

Voting ensemble
→ same model are different
Ex: SVM LOR DT

Boosting

① Ada Boost

Stacking

→ same as voting but here we assign some weight for op of each model

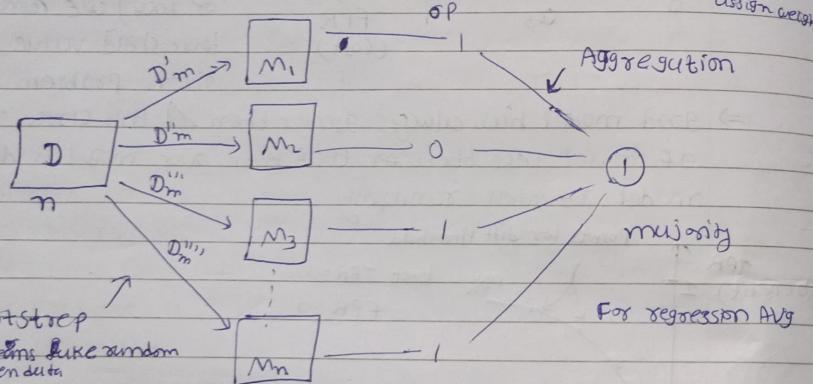
② gradient boosting model
③ XGBoost

Random forest

* Bagging

Raw sample with replacement

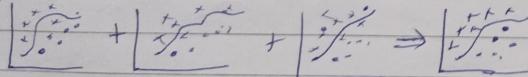
$m < n$



training

→ After ~~bootstrapping~~ is done from test Dataset If we get data to those models and take majority of it's op us final op

⇒ why Ensemble Techniques work



⇒ Voting & boosting can have diff weights but Bagging has same

* Random Forest

$D' \subset D$

$m < n$

$d' < d$

$RSTFS$

DT_1

M_1

I

M_2

I

M_3

I

D'

D_0

I

M_4

I

M_5

I

M_6

I

M_7

I

Row sample & Feature

sample with replacement

Decision trees

Decision

⇒ ~~Decision Tree~~ has low bias and high variance

⇒ Here we combine Decision trees so high variance convert into low variance

⇒ In regression problem we take op us Avg or mean or of all ops

- In case of Bagging we take equal importance of each vote
- but for AdaBoost we don't
- AdaBoost is made by adding multiple weak learners

PAGE NO.:

* ADA BOOST (Boosting technique)

$$f = \alpha_1 H_1 + \alpha_2 H_2 + \dots + \alpha_m H_m$$

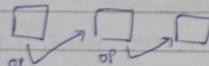
⇒ assign weights

- 2) Create decision stump for every feature column
- 3) Best decision stump (using entropy)
- 4) Find α (tree with only one split)
(max depth = 1)
- 5) ↑ weight whichever point incorrectly classified
and ↓ weight ↓ correctly classified
- 6) Same for $\alpha_2, \alpha_3, \dots, \alpha_m$
- 7) After getting eqn $f = \alpha_1 H_1 + \alpha_2 H_2 + \dots + \alpha_m H_m$ we can find OP directly by putting inputs in eqn

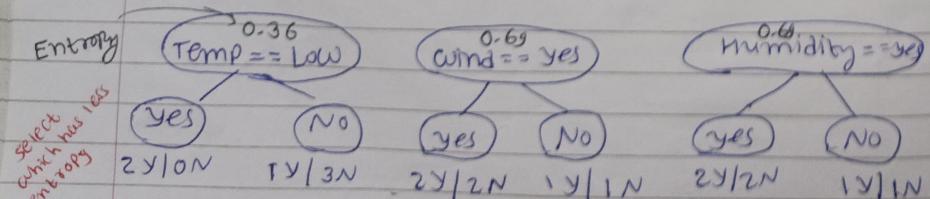
1) weak learner (accuracy is just over 50%)

2) weight

3) Dependency (give 1st op into 2nd, 2nd op to 3rd...)

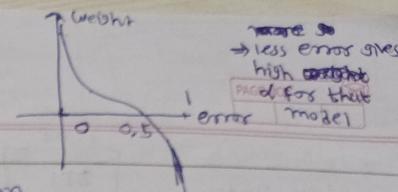


| Temp | Wind | Humidity | Play | sample weight | Updated weight | Normalized weight |
|------|------|----------|------|---------------|----------------|-------------------|
| High | No | Yes | N | 1/6 | 0.07 | 0.09 [0 - 0.09] |
| Low | Yes | No | Y | 1/6 | 0.07 | 0.09 |
| High | Yes | Yes | N | 1/6 | 0.07 | 0.09 |
| High | No | Yes | Y | 1/6 | 0.37 | 0.51 |
| Low | Yes | Yes | Y | 1/6 | 0.07 | 0.09 |
| High | Yes | No | N | 1/6 | 0.07 | 0.09 |



⇒ Now we select column which has low entropy as a weak learner and it's called base learner

$$f = \alpha_1 H_1 + \alpha_2 H_2 + \alpha_3 H_3$$



⇒ So we select Temp column

$$\text{Total error} = \frac{1}{6} \quad (\because \text{mode no has 1 yes and weight of that yes } 1/6)$$

$$\text{Performance of Stump} = \frac{1}{2} \log_e \left(\frac{1 - \text{TE}}{\text{TE}} \right)$$

$$\alpha_1 = \frac{1}{2} \log_e (5) = 0.8$$

⇒ now update the weight

Increase weights whose guess is incorrect and decrease all others weight.

$$\text{increase weight} \times e^{\text{perf}_{\text{stump}}} = \frac{1}{6} \times e^{0.8} = 0.37$$

$$\text{decrease} = \frac{1}{6} \times e^{-0.8} = 0.07$$

⇒ select random data from original dataset (using probability)

T W H P
H Y N N 1/6

H N Y Y 1/6

H N Y Y 1/6

L Y N Y 1/6

H N Y Y 1/6

L Y Y Y 1/6

zampling weight

[0.09 - 0.18]

[0.18 - 0.27]

[0.27 - 0.36]

[0.36 - 0.45]

[0.45 - 0.54]

[0.54 - 0.63]

[0.63 - 0.72]

[0.72 - 0.81]

[0.81 - 0.90]

[0.90 - 1]

T = L E = 0.36

W = Y E = 0.31

H = Y E = 0.22

YES NO

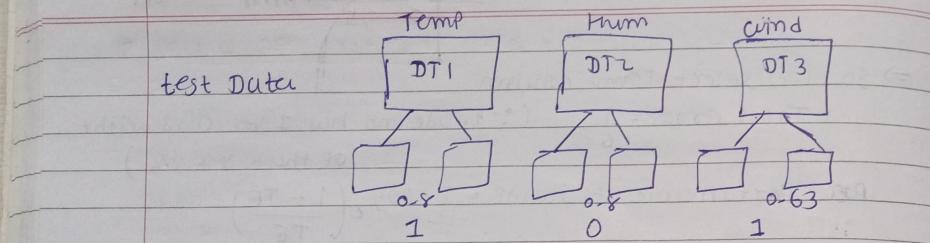
⇒ now we select Humidity column (\because low entropy)

$$\text{Total error} = 1/6$$

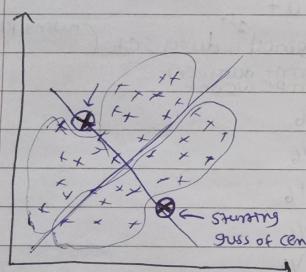
$$\alpha_2 = 0.8$$

⇒ Update weights and sume for this also

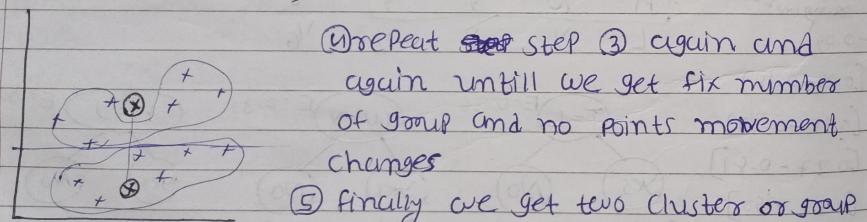
⇒ let assume we final get $\alpha_2 = 0.63$ so we get our final model $F = 0.8 H_1 + 0.8 H_2 + 0.63 H_3$



* K-means Clustering (unsupervised technique)
 ⇒ K-means clustering finds the similarity between the points and group them into cluster



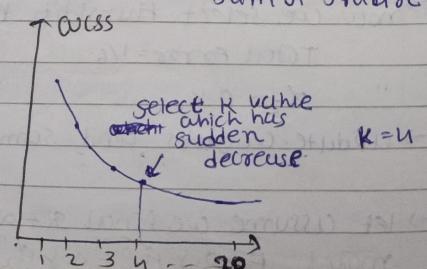
- ① K values \Rightarrow Centroids $K=2$
 $(\text{If } K=2 \text{ then } 2 \text{ clusters})$
- ② initialize two centroid randomly
- ③ select group & avg and update centroid to new avg value



