

Date: 30-04-2025

1 Test 1 - EDA
2 Test 2 - Stats

To do:

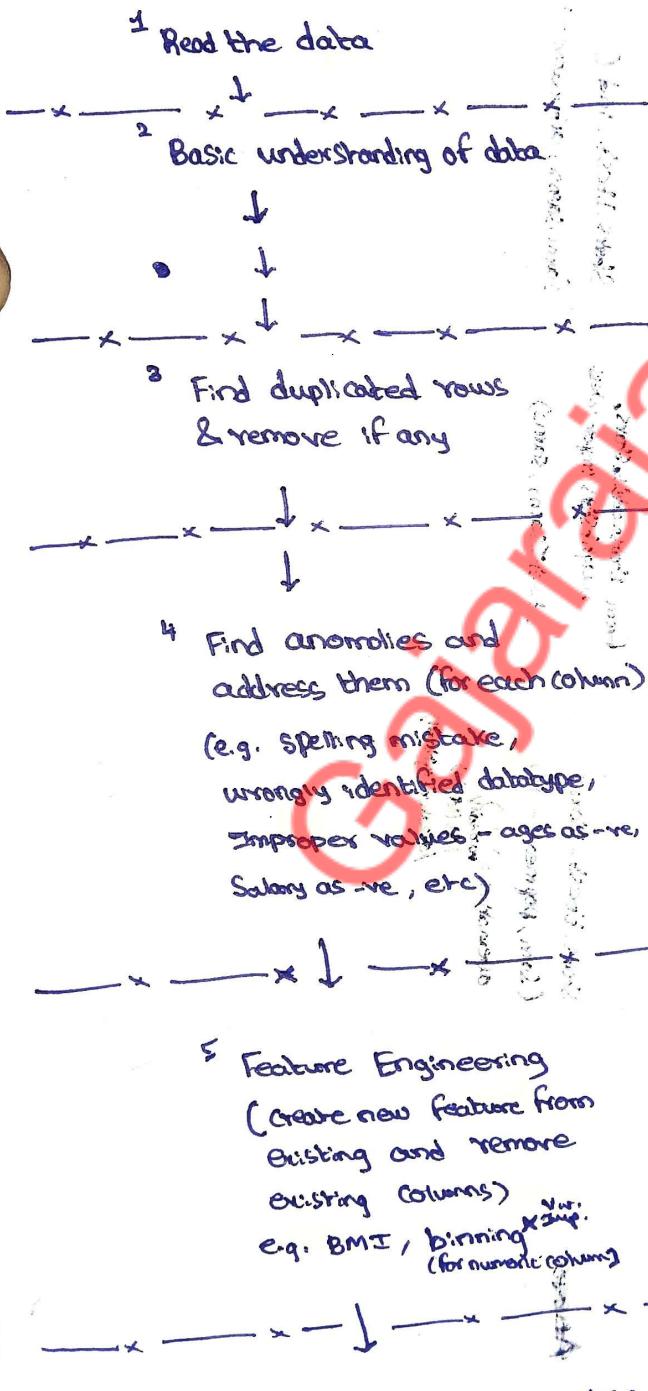
- Stats sumup (day-2)
- ✓ EDA sumup
- ✓ Stats imp. Python snippets
- ✓ EDA imp. Python snippets

- 5 Case study in class & take-home
- 6 Interview Qns - Stats.
- 7 Interview Qns - EDA
- 8 Stats other concepts (day 1/2)

a
Replace
remove
binning
date func
drop
delele
stringfunc

EDA Sum up

Flow of process



Python snippets for each

pd.read_csv()

df.info(), df.shape, df.size

[Note: pandas is on top of numpy]

↳ so attributes & methods.

a df.duplicated().sum()

b df[~df.duplicated()]

c df[col].str.replace(pattern, replace)

↳ old(key) → new(value)

d df[col].astype(int)

e df[col].astype(str)

f df[col].astype('category')

g pd.to_numeric(df[col], errors='coerce')

h pd.to_datetime(df[col], errors='coerce')

i df[col].astype(float)

j pd.cut(df[col], bins=, labels=, include_lowest=True)

k df['Salary'].transform(lambda x: x/1000)

l df['isometry'].apply(lambda x: x+1.1)

m df['dept'].map(dict)

