

git clone https://github.com/username/repo.git
git init
git add .
git commit -m "First commit"
git branch -M main
git push -u origin main

→ make file in git hub

→ clone to folder

wing# git clone https://github.com/username/repo.git

open folder and open git bash right click then run in.

→ git init initialize empty git repository inside git repository

→ git status

→ git add . This command add whatever we create into git repo

→ git commit -m "First commit adding.html"

→ push this change to git hub so every person can see the changes

→ git branch -M main

git push -u origin main

Git & GitHub

- untracked file: git does not know about that file
- we sent our project (which in our system) to git repository
- we can't send direct to git
 - so first project to staging index
 - and then staging index to git

let we have two files in project index.html and contact.html

\$ git add contact.html → sent contact.html to
touched file stage git add . will file to stage

So at this point of time we have 1 file in stage and
to file at project and 0 files at git repo, index.html
untouched file

\$ git add index.html
2 file stage, 2 file project, 0 in git repo

- to commit file
- \$ git commit -m "First commit 7 June 23"
all the files push to git repo (not GitHub)
- \$ git commit -a -m "First commit" add & modified in one line
- If we don't want to track some file then make new
file .gitignore and write name of those files.

- \$ git log give history of all commit Author & Date

Cmd also give SHA (Inception algo) in every project it's unique

\$ git log --oneline give history in one line

\$ git show {SHAkey} show that version of file

→ \$ git config --global user.name user.email

\$ ls \$ ls -a give list in folder

\$ mkdir make directory

\$ cat {filename} can check what change we made

\$ git branch -M master

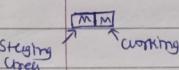
\$ git diff compare working tree to staging area

\$ git diff --staged compare staging area to last commit

\$ touch (unst.html) to make unst.html file

\$ git rm --cached unst.html remove file from staging area
(untracked)
\$ git rm unst.html delete file from staging area & working directory
permanent delete

\$ git status -s



\$

\$ git checkout {filename} match file to previous commit
\$ git checkout ~~the command~~ -f > all files ..

⇒ git ignore

→ If we write file name inside .gitignore then git ignore that files anywhere in root node

→ *.html = ignore all html files

→ /filename = ignore from root node only (not root → lib → filename)

→ foldername/ = ignore folder name
↑ indicate that ^{this is} folder

http link = origin
\$ git remote

→ branch

\$ git branch gives all the branch

\$ git branch slave make new branch name 'slave'

\$ git checkout slave switched branch to 'slave'

\$ git merge slave merge slave to master (∴ run from master)

\$ git branch -d slave delete slave

\$ git branch -M master change branch name to master

\$ git checkout -b New make 'new' branch and switched to that branch

→ to push your programs (projects) to github
new make file in github

\$ git clone {https link of github file}
add file into your computer

If you want to update your changes in github then you need to push that on github

\$ git remote, \$ git remote -v

OP: origin ← origin is http link which we got from github

\$ git push origin master Push project to github
↑ name of branch that you want to push

existing can be any name

\$ git remote add origin {https link}

\$ git branch -M master first commit file then push

\$ git push -u origin master

★ ML Revision 27 July 23

$$\Rightarrow \text{F Factor} = \frac{(1+\beta)^2}{(\beta^2)FR} \frac{PR}{\text{Type I}} \quad \text{FP Impr } \downarrow \beta, \text{ FN } \uparrow \beta, \text{ Type I}$$

\Rightarrow Maximum likelihood

TL Prob of correctly classified

Cross entropy take -log of max liklihood.

$$\Rightarrow \text{Huber loss} = \begin{cases} \frac{1}{2}(y-\hat{y})^2, & |y-\hat{y}| \leq \delta \\ \delta|y-\hat{y}| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases} \text{ hyper parameters}$$

\Rightarrow multi class cross entropy

autogorical cross entropy

$$\text{LOSS} = - \sum_{i=1}^n y_{ij} \log(\hat{y}_{ij})$$

For ith row $\sum_{j=1}^c$ y_{ij} \hat{y}_{ij}
 actual ↑ pred

$$\text{Total LOSS} = - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c y_{ij} \log(\hat{y}_{ij})$$

sparse categorical cross entropy

instead of one hot encoding

we do normalizing to each op

Ex: $y \rightarrow 1$
 $y_1 \rightarrow 2$ and calculate
 $y_2 \rightarrow 3$ only for that y

after updating weights for every model for test data

we get OP

take soft max of that and max prob gives OP

\Rightarrow Entropy & gini impurity (For purity check)

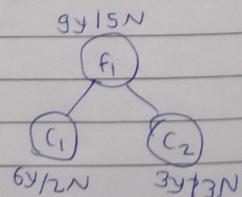
$$H(S) = -P_+ \log P_+ - P_- \log P_-$$

$$GI = 1 - \sum_{i=1}^c (P_i)^2$$

\Rightarrow IG (Feature selection)

$$\text{Gain} = H(S) - \sum_{v=1}^k \frac{|S_v|}{|S|} H(S_v)$$

↑ Parent
↑ Total sample
↑ sample in child node



* DL Revision 28 July 23

- ⇒ ML use statistic technique for this find relation b/w input & output while DL follow logical structure called Neural Networks
- ⇒ $Z_1 = W_1 A_0 + b_1$
 $A_1 = \sigma(Z_1)$
- ⇒ $Z_2 = W_2 A_1 + b_2$
 $A_2 = \sigma(Z_2)$

→ regularization &
→ drop out layers apply for last layers if it work then try
first & second last so on--