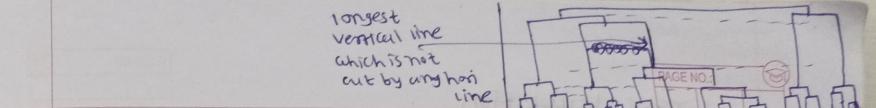
Elbow method

$$K = 1 \text{ to } 20$$

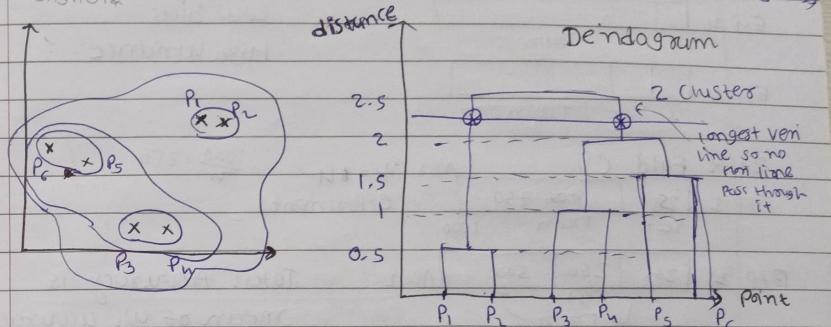
\Rightarrow WCSS is square sum of each point to its nearest centroid



WCSS = Within Cluster Sum of Square



* Hierarchical clustering intuition (unsupervised ml)



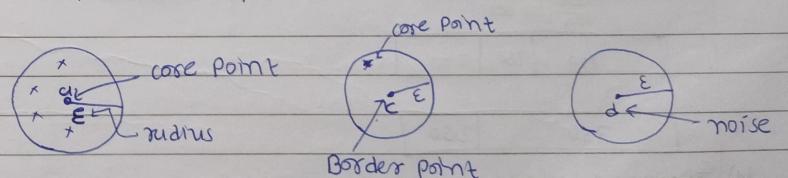
\Rightarrow This method use to know how many number of cluster should be use in order to classify points

* DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Core Point: at least $M_{in} Pts$ should be in range of core point

Border Point: does not satisfy min point condition but atleast one core point present in range of Border point

Noise Point: no point Present



Advantages: Is great at separating high density versus low density cluster, & great with handling outliers

Disadvantages: does not work well when clusters of varying density

\Rightarrow not work well with clusters of similar density

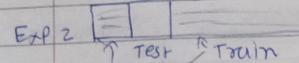
\Rightarrow struggles with high dimension data

Cross Validation

1) Leave one out cross validation (LOOCV)

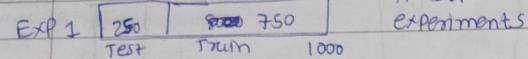


Low bias
High Variance



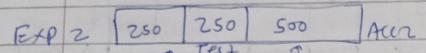
High Variance

2) K Fold CV

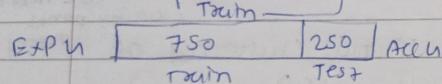
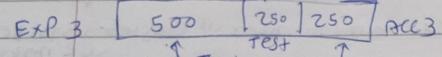


$$1000 \approx 250$$

$K=4$



Acc 2
Total Accuracy is
mean of all accuracy



3) Stratified CV

$$\frac{150}{1000} \text{ Yes } \frac{850}{1000} \text{ No}$$

$$\frac{250}{1000} \times 150 \text{ Yes } \frac{750}{1000} \times 150 \text{ Yes}$$

$$\boxed{250 \quad 750}$$

$$- \boxed{250 \quad 750 -}$$

$$\frac{250}{1000} \times 150 \text{ Yes}$$

$$\boxed{750 \quad 250 -}$$

$$\boxed{750 \quad 250}$$

\Rightarrow every time the test data
and train data contain
same ratio of numbers
of instance

4) Time series CV

Day 1 IN

Day 2 Day 1 D2 D3 D4 D5

Day 3

Day 4

Day 6 D1 D2 D3 D4 D5 D6 D7

Day 7

D1 D2 D3 D4 D5 D6 D7

O/P

D6

D7

DY

let
we know
need
Day 6-7
O/P

Bayes' Theorem

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

Probability of A when B is given

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

* Naive Baye's classifier

$$\begin{aligned} P(y|x_1, x_2, \dots, x_n) &= \frac{P(x_1|y) P(x_2|y) \dots P(x_n|y) P(y)}{P(x_1) P(x_2) \dots P(x_n)} \\ &= \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(\prod_{i=1}^n P(x_i))} \end{aligned}$$

$$P(y|x_1, x_2, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Example

outlook

| | yes | No | $P(y)$ | $P(N)$ |
|----------|-----|----|--------|--------|
| sunny | 2 | 3 | $2/5$ | $3/5$ |
| overcast | 4 | 0 | $4/5$ | $0/5$ |
| rainy | 3 | 2 | $3/5$ | $2/5$ |
| Total | 9 | 5 | 100% | 100% |

Temp

| | yes | No | $P(y)$ | $P(N)$ |
|-------|-----|----|--------|--------|
| HOT | 2 | 2 | $2/4$ | $2/4$ |
| MILD | 4 | 2 | $4/6$ | $2/6$ |
| COLD | 3 | 1 | $3/5$ | $1/5$ |
| Total | 9 | 5 | 100% | 100% |

$$\begin{aligned} P(y) &\text{ & } P(N) \quad \text{Assume Today(sunny, HOT) ref find prob} \\ \text{yes} & \quad 9/14 \quad P(\text{yes}|\text{today}) = P(\text{sunny}/\text{yes}) P(\text{hot}/\text{yes}) \\ \text{No} & \quad 5/14 \quad P(\text{Today}) \leftarrow \\ \text{Total} & \quad 14 \quad 100\% \quad P(\text{yes}|\text{Today}) = \frac{2}{9} \cdot \frac{2}{4} \cdot \frac{9}{14} \xrightarrow{\text{const}} 0.031 \end{aligned}$$

$$P(\text{No}|\text{Today}) = 3/5 \cdot 2/5 \cdot 5/14 = 0.08571$$

$$\begin{aligned} \text{normalize } P(\text{yes}|\text{Today}) &= \frac{0.031}{0.031 + 0.08571} = 0.27 \\ P(\text{No}|\text{Today}) &= 1 - 0.27 \\ &= 0.73 \end{aligned}$$

So, Today person will Not play

How to apply Naive Baye's classifiers on Text data (NLP)

$$f_1 \ f_2 \ f_3 \ f_4 \ 0/p$$

The food Delicious Bud

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 1 |
| 1 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 0 | 0 | 0 | 1 | 0 |

In sentence many words $[x_1, x_2, x_3, \dots, x_n]$

$$\begin{aligned} P(y=\text{yes} \mid \text{sentence}) &= P(y=\text{yes} / x_1, x_2, \dots, x_n) \\ &\propto P(y) \prod_{i=1}^{n=3} P(x_i | y=\text{yes}) \cdot P(x_3 | y=\text{yes}) \\ &\propto 2/5 \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \\ &\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \\ &P(y) \quad P(x_1 | y=\text{yes}) \quad P(x_2 | y=\text{yes}) \\ &= 1/10 \\ &= 0.1 \end{aligned}$$

$$P(\text{No} | y=\text{No} / x_1, x_2, x_3)$$

→ For Numerical data

height weight gender

Find Prob of being male given { $h=185, w=170$ }

we have to

$$P(M | h=185, w=170) \neq$$

$P(F | h=185, w=170)$ and which ever is large (\therefore we ignore denominator because we need to calculate both prob)

| | | |
|-----|-----|---|
| 172 | 150 | M |
| 180 | 180 | M |
| 165 | 165 | M |

| | | |
|-----|-----|---|
| 139 | 110 | F |
| 145 | 150 | F |
| 160 | 130 | F |

$$P(F | h=185, w=170) = P(h=185 | M) P(w=170 | M) P(M)$$

assume normal distribution so we find like of all males H & W and by putting $x = H$ or W of given que. and find prob $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$

base model + R_n
for base model
more precise
For DT 1
For DT 2

SUM kernels

For degree d polynomials the polynomial kernel is

$$K(x,y) = (x^T y + c)^d$$

when $c=0$ kernel is called homogeneous

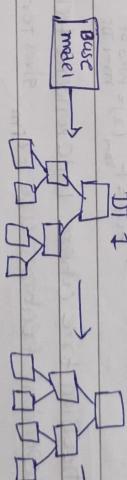
features
 x_1
 x_2
 y
dip

$$f(x_1, x_2) = (x_1^T x_2 + c)^d$$

$$= \begin{bmatrix} x_1^T & x_2^T \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$\text{Base model} = \frac{50+70+80+100}{4} = 75$$

Base model give dlp for all inputs



Step 1) Base model

- 2) compute Residual, Error, Pseudo Residual
- 3) Construct DT input Degree, exp & op (not salary)

$$F(x) = h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_m h_m(x)$$

$$\text{Op} \Rightarrow \text{Base model} + \sum_{i=1}^m \text{DT}_i + \text{hyper parameter}$$

Salary (base model) (not salary) \rightarrow 0.12 (0.61) 0.01 (0.01)
Degree (base model) (not salary) \rightarrow 1.00 (0.91) 0.71 (0.71)
Exp (base model) (not salary) \rightarrow 0.00 (0.00) 0.00 (0.00)

