

Platform : Google Colab  
Pg 3/11

**(EDA)** End-to-End Exploratory Data Analysis  
process to perform investigation on Data to discover patterns, some other things to know your Data well

- ① Understand the Data
- ② Data Cleaning
- ③ Feature Selection
- ④ Assumption Checking

### Library Required :

① Pandas :  
i) python library for data manipulation and Analysis  
ii) Data structures = Dataframes & Tabular Data

② Seaborn :  
i) statistical Data Visualization library.  
ii) provides informative blueprints  
couplets histograms barplots ...

```
import pandas as pd
import numpy as np
import Seaborn as sns
import matplotlib.pyplot as plt

sns.set(style='whitegrid')
```

### Step 2 : Data Loading & initial Inspection

```
!git clone 'https://...'
```

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

**DATA SET**

Pclass	- Passenger class (1=1st, 2=2nd, 3=3rd)
Survived	- (0=No, 1=Yes)
Name	- Name
Sex	- Sex
Age	- Age
SibSp	- Number of Siblings/Spouse
Parch	- Number of Parents/Children
Ticket	- Ticket Number
fare	- fare
Cabin	- Cabin
Embarked	- C=Cherbourg, S=Southampton, Q=Queenstown
boat	- Lifeboat
body	- Body Identification number
home.dest	- Home Destination

```
df.head() // first 5 entries
df.tail() // last 5 entries
df.shape // (Rows, Columns)
df.info() // Attribute
```

Problem: Some values are missing. How to fill those values?

### Column Manipulation

```
create new column like Has.Cabin (Yes/No);
```

```
print("Descriptive Statistics:")
df.describe()
```

```
Count mean std min 25% 50% 75% max
```

### How To Filter out Data.

```
Column wise df['Age'] // Display only column
df['Age'] > 18 // Display True or False
df[df['Age'] > 18] // df:
```

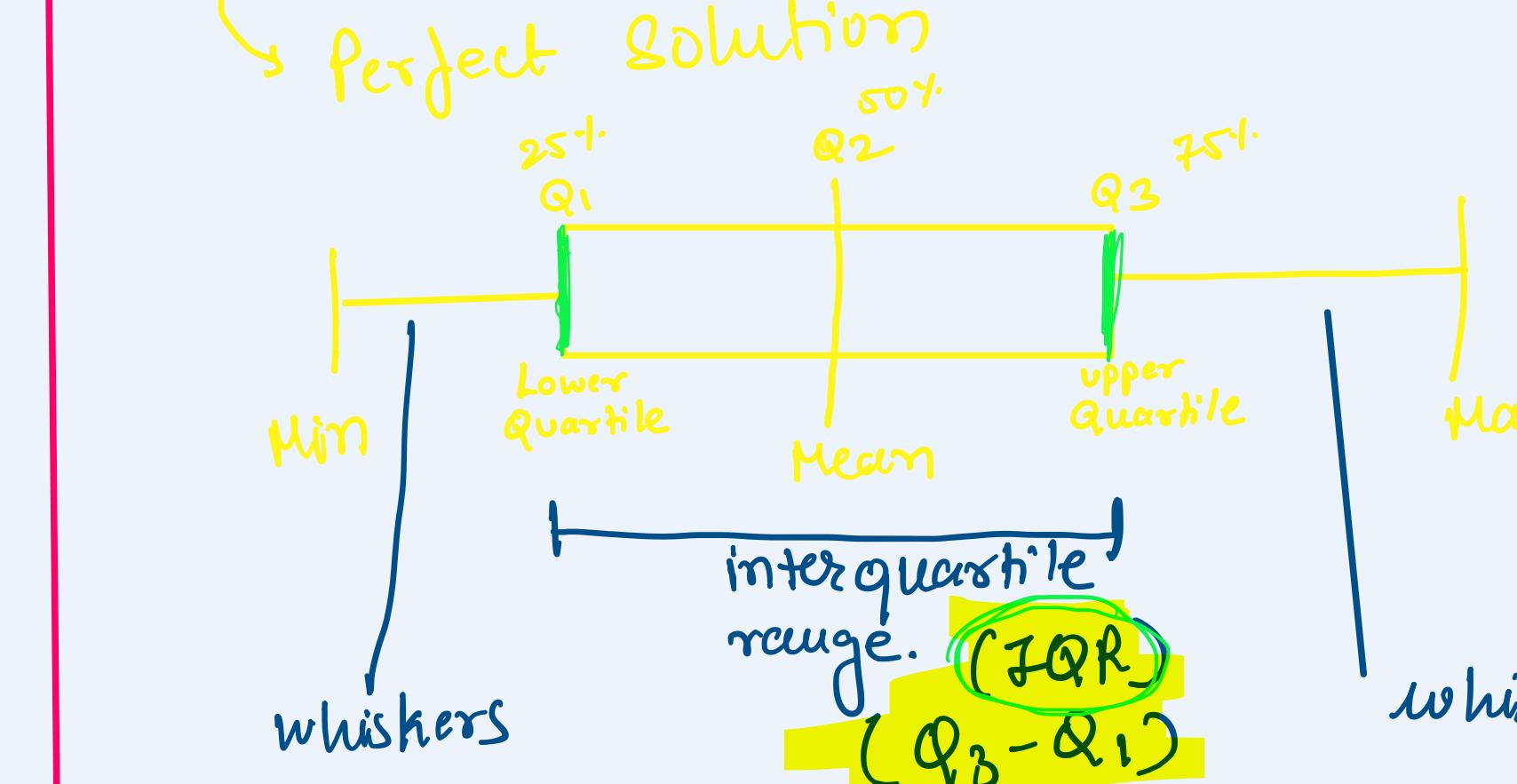
```
df[df['Embarked'] == 'C']
```

```
df[(df['Sex'] == 'female') &
    (df['Embarked'] == 'C')]
```