CNN best result = $P \approx 0.91 - 0.95$

RNN = $P \approx 0.2 - 0.3$

ANN $P \approx 0.1 - 0.5$

→ This method is like RF we use NN instead of Tree

→ For every epoch remove P fraction of neuron in each HL

→ For O/P activate all features and all neurons and multiply $(1-P)$ with every weight.

dead neuron or dead activation prob

$$\Rightarrow \text{ReLU} = \begin{cases} z, & z > 0 \\ 0, & z \leq 0 \end{cases}, \text{Leaky ReLU} = \begin{cases} \alpha, & z > 0 \\ \alpha z, & z \leq 0 \end{cases}$$

$$\text{ELU} = \begin{cases} z, & \text{if } z > 0 \\ \alpha(e^z - 1), & z \leq 0 \end{cases}$$

$$\text{mean of O/P close to 0} \quad \mu \approx 0$$

$$\text{SELU} = \lambda \text{ELU} \quad \lambda \approx 1.6732 \dots$$

$$\text{self normalizing} \quad \lambda = 1.0507 \dots$$

$$\text{O/P is } \sigma = 1, \mu = 0$$

σ = vanishing grad prob

tanh = vanishing

$$\mu \neq 0$$

$$\mu = 0$$

skewish = used when NN > 10 layers

$$\propto \sigma$$

$$\text{softplus} = \ln(1 + e^x)$$

$$\text{softmax} = e^{x_i}$$

$$\sum_{i=1}^n e^{x_i}$$

⇒ Baye's $P(A \mid B) = \frac{P(B \mid A) P(A)}{P(B)}$

⇒ Translation covariance

After applying convolution layer feature which is find by CNN is dependent on location this is called translation variance.

We can avoid this by using pooling, pooling is a way to down sample feature map so features become independent of location and it called translation invariance.

Image → ^{conv} + ReLU + ^{conv} + Pooling → flatten → Fully connected → O/P layers

For CNN Trainable parameters change only depends on filters size and number of filters

(3, 3, 3)

So filter has $(3 \times 3 \times 3 \times 50)$ + 50 learnable parameters

$$\Rightarrow V_t = \beta V_{t-1} + (1-\beta) S_t \quad 0 < \beta < 1$$

$$V_t \text{ approximating avg over } \frac{1}{1-\beta}$$

$$V_t = (1-\beta) (\beta^{n-1} u_1 + \beta^{n-2} u_2 + \dots + u_m)$$

⇒ For Deep Learning L2 (Ridge Regularization) gives best result

$$\Delta w = w_{old} - \eta \frac{\partial L}{\partial w_{old}}$$

$$L = L^0 + \frac{\lambda}{2} \sum_{i=1}^n \|w_i\|^2 =$$

↑
without regu.

$$\frac{\partial L}{\partial w_{old}} = \frac{\partial L^0}{\partial w_{old}} + \frac{\lambda}{2} (2w_{old})$$

$$w_{new} = (1-\lambda n) w_{old} - \eta \frac{\partial L}{\partial w_{old}}$$

↓
PROM

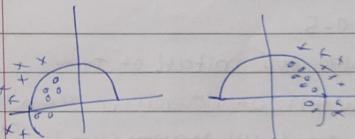
$$1-\lambda n < 1$$

$$0 < \lambda < \infty$$

$$0 < n < \infty$$

(LR)

⇒ Covariate shift



Input distribution and O/P distribution is diff but O/P is P rel. measure
(in figure given by some curve)

⇒ Batch norm

⇒ change in activations due to the change in network parameters during training is called internal covariate shift.

⇒ In Batch normalization we convert O/P of ~~ML~~ HL into Gaussian distribution ($\mu=1, \sigma^2=0$) so internal covariate shift reduce and training will be fast and stable.

⇒ If we don't use ~~normal~~ Batch norm. LR should be low and

carefully ini param. otherwise training will be unstable.

* K-Means Clustering

⇒ From sklearn.cluster import KMeans

Km = KMeans(n_clusters=5)

Km.fit_predict(df)

⇒ Km.inertia_ gives WCSS = within cluster sum of squares

⇒ Km.fit_predict(x) : fit and predict clusters for each row in dataset x

* Stacking Classifier

sklearn.ensemble.StackingClassifier

⇒ clf = StackingClassifier(

estimators=estimators

final_estimator=LogisticRegression(),

CV=10

* Agglomerative Clustering

⇒ From sklearn.cluster import AgglomerativeClustering

cluster = AgglomerativeClustering(n_clusters=5,

affinity='euclidean',

linkage='ward')

labels = cluster.fit_predict(data)

★ Decision Trees

Plot tree

```
from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier()
clf.fit(X_train, y_train)
plot_tree(clf)
```

★ Voting ensemble

```
From sklearn.ensemble import VotingClassifier
estimators = [('rf', clf1), ('rf', clf2), ('knn', knn)]
vc = VotingClassifier(estimators=estimators, voting='hard')
list of tuples
curr also assign weights
weights = [1, 1, 3]
```

