

#EDA#

- ① Numerical Features
- ② Categorical Features
- ③ Missing Value
- ④ Outliers
- ⑤ Temporal Variable
- ⑥ Cleaning

#Feature Engineering#

- ① Handling Missing value
- ② Handling Imbalanced Dataset
- ③ Treating the Outliers
- ④ Scaling the data {Standardization down, Normalization}
- ⑤ Converting the categorical feature into numerical features.

#Feature Selection

- ① Correlation
- ② K Neighbour
- ③ Chi-square
- ④ Genetic Algorithm
- ⑤ Feature Importance {Extra Tree Classifier}

- ① Handling Missing Value

Type of Encoding

①

① Handling Missing Value.

- ① Numerical feature → If Outliers → then replace with Median.
→ NO Outliers → then replace with Mean.

- ② Categorical feature → Replace with Mode.

Numerical Feature.

- ① Drop the entire row which has missing value. {for huge data only 1 million}
② Create a separate model to handle missing value. {much more time}
③ Statistical method, Mean, Median, Mode. {Best}

Categorical Feature.

- ① Drop the entire row which has missing value {for huge data only 1 million}
② Replace with Mode.
③ Apply classifier algorithm to predict. {Most Efficient} {KNN}
④ Apply Unsupervised ML. {Clustering}

Type of Encoding :-

- ① Nominal Encoding {Order not Important} → One Hot Encoding
Ex:- Gender → Male
 Female
- One Hot Encoding with many categorical {KOD Orange}
→ Meam Encoding {Pincode}
- ② Ordinal Encoding {Order Importamt} → Label Encoding
Ex:- Education.
- Label Rank
P.hd 4
Master 3
12th 1
B.com 2
- Target Guided Ordinal Encoding

Label	Rank
P.hd	4
Master	3
12 th	1
B.com	2

(3)

Day 1:- Feature Engineering {Missing Value}

Lifecycle of Data Science Projects.

1. Data Collection Strategy - from company site, 3rd party API, Surveys.
2. Feature Engineering - Handling missing value.

Why are there Missing Value? {Survey → Depression Survey}

1. They hesitate to put down the information.
2. Survey informations are not that valid.
3. Men → Salary {don't say} salary.
4. Women → Age {don't say} age.
5. People may have died

Note:- Dataset should be collected from multiple sources.

What are the different types of missing Data?

① Missing Completely at Random (MCAR) :- A variable is missing completely at random (MCAR) if the probability of being missing is the same for all the observations. When data is MCAR, there is absolutely no relationship between the missing data and any other value, observed or missing, within the dataset. In other words, those missing data points are a random subset of the data. There is nothing systematic going on that makes some data more likely to be missing than others.

Inshort :- No relationship between missing value with any feature.

Example :- In [titanic] dataset, the column [Embarked] has a missing value which has no relationship between dependent and other feature.

② Missing Data Not at Random (MNAR) :- Systematic Missing Value there is absolutely some relationship between the data missing and any other values, observed or missing, within the dataset.

Inshort :- Relationship between missing data with other feature.

Example
on miss
which

③ Missin

IMP.

AU +

① Mean

values

Imput

① mea

or if

Advan

Discr

Day

② Lan

impul

we

Adv

Discr

③ Co

m

Adv

Discr

Eg:

(3)

(4)

Example :- In [titanic] dataset, the [Age] column and [Cabin] column has some missing values. which has relationship between each other. which means if person is died then we cannot find its cabin and age.

③ Missing at Random :- Men { Hide Salary } women { Hide Age }

IMP.

All the rep technique of handling Missing Value?

- ① Mean / Median / Mode
- ② Random Sample Imputation
- ③ Capturing NAN values with a new feature
- ④ End of Distribution Imputation
- ⑤ Arbitrary Imputation
- ⑥ Frequent Category Imputation.

① mean / Median / Mode :- Data should be MCAR. If no outlier use mean or if outlier then use median for categorical variable use mode.

Advantages :- ① Easy to implement (Robust to outlier).

② Faster way to obtain the complete dataset.

Disadvantage :- ① Change or Distortion in the original variance ↑ ↘
 ② Impacts Correlation.

Day 2 :-

② Random Sample Imputation :- Data should be MCAR. Random Sample Imputation consists of taking random observation from the dataset and we use this observation to replace null value.

Advantages :- ① Easy to Implement

② There is less distortion in variance ↗ ↘

Disadvantage :- ① Every situation randomness won't work.

③ Capturing NAN value with a new feature :- It works well if the data are not missing completely at random MNAR.

Advantages :- ① Captures the importance of missing values.

Disadvantage :- ① Creating Addition feature (Curse of Dimensionality)

	Age	Survived	Age-new	Age	Survived	Age-new
1	1	0	0	1	0	8
2	2	1	0	2	1	0
3	3	0	0	3	0	1
NAN		1	0	4	0	0
4	4	0	0			

Replace with Mean / Median

(4) End of Distribution Imputation :- In this technique we will replace NaN value with far end of the distribution. $[df.Age.mean() + df.Age.std() * 3]$

$[df.Age.mean() + df.Age.std() * 3] \leftarrow \text{same}$

Eg:- Age survived

22	0
38	0
26	1
NaN	0
30	0

Age survived

22	0
38	0
26	1
60	0
30	0

far end distribution

Eg:- Age Survived		Age Survived		Age - End-Distribution
22	0	22	0	22
38	1	38	1	38
26	0	26	0	26
NaN	1	23	1	60
30	0	30	0	30

median

far end distribution

Advantages:-

- ① Easy to implement.