↳ old(key) → new value

n df[age].apply(lambda x: abs(x) if x<0 else x)

o df[age].apply(lambda x: abs(x) if x<0 else x)

p df.groupby('dept')[['age']].transform(lambda x: x.fillna(x.median()))

q df.dropna() → drop rows with any missing value

r df.drop(['column1', 'cols2=1]) → specific column drop

s df.dropna(subset=[colname])

6 Select numeric and categorical columns separately

a df.select_dtypes(include=np.number)

• columns.to_list()

b df.select_dtypes(exclude=np.number)

• columns.to_list()

7 Basic descriptive statistics (five point summary)

a df[num].describe()

b df[cat].describe()

8 Treat missing values

(i) Single value treatment -

mean, median [num] df[col3].ffill()

mode [cat]

(ii) Multi value treatment - bfill, ffill

df['age'].bfill(), df['age'].ffill()

(iii) Logical imputation

df['age'].apply(lambda x: df['age'].median() if x < 0 else x).fillna(df.groupby('dept')[['age']].transform('median'))

(iv) Remove the observation / row / column

list of num. col names

9 Treat "outliers" [only for numeric columns]

[as cat. only imbalance]

identify {

[entire row is outlier]

{ iqr

[if only one value in it is outlier]

{ z-score

Q1 = df[num].quantile(0.25)

Q3 = df[num].quantile(0.75)

iqr = Q3 - Q1

uw = Q3 + 1.5 * iqr

lw = Q1 - 1.5 * iqr

df[(df[num] < lw) | (df[num]

> uw), anyaxis=1]]

to identify the outliers (rows)

10 Advance descriptive statistics [To get claims for sample]

(i) Univariate analysis

} for each metric & viz.

(ii) Bivariate analysis

(iii) Multivariate analysis

11 Inferential Statistics [hypo. testing & estimation]

(i) One sample testing - z/t test

(ii) Two sample testing - z/t test, matched pair t-test

(iii) more than two sample testing - Anova, Tukey test, Kruskal wallis

(iv) Categorical columns - Chi-square

(v) Pearson r

Piechart



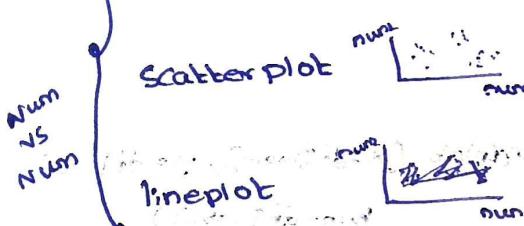
count = df['cat'].value_counts()

Count.plot(kind='pie', labels=Count, index=False, autopct='%1.1f%%')

Bivariate analysis

Metric / plot

Correlation



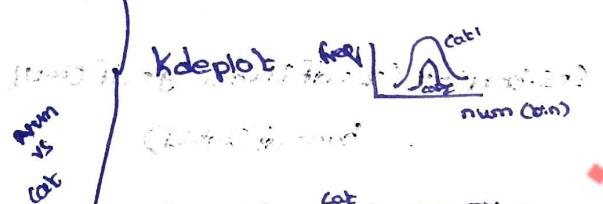
Python code:

df['num1'].corr(df['num2'])

Sns.scatterplot(x=df['num1'], y=df['num2'])

plt.plot(df['num1'], df['num2'], 'xg', linestyle='--')

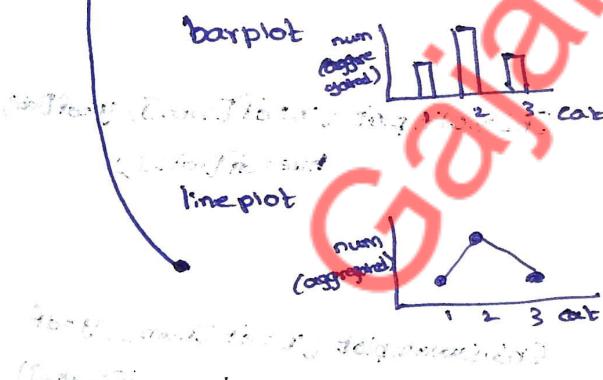
df.groupby('cat')[num].mean()



Sns.kdeplot(x=df['num'], hue=df['cat'])