★ Bagging

```
from sklearn.ensemble import BaggingClassifier
bag = BaggingClassifier(
    base_estimator=DecisionTreeClassifier(),
    n_estimators=500,
    max_samples=0.5, max_features=0.5
    bootstrap=True, bootstrap_features=True
    random_state=42,
    verbose=1,
    n_jobs=-1)
```

taking 50% of data to every DT

Sample with replacement

divide task in CPU

Plot tree

→ bag. estimators - samples

list of array each array contain index number of which assign to test that DT

→ bag. estimators - features

random to each feature

★ Random Forest

rf = RandomForestClassifier(max_features=2, n_estimators=100, max_features=5, auto, sqrt, log2, 0.3, None, bootstrap=True, max_sample=100, max_depth=5, criterion={gini, entropy}, min_sample_split=5, 0.1, min_leaf_nodes=, min_impurity_decrease=)

hyper DT parameters

How many DT
100 Rows

★ AdaBoost

ada = AdaBoostClassifier(n_estimators=1500, learning_rate=0.1, default=DT with height=1)

base_estimator=,

algorithm='SAMME.R', 'SAMME'

default

★ Date time

```
pd.to_datetime(df['date'])  
df['date'].dt.year / df['date'].dt.day / df['date'].dt.week  
- df['month'].month_name() / df['day'].name()  
- df['week'].dt.week / df['month'].dt
```

★ Missing indicator

$\Rightarrow \text{mi} = \text{MissingIndicator}(\text{x_train}[:, \text{col1}])$

$\text{mi}.fit(\text{x_train}[:, \text{col1}])$

$\text{mi}.transform(\text{x_train}) \rightarrow \text{OP} = \text{array of } (n, 1) \text{ with T/F}$
 $(\text{x_train}[:, \text{col1}])$

$[[\text{False}],$
 $[\text{True}],$
]

\Rightarrow simple imputer also contain missing indicator so no need to use superlative

```
SimpleImputer(add_indicator=True)
```

$\Rightarrow \text{Knn} = \text{KNNImputer}(n_neighbors=3, \text{weights}='distance')$

'uniform'

for arbitrary value imputer

$\Rightarrow \text{SimpleImputer}(\text{strategy}='mean' / \text{median})$

most_frequent / constant / add_indicator

$\Rightarrow \text{Si} = \text{SimpleImputer}(\text{strategy}='constant', \text{fill_value}=1000)$

(MICE)

imputer

$\Rightarrow \text{ite} = \text{IterativeImputer}()$

$x_{\text{train}} = \text{ite}.fit(\text{x}_{\text{train}})$

$\text{tol} = 10^{-3}$

max_iter

add_indicator

initial_strategy=

★ PCA

```
from sklearn.decomposition import PCA
```

~~PCA~~ $\text{Pca} = \text{PCA}(n_components=100)$

$\text{Pca}.fit(\text{x}_{\text{train}})$

$\text{Pca}.transform(\text{x}_{\text{test}})$

$\text{Pca}.explained_variance = \text{Eigen values}$

$\text{Pca}.components = \text{Eigen vectors}$

$\text{Pca}.explained_variance_ratio$

★ Polynomial logistic regression (non linear logistic regression)

\Rightarrow from sklearn.preprocessing import PolynomialFeatures

$\text{Poly} = \text{PolynomialFeatures}(\text{degree}=3, \text{include_bias=False})$

$\text{x_trf} = \text{Poly}.fit(\text{x})$

after this $\text{x} \rightarrow x_1 | x_1^2 | x_1^3 | x_2 | x_2^2 | x_2^3 | \dots$

and now we can apply Logistic Regression on this modified x_{trf}

Pundus Profiling

```
⇒ !pip install Pundus-Profiling
from Pundus_Profiling import ProfileReport
prof = ProfileReport(df)
prof.to_file(output_file='output.html')
```

* one hot encoding using Pandas

```
⇒ pd.get_dummies(df, columns=['fuel', 'owner'], drop_first=True)
```

* one hot encoding using sklearn

```
⇒ from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder(drop='first', sparse=False, dtype=int32)
ohe.fit(x_train[['col1', 'col2']])
ohe.transform([[1, 2]])
```

give sparse array or gives sparse matrix

```
trf = FunctionTransformer(func=np.log1p, kwds={'method': 'Box-Cox'})
```

```
⇒ SimpleImputer()
ordinalEncoder(categories=[list]) si.transform(series)
OneHotEncoder(drop='First', ...)
```

gives np array

* Column transfer make_column_transformer

```
⇒ from sklearn.compose import ColumnTransformer
```

```
transformer = ColumnTransformer(transformer=[('num1', SimpleImputer(), ['Fever']), ('num2', ordinalEncoder(categories=[list]), ['(cough)'])], remainder='passthrough')
```

↑ list or
↑ list or
↑ column

```
transformer.fit_transform(x_train) & transformer.transform(x_test)
```

name! trtf.named_transformers_

From sklearn import set_config
set_config(display='diagram') gives diagram of pipeline

* Pipeline

```
Pipe = Pipeline([(t1, trf1), (t2, trf2)])
```

```
Pipe = make_pipeline(trf1, trf2, trf3)
```

Pipe.fit(x_train, y_train) If model is inside pipeline other wise fit_transform & train

Pipe.predict(x-test)

Pipe.named_steps['trf1'].transforms_[0][1].statistics_

⇒ cross_val_score(Pipe, x_train, y_train, cv=5, scoring='accuracy').mean()

param = {trf1__max_depth: [1, 2, 3, 4, 5, None]} ← dict

grid = GridSearchCV(Pipe, param, cv=5, scoring='accuracy') ← run more than one items

grid.fit(x_train, y_train)

grid.best_score_ / grid.cv_results_

grid.best_params_

threshold = n) RandomSearchCV
bi = Binarizer(copy=False) take random from
bi.fit() param ← param &
bi.transform() global column itself
convert continuous values to binary 0 or 1
value < threshold or > threshold by default threshold = 1 or 'onehot'