- ② Captures the importance of missings if there is one.

Disadvantages:-

- ① Distorts the original distribution of the variable.

- ② If missings is not important, it may mask the predictive power of the original variable by distorting its distribution.

- ③ If the number of NA is big, it will mask true outliers in the distribution.

- ④ If the number of NA is small, the replaced NA may be considered an outlier and pre-processed in a subsequent step of feature engineering.

(5) Arbitrary Value Imputation :- This technique was derived from Kaggle competition it consists of replacing NaN with arbitrary value based on random choice or personal whim, rather than any reason or system. In these technique we will replace nan value with minimum or maximum value i.e. outlier value.

Eg:- Age survived

10	0
NaN	1
2	0
NaN	1
13	0
15	0
100	1
Ymax	1

Age survived

18	0
11	1
2	0
13	0
16	0
100	1

Age min

10
11
2
2
13
16
100

Age max

10
11
2
2
13
16
100

Day 3

Advan

Disadv

② If m

varia

③ Hard

Ham

① Fre

elem

Advan

Disadv

use

② I

② Capt

Note

Ham

② O

③ m

① L

②

No

① One

② On

③ M

(5)

replace NaN
[pre:std]

Day 3

- Advantages :- ① Capture the importance of missingness if there is one.
Disadvantages :- ① Distorts the original distribution of the variable 
② If missingness is not important, it may mask the predictive power of the original variable by distorting its distribution.
③ Hard to decide which value to use.

H

Handling Missing Value in Categorical features :-

- ① Frequent Category Imputation - Replacing NaN value with most frequent element i.e. mode. $\{ \text{for data missing } 10-15\% \}$

Advantages :- ① Faster way to implement

Disadvantages :- ① Since we are using the more frequent labels, it may use them in an over represented way, if there are many nans.
② It distorts the relation of the most frequent label.

- ② Capturing NaN value with a new feature - } Refer page no. 4 } { 1, 0 }

Note :- Suppose if we have more frequent categories, we will replace NaN with a new category i.e. $\{ \text{Missing} \}$ $\{ \text{for data missing more than } 50\% \}$

Handling Categorical Variable :-

- ① One Hot Encoding
- ② One Hot Encoding with many category $\{ \text{KDD} \text{ or Orange Cup} \}$
- ③ Mean Encoding $\{ \text{Pincode} \}$

① Label Encoding

② Target Guided Ordinal Encoding

{ **Ordinal Encoding** order $\{ \text{no imp} \}$ }

Nominal Encoding

Nominal Encoding $\{ \text{Order not Important} \}$

① One Hot Encoding

② One Hot Encoding with many category (KDD)

③ Mean Encoding $\{ \text{Pincode} \}$

P.T.O

Ordinal Encoding $\{ \text{Order Imp} \}$

① Label Encoding

② Target Guided Ordinal Encoding

(7)

Gender	Gender-Male	Gender-Female	Ex. Gender
Male	1	0	male
Female	0	1	Female
Female	0	1	
Male	1	0	

pd.get_dummies(dataset).head()
drop_first = True.

- ② One Hot Encoding with many category in a feature {Mercedes.csv}
- In this type of encoding we take first 10=now or more frequent element from each categorical column. {KDD Orange}
 - In mercedes dataset, we have column like [$x_0, x_1, x_2, x_3, x_4, x_5, x_6$] with many category value. {Highly Categorical}
 - In this only top 10 will replace with One Hot Encoding and rest of the value will become 0 which are not more frequent.

{Ordinal}

- ③ Ordinal Number Encoding:- In this type of encoding, we assign rank to our data based upon priority.

Ex. Education	Rank	Day	Rank
12 th	4	Monday	1
Graduation	3	Tuesday	2
PhD	2	Wednesday	3
Master	1	Thursday	4
		Fridgy	5
		Saturday	6
		Sunday	7

{Nominal}

④ Count or Frequency Encoding :- In these encoding, we will replace the value with its count.

Eg. Country	Country
India	3
USA	1
US	1
India	3
India	3
Cuba	1

Advantages:- ① Not increase feature space.

Disadvantage:- ① It will provide same weight if the frequency are same.

i.e. If India has count = 3 and Cuba has count 3 then the data become inappropriate.

{5} Target Guided {Ordinal Encoding} :- ① Ordering the labels according to the target.

② Replace the labels by the joint probability of being 1 or 0.

In titanic dataset, cabin column had value as E8S, D2S, A22, B33 and NAn.

→ we will replace NAn with "Missing".

→ Then we will take first letter of every column with [E, D, A, B, M]. mean → missing

→ Then sort the value [A, E, M, D, B] then assign rank to [A:1, E:2 ...]