\*\* Descriptive Statistics (contd.)

### Step 3 : Data Cleaning

**Outliers** → An entity which does not belongs to that range



### Missing Value Imputation

- ① Numerical Data (Age / Fare) → Median
- ② Categorical Data → Mode (most frequent value)
- ③ High Cardinality / Too many Missing (cabin) → Drop the column or Create separate column

```
df.isnull().sum() // shows which columns has how many null values
```

```
# To handle the missing value values
Calculate, print, & fill that value
```

```
median_age = df['Age'].median()
print(median)
df['Age'].fillna(median_age) // df['Age']
```

### To find mode with index, mapping

```
mode_embarked = df['Embarked'].mode()[0]
df['Embarked'].fillna(mode_embarked)
```

### To handle Cabin values.

```
df['Has.Cabin'] = df['Cabin'].notna().astype(int)
df.drop(['Cabin'], axis=1, inplace=True)
```

check with:

```
df
df['Has.Cabin'].value_counts()
df.isnull().sum()
```

### STEP : UNIVARIATE Analysis

> Data is Analyzed contains only one Variable.

> To find patterns

- ① Categorical variable → Frequency Table, Bar charts (count plot), Pie charts
- ② Numerical Variable → Histogram, Kernel Density Plot, Box Plot

Point ("Analyzing Categorical features")

```
# Set up
fig, axes = plt.subplots(2, 3, figsize=(18, 12))
fig.suptitle('Univariate Analysis', fontsize=16)
```

```
# plotting
sns.countplot(ax=axes[0, 0], x='Survived', data=df)
set_title('Survived Distribution')
```

```
sns.countplot(ax=axes[0, 1], x='Pclass', data=df)
sns.countplot(ax=axes[0, 2], x='Gender', ...)
```

```
sns.countplot(ax=axes[1, 0], x='Sex', ...)
```

```
sns.countplot(ax=axes[1, 1], x='Embarked', ...)
```

```
sns.countplot(ax=axes[1, 2], x='...', ...)
```

```
plt.tight_layout(rect=[0, 0, 1, 0.95])
plt.show()
```