⇒ Kbin = KBinsDiscretizer(n_bins=10, encode='ordinal',
 strategy='quantile')

Kbin.n_bins_

-bin-edges-

new column edge table can get by

1-10
10-20
pd.cut(x=x_train['Age'], bins=list of bin edges)

interview questions

Q) Diff b/w parametric & non parametric ML algo

parametric

i) Parametric ML algo: we make assumption of nature of the function Ex Linear regression

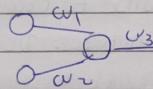
→ and it has fix number of parameter It does not depends on no. of data

linear reg, linear sum, Naive Bayes

$$\text{Ex } y = mx + c$$



or



↑
assumption
all independent
comes one
independent
simple neural netw

ii) Non parametric ML algo: No assumption

→ No fix para grow with ~~no.~~ ^{number} of row

Ex Trees

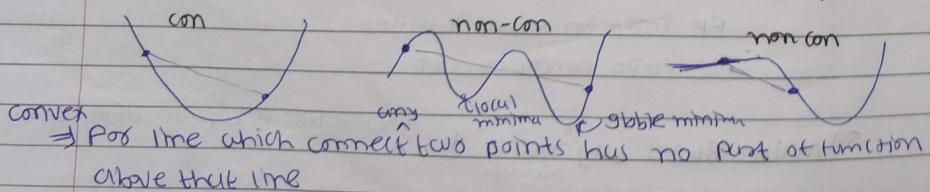
Other SVM with diff kernel

KNN, DT, RF, XGBoost, Ada...

Complex neural networks

Param	Non-param
1) Simple func	1) Complex
2) less data	2) more data
3) underfit	3) overfit

Q) what is difference b/w convex and non convex loss fun, what happens when we have a non convex function?



⇒ Strictly convex fun has only one minima

non-convex has more than one minima

⇒ when we have non convex function then some time solution stuck in local minima

6) give example where medium is more imp than mean?

Eg:
salary 3L, 5L, 8L, 2L

This kind of situation take medium

→ when data has outliers take median

7) what do you mean by unpredictable effectiveness of

data in ML?

⇒ as data increases algo does not make simple and work
algo also give some performance in compare to complex and

powerful algo.

8) what is lazy learning algorithm?

how it is diff from eager learning?

why is KNN A lazy learning ML algo?

two type → lazy: generalization of training data in theory, delayed

→ eager: until query is made by system.
where system try to generalize the training data before receiving queries.

⇒ Lazy Learning store training data because it need during prediction
Prediction case show

→ eager learning not store data, just prediction

9) what is semi supervised ML?

→ supervised has labeled op

types
unsupervised
or ml
reinforcement
semi supervised

→ unsupervised has no labeled op

→ semi supervised
Eg: google photos cluster different

photos in ~~one~~ group and one
need to keep only ~~one~~ group.
each photo in group has some label

3) when use Deep learning?

when need high performance

⇒ have large numbers of data

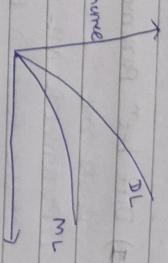
⇒ It problem is very complex and we don't have domain knowledge

about it then use DL algo

Eg: cut-dog classification

⇒ it is not possible to create rules for this problem and ~~not~~

create features. DL algo create features & rules



against

⇒ data is small then we ml

⇒ if you don't have cost measure GPU

⇒ interpretability or explainability

Eg: block account using DL can't explain why block account

y) give diff between FP & FN give example when important

Actual	0	1	
Predicted	0	FP	FN
1	TP	FN	TP

FP: impo in spam classifier

FN: impo cancer or not

5) Naive Bayes why "naive"? (NaiveBayes is base on bayes theorem)

ignorance

⇒ Naive Bayes algo consider that data has no relationship ~~but~~ PLANB = P(A|B)P(B)

↓ (ignores dependencies between variables)
assuming they are independent

$$P(Y|X) = P(X \cap Y)$$

$$= P(X_1 \cap X_2 \cap \dots \cap X_n \cap Y) = P(X_1|Y)P(X_2|Y) \dots P(X_n|Y)$$

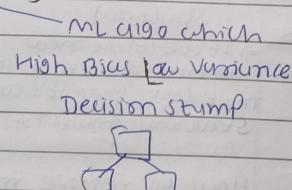
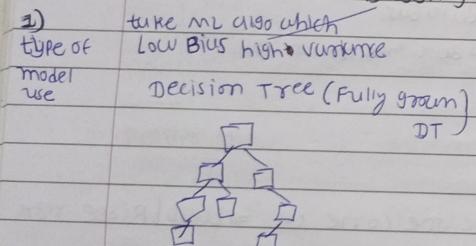
1) what is 'no free lunch' in ML?

→ without any assumption while making ML model, we can't say which model gives good result on particular dataset.

2) difference b/w structured & unstructured data?

structured	unstructured (qualitative)
→ tabular data	→ non tabular EX social media post, text file, sound file,
→ organized	→ unorganized
→ easy to apply ML algo	→ apply DL for unstructured data true in NLP, text mining, computer vision

3) difference b/w Bagging vs Boosting



2) Sequential vs Parallel Learning

Bagging → Parallel

Boosting → Sequential

3) weightage of base learner

Bagging → same weightage

Boosting → different weightage

4) Assumptions of linear regression

1) Linear Relationship b/w input & op

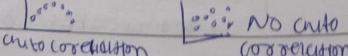
2) no multi collinearity

3) Normality of Residuals $\rightarrow \text{Residual} \sim N(0, \sigma^2)$ for every point distributed normally

↳ Homoscedasticity $\rightarrow \text{Residual} \sim N(0, \sigma^2)$ scatter is uniform in X axis

5) NO AutoCorrelation of errors

↳ There is no pattern on above graph



1) what is OOB error and how is it useful?