→ Then replace in dataset [0, 1, 2, 3, 4, 5]

{6} Mean Encoding {Nominal} :- we will replace with mean.

m = dataset.groupby(['cabin'])['survived'].mean().to_dict()

print(m)

{'A': 0.46, 'B': 0.74, 'C': 0.59, 'D': 0.75}

Now, we will replace in our dataset.

df['mean-ordinal-encode'] = df['cabin'].map(m)

Survived cabin mean-ordinal-encode ① capturing information

0 M 0.59

within the label therefore rendering more predictive feature

1 C 0.59

② Monotonic relationship between variable and target.

1 M 0.69

0 C 0.59

① Prone to overfitting leads to ↑

1 A 0.46

0 B 0.74

Note :-

- ① Probability
- ② Probability
- ③ $\text{pr}(\text{survived})$
- ④ Dictionary
- ⑤ Replace

⑦ Probability Ratio Encoding :-

as "cabin" & "Survived".

Survived	Cabin
0	Nan
1	C85
1	Nan
1	C123
0	Nan

In titanic dataset, we will use column

we will replace Nan value with "missing".

`df['Cabin'].unique()`
[Missing, C85, C123, C4, A16, A27 ...]

`df['Cabin'].astype(str).str[0] = df['Cabin']`
`df.head()`

Survived	Cabin
0	M
1	C
1	M
1	C
0	M

`df.Cabin.unique()`
['M', 'C', 'E', 'G', 'D', 'A', 'B', 'F', 'T']

`pob=df.groupby(['Cabin'])['Survived'].mean()`

A 0.04 - B 0.466

B 0.744

C 0.75

⋮

T 0.00.00.

`pob = pd.DataFrame(pob)`

cabin	survived
A	0.466
B	0.744
C	0.75
⋮	⋮
T	0.00

`pob['Pob-ratio'] = pob['Survived'] / pob['Died']`
`pob.head()`

Cabin	Survived	Died	Probability-ratio
A	0.46	0.53	0.875
B	0.744	0.25	2.91
C	0.59	0.40	1.45
⋮	⋮	⋮	⋮
T	0.00	0.00	0.00

`pob['Pob-ratio'].to_dict() = encoded`
{'A': 0.875, 'B': 2.91, 'C': 1.45...}

`pob['Died'] = 1 - pob['Survived']`

Cabin	Survived	Died
A	0.46	0.53
B	0.744	0.25
C	0.59	0.40
⋮	⋮	⋮
T	0	1

`df['Cabin-encode'] = df['Cabin'].map(encoded)`
`df.head()`

Survived	Cabin	Cabin-encode
0	M	0.48
1	C	1.45
1	M	0.48
0	M	1.45
⋮	⋮	⋮

Type

① Nominal

③ Scale

Transform

① Standard

Scale

df = Standardize(df)

9

Note :-

- ① Probability of survived based on cabin \rightarrow Categorical Feature
- ② Probability of Not survived $\rightarrow 1 - \text{pr}(\text{Survived})$
- ③ $\text{pr}(\text{Survived}) / \text{pr}(\text{Not survived})$
- ④ Dictionary to map cabin with probability
- ⑤ Replace with the categorical feature.

Transformation of Features.

Why transformation of features are required? To scale in a similar unit
 To scale down the data and make calculation fast. Machine Learning algoth
 like Linear Regression, Logistic Regression, KNN, Kmeans, Hierarchical Clustering and
 Deep Learning Technique like ANN, CNN and RNN we required feature transformation.
 Note:- we don't apply transformation to all ML algo only above ~~mention~~ mention
 algo ~~we require transformation whereas in DL we require transformation.~~

Types of Transformation:-

- ① Normalization & Standardization ② Scaling to Min & Max value.
- ③ Scaling To Median & Quantiles ④ Gaussian Transformation { Logarithmic Transformation } { Reciprocal } { Square Root } { Exponential (exponent) } { Boxcox }

① Standardization :- we try to bring all the variables or features to a similar scale. Standardization means centering the variable at zero. $Z = \frac{x - \mu}{\sigma}$

df	Survived	Pclass	Age	Fare
0	3	22	7.25	
1	1	38	71.28	
1	3	26	7.92	
1	1	35	53.10	
0	3	35	8.05	

from sklearn.preprocessing import StandardScaler
 scalar = StandardScaler()

df_scaled = scalar.fit_transform(df)

[-0.789, 0.827, ...]

plt.hist(df_scaled['Survived'], bins=20)...

② Normalization :- {Mostly used in CNN DL}
 → Min Max scaling scales the value between 0 to 1.

$$x_{\text{scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

from sklearn.preprocessing import MinMaxScaler

min_max = MinMaxScaler()

df_minmax = pd.DataFrame(min_max.fit_transform(df), columns=df.columns)

Survived	Pclass	Age	Fare	plt.hist(df_minmax['Pclass'], bins=20)
0.0	1.0	0.27	0.014	
1.0	0.0	0.47	0.139	
1.0	1.0	0.32	0.015	
1.0	0.0	0.43	0.103	
0.0	1.0	0.43	0.015	

③ Robust Scaler :- It is used to scale the feature to median & quantiles

Scaling using median and quantiles consists of subtracting the median to all the observations, and then dividing by the interquartile difference. The interquartile difference is the difference between the 75% & 25% quantile.