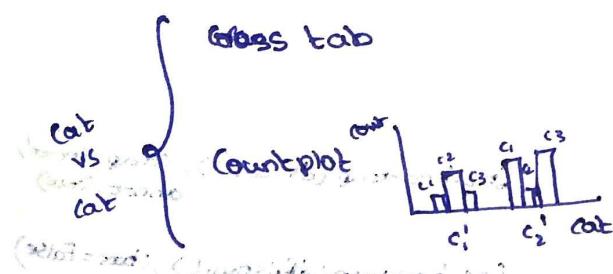
Sns.boxplot(x=df['num'], y=df['cat'])

Sns.boxplot(x=df['cat'], y=df['num'])



Sns.lineplot(x=df['cat'], y=df['num'])

Pd.crosstab(df['cat1'], df['cat2'])



Sns.countplot(x=df['cat1'], hue=df['cat2'])

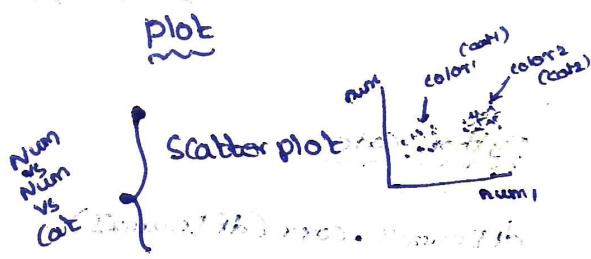
Logistic

and etc

etc

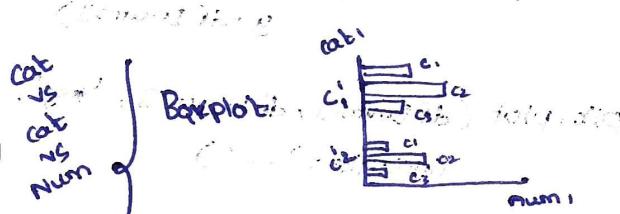
Decision Tree

Multivariate Analysis

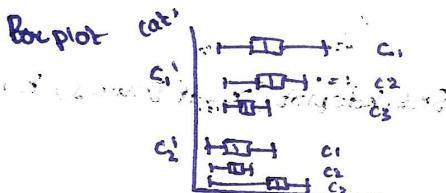


Python code

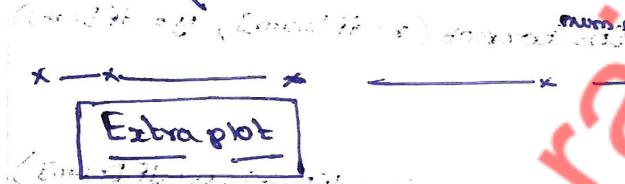
```
Sns.scatterplot(x=df['num1'], y=df['num2'],
                 hue=df['cat1'])
```



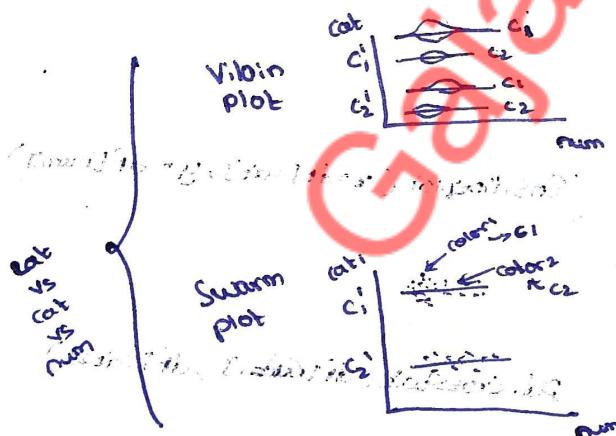
```
Sns.boxplot(x=df['num1'], y=df['cat1'],
             hue=df['cat2'])
```



```
Sns.buplot(x=df['num1'], y=df['cat1'],
            hue=df['cat2'])
```



```
Sns.violinplot(x=df['num1'], y=df['cat1'],
                hue=df['cat2'])
```



```
Sns.swarmplot(x=df['num1'], y=df['cat1'],
                hue=df['cat2'])
```