OOB = out of bag Evaluation

→ take ex of RF some samples are not selected for training ($\approx 37\%$) we can test over model on this data (as validation set)

2) when you prefer DT over RF?

→ explainability

→ low computation power

→ If we know some feature is very important and training should contain those feature while training ...

3) why logistic regression not logistic classification?

Logistic regression is ~~linear~~ ^{non-linear} regression with S function. And we set some threshold if $S(p) > \text{threshold}$ then $\rightarrow 1$

→ Logistic regression called ~~soft~~ ^{soft} regression because it's give continuous value so...

4) what is online ml?

and In what kind of scenario it is useful?

→ is technique in which train model on server (online).

→ This is called incremental learning

$$m \rightarrow m' \rightarrow m'' \rightarrow m'''$$

We ~~can~~ ^{only} train model not every model support but some algo like SVM, SGD regression support.

advantage: less costly

faster training

→ Ex: chatbots.

→ Swift Key → Keyword Android

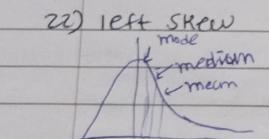
→ YouTube → once you watch video feed below their video change

→ Concept drift ~~Ex:~~ recommendation sys $\xrightarrow{\text{sales}} \text{festival}$ $\xrightarrow{\text{season}}$

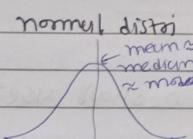
→ since we don't have all house data we can't find $f(x)$ but find $f'(x)$ which is close to $f(x)$
 $f(x) - f'(x) = \text{reducible error}$
 This error reduce as we get more data or using different models

2) For which case outliers affect more?

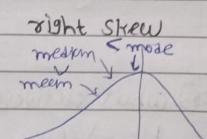
- outliers are mostly affect on weight base algo like linear logistic, AdaBoost, Deep learning
- not much effect on tree base algo like RF, XGBoost, Gradient boosting



Ex wealth distn



Symetric distn
IRIS dataset

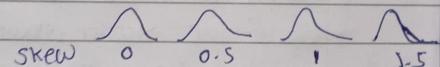


Negative skew

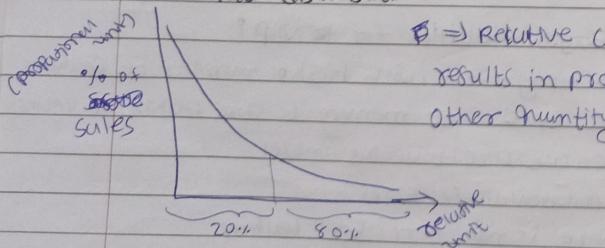
life span of human being

mode = most frequent value

$$\text{Skew} = \frac{n}{(n-1)(n-2)} \sum \left[\frac{(x_i - \bar{x})}{s} \right]^3$$



23) power law distribution



Ex 1) 80% sales come commig from 20% of overal product

- 80% of windows crash due to 20% bugs
- 80% data scientist are 20% software

⇒ because of multi collinearity can't explain importance of feature in output

Ex some invstg done by scientist only two scientist is from same background math so it is very difficult to find contribution of each those two.

⇒ some time we want to find relation b/w feature and op so we use linear regression

⇒ Is multicollinearity always bad?

If we can't find feature impo then yes

If only prediction then No, for only pred we better algo than linear regression

type of multi collin

Structural

→ If we create new features and it contains multi collin then it's called structural multi collin

data base

→ data has multi collin

How to remove multi collin

→ increase data → remove one corrone → Lasso/Ridge regns

⇒ PCA

19) Does multicollinearity affects all ML algo?
mostly impact on parametric algo

20) Reducible vs irreducible errors

→ data has some noise so we add bias in $y = mx + b$ even if we have x can't predict y so we add bias this error can't be reduce so it's called irreducible error

⇒ Reducible error: $y = f(x) + \text{error (bias)}$

let we has 1000 house data to predict price but we (sample data)
want to mettle formula for population (true houses)

Log normal distribution: is continuous prob distribution of random variable whose logarithm is normal distribution

PAGE NO.:

PAGE NO.:

23) what is assumption for SVM?

⇒ No assumption

24) Decision Tree

⇒ No assumption

⇒ Feature scaling not required (\because DT is one tree approach instead of all trees)

⇒ Can handle continuous and categorical variable

⇒ Use for class & regression

⇒ No linear parameter not affect performance of DT. (DT is not sensitive to scale changes)

⇒ Can handle missing value automatically (values not missing in splits)

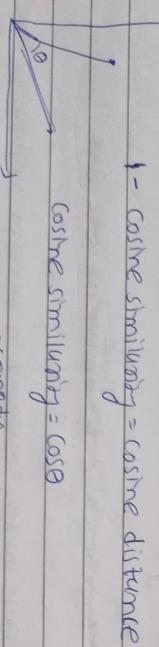
⇒ Less robust to outliers and can handle automatically.

⇒ Less training time compared to RF (\because train only one tree not forest)

DiscuN: over fitting, high variance, unstable: adding a new data point can

Not suitable for large dataset lead to regenerate overall tree

29) cosine similarity and cosine distance



1 - cosine similarity = cosine distance
 \Rightarrow all when we use RNN we have weight in matrix
 \Leftarrow use Xavier glorot or orthogonal for O & softmax both in one model

Regression = linear

Classification = binary
 \Leftarrow min = softmax

ML → RNN

output depends on problem

fact

classification = binary
 \Leftarrow min = softmax

25) why use Naive Bayes classifier for NLP?