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

$$x_{\text{scaled}} = (x - x_{\text{median}}) / IQR$$

df_robust_scaler = pd.DataFrame(scaler.fit_transform(df), columns=df.columns)
 df_robust_scaler.head()

Survived	Pclass	Age	Fare	plt.hist(df_robust_scaler['Age'], bins=20)
0.0	0.0	-0.46	-0.31	
1.0	-2.0	0.789	2.46	
1.0	0.0	-0.15	-0.28	
1.0	-2.0	0.53	1.67	

(11)

Gaussian Transformation :-

`df = pd.read_csv('titanic.csv', usecols=['Age', 'Cabin', 'Survived'])
df['Age'] = df['Age'].fillna(df['Age'].median())`

`import sys
scipy.stats as stat
import pylab`

Note:- If we want to check whether feature is gaussian or normal distribution use Q-Q plot.

~~def plot-data~~ def plot-data(df, feature):
plt.figure(figsize=(10, 6))
plt.subplot(1, 2, 1)
df[feature].hist()
plt.subplot(1, 2, 2)
stat.probplot(df[feature], dist='norm', plot=pylab)
plt.show()

Logarithmic Transformation

`import numpy as np
df['Age-Log'] = np.log(df['Age'])
plot-data(df, 'Age-Log')`

Exponential Transformation

`df['Age-Exponential'] = df.Age ** (1/1.2)
plot-data(df, 'Age-exponential')`

Reciprocal Transformation

`df['Age-Reciprocal'] = 1 / df.Age
plot-data(df, 'Age-reciprocal')`

Box-Cox Transformation
`stat.boxcox(df['Age'])` P70

Square root Transformation

`df['Age-square'] = df.Age ** (1/2)
plot-data(df, 'Age-square')`

(13)

Box Cox Transformation :- The Box Cox transformation is defined as:

$$T(y) = \frac{(y^{\lambda} - 1)}{\lambda}$$

where y is the response variable and λ is the transformation parameter.
 λ varies from -5 to 5 . In the transformation, all values of λ are considered and the optimal value for a given variable is selected.

`df['Age_Boxcox'] = parameters = stat.boxcox(df['Age'])`

`print(parameters)`

0.79

`plot_data(df, 'Age_Boxcox')`

(13)

Handling Imbalanced Dataset

{pip install imbalanced-learn}

we will use CreditCard Fraud Detection which has no missing value and output feature is "Highly Imbalanced." Class 0 - 284315 1 - 492

`x = df.drop("Class", axis=1)`

`y = df['Class']`

{for Imbalance dataset it is not worth}

{we will apply logistic regression to check what happen if our data is highly Imbalance? }

from sklearn.linear_model import LogisticRegression

from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

from sklearn.model_selection import KFold DO not check accuracy IF data

is highly Imbalanced.

`import numpy as np`

from sklearn.model_selection import GridSearchCV

`log-class = LogisticRegression()`

`grid = {'C': 10.0 ** np.arange(-2, 3), 'penalty': ['l1', 'l2']}`

`cv = KFold(n_splits=5, random_state=None, shuffle=False)`

from sklearn.model_selection import train_test_split

`x-train, x-test, y-train, y-test = train_test_split(x, y, train_size=0.7)`

`clf = GridSearchCV(log-class, grid, cv=cv, n_jobs=-1, scoring='f1-macro')`

`clf.fit(x-train, y-train)`

*use core of computer

`y-pred = clf.predict(x-test)`

`print(confusion_matrix(y-test, y-pred))`

`print(accuracy_score(y-test, y-pred))`

`print(classification_report(y-test, y-pred))`

`[[85251 47]`

`[41 104]]`

0.995

Ignore
bad model

precision recall f1 score support

:

:

:

:

:

ensemble technique work
well for Imbalance Dataset. (15)

i will also be doing
be always moving
working; will go with
oversampling

```
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier(class_weight = class_weight)  
classifier.fit(X_train, y_train)  
  
y_pred = classifier.predict(X_test)  
print(confusion_matrix(y_test, y_pred))  
print(accuracy_score(y_test, y_pred))  
print(classification_report(y_test, y_pred))
```

Under Sampling → Best for small dataset worst for Big Dataset.

```
from collections import Counter  
Counter(y_train)
```

```
from collections import Counter  
from imblearn.under_sampling import NearMiss  
nm = NearMiss(0.8)
```

```
X_train_ns, y_train_ns = nm.fit_sample(X_train, y_train)  
print("The number of classes before fit {} ".format(Counter(y_train)))  
print("The number of classes after fit {} ".format(Counter(y_train_ns)))
```

```
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier()  
classifier.fit(X_train_ns, y_train_ns)
```

```
y_pred = classifier.predict(X_test)  
print(confusion_matrix(y_test, y_pred))  
print(accuracy_score(y_test, y_pred))  
print(classification_report(y_test, y_pred))
```