Algorithm Input: (x_i, y_i) , $L(y, F(x))$, No of Trees

Step 2) Initialize model with constant value

$$F_0(x) = \arg \min_V \sum_{i=1}^n L(y_i, V)$$

Loss function
 $\frac{1}{2}(y - \hat{y})^2$

We have to find $V = \hat{y}$ so that loss is min

$$f(x) = f_0(x) + f_1(x) + \dots$$

(2)

Step 2) Iterate $M=1$ to M :

i) compute so-called pseudo-residuals

$$\hat{x}_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right] \quad \text{for } i=1, \dots, n$$

$$F(x) = F_{m-1}(x) + \text{path of } \frac{\partial L}{\partial F}$$

$$\hat{x}_{im} = -2 \frac{(y_i - \hat{y}_i)(-1)}{2} = y_i - \hat{y}_i$$

loss = $\sum \frac{1}{2}(y_i - \hat{y}_i)^2$
 now create decision tree where independent features have to find less min is x_i and dependent features are \hat{x}_{im} gives terminal regions

$$V = \frac{\sum y_i}{n}$$

for less = more

more

less

$$F_m(x) = F_{m-1}(x) + \arg \min_{V_m} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + V_m)$$

Find V such that V_m is $\arg \min_{V_m} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + V_m)$

terminal regions

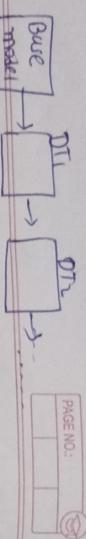
removal

$$F_m(x) = F_{m-1}(x) + V_m h_m(x)$$

old value
 value from DT

Step 3) output $F_m(x)$

$$f_0(x) + f_1(x) + \dots + f_m(x)$$



PAGE NO.:

PAGE NO.:

XG BOOST
dp=0

New residual
probabi
0.6
0.4

| Classification | Salary | Credit Approval | Residual | probabi | residual |
|----------------|---------|-----------------|----------|-----------------|-----------------|
| ≤ 50K | Good | Bad | 0 | 0 - $p_x = 0.5$ | 0.6 - 0.5 = 0.1 |
| ≤ 50K | 6 | 1 | 0.5 | 0.3 | 0.7 - 0.3 = 0.4 |
| > 50K | B | 0 | -0.5 | 0.7 | -0.7 |
| > 50K | 6 | 1 | 0.5 | 0.2 | 0.8 |
| ≤ 50K | Nominal | 1 | 0.5 | 0.1 | 0.9 |
| ≤ 50K | N | 0 | -0.5 | 0 | 0 |

Op is either 0 or 1 and since this is classification problem probability of 0 or 1 $p_x = 0.5$

[-0.5, 0.5, 0.5, -0.5, 0.5, -0.5]

base model
 $p_x = 0.5$

[≤ 50K Salary > 50K]
[-0.5, 0.5, 0.5] [-0.5, 0.5, 0.5]

⇒ For xgboost we need to classify into binary branch if there are two then more category in salary we need to classify into two child only

$$\text{Similarity weight} = \frac{\sum_{i=1}^n \text{Pr}(C_i \mid \text{Residual}_i)^2}{\sum_{i=1}^n \text{Pr}(C_i \mid \text{Residual}_i)}$$

In this problem let

similarity weight

$$\text{for left side} = \frac{[-0.5 + 0.5 + 0.5 - 0.5]^2}{[0.5(0.5) + 0.5(0.5) + 0.5(0.5) + 0.5(0.5)]} = \frac{0}{1} = 0$$

≤ 50K

$$\text{similarity weight} = \frac{(-0.5 + 0.5 + 0.5)^2}{[0.5(0.5) + 0.5(0.5) + 0.5(0.5)]} = \frac{0.25}{3} = \frac{1}{3}$$

for right side

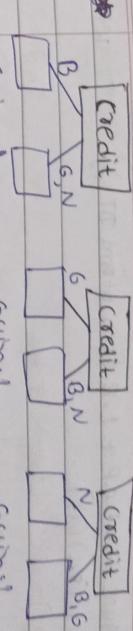
$$= \frac{[0.5(0.5) + 0.5(0.5) + 0.5(0.5)]}{0.75}$$

> 50K

$$\text{Sm for parent} = \frac{0.25}{7(0.25)} = 0.14$$

salary

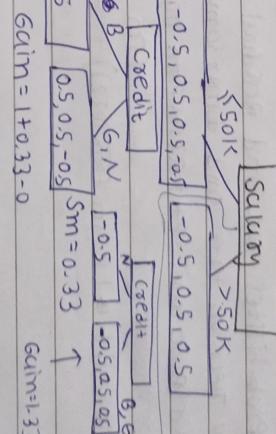
$$= \log \left(\frac{0.5}{1-0.5} \right) = 0$$



Credit
B
G,N
B,N
N
G,C
B,C

Gain=1
Gain=0

Let assume salary give highest gain among salary and credit so salary is best way of constructing tree now we are doing split w.r.t. credit



Gain=1.33
Gain=0.33

So we can select any of one



Post Prunning

Cover value = $Pr(1-Pr)$

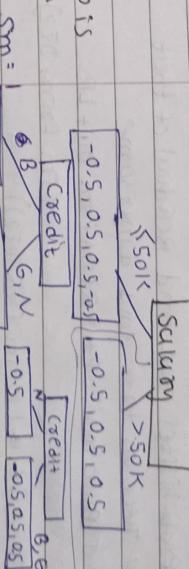
In this problem $Pr = 0.5$
hyperparameter $\lambda = 0$

$$\text{Sm} = \frac{1}{5} \quad \begin{cases} B \\ N \end{cases}$$

Sm = 1

$$\begin{cases} \text{Credit} \\ \text{Credit} \end{cases}$$

$$\begin{cases} G,N \\ B,N \end{cases}$$



Post Prunning
Cover value = $Pr(1-Pr)$

Cover value = 0.25

In this problem $Pr = 0.5$
hyperparameter $\lambda = 0$

$$\begin{cases} \text{Credit} \\ \text{Credit} \end{cases}$$

$$\begin{cases} G,N \\ B,N \end{cases}$$

If gain is less then

Cover value then we cut

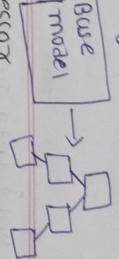
Bgain = Sm Left child + Sm right child - Sm of parent

$$Guin = 0 + 0.33 - 0.14 = 0.21$$

More we traverse tree according to
Input data and find similarity weight
at current node

| |
|-----------|
| Page No.: |
| |
| |

$$Op = \text{avg} = 51K$$



* Regressor

| | Exp | Gdp | Salary | Residual | Op | Residual |
|-----|-----|-----|--------|----------|-------|----------|
| 2 | Yes | 60K | -11 | 6 | 6 | -6 |
| 2.5 | Yes | 62K | -9 | 6 | 6 | -6 |
| 3 | No | 52K | 15.5 | 53.5 | 52.5 | -1.5 |
| 4 | No | 60K | 9 | 62 | 60.5 | -2 |
| | | | | | 62-63 | -1 |
| 5.5 | Yes | 62K | 11 | 63 | - | - |

After finding all Probabilities we find residual of that

and creating new decision tree and doing same

now $Pr = \sigma(\text{base model } Op + LR(\text{sm from DT1}) + LR(\text{sm from DT2}))$

$$Pr = \sigma(0 + 0.1(1) + LR(\text{sm from DT2}))$$

$$Sw = \frac{110}{14} \quad \text{let } \lambda = 1 \quad Sw = (-9 + 1, 9 + 11)^2 = 28.8$$

$$\text{Gain} = 65.5 + 28.8 - 0.16$$

$$\boxed{\text{Gain} = 93.84}$$

$$\boxed{\text{Exp}} \quad Sm = 0.166$$

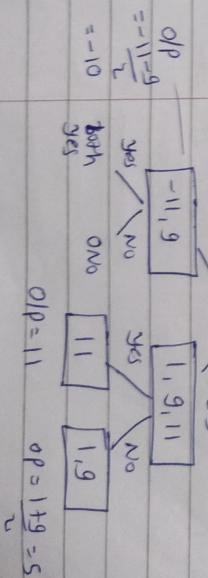
$$\boxed{\text{Exp}} \quad > 2.5$$

$$Sm = 133.33 \quad \boxed{-11, -9} \quad Sm = 110.25$$

$$\boxed{\text{Gain} = 113.42}$$

similarly we find gain for all and choose max gain split

Let assume above tree ($\leq 2.5, > 2.5$) is best



| |
|-----------|
| Page No.: |
| |
| |

| |
|--|
| |
| |
| |

Tensors

Now, we pass the data {2, yes} into Decision Tree

and get value of output

$$* \quad OLP_2 = 51 + (0.5) [-10] = 46$$

$$\begin{matrix} \text{Base} \\ \uparrow \\ LR \end{matrix} \quad \begin{matrix} \text{Op from} \\ \text{DT}_1 \end{matrix}$$

After finding OLP_2 column we find residual 2 and make corner tree DT_2

$$OLP_2 = 51 + LR \bar{DT}_1 + LR \bar{DT}_2$$

$$\begin{matrix} \text{Base model} \\ \uparrow \\ \text{Op from DT}_1 \end{matrix} \quad \begin{matrix} \text{Op from} \\ \text{DT}_2 \end{matrix}$$