⇒ because it work very well with high number of features also less word 2 vec take huge memory to store feature in vector.

⇒ converge fast when we use training more

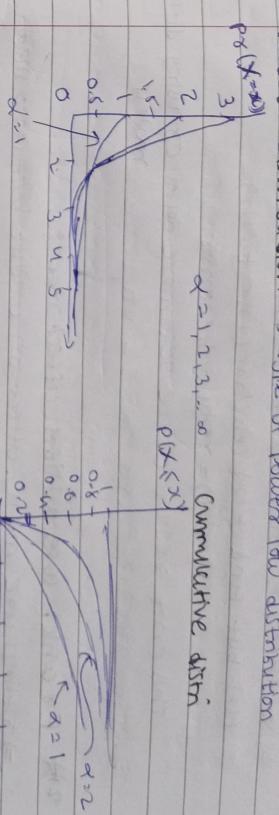
⇒ work well with categorical features

DiscuN: assume there is not good best features

Adv: can handle missing value dataset ($\because \sum_{i=1}^n p(x_i/y) = 1$)

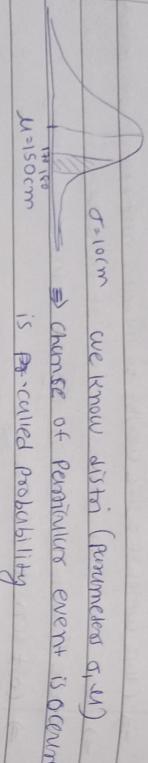
⇒ robust to outliers just skip & which we don't know (\because take our prob for that feature)

⇒ Pareto distribution is type of power law distribution



Let take example of normal (Gaussian distribution)

\Rightarrow Normal dist is continuous distribution



\Rightarrow chance of particular event is occurs

is \rightarrow curved probability

$P(170 \text{ to } 180) = \text{Area under } 170 \text{ to } 180$

$$P(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$Z = \frac{x-\mu}{\sigma}$$

$$P(170 \text{ to } 180) = 0.02$$

\Rightarrow now likelihood is

so let we take one point and its 100cm then what is

likelihood that it follows normal dist $\mu=150 \text{ cm}$

$$L(\mu, \sigma | x = 100) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\downarrow \sigma = 10$

$$\text{Very small} \rightarrow = 1.17 \times 10^{-7}$$

$L(\mu, \sigma | x = 130) = 0.065$

$$L(\mu, \sigma | x = 100) = 0.024$$

$$L(\mu, \sigma | x = 100) = 1.17 \times 10^{-7}$$

\Rightarrow probability that certain event occurs

out of all possible event $0 < P \leq 1$

likelihood = in statistical context, likelihood is a function that

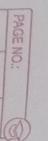
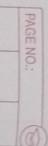
measures the plausibility of a particular parameter value given

some observed data. It quantifies how well a specific outcome

support specific parameter values.

It quantifies how good one's model is, given set of data that has been observed

probability describe outcome while likelihood describe models



30) Data Engineer vs Data Analyst vs Data Scientist

Data warehouse

DATA ENGINEERS: ensure that data will be stored in much more efficient way

Data Scientist $\left\{ \begin{array}{l} \text{Exploratory Data Analysis} \\ \text{and Statistical analysis on data and from data} \\ \rightarrow \text{Discover data using power BI, tableau} \end{array} \right.$

DATA ANALYSTS: Perform some EDA and Statistical analysis on data and understand some more information

Model creation

model deployment

3) Probability vs Likelihood

Parameter \rightarrow event

event \rightarrow Parameter

0	0	0
0	0	0
0	0	0

\leftarrow Blue or Black

B = Black
Bl = Blue

$$P(B) = \frac{3}{5}$$

Bernoulli distribution

\Rightarrow we take 5 ball and let those all are Blue

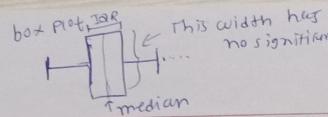
Bl Bl Bl Bl Bl

Base on this observe data or event we have to find the possibility of give parameter ($P(B) = 3/5 \neq P(B)$)

So possibility $\rightarrow \left(\frac{3}{5}\right)^5$

that means base on given data (all ball are blue) probability of blue ball is green is $\frac{3}{5}$ is very low ($\frac{3}{5}$)⁵

K.F.I



PAGE NO.:

34) Time series cross validation

Fold 1 1, 2, 3
validation 4
(using rolling window approach)

Fold 2 1, 2, 3, 4

" 5

Fold 3 1, 2, 3, 4, 5

" 6

walk forward cross validation

- 1
- 2
- 1, 2
- 3
- 1, 2, 3
- 4

35) why you take σ not σ^2 ?

⇒ because σ^2 change the unit from original data so we prefer σ instead of variance

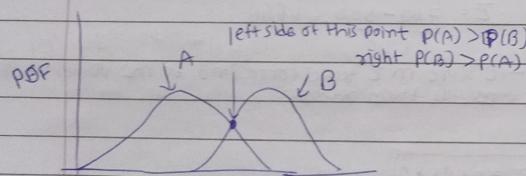
SD

36) why variance prefers instead of Mean abs deviation?