Note:- Interview how do you handle Imbalanced dataset? → If the dataset is small i will definitely go with Under Sampling but i will be focusing on all the performance matrix recall precision f1 score apart from that i will also be focusing on based upon domain knowledge ~~PP~~ whehter i should reduce with reducing FP with FN based i will be selecting ROC Score that is the probability value 0.6 0.7 sometimes like that

```
from imblearn.OverSampling  
os = RandomOverSampler()  
X_train_ns, y_train_ns = os.fit_resample(X_train, y_train)  
print("X-train_ns, y-train_ns")  
print(X_train_ns)  
print(y_train_ns)  
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier()  
classifier.fit(X_train_ns, y_train_ns)
```

SMOTE Tomek

```
from imblearn.over_sampling import SMOTETomek  
os = SMOTETomek()  
X_train_ns, y_train_ns = os.fit_resample(X_train, y_train)  
print(X_train_ns)  
print(y_train_ns)  
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier()  
classifier.fit(X_train_ns, y_train_ns)
```

Ensemble - 1

```
from imblearn.ensemble import EasyEnsemble  
easy = EasyEnsemble()  
easy.fit(X_train_ns, y_train_ns)  
y_pred = easy.predict(X_test)  
print(y_pred)
```

15
i will also be doing SMOTE technique i will be oversampling technique ; will be always moving around understanding performance matrix finally it is not working ; will go with DT, Randomforest, XGBoost where I can also provide class weight.

16

Oversampling

```
from imblearn.over_sampling import RandomOverSampler  
os = RandomOverSampler(0.75)  
X_train_ns, y_train_ns = os.fit_sample(X_train, y_train)  
print(-11 - 11 - 11)  
print(-11 - 11 - 11)  
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier()  
classifier.fit(X_train_ns, y_train_ns)  
→ 4 line
```

SMOTETomek

```
from imblearn.combine import SMOTETomek  
os = SMOTETomek(0.75)  
X_train_ns, y_train_ns = os.fit_sample(X_train, y_train)  
print(-11 - 11)  
print(-11 - 11)  
from sklearn.ensemble import RandomForestClassifier  
classifier = RandomForestClassifier()  
classifier.fit(X_train_ns, y_train_ns)  
→ 4 line
```

Ensemble technique

```
from imblearn.ensemble import EasyEnsembleClassifier  
easy = EasyEnsembleClassifier()  
easy.(X_train, y_train)  
easy.  
y_pred = easy.predict(X_test)  
→ 3 line
```

(17)

Outlier Handlings -

- ① Naive Bayes Classifier - Not sensitive to outliers
- ② SVM - Not sensitive to outliers
- ③ Linear Regression - Sensitive to outliers
- ④ Logistic Regression - Sensitive to outliers
- ⑤ Decision Tree Regressor or classifier - Not sensitive to outliers.
- ⑥ Ensemble (RF, XGBoost, GB) - Not sensitive to outliers.
- ⑦ KNN - Not sensitive to outliers.
- ⑧ Kmeans - Sensitive to outliers.
- ⑨ Hierarchical - Sensitive to outliers.
- ⑩ PCA - Sensitive to outliers.
- ⑪ Neural Network - Sensitive to outliers.
- ⑫ DBScan - Sensitive to outliers.

```
import pandas as pd
```

```
df = pd.read_csv('titanic.csv')
```

```
df.head()
```

```
df['Age'].isnull().sum() OP - 177
```

```
import seaborn as sns
```

```
sns.distplot(df['Age'].dropna())
```

Now, we will add outlier in age column for demonstration.

```
sns.distplot(df['Age'].fillna(100))
```

Note: If the data is normally distributed we will use IQR to remove outliers.

Gaussian Distributed

```
figure = df.Age.hist(bins=50)
```

```
figure.set_title('Age')
```

```
figure.set_xlabel('Age')
```

```
figure.set_ylabel('No of Passenger')
```