\Rightarrow There is hyper parameter r has some value node which has Information < r we pass down (cut that sub tree)

- * 0-D Tensor / scalar

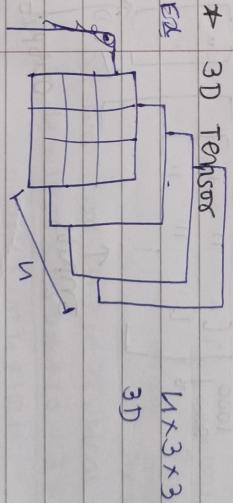
Ex 2

* 1D tensor / vector
 $\begin{bmatrix} 1, 2, 3, 4, 5 \end{bmatrix} \rightarrow 1D \text{ tensor}$
 Collection of 0D
 Tensors / scalars 1D array / matrix

Ex

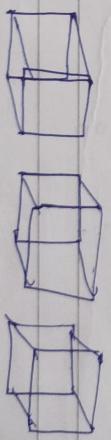
* 2D tensor / matrices
 $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ collection of 1D tensors / vectors

* 3D tensor
 $n \times 3 \times 3$
 3D

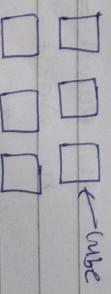


Collection of 2D tensor

* nD tensor



* 5D tensor

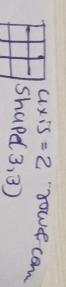


collection of 3D tensor

collection of nD tensor

\Rightarrow In machine learning there are mostly 0-5D tensor used

No of axis = rank = No of dimension



process command to render into Blue (root) environment
conk, but deactivate → A spider can run F5
conda activate envname: enter into `envname` environment

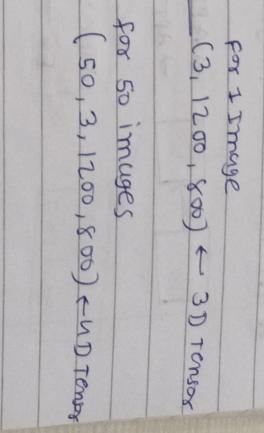
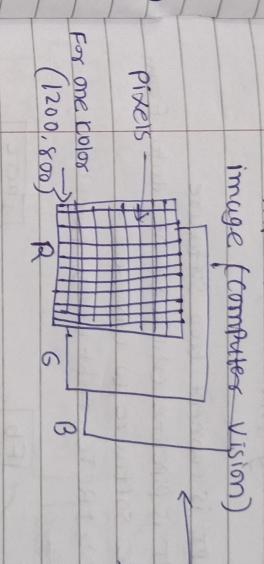
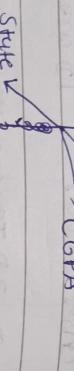
PAGE NO:

DATE:

Ex of 1D tensors

| CGPA | ID | state | placement |
|------|-----|-----------------|-----------|
| 8.0 | g1 | w ^B | 1 |
| 7.2 | 102 | K ^{S1} | 0 |
| | | | ? |

each Row is 1D tensor

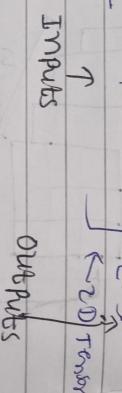


Example of 5D Tensor

Videos = collections of images
Let assume we have video of 60 sec

30fps

WxHxP(WxHxT x C)
1800x1200x30x3x3



Ex of 3D tensors

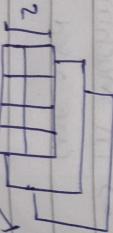
NLP = natural language processing

Hi Nitish | Ruhul | Ankit
Hi Nitish ⇒
Hi Ruhul ⇒
Hi Ankit

$$\left[\begin{bmatrix} [1, 0, 0, 0, 0] \\ [1, 0, 0, 0] \\ [1, 0, 0, 0] \end{bmatrix}, \begin{bmatrix} [0, 1, 0, 0] \\ [0, 1, 0, 0] \\ [0, 1, 0, 0] \end{bmatrix} \right]^*$$

time series data

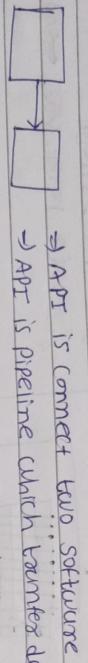
highest lowest



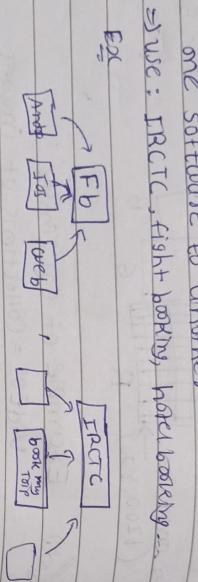
| Day 1 | - | - | For 1 year [365, 2] → 2D tensor |
|---------|---|---|---------------------------------|
| Day 2 | - | - | For 1 year [365, 2] → 2D tensor |
| Day 365 | - | - | For 1 year we have 102D tensor |

3D tensor

★ API (Application programming interface)



⇒ API is connect two software
⇒ API is pipeline which transfer data from one software to another

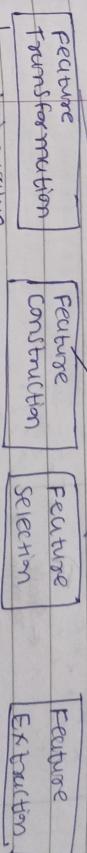


⇒ use: IRCTC, flight booking, hotel booking...

Ex:

| | |
|-----------|-------|
| PAGE NO.: | _____ |
|-----------|-------|

★ Feature Engineering



- ⇒ missing value handling
- ⇒ categorical features
- ⇒ outlier detection
- ⇒ function transformer
- ⇒ Feature Scaling

* Feature scaling - Standardization

⇒ Standardization

$$Age(x_i)$$

$$x'_i = \frac{x_i - \bar{x}}{\sigma} \quad \text{mean}$$

$$\bar{x} \leftarrow x_1$$

$$15 \leftarrow x_2$$

$$33 \leftarrow x_3$$

$$63 \leftarrow x_4$$

$$90 \leftarrow x_5$$

mean
of salary

$\sigma = 5$

$\sigma = 1$

500 values

which to use

⇒ K means

⇒ K nearest neighbours

⇒ PCA

⇒ Artificial neural network Share is same

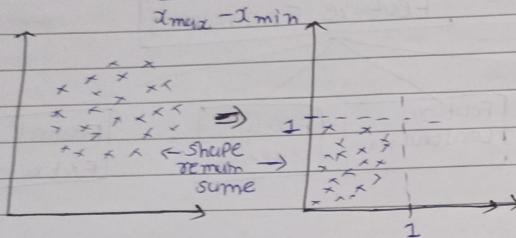
⇒ Gradient descent

| | |
|-----------|-------|
| PAGE NO.: | _____ |
|-----------|-------|

* Features Scaling - Normalization

⇒ Min-Max scaling Normalization

$$x'_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \quad 0 < x'_i < 1$$



⇒ Mean normalization

$$x'_i = \frac{x_i - \bar{x}_{\text{mean}}}{x_{\max} - x_{\min}} \quad -1 < x'_i < 1$$

mean centering

⇒ Not given direct class/function
in SK learn

⇒ used when we need centered data

⇒ Max Abs scaling

$$x'_i = \frac{x_i}{|x_{\max}|} \quad \text{SK learn: MaxAbsScaler}$$

⇒ used when we have sparse data (data which contain more zeros)

⇒ Robust scaling

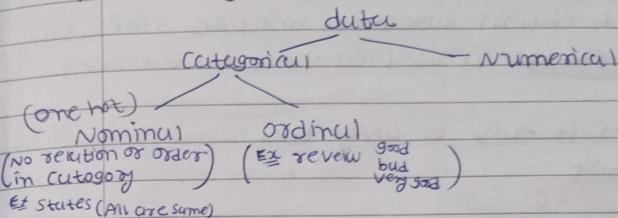
$$x'_i = \frac{x_i - \bar{x}_{\text{medium}}}{IQR}$$

$IQR = 75^{\text{th}} \text{ value} - 25^{\text{th}} \text{ value}$

⇒ used when data contain more outliers

Lable encoding: work same as Ordinal encoder but we when o/p lable(y) are given as category

* Encoding Categorical variables



⇒ ordinal Encoding (For ordinal data)

| | Education | |
|----------------|-----------|--------------|
| High school | HS → 0 | PG > UG > HS |
| under graduate | UG → 1 | ↑ ↑ ↑ |
| post gradu. | PG → 2 | PG → 2 |
| | PG → 2 | UG → 1 |
| | | HS → 0 |

⇒ one hot Encoding (For nominal data)

| Color | Target | Color_Y | Color_B | Color_R | Target |
|-------|--------|---------|---------|---------|--------|
| Y | 0 | 1 | 0 | 0 | 0 |
| PY | 1 | 1 | 0 | 0 | 1 |
| B | 1 | 0 | 1 | 0 | 1 |
| Y | 1 | 1 | 0 | 0 | 1 |
| R | 1 | 0 | 0 | 1 | 1 |
| Y | 0 | 1 | 0 | 0 | 0 |

[1, 0, 0] → yellow

[0, 1, 0] → blue

[0, 0, 1] → red

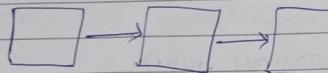
Multicollinearity: There should be no relation between Independent variables.