### Numerical Distribution

- ① Histogram
- ② K.D.

```
(as line over histogram)
```

```
fig, axes = plt.subplots(1, 2, figsize=(16, 8))
sns.kwplot(ax=axes[0], x=df['Age'],
            x_kde=True,
            bins=30).set_title('Age Distrb')
```

```
plt.show()
```

### STEP 5 : Bivariate Analysis

① Target : Survived → primary

- Types:  
 ① Categorical vs Numerical : Barplots, Box plots, Violin Plots  
 ② Categorical vs Categorical : Crosstabs  
 ③ Numerical vs Numerical : Scatter plots

```
fig, axes = plt.subplots(2, 2, figsize=(16, 12))
fig.suptitle('Bivariate Analysis', fontsize=16)
```

```
sns.barplot(ax=axes[0, 0], x=df['Pclass'], y=df['Survived'], data=df)
-- Sex vs Survival
Embarked vs Survival
Mas.Cabin vs Survival
```

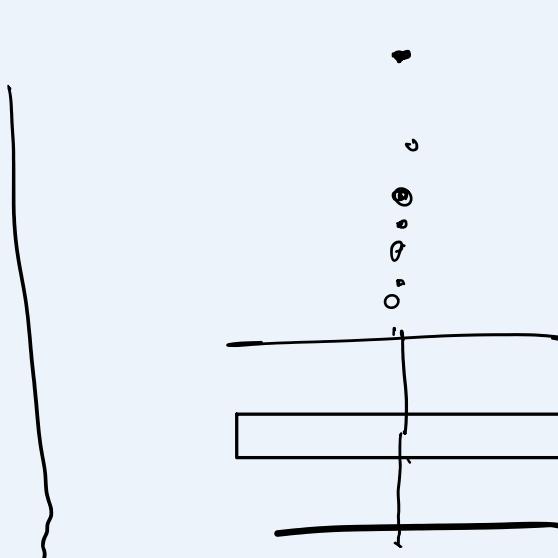
```
plt.tight_layout(rect=[0, 0, 1, 0.95])
plt.show()
```

### Categorical vs Numerical [Age]

```
g = sns.FacetGrid(df, col='Survived', height=6)
g.map(sns.histplot, 'Age', bins=25, kde=True)
plt.title('Survival vs Age', y=102)
plt.show()
```

### # Deeper Dive - Outliers Analysis

```
plt.figure(figsize=(10, 8))
sns.boxplot(y='Fare', data=df)
plt.title('Box plot for Ticket Price')
plt.ylabel('Fare')
plt.show()
```



### # Feature Engineering

→ create new features from existing one to uncover & get useful information.

Techniques:

- ① Combining features (Sibling + Parents = Family)
- ② Extracting from Text (Extracting title from Name)
- ③ Binning (Binning Age into 'child', 'Adult', 'Senior')

### # Create a 'Family Size' Column

```
df['FamilySize'] = df['SibSp'] + df['Parch'] + 1
```

### # Creating a 'IsAlone' Features

```
df['IsAlone'] = 0
df.loc[df['FamilySize'] == 1, 'IsAlone'] = 1
```

```
print("Created 'FamilySize' & 'IsAlone' feature")
df[['FamilySize', 'IsAlone']].head()
```

Now create barplot using these features

```
# 3. extract 'Title' from 'Name'
df['Title'] = df['Name'].str.extract('([A-Za-z]+[.])', expand=False)
```

# Let's see the different titles

```
print("Extracted Titles:")
df['Title'].value_counts()
```

### # TO Replace Some Titles with Something

```
df['Title'] = df['Title'].replace('Lady', 'Countess', ...)
```

```
Mlle → Miss
```

```
Mrs → Mrs
```

```
Master → Mose
```