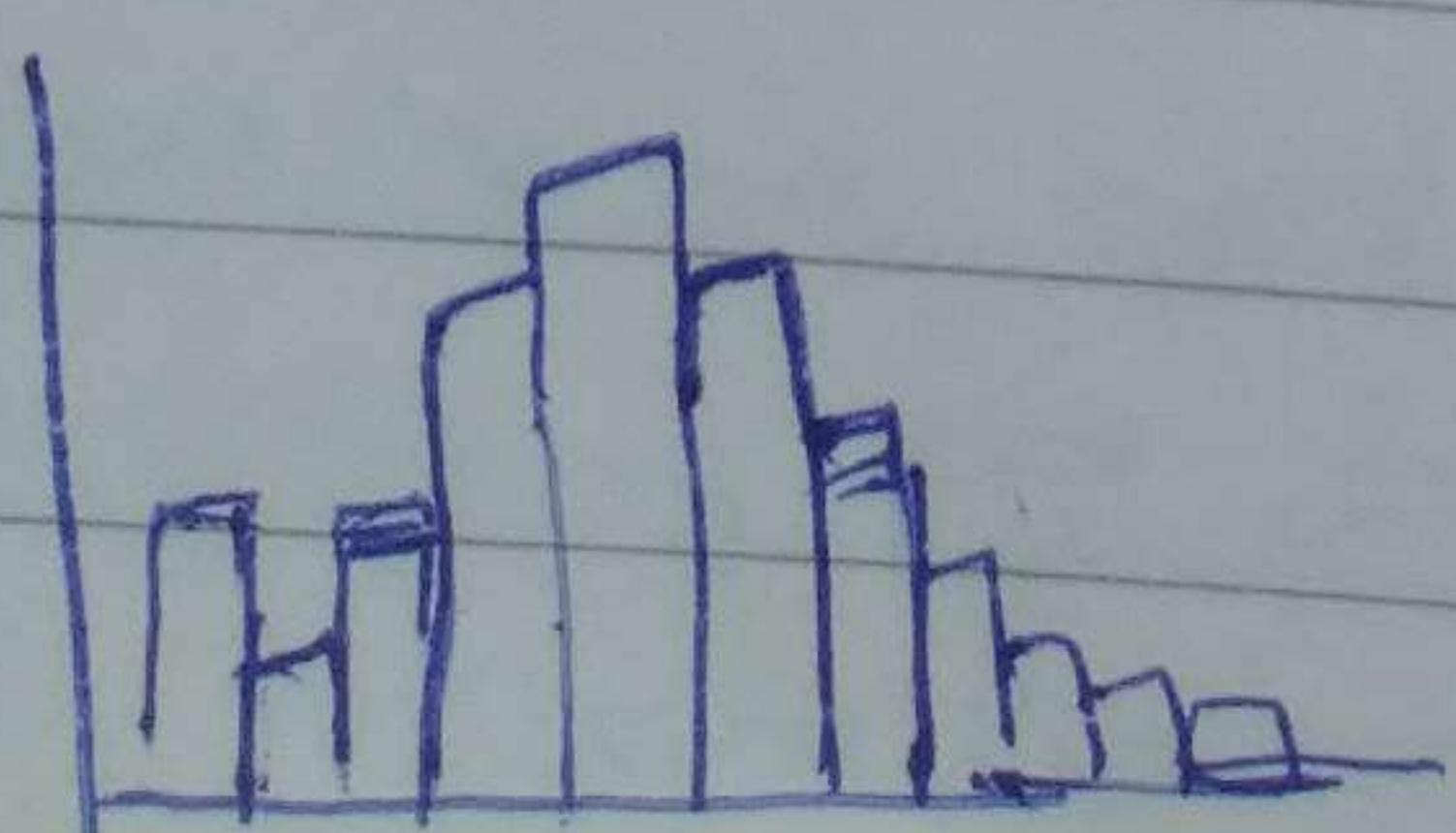


figure = df.
df['Age'].de

Assuming the
the bounda

upper-bound
lower-bound
print(upper-
-13.88, 73

IQR = df.

17.875

lower-tence
higher-fen

lower-brif
higher-fend

Skewed

figure = d

figure.set

figure.set

figure.set

df.boxplot

df['Fam

(17)

```
figure = df.boxplot(column = 'Age')
df['Age'].describe()
```

Assuming the age follows a Gaussian Distribution we will calculate the boundaries which differentiates the outliers.

$$\text{upper-boundary} = \text{df}['Age'].mean + 3 * \text{df}['Age'].std()$$

$$\text{lower-boundary} = \text{df}['Age'].mean - 3 * \text{df}['Age'].std()$$

```
print(upper-boundary), print(lower-boundary), print(df['Age'].mean())
-13.88, 73.27, 29.69
```

$$IQR = \text{df}.Age.quantile(0.75) - \text{df}.Age.quantile(0.25)$$

$$17.875$$

$$\text{lower-fence} = \text{df}['Age'].quantile(0.25) - (IQR * 1.5) \quad \{-6.68\}$$

$$\text{higher-fence} = \text{df}['Age'].quantile(0.75) + (IQR * 1.5) \quad \{64.81\}$$

$$\text{lower-extreme} = \text{df}['Age'].quantile(0.25) - (IQR * 3) \quad \{-33.5\} \quad \{\text{For Extreme}\}$$

$$\text{higher-extreme} = \text{df}['Age'].quantile(0.75) + (IQR * 3) \quad \{91.6\} \quad \{\text{outlier}\}$$