→ So we remove one of column

* Pipelines

⇒ Pipelines which chains together multiple steps so that O/P of each step is used as input to the next step.

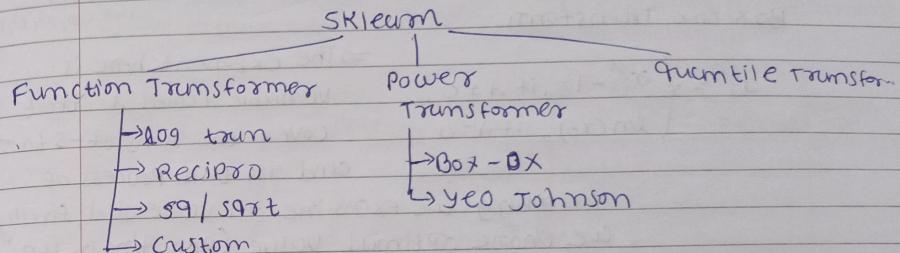
⇒ Pipeline makes it easy to apply the same preprocessing to train and test



⇒ Transformation is used for linear Regression, logistic regression
⇒ Decision Tree have no effect on Transf.

* Function Transformer

→ Using transformer we can convert our data into normal distribution



How to find data is normal or not?

use: sns.distplot OR pd.skew() OR QQ Plot

log transforms

Age

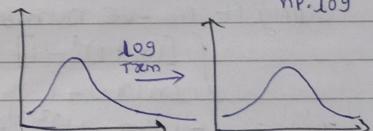
$$21 \rightarrow \log 21$$

$$35 \rightarrow \log \frac{10}{35}$$

80

⇒ log transform use when data is right skew

np.log



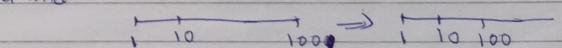
⇒ when data contains zero values

then we use np.log1p

which add 1 to data and ⇒

then take log.

$$\log(1+x)$$



Reciprocal Transforms

$$\frac{1}{x}$$

⇒ convert small values to big and big to small

Sq Transforms x^2
use for left skewed

$$\sqrt{x}$$

★ Power Transformer

Box Cox Transform

$$x_i^{(\lambda)} = \begin{cases} x_i^\lambda - 1, & \text{if } \lambda \neq 0, \\ \ln(x_i), & \text{if } \lambda = 0 \end{cases}$$

This is general transform for all
like $x_1^1, x_1^2, x_1^3, \dots$

\Rightarrow The exponent here is

variable called λ that
cover the range of -5 to 5
and in this process of
searching, we examine all value of λ , finally
we choose optimal value (resulting in the
best upto a normal distn) for your variable

\Rightarrow Applicable only for $x_i > 0, \forall i$

Yeo-Johnson Transform

\Rightarrow This transformation is somewhat somewhat of an
adjustment to Box-Cox transformation, by which we can
apply it to -ve numbers.

$$x_i^{(\lambda)} = \begin{cases} [(x_i + 1)^{\lambda} - 1] / \lambda & \text{if } \lambda \neq 0, x_i > 0 \\ \ln(x_i + 1) & \text{if } \lambda = 0, x_i > 0 \\ -[(-(x_i + 1)^{2-\lambda} - 1)] / (2-\lambda) & \text{if } \lambda \neq 2, x_i < 0 \\ -\ln(-x_i + 1) & \text{if } \lambda = 2, x_i < 0 \end{cases}$$

* Principle component Analysis (PCA) Features extraction algorithm

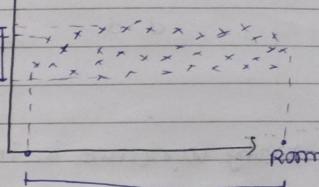
\Rightarrow PCA used to convert higher dimensional data into
lower dimensions. Benefits: faster execution of algorithms
& easy to visualize

\Rightarrow Geometric intuition

| No of Rooms | No of grocery shop | Price |
|-------------|--------------------|-------|
| 3 | 2 | 60 |
| 4 | 6 | 130 |
| 5 | 6 | 180 |
| 7 | 10 | 90 |

grocery shops

feature selection
(not PCA)



$d > d'$ so we select x-axis

which has more variance
Here Rooms has more variance

| No of rooms | No of washrooms | Price |
|-------------|-----------------|-------|
| 1 | 1 | 60 |
| 2 | 1 | 130 |
| 3 | 1 | 180 |
| 4 | 2 | 90 |
| 5 | 2 | 100 |
| 6 | 2 | 110 |
| 7 | 3 | 120 |
| 8 | 3 | 130 |
| 9 | 3 | 140 |
| 10 | 4 | 150 |

\Rightarrow Here rooms and washrooms
are equally important feature
so here feature selection is
not important

\Rightarrow So we can not decide
(mathematically also) which
is important feature

\Rightarrow In this kind of case
feature extraction is used

\Rightarrow Since both features are equally
important we create new feature
called size which tells us about rooms and washrooms.

\Rightarrow So, we use size instead of rooms & washrooms. This is
what PCA does

\Rightarrow PCA makes new set of features from existing features and
select subset from new features which is more important

$d \approx d'$

\Rightarrow So, we use size instead of rooms & washrooms. This is
what PCA does

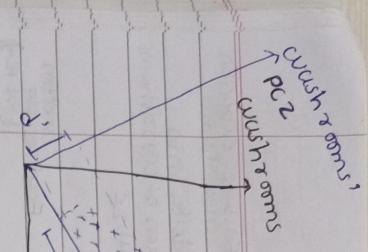
\Rightarrow PCA makes new set of features from existing features and
select subset from new features which is more important

$$\text{proj}_{\vec{a}} \vec{b} = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}|^2} \vec{a}$$

$$\text{proj}_{\vec{a}} \vec{b} = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}|^2} \vec{a}$$

PAGE NO.:

PAGE NO.:



$d > d'$

rooms

\Rightarrow Why variance is important

$$\text{mean} = -5 + 0 + 5 = 0$$

$$\sigma^2 = \frac{(-5 - 0)^2 + (0 - 0)^2 + (5 - 0)^2}{3} = \frac{50}{3}$$

$$\sigma = \sqrt{\frac{50}{3}}$$

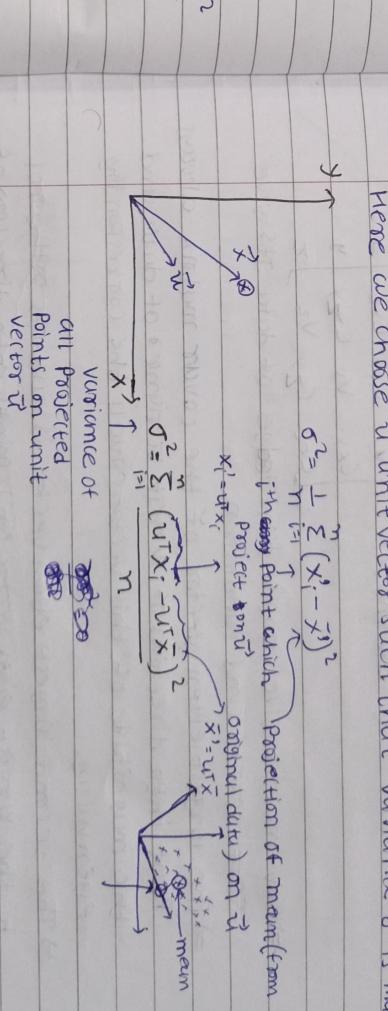
$$\text{mean} = 0$$

$$\text{variance} = \frac{25 + 0 + 25}{3} = \frac{50}{3}$$

$$\sigma^2 = \frac{50}{3}$$

$$\text{mean} = 0$$

$$\begin{array}{c} \text{mean} = 0 \\ \sigma^2 = 200/3 \\ \sigma = \sqrt{200/3} \\ \sigma^2 \text{ is not spread but} \\ \sigma^2 \text{ is spread} \end{array}$$



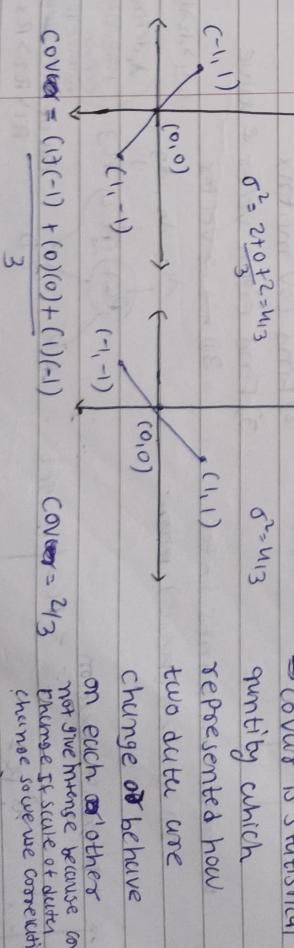
Here we choose \vec{u} unit vector such that variance σ^2 is max

\Rightarrow PCA select \vec{u} such that σ^2 is maximum
 \Rightarrow variance is only measured in 1D. while covariance give us 2D relations. (relation between x & y)

Covariance and Covariance matrix

$$\text{mean abs deviation} = \frac{\sum |x_i|}{n} \quad \text{optimization algo in PCA.}$$

\Rightarrow while square function is differentiable so we prefer σ^2 instead of MAD. to denote spread of data



$$\text{Cov}(x,y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Pearson Correlation} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} = -1.601$$

carrying this linear transformation of vector V

$$\text{cov}(x_1, x_1) = \text{var}(x_1)$$

$$\text{cov}(x_1, x_2) = \text{cov}(x_2, x_1)$$

PAGE NO.:

⇒ eigen value define how that particular one principle component define or explain variance of whole data
 $\Rightarrow \lambda$ is also called explain variance.
 $\lambda_1 \times 100\%$ how many percent....