### # Multi-variate Analysis

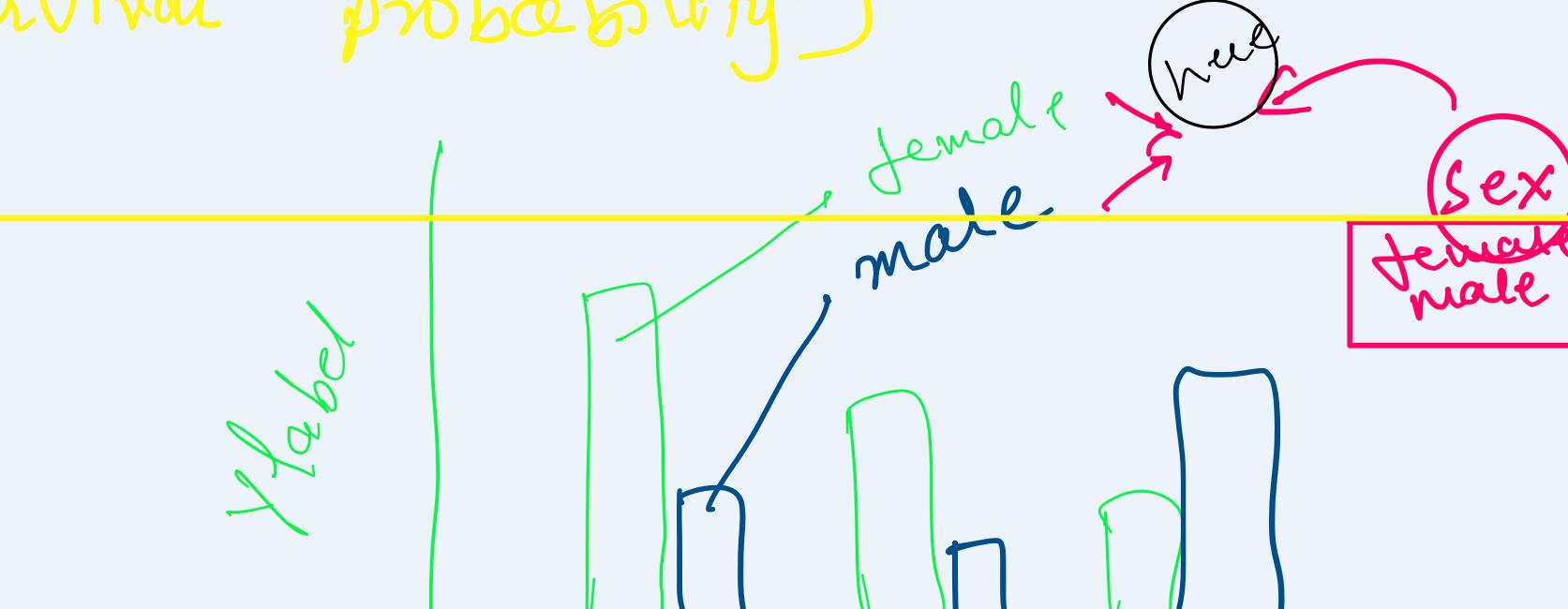
interaction between Diff Variables Simultaneously

```
sns.catplot(x='Pclass', y='Survived', hue='Sex', data=df,
            kind='bar', height=6, aspect=1.5)
```

```
plt.title('Survival Rate by Pclass & sex')
```

```
plt.ylabel('Survival probability')
```

```
plt.show()
```