Skewed feature:

```
figure = df.Fare.hist(bins=50)
```

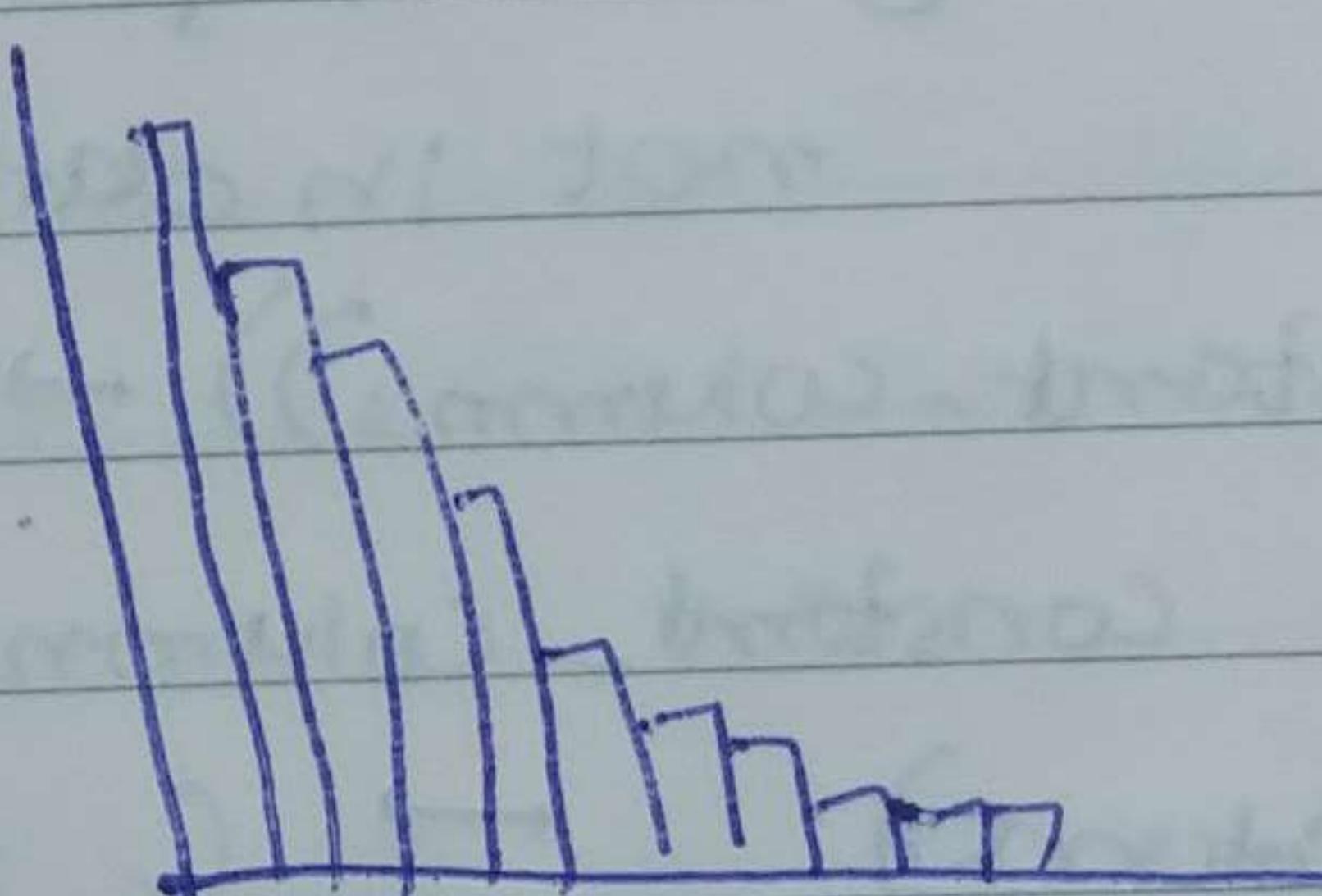
```
figure.set_title('Fare')
```

```
figure.set_xlabel('Fare')
```

```
figure.set_ylabel('No of passenger')
```

```
df.boxplot(column="Fare")
```

```
df['Fare'].describe()
```



(18)

(1)

Feature selection

Feature selection - {Dropping constant feature} - In this step we will be removing the features which have constant features which are actually not important for solving the problem statements.

```
import pandas as pd
data = pd.DataFrame({ "A": [1, 2, 4, 1, 2, 4], "B": [4, 5, 6, 7, 8, 9],
                      "C": [0, 0, 0, 0, 0, 0], "D": [1, 1, 1, 1, 1, 1] })
```

① Variance Threshold :- Feature selector that removes all low-variance feature. This feature selection algorithm looks only at the feature (X), not the desired outputs (y), and can thus be used for unsupervised learning.

It will remove zero variance feature.

```
from sklearn.feature_selection import VarianceThreshold
var_thresh = VarianceThreshold(threshold=0)
var_thresh.fit(data) → VarianceThreshold(threshold=0)
```

```
var_thresh.get_support() → array([True, True, False, False])
data.columns[var_thresh.get_support()] → Index(['A', 'B'], dtype='object')
```

```
constant_columns = [column for column in data.columns if column
                    not in data.columns[var_thresh.get_support()]]
print(len(constant_columns)) → 2
```

```
for feature in constant_columns:
    print(feature) → C D
```

```
data.drop(constant_columns, axis=1)
```

A	B
1	4
2	5
4	6
1	7
2	8
4	9

② Feature selection
removing the
we will take b
test split.
x_train, corr()
import seaborn as
plt.figure(figsize
corr = x_train.c
sns.heatmap(corr
plt.show()

now, we will rem

Note:- If indep
do not drop
with other indep
def correlation

```
col_corr =  
corr_matrix =  
for i in  
    for j in  
        if
```

```
corr_features  
len(set(corr  
corr_features
```

```
X_train_dro  
X_test_dro
```

(1)

will be actually
 ② Feature selection - with Correlation :- In this step we will be removing the features which are highly correlated.
 we will take boston dataset and split feature in x & y then train test split.