2D transformation

$$(1000, 3) \cdot \begin{pmatrix} 3 & 2 \\ 2 & 3 \end{pmatrix} \rightarrow (1000, 2)$$

PAGE NO.:

$$\left[\begin{array}{c|cc} x_1 & x_2 & \\ \hline x_1 & & \\ x_2 & & \\ \vdots & \vdots & \end{array} \right] \cdot V = \lambda \cdot V$$

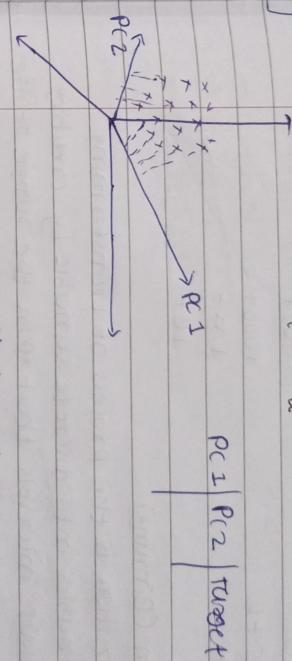
↑ even value

$$\text{cov matrix} = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix}$$

$$= \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & \text{var}(x_2) \end{bmatrix}$$

For 3 column

$$\text{cov matrix} = \begin{bmatrix} V_x & C_{xy} & C_{xz} \\ C_{xy} & V_y & C_{yz} \\ C_{xz} & C_{yz} & V_z \end{bmatrix} Z$$



⇒ cov matrix give complete knowledge of data. It's give spread and orientation.

⇒ The largest eigen vector of the covar matrix always

point into direction of largest variance of data and the magnitude of this vector equals the corresponding eigenvalue.

⇒ The second largest eigen vector is always orthogonal to the largest eigen vector and points into direction of the second largest spread of the data.

⇒ step by step solution

$$f_1 \quad f_2 \quad f_3 \quad \text{first 3 column}$$

of features

1) mean centering

2) Find cov matrix

3) Find eigen vector & value

v) transform points into

1D or 2D (here) \bullet dot product

$$1D \rightarrow 3 \text{ vector}$$

$$A_1, A_2, A_3 \quad (\lambda_1 > \lambda_2 > \lambda_3)$$

↑ 1D transformation

$$\text{PC1} \quad \left| \begin{array}{c} \text{target} \\ \text{1D} \\ \text{2D} \end{array} \right. \quad \text{dot Product}$$

$$\text{PC1} \quad \left| \begin{array}{c} \text{target} \\ \text{1D} \\ \text{2D} \end{array} \right. \quad \text{result}$$

cov of

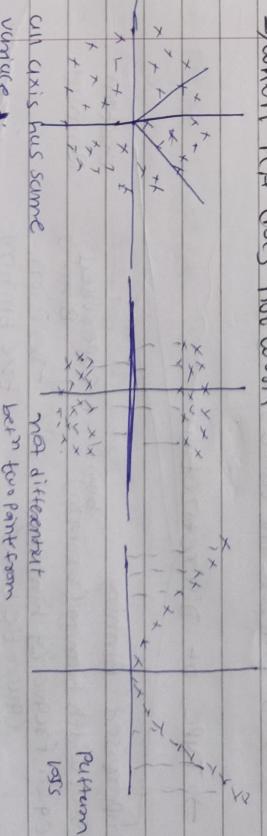
$$\text{PC1}, \text{PC2}, \text{PC3} \quad \text{explain variance}$$

data

$x_i^T \cdot x_i, \forall i$

$$(1000, 3) \cdot (3, 1) \rightarrow (1000, 1)$$

shape(1,3) (1000, 3)



Binning (Discretization) and Binarization

⇒ Encoding Numerical Features

No of downloads

23

1734574

9927

bins
100+
1000+
1M+
1B+

★ Discretization (Binning)

Discretization is the process of transforming continuous variables into discrete variable by creating set of contiguous intervals that span the range of the variable's values. Discretization is also called binning where bin is an alternative name of interval.

why use Discretization

⇒ To handle outliers

⇒ To improve the value spread

⇒ Type of Discretization

unsupervised Binning
equal width binning (uniform)
equal frequency binning (quantile binning)
K mean binning

supervised Binning
Decision tree binning

custom Binning

i) equal width binning / uniform binning

$$\frac{\text{max} - \text{min}}{\text{bins}}$$

⇒ Can handle outliers

⇒ No change in spread of data

ii) Equal frequency / quantile Binning

⇒ If intervals is 10 then each interval contain 10% of total observations

⇒ width of each intervals are may different

⇒ make value spread uniform

⇒ can handle outliers

iii) K means binning



⇒ also used is same as K mean clustering

⇒ This binning is used when data is spna in form of clusters

⇒ Encoding and Discretized variable

SKlearn

↳ KBinsDiscretizer()

bins=?
strategy
uniform
quantile
K means

encoding
ordinal
one hot encoding

iv) custom / Domain Based Binning

base on your domain knowledge you can do binning

[0-18] → Kids

[18-60] →

[60-80] →

★ Binarization (special case of binning)

→ Ex For black and white pixel $0 - 127.5 = 0$

$$127.5 - 255 = 1$$

→ Ex salary < 6L then 0

> 6L then 1

Mix Data

⇒ Ex 1 cabin

B5 ⇒ category numbers

C23 B 5

D41 C 23

D 41

⇒ If we assume cabin column is categorical then there are so many categories so we need to split it into category and numbers

⇒ Ex 2 num category

7 7 NA

3 3 NA

A NA A

C NA C

Handling Missing Data

PAGE NO.:

PAGE NO.:

Missing Values

Remove them
numerical
categorical
mean/median
random
end of distribution
Arbitrary
 $Q_1 - 2.5 \times IQR$
 $Q_3 + 2.5 \times IQR$

⇒ can do only help
of that column
or no need other
columns

Simple imputer()
Impute
univariate
multivariate
KNN
imputer
(MICE)

⇒ Arbitrarily value imputer:
Replace fill missing data with arbitrary random values
like 1,99,1009, missing, 0,-1,5,8

⇒ by giving random values we create difference betw where
data is available and where not.

⇒ when data is not missing randomly and we cant use CCA
then we can use Arbitrarily imputer

⇒ discarding observations where values in any of the variables

(row) are missing

⇒ complete case analysis (CCA)
CCA also called "list-wise deletion" of cases consist in

disregarding observations where values in any of the variables

(column)

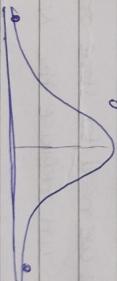
Observation for which there is information in all the
variables in data set

Assumption for CCA: missing Completely random

→ < 5% data missing

⇒ you have to ensure that before and after data removing
distribution of data is same.

⇒ when data is not missing randomly



$Q_1 = 25^{th} \text{ percentile}$

a → mean + 3σ

c → $Q_1 - 1.5 \times IQR$

$IQR = Q_3 - Q_1$

b → mean - 3σ

d → $Q_3 + 1.5 \times IQR$

skewed right

same as missing values
not random

⇒ when data is not missing randomly

⇒ Mean / Median ADV: simple DISADV: changing distribution
→ outliers come in picture (mean/median)
→ covariance or (cov) with other
column change

→ use only when data missing
completely at random & < 5% data missing

★ Handling Categorical missing data

⇒ most frequent value imputation

(mode)
fill NA with most frequency value in that column

Assumption: missing completely at random (MCAR)

mode should be much more w.r.t others
missing < 5%

⇒ missing category imputer (some like arbitrary then can use values so we use 'missing' to indicate that this is missing)
fill NA with 'missing' so that machine learning algo know that those data are missing

★ Random Imputation (Numerical & Categorical)

⇒ select random number from data and fill it with it

⇒ It can be applied for both numerical and categorical

⇒ distribution and variance is not change much

⇒ covariance matrix is change and this change is random
so it's affect on covariance matrix

⇒ memory heavy for deployment, as we need to store original training set to extract value from it and replace the NA in coming observation

⇒ well suited for linear models as it does not distort the distribution, regardless of % of NA