⇒ variance give smooth function

⇒ when outliers are more then use MAD

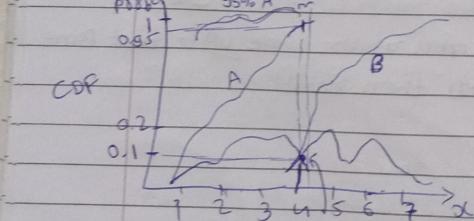
37) what is importance of PDF & CDF



If $x < \bar{x}$ then A \leftarrow 95% correct

else If $x > \bar{x}$ then B

10% wrong
90% correct.



32) How handle imbalance datasets?

→ number of sample which is less
(is very) less - which more

⇒ Class weights & give weight to increase importance of minority class

⇒ use evolution metrics instead of accuracy

⇒ use Ensemble method E.g RF, Adaboost, bagging

⇒ cost sensitive Learning:

Some algo have build-in options for handling imbalance data.

For E.g. SVM has option for adding weights to penalize misclassifying the minority class.

⇒ Data augmentation in minority class

⇒ collect more data

⇒ use stratified cross validation for find accuracy

33) how to make model robust to outliers?

⇒ remove outliers using z-score, IQR, visual inspection using box plot

⇒ Data transformation (using fun transformers)

⇒ use loss function which less sensitive to outliers E.g. MAE instead of MSE

⇒ Ensemble techniques reduce effect of outliers

⇒ data segmentation: If delta contains mix outliers rich and outliers free regions build separate model for each segment.

⇒ Use median instead of mean

⇒ Normalize / scale data

⇒ Cross-validation: Employ robust cross-validation like k-fold cross validation to evaluate the model performance and ensure that generalize well.

u) How to handle imbalance dataset in deep learning?

→ use weighted neural network

weights = {0:1, 1:550} If 0 is 550 times more than 1 then assign 1 to 0 & 550 to 1

model.fit(x_train, y_train, class_weight=weights, epochs=10)

w) How to do hyper parameter tuning in ML & DL?

→ ML → GridSearchCV, RandomSearchCV

dict = {'c': [0.1, 0.5, 1, 10],
'B': ['linear', 'rbf']}

Hyper parameters
tuning using Randomized
Search CV

grid = GridSearchCV(estimator=SVM, dict, cv=5, scoring='accuracy')
grid.fit(x_train, y_train)

n_jobs = -1

to make process
faster

u2) How can you make process faster while using hyper
parameter tuning?

using n_jobs = -1

u3) When use normalization and when standardization

→ Normalise data from range 0 to 1

→ use normal when features are in diff range and you want to
bring them all to a comparable range.

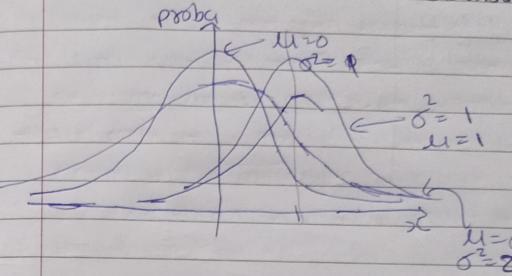
→ It's specially useful when algo relies on the magnitude of
features and works well with features that have meaningful
minimum & maximum

→ Stem = $\mu=0, \sigma^2=1$

→ use standardization when data has varying distributions and
you want to make sure all features have comparable mean & σ^2

→ Stem is commonly used in ML algo (SVM, K-nearest neighbor) relies on the
distance between data points. As it centers the data at zero and
gives equal impo of +ve & -ve deviations from mean.

38) How σ and μ affect distribution graph



σ^2 = spread of data
 μ = centrality of model

For $\sigma^2 = 1$ 44%

-15 to 15 55%

-25 to 25 95% (68+27)

-35 to 35 99% (68+27+4)

39) Imagine a group of 300 students took a test. Rohit scored 700 out of 1000. The avg score was 500 and standard deviation was 120. Assuming test scores follow normal distribution find out how well rohit performed in comparison to his peers.

$\sigma = 120$

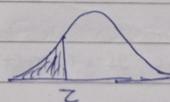
$\mu = 500$

Z for rohit is $Z = \frac{x-\mu}{\sigma}$

$$Z = \frac{700 - 500}{120} = 1.666$$

From Z table

(-ve score in Z table correspond to the values which
are less than mean)



$$Z = 0.95 + 5$$

Z score gives area under curve up to that point

so 95% student score less than rohit

- Q5) In what case we use SGD regression over linear regression
- ⇒ linear regressor gives exact value but it uses matrix inversion so it is very costly computation $(x^T x)^{-1} x^T y$
 - ⇒ SGD regressor is useful when data is very large
 - ⇒ choice is dependent on size of data and computational resources available.

- Q6) Is Ridge Regression use Matrix inversion?

No,

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

while the matrix inversion appears in the eqn for optimal β
modern implementation of ridge regression in libraries like
Scikit-learn use more efficient numerical methods, such as
SVD (singular value decomposition) to solve this problem
without explicitly performing matrix inversion

- Q7) Kernel Function (linear, rbf, poly, sigmoid, precomputed)

RBF: Radial Basis function:

$$y = e^{-\frac{(x+b)^2}{2}}$$

$y = e^{-\frac{x^2}{2}}$ ← x is vector sum of features

$$x_1, x_2 \rightarrow z = e^{-(x_1^2 + x_2^2)} + e^{-x_1^2}$$

1	1
1	1
1	1
1	1

- ⇒ In summary normal use when you want feature $[0, 1]$ & stand use when you want $\sigma^2=1$, $E=0$
- ⇒ Stand can be more robust since it's less influence by extreme value compare to normal
- ⇒ Although it is not fix which use when some algo might perform better with specific scaling method
- It's good to try both and select best