$x_train, corr()$

import seaborn as sns

plt.figure(figsize=(12, 10))

corr = $x_train, corr()$

sns.heatmap(corr, annot=True, cmap=plt.cm.coolwarm_r)

plt.show()

Now, we will remove feature which are highly correlated.

Note:- If independent feature is highly correlated with dependent feature do not drop that column. If one independent feature is highly correlated with other independent feature then we can drop one column.

def correlation(dataset, threshold):

col_corr = set()

corr_matrix = dataset.corr()

for i in range(len(corr_matrix.columns)):

for j in range(i+1):

if abs(corr_matrix.iat[i, j]) > threshold:

colname = corr_matrix.columns[i]

col_corr.add(colname)

return col_corr

corr_features = correlation(x_train, 0.7)

len(set(corr_features))

corr_features

$x_train.drop(corr_feature, axis=1)$

$x_test.drop(corr_feature, axis=1)$

(2)

③ Feature Selection - Information Gain - Mutual Information in Classification Problem Statements.

Mutual Information estimate mutual information for a discrete target variable. Mutual Information between two random variable is a non-negative value, which measure the dependency between the variable. It is equal to zero if and only if two random variables are independent, and higher value mean higher dependency. { use wine data, $X_f Y$, train test split. }

```
from sklearn.feature_selection import mutual_info_classif {0 to 1}
mutual_info = mutual_info_classif(X_train, y_train)
mutual_info
```

```
from sklearn.feature_selection import SelectKBest
sel_five_cols = SelectKBest(mutual_info_classif, k=5)
sel_five_cols.fit(X_train, y_train)
X_train.columns[sel_five_cols.get_support()]
```

Difference Between Information Gain and Mutual Information

$$I(X;Y) = H(X) - H(X|Y) \text{ and } IG(s,a) = H(s) - H(s|a)$$

As such, mutual information is sometimes used as a synonym for information gain. Technically, they calculate the same quantity if applied to the same data.

③

④ Feature Selection Problem Statement

Mutual Information

we will use

from sklearn.feature

mutual_info = m

mutual_info

mutual_info = p

mutual_info = i

mutual_info = s

from sklearn.feature

selected_top_c

selected_top_u

selected_top_p

selected_top_s

selected_top_r

selected_top_d

selected_top_m

selected_top_n

selected_top_l

selected_top_k

selected_top_j

selected_top_i

selected_top_h

selected_top_g

selected_top_f

selected_top_e

selected_top_d

selected_top_c

selected_top_b

selected_top_a

selected_top_z

selected_top_y

selected_top_x

selected_top_w

selected_top_v

selected_top_u

selected_top_t

selected_top_s

selected_top_r

③

③ Feature Selection - Information Gain - mutual information in Regression

information in
a discrete target
variable is a non-negative
value. It is equal to zero
and higher value
put.

if $\{0 \text{ to } 1\}$

mutual information estimate mutual information for a continuous target variable.
we will use housing dataset in that only numerical feature x_{4Y} , train test.

from sklearn.feature_selection import mutual_info_regression

mutual_info = mutual_info_regression(x_train.fillna(0), y_train)

mutual_info

mutual_info = pd.Series(mutual_info)

mutual_info.index = x_train.columns

mutual_info.sort_values(ascending=False)

from sklearn.feature_selection import selectPercentile

selected_top_columns = selectPercentile(mutual_info_regression, percentile=20)

selected_top_columns.fit(x_train.fillna(0), y_train)

selected_top_columns.get_support()

x_train.columns[selected_top_columns.get_support()]

sym for
quantity if \rightarrow

④ Chi-square Test for Feature Selection:

Compare chi-squared stats between each non-negative feature & class.
This score can be used to evaluate categorical variable in a classification test
 \rightarrow This score can be used to select the n-feature features with the highest
values for the test. chi-squared statistic from x, which must contain only
non-negative features such as booleans or frequencies (e.g. terms counts in
document classification), relative to the class.

\rightarrow Recall that the chi-squared test measure dependence between stochastic
variables, so using this function "weed out" the feature that are the most likely
to be independent of class and therefore irrelevant for classification. The Chi
square statistic is commonly used for testing relationship between categorical
variables.

- It compares the observed distribution classes of target Y among the different category of the feature, against the expected distribution of the target class, regardless of the feature categories.
- We will use titanic dataset categorical feature with feature encoding of categorical data: X4Y, test-train split.

```
from sklearn.feature_selection import chi2
```

```
f_p-value = chi2(X_train, Y_train)
```

```
f_p-values ([...]), ([...])
```

```
import pandas as pd
```

```
p-values = pd.Series(f_p-value[1])
```

```
p-values.index = X_train.columns
```

```
p-values
```

```
p-values.sort_index(ascending=False)
```

Observation

Sex column is the most important column when compared to the output feature Survived.