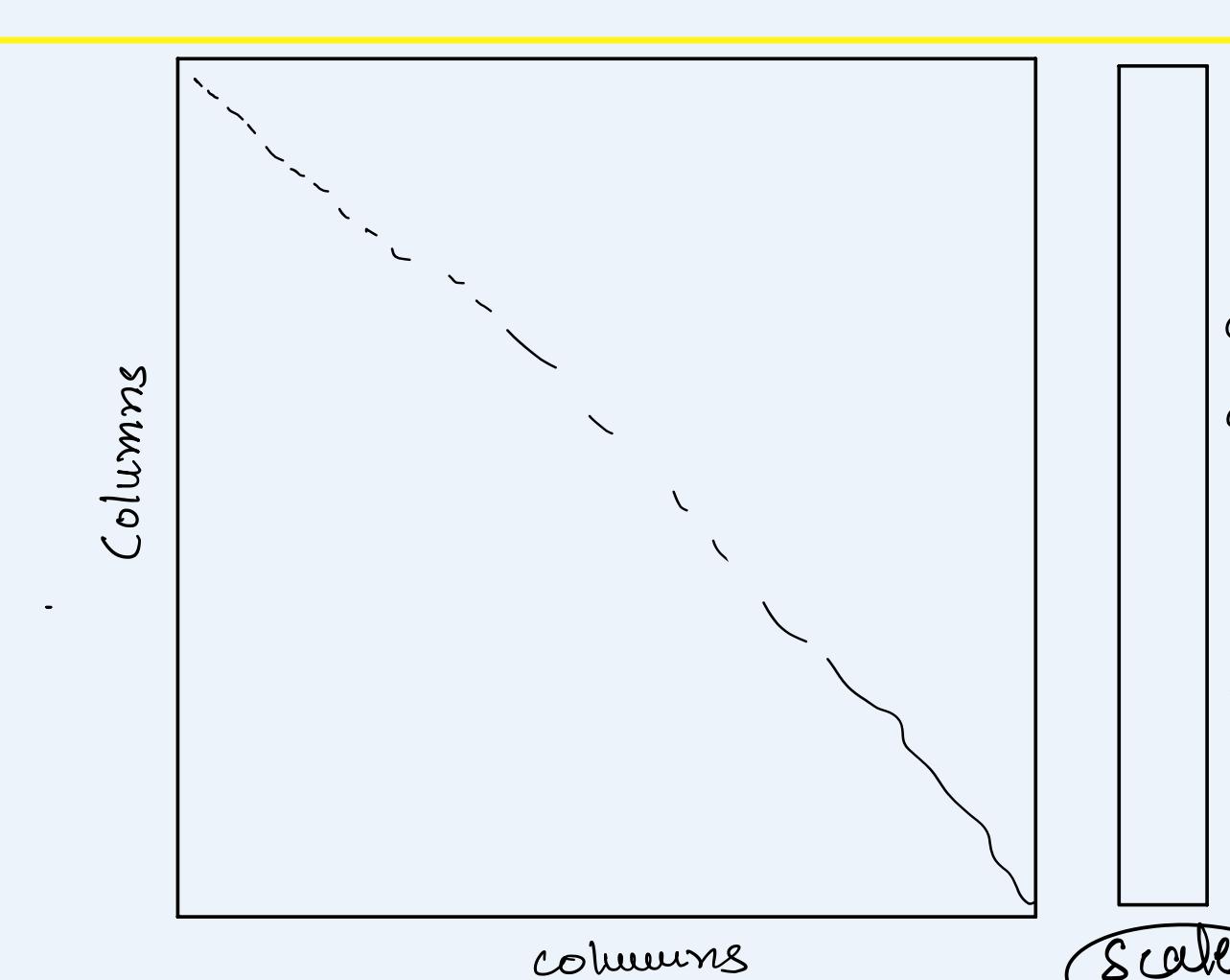


### # Correlation Analysis

How one value acts on other value dependency.

```
plt.figure(figsize=(14, 10))
numeric_cols = df.select_dtypes(include='mpy.number')
Correlation Matrix = numeric_cols.corr()
```

```
sns.heatmap(Correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Matrix of Numerical Features')
plt.show()
```



### Profiling Summary

### Summary

### Violin plot

```
plt.figure(figsize=(14, 8))
sns.violinplot(x='Sex', y='Age', hue='Survived', data=df, split=True,
                palette={0: 'blue', 1: 'orange'})
plt.title('---')
plt.show()
```