★ Missing Indicators

⇒ Ex Age Fare Age-NA

| | | |
|----|----|---|
| 18 | 32 | F |
| 20 | 35 | F |
| 19 | 32 | F |
| NA | 35 | T |

⇒ adding Age-NA can

help model to differentiate missing values now so model performance is increase

→ This method is not work every time but it's work sometimes

★ KNN Imputer

→ 1) find K nearest neighbour

2) find the value

| S.NO | Feature 1 | F2 | F3 | F4 |
|------|-----------|----|----|----|
| 1 | 33 | — | 67 | 21 |
| 2 | — | 45 | 68 | 12 |
| 3 | 23 | 51 | 71 | 18 |
| 4 | 40 | — | 81 | — |

to find value in S.NO 2 let consider S.NO 2 & 3

$$d = \sqrt{w((45-51)^2 + (68-71)^2 + (12-18)^2)}$$

Total number of coordinates
number of present coordinates

let S.NO 2 & 4

$$d = \sqrt{w((45-51)^2 + 67^2)}$$

$$d = \sqrt{w[(68-81)^2]}$$

MCAR: missing completely at random
 MAR: missing value can fill by other column (there is relation)
 MNAR: " " cannot "

For now no 2 d from point 1 is $d_1 = 11.29$
 from point 3 $d_3 = 9$
 $d_5 = 5.19$

$$d_5 = 5.19$$

Let consider K=2 so for K=2 we ~~find~~ closest 2 Point
 1st column for value of 2nd raw 1st column value

$$\frac{23+10}{2} = 31.5$$

$$\text{If } K=3 \text{ then } \frac{23+10+33}{3} = 96/3$$

Adv: more accurate

Disadv: more calculation

training set is on server so memory problem

$$\Rightarrow \text{for uniform weight} \cdot w_2 = \frac{23+10}{2}$$

Step 1: Fill all NaN with mean of that column

| | | |
|-----|-----|-----------|
| R&D | Adm | Marketing |
| 8 | 15 | 30 |
| Nan | 5 | 20 |
| 15 | 10 | 11 |
| 12 | Nan | 26 |
| 2 | 15 | Nan |

Step 2: Remove all column 1

REDA Adm Marketing Step 3: find missing values here 9.25
 8 15 30 missing values here 9.25
 15 10 11 Step 3: find missing values here 9.25
 12 11.25 26 (left to right)
 2 15 29.25

REDA Adm Marketing Step 3: find missing values here 9.25
 8 15 30 missing values here 9.25
 15 10 11 Step 3: find missing values here 9.25
 12 11.25 26 (left to right)
 2 15 29.25

| | | |
|----|-------|-------|
| 8 | 15 | 30 |
| 15 | 10 | 11 |
| 12 | 11.25 | 26 |
| 2 | 15 | 29.25 |

x = training dataset input
 y = testing dataset
 & predict on

Outliers

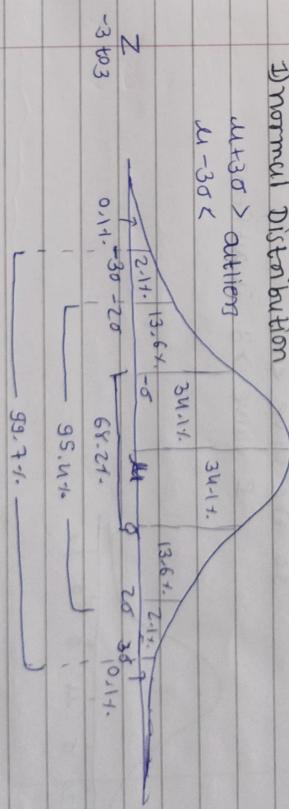
\Rightarrow outliers use mostly affect on weight base algo like
Linear, Logistic, Adaboost, Deep learning
 \Rightarrow not much impact on Tree based random forest, XGBoost
 gradient boosting

| Iteration 0 | Iteration 1 | | | | | | Diff |
|-------------|-------------|-------|-----|-------|-------|-----|-------|
| | Red | Adm | Mur | Red | Adm | Mur | |
| ✓ 9.25 | 15 | 30 | 8 | 15 | 30 | 0 | 0 |
| ✓ 15 | 10 | 11 | 15 | 10 | 11 | 0 | 0 |
| ✓ 12 | 11.25 | 26 | 12 | 11.06 | 26 | 0 | -0.19 |
| ✓ 2 | 15 | 39.25 | 2 | 15 | 31.56 | 0 | 8.23 |

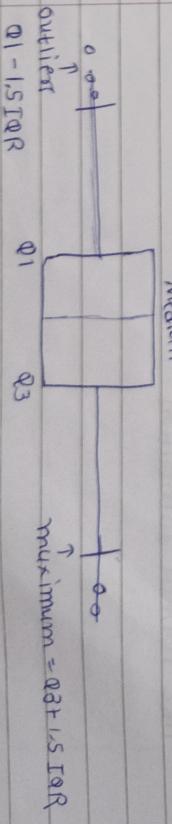
doing same process with Iteration 1 as taking * until convergence occurs.

How to detect outliers

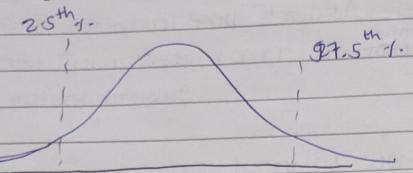
Normal Distribution



2) Skew Distribution



3) Other Distribution



How to deal with outliers

i) outlier removal using z-score (Normal dist)

Capping: assign outliers as $\mu + 3\sigma$ or $\mu - 3\sigma$ whichever is closer

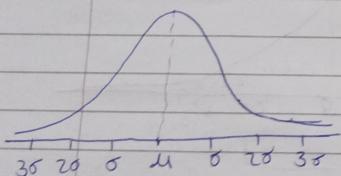
$z \text{ score } z = \frac{x - \mu}{\sigma}$

$$z_i = \frac{x_i - \mu}{\sigma}$$

For like normal dist $\mu \approx 0$

so the value < -3

& value > 3 is outliers



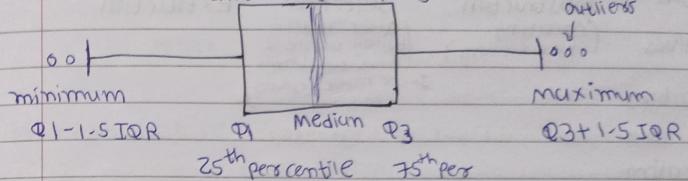
Trimming: remove those row which has outliers

⇒ This method is used when distribution of data is ~~skewed~~ like normal distribution

ii) outlier removing using IQR method

⇒ This method is used when data distribution is skewed

Interquartile Range (IQR)

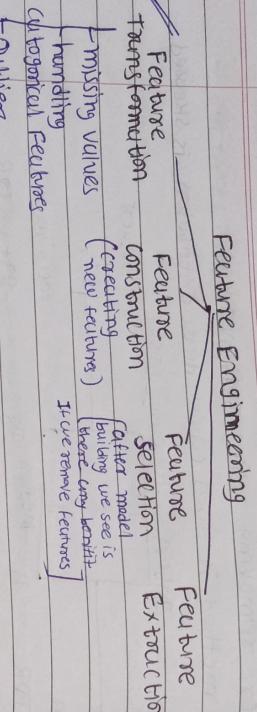


iii) outlier removing using $\pm k$ percentile method

⇒ we can also choose any values like 1st percentile & 99th per.. and can do capping & trimming

Course of dimensionality

⇒ Performance of model is increasing as we add features but after certain feature if we add more then performance is decreases. Ex



Outliers
Feature scaling

⇒ we have completed Feature transformation now we covers Feature construction where we convert features to new features

* Feature construction

There is no fix process to do base on domain knowledge

only you can do Feature construction

⇒ Ex Sibsp | Parenthetically Family type

0 → alone
1 → small
2 → big

⇒ Feature splitting

Ex Name
Mr Ankit → Mr Ankit

⇒ Multiclass Precision & Recall

| | predicted | | | |
|--------|-----------|-----|--------|-------|
| | Dog | Cat | Rabbit | Total |
| Dog | 25 | 5 | 10 | 40 |
| Cat | 0 | 30 | 5 | 35 |
| Rabbit | 5 | 10 | 20 | 35 |
| Total | 29 | 45 | 35 | 109 |

F1 Score lies near to low values

$$F_1 = \frac{2PR}{P+R}$$

P → precision
R → recall

Am → harmonic mean
Hm → geometric mean

Ex P=R=80 P=60 R=100 Am=80 Hm=80

$$Am = \frac{80+80}{2} = 80$$

$$Hm = \sqrt{60 \times 100} = 77.46$$

$$F_1 \text{ score} = \frac{2 \times 80 \times 80}{160} = 75$$

Precision: what proportion of predicted Positive is truly Positive

$$P_{Dg} = \frac{25}{29} = 0.86$$

$$P_{Cg} = \frac{30}{45} = 0.66$$

$$P_R = \frac{20}{34} = 0.58$$

Macro Precision = avg of all precision

$$\text{Macro P} = \frac{0.86 + 0.66 + 0.58}{3} = 0.70$$

Weighted precision = In this we multiple all precision value with its weight

weight of class = actual Number of that class in data / total data

$$\text{Weighted P} = \frac{10}{108} \times 0.86 + \frac{34}{108} \times 0.66 + \frac{84}{108} \times 0.58 = 0.71$$