Heatmap



Correlation



nullvalues



with hue = categorical



w/o hue = categorical

```
Sns.heatmap(df['corr'], cmap='Greens',
            annot=True)
```

```
Sns.heatmap(df['isnull'], cmap=False)
```

```
Sns.pairplot(df, hue=df['cat1'])
```

```
Sns.pairplot(df)
```

B/w
numerical
values ← pairplot

Inferential Statistics

Num vs Num

① Pearson r

H_0 : not related variables / samples

Prerequisite: a. Data normal
b. Equal variance H_1 : Realized Sample
(so different units)

as num vs num

then if it's not met → p-value

Python code

```
Stats.pearsonr(Sam1, Sam2)
```

Num vs Cat

Master: If sample size for both is small

① If cat = 2(samples)

a) Two Sample T-test

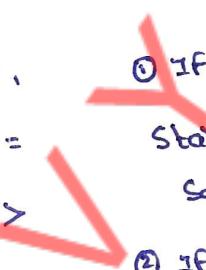
H_0

H_1

Prerequisite: a. Independent Samples
b. Data normal
c. Equal variance or not

if not met then → Welch's t-test, F-test, ...
Monotonicity

[on Num vs Mo]



② If equal variance

```
Stats.ttest_ind(Sam1, Sam2, alternative="...")
```

③ If unequal variance

```
Stats.ttest_ind(Sam1, Sam2, alternative="...", equal_var=False)
```

② If cat > 2(samples)

a) ANOVA

Prerequisite: a. Independent Samples
b. Data normal
c. Equal variance

b) Tukey HSD (after ANOVA)

H_0 : =

H_1 : ≠

Cat vs Cat

① Chi-square (contingency)

H_0 : not related

H_1 : realized

Finality:

1. Reject H_0

2. fail to reject H_0

```
import statsmodels.stats.multicomp as mc
mc.friedman_test(*obs)
report=mc.pairwise_tukeyhsd
(dft[cat], df[cat])
report.summary()
```

Stats.chi2_contingency
(observed data)
↳ is crosstab

reject - reject

Fail

Feature Engineering

Encoding

Ordinal Encoding

Python code

```
from sklearn.preprocessing import OrdinalEncoder  
o = OrdinalEncoder(categories=[['low', 'med', 'high']])  
df[['col1']] = o.fit_transform(df[['col1']])
```

One Hot Encoding

Python code

```
df_new = pd.get_dummies(df, columns=[  
    ['col1', 'col2', ...]]).dtypes.int
```

(N-1) dummies

```
df_new = pd.get_dummies(df, columns=[...],  
    drop_first=True)
```

freq. encoding

```
freq = df['col1'].value_counts(normalize=True)  
df['col1'] = df['col1'].replace(freq)
```

label Encoding

```
from sklearn.preprocessing import LabelEncoder  
l = LabelEncoder()
```

```
df['col1'] = l.fit_transform(df['col1'])
```

Scaling / Data normalization

Standard Scaler

```
from sklearn.preprocessing import StandardScaler  
s = StandardScaler()
```

```
df['col1'] = s.fit_transform(df[['col1']])
```

Min-Max

```
m = MinMaxScaler()
```

```
df['col1'] = m.fit_transform(df[['col1']])
```

Robust Scaler

X-median
IQR

```
from sklearn.preprocessing import RobustScaler
r = RobustScaler()
df[col] = r.fit_transform(df[col])
```

Data Transformation

Box-Cox
(> 0)

```
from sklearn.preprocessing import PowerTransformer
P = PowerTransformer(method='box-cox')
df[col] = P.fit_transform(df[col])
```

Yeo-Johnson
(+, -, 0)

```
from sklearn.preprocessing import PowerTransformer
P = PowerTransformer()
df[col] = P.fit_transform(df[col])
```

Date handling

* pd.to_datetime():

- ↳ expected 'ip' format: /./m./y.2 /./y ./.H : %M
- ↳ format = mixed ("if or /" as months in diff. rows)
- ↳ errors='coerce'
- ↳ if ip: Null, String: then OLP: NAT (Not Available Time Stamp)

* date & time extract

- ↳ df[col].dt.date
- ↳ df[col].dt.time

String handling

e.g.: df[col].str.upper()

* df[col].str[3]

* df[col].str.contains('Tech', case=False, na=False)

* df[col].str.replace('Human', 'People', case=False)