- Q8) What is diff betn hard & soft margin?

$$\text{hard margin} = \arg \max_{w^*, b^*} \frac{z}{\|w\|} \text{ such that } y_i(w^T x_i + b) \geq 0$$

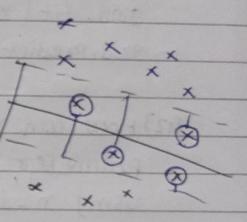
$$= \arg \min_{w^*, b^*} \frac{\|w\|}{2}$$

Soft margin: allow outliers in model

$$\arg \min_{(w^*, b^*)} \frac{\|w\|}{2} + C \sum_{i=1}^n E_i \leftarrow \begin{array}{l} \text{hinge loss} \\ \text{hyperplane} \\ \text{margin error} \end{array}$$

→ If hyperparameter $C \uparrow$ give more focus to reduce misclassified point

→ $C \downarrow$ more focus to give large margin with some misclassified points



$$E_i = \text{loss} = f_{\text{func}}[y_i, f(w^T x_i + b)]$$

Same analogy as Logistic regression

$$\text{Loss term} \rightarrow \text{logistic loss} + \lambda \|w\| \leftarrow \text{regularization term}$$

$$\text{Similar to } \left(\arg \frac{\|w\|}{2} \right) \text{ similar to } \text{hinge loss } E$$

$$C \propto \frac{1}{\lambda}$$

Mutually Exclusive events: A & B never occurs together

A, B $P(A \cap B) = 0$

Independent event $P(A \cap B) = P(A)P(B)$

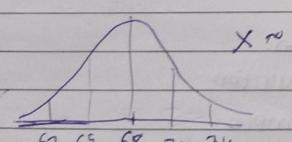
If event does not depend on 2nd event

- 51) why normal distribution is very important?
- ⇒ Commonly in nature many natural phenomena follow a normal distribution such as height of people, weights of objects.
- IQ score of population. and many more
- ⇒ Normal distribution provide a convenient way to model and analyse such data.
- ⇒ Now we have enough study and information about normal distribution so we can't dataset to be normally distributed.

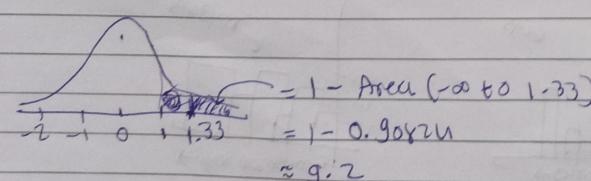
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

Area of $e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} = \sigma\sqrt{2\pi}$ so we ~~divide~~ by $\sigma\sqrt{2\pi}$ to get Area = 1

- 52) Suppose the heights of adult males in certain population follow a normal distn which $\mu=68$ inches & $\sigma=3$ inch. What is prob that randomly selected adult male from the population is taller than 72 inches?

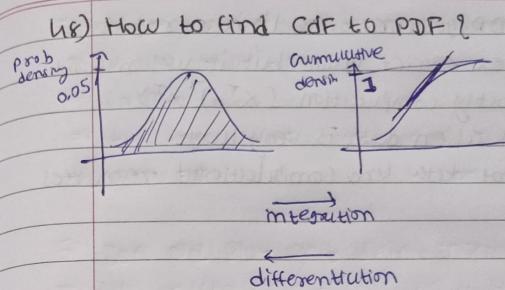


$$Z = \frac{72 - 68}{3} = 1.33$$



$$= 1 - \text{Area } (-\infty \text{ to } 1.33)$$

$$\approx 0.90821$$



- 69) What is advantage of Normal distribution?

⇒ Using standard normal distribution allow to compare different distribution with each other.

Normal distn $X \sim (\mu, \sigma)$

standard normal dist $Z \sim (0, 1)$

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$

How to convert Normal distn to standard normal distn.
Convert each column to $Z = \frac{X - \mu}{\sigma}$

so final data has standard normal distn $(\mu=0, \sigma=1)$

- 50) F1 score (F1 is for multi class classification)

F1 is used evaluate the performance of binary classification model. It combine precision and recall both.

⇒ Combine both P, R in one and use for imbalance betⁿ classes.

⇒ F1=1 means perfect P&R

F1=0 means poor P&R

T-test work good with smaller sample size

PAGE NO.:	1
	2
	3

59) what is concept drift?

some
when dependent variable change after time or build model
data drift & independent " " "

Q-Q Plot: A normal Prob Plot Plots the observed data against the expected values of normal distri. If data fall along a straight line the distribution is likely to be normal

PAGE NO.:	1
	2
	3

53) How to find data is normally distri or not?

- ⇒ Visual inspection
- ⇒ Q-Q plot
- ⇒ Statistical test like Shapiro-Wilk, Anderson-Darling test
- ⇒ $\text{Skew} = 0$ Kolmogorov-Smirnov test

54) Does Q-Q Plot only detect normal distribution?

- ⇒ No, it can detect any use y -quantile diff

55) How to check distribution is Pareto?

$$y = \frac{\alpha x_m^{\alpha}}{x^{\alpha+1}}$$

take log both side and check graph If it is \downarrow then distribution y is pareto

→ OR we can ~~can~~ check by using Q-Q plot

56) Proof of CLT

No, Proof ^{its} exist in nature

57) diff betⁿ parameter & statistics.

- ⇒ Parameter characteristic of Population
- ⇒ Statistic " " " sample

58) diff betⁿ bar graph and histogram