\Rightarrow when all classes ~~not~~ have same number then use Macro otherwise weighted

same for recall also
& F1 score

$$F1_D = \frac{2P_D R_D}{P_D + R_D}$$

$$F1_C = \frac{2P_C R_C}{P_C + R_C}$$

$$F1_R = \frac{2P_R R_R}{P_R + R_R}$$

Cum find macro & weight F1 score

Projects

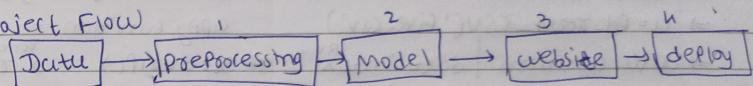
* Recommended system

Content base
recommendation base
of similarity of content
(tags)

Collaborative filtering
recommend content
on base of similarity
of users

Popularity base
Hybrid
[combine both]

Project Flow



CUR Prediction

- `str.isnumeric()` or `not str.isalpha()` gives num or unique values
- `unique()` / `unique(keep='first')` gives all unique values
- `backup = cur.copy()` copy dataset into backup
- `cur.drop('col1', inplace=True, axis=1)` except:
- `drop col1`
- `c = pd.read_csv('file.csv')`
- `data.columns` | `data.index` list of all column's name
- `c['col1'].value_counts()` gives count of unique values
- `series.str.split(' ').str.get(0)` remove abc in x
- `split string with ' ' and give 0th element in every row`
- `series.str.replace(',', '')`
- `cur[~cur['final_type'].isna()]` excluding those rows
- `series.str.slice(0, 3)` gives first 3 element in list
- `.Join(' ') / ' '.join(list)` Join list with ' '
- `df['brand'].replace('rel', 'uncommon')`
- `cur.reset_index(drop=True)`
- `cur.to_csv('clean_data')` save data frame in to csv
- `flights.groupby('year').sum().reset_index()`

Flights

| Year | Month | Passenger | Year | Passenger |
|--------|-------|-----------|--------|-----------|
| 0 1980 | Jan | 112 | 0 1980 | 1120 |
| 1 1980 | May | 98 | 1 1980 | 1180 |
| 2 1980 | April | 105 | 2 1981 | 1250 |
| | | | 3 1982 | 1800 |

House Price

- `data.isna().sum()` missing values in every row
- `data.drop(columns=['a', 'b'])`
- `try:`
- `data.drop('col1', axis=1)`
- `data['col1'].apply(fun)` apply function on column
- `x.strip()` remove if it is whitespace
- `x.strip('q')` remove Leading and trailing whitespace
- `x.strip('abc')`
- `x2 If x is in array1 else x`
- `x2 If x is in array1 else return x`

Book Recommendation System

- `books = pd.read_csv('book.csv', sep=';', error_bad_lines=False)`
- `books.rename(columns={'a': 'b', 'x': 'y'})` rename the columns of dataset
- `ratings['user_id'].isin(y)` True: If ~~that~~ that 'user_id' is in y
else False
- `data.groupby('title').count()` type = Pandas Series
- | title | f1 | f2 |
|-------|----|----|
| "1 | 2 | 5 |
| "2 | 1 | 6 |

 occurrence of feature f2 in title
- `data.groupby('title').f1.count()` type: Pandas Series
when we do reset_index() then it becomes Pandas data frame with index 1, 2, ... and series Book name & f1
- `data.drop_duplicates(['col_1', 'col_2'])` drop row which has same col1 → col2 pair

- `data.pivot_table(columns='col', index='ind', values='val')`

$\xrightarrow{\text{col}}$
 $\downarrow \text{index}$ $\xrightarrow{\text{values}}$
 $\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$ $\xrightarrow{\text{data.fillna(0)}}$

- `from scipy.sparse import csr_matrix`
- `csr_matrix(data) ← type = SciPy-Sparse, CSR-Matrix`
- `array of 0's calculations used when sparse matrix`
- `data.iloc[2, :]` row of df at index 2. series
- `data.iloc[2, :].values` numpy.ndarray shape = ~~(n, 1)~~
- `data.iloc[2, :].values.reshape(1, -1)` shape = $(1, m)$
- `data.iloc[:, :-1]`
- `distance, suggestions = model.kneighbors(row, n_neighbors)`
array of dis array of suggestions
- `data.index[0]` index of row 0th
- `data.index` all the rows index

Numpy (Numeric Python)

$\Rightarrow u = np.array([1, 2, 3])$ size = 3

$u = u[:, np.newaxis]$

$size = (3, 1)$ row column [3]

$a = u[improvedaxis, :, :]$

$size(1, 3) [1, 2, 3]$ row column

$\Rightarrow x = np.array([1, 2, 3], [4, 5, 6])$

$np.sum(x)$ sum all elements

$\Rightarrow x[1:, ::2]$ all the elements from row 1 to last & column with step of 2

$\Rightarrow x.sum(axis=0)$ sum of all columns

$\Rightarrow np.count_nonzero(x)$ number of non-zero elements in x

$\Rightarrow np.hstack((arr1, arr2))$ combine 2 arrays horizontally

$\Rightarrow np.concatenate((array1, array2), axis=1)$ axis = 1

$\Rightarrow arr1 = (20, 1)$ arr2 = (10, 1) concat to 1

$\Rightarrow x[[1, 2], [3, 4], then x[[y=3], [1, 2]] =$ number you take only those data which has y=3 and take its 1st column(index) all rows

$price$

$\begin{bmatrix} 10 \\ 10 \\ 10 \end{bmatrix}$ avg price

company

$\Rightarrow df['col1'].apply(lambda x: x)$

$\Rightarrow df['col1'].str.split(')$

Randu Series and their column contain list of elements

series contain list of elements

$\Rightarrow df['col1'].str.split('), expand=True)$

Randu Series contain

series (column) of repeated

element (with ')

Pandas

$\Rightarrow df = pd.DataFrame(dict1)$
 $\Rightarrow df.to_csv('friends.csv')$
 $\Rightarrow df.index = [1, 2, 3]$
 change index to given name
 $\Rightarrow df['col1'].sum()$ sum of that series
 $\Rightarrow df['col1'].values$ convert array from DF
 $\Rightarrow df = pd.DataFrame(np.random.rand(3, 5), index=np.arange(3))$
 $\Rightarrow df.dtype$ data type of each column (334, 5), index=np.arange(334)
 $\Rightarrow df.to_numpy()$

convert df to numpy array

$\Rightarrow df1 = df$

df1 is pointing to df

df1 is not copy of df

For copy $df1 = df.copy()$

$\Rightarrow df.loc[df['A'] < 5]$

loc used when col name
of row name

$\Rightarrow df.iloc[1:3, 2:7]$ col 2 to 7
 $\Rightarrow df.iloc[[0, 5], [1, 2]]$ row=0, 5 col=1, 2

$\Rightarrow df.drop$ by default how='any'
axis=0 (row)

$\Rightarrow df.dropna()$

any=drop If any ele is na

All=drop only If all ele na

$\Rightarrow df.drop_duplicates(subset=['name'])$

drop dup in column name

keep='last', False

keep='first' drop un dup

\Rightarrow list

$\Rightarrow l = m.copy()$ l=m change $\Rightarrow dic.get('name')$

l If we change m insert num at position 1
 $\Rightarrow l.insert(1, num)$ is not found it return none
 $\Rightarrow l.append(m)$ l → append m
 while dic['name'] throw error if not found

$\Rightarrow dic = \{'name': 'Kurum',$

'age': 19}

key value

dict.keys() dict_keys['name', 'age']

dict.items()

[('name', 'Kurum'), ('age', 19)]

$\Rightarrow dic.update(new_dic)$ add

$\Rightarrow dic.clear()$ empty

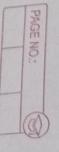
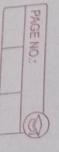
$\Rightarrow dic.pop('name')$ give value of name remove name

$\Rightarrow del dic['name']$ change name as

$\Rightarrow def fum(**name)$ dict
print(name['fname']),
name['mname'])

name(fname="Bucham", mname='Hi')

$\Rightarrow empty_dic = dict()$ or {}

Plots

⇒ sets
→ hi = set()
→ hi = set() empty set

⇒ Tuples tuples can't be changed
→ tup = (1,) need comma for single
→ tup3 = tup + tup1 This is

Possible because here we are not
changing tuple but making new

→ def fun(* nume)
tuple argument

→ titanic.groupby('Pclass').mean()

→ sns.pairplot(iris, hue='Species')

→ sns.lineplot(new['Year'], new['Passenger'])

⇒ sns.scatterplot(tips['bill'], tips['tip'], hue=['Sex'], style=tips['Smoker'], size=tips['Family-size'])

⇒ sns.boxplot(titanic['Pclass'], titanic['Age'], hue=titanic['Sex'])

⇒ sns.distplot(titanic['Age'])

⇒ pd.crosstab(titanic['Pclass'], titanic['Survived'])

⇒ sns.heatmap(

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

↓)

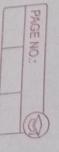
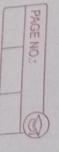
↓)

↓)

↓)

↓)

↓)

Plots

⇒ sns.kdeplot(data['Age'], density=True, rug=True, kde=True, kernel='gaussian', bw_adjust=5) in PDF (prob density function)

Putte kde = color_dict

color_dict = {1: 'red', 2: 'green', 3: 'blue'}

df['col'].plot.density(color='red')

df['col'].hist(values, bins=100)

sns.kdeplot(data['Age'], density=True, rug=True, kde=True, kernel='gaussian', bw_adjust=5)

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑

↑