



EDA HandWritten Notes

- Darshan R M

"Learn, Share,
and
Collaborate"

Week - 14

Exploratory data analysis

EDA

T

→ ~~it's~~ done to understand specific data set.

Wine

* `summary → df.info()`

* ~~descriptive~~ ~~statistics~~ → `df.describe()`

* List down column name → `df.columns`

* ~~shape~~ → `df.shape`

* `df[~].unique() → unique possible value`

* null values → `df.isnull(), sum()`

* duplicate values → `df.duplicated()`
record

↓ drop

`df.drop_duplicates()`

* `sns. heatmap(df.corr())`

↓

Correlation

`annot = True`

* `df[''].value_counts()`

↓ plot

> `df[''].value_counts().plot(kind='bar')`

* pairplot

for univariate, bivariate analysis

> `sns.pairplot(df)`

To see how 1 feature is related to other and that feature with other feature.

* `sns.catplot(x=' ', y=' ', data=df, kind='scatter')`

* Data checks

waste
weights
for
each

* ① check missing values.

② check duplicates.

③ check data type.

④ check the number of unique values

⑤ check statistics of each column.

⑥ check various categories present in different categorical columns.

* missing values. → replace, drop
based on amount of data.
(Feature engineering)

* ~~drop~~ duplicate value → should be removed
~~as~~ as ml need different kind of data.

df.nunique() → Each column unique values

* separate numerical & categorical feature

Categor = [feat for feat in df.columns
if df[feat].dtype == 'object']

numerical → 'int64'

fig, axis = plt.subplots(1, 2, figsize=(-, -))
row, column

> plt.subplot(1, 2, 1) + 1st row 1st column
1st figure

> sns.histplot(data=df, x='Coverage', bins=50,
kde=True, color='g')
hue = 'gender'

based on
gender

~~sns. heatmap~~
 df[gender == 'male']
 df[gender == 'female']
 n = 'average', palette = True, hue = 'lunch')

df = pd.read('---', header = 1)

↳ making 1st row as header.

Forest Fire

① Load dataset

② check df.info()

→ data types not correct
to find it no string
should be there

③ remove string and nan values.

→ df[df.isnull().any(axis = 1)]

Creating region column

df.loc[122, 'Region'] = 2

df[['region']] = df[['Region']].astype(int)

~~Handling~~.

df = df.dropna().reset_index(drop=True)

drop(122), -11-

- ④ 2. removing extra space present
in column name

→ df.columns.str.strip()

- ⑤ change datatype to corresponding
datatype.

⑥ df.describe()

- ⑦ removing extra

→ storing cleaned data in new
file.

EDA → encoding: 'df['class']' in np.where(df['class'])
= 'not first, 0, 1'

density plot of all features

> plt.style.use('seaborn')
df.hist(bins=30, figsize=(10, 15), ec='b')
plt.show()

Flight

Date

12/6/2023

→ need to convert
into separate date
month & year
columns

object

object

$df['date'] = df['date.of.year'].str.split('/').str[0]$

0th index

month

[1]

year

[2]

→ convert into numerical

↓ typecast

$df['Date'] = df['Date'].astype(int)$

→ drop

$df.drop(['-'], axis=1, inplace=True)$

column

arrival time

21:35 10 march



$df['arrival.hour'] = df['-'].str.split(':').str[0]$

$df['arrival.minute'] = df['-'].str.split(':').str[-1][1:-1]$

object

10 march

$df['arr_min'] = df['arrival.minute'].str.split(' - ').str[0]$

integer



$df['-'] = df['-'].astype(int) \rightarrow drop$

10:32

* Similarly, "dep-time" feature

* Total stops - Categorical

> $df[{}^{\circ} \rightarrow {}^{\prime}]$.unique

Convert into
ordinal label

> $df[{}^{\circ} \rightarrow {}^{\prime}] = df[{}^{\circ} \rightarrow {}^{\prime}]$.map({
 'non-stop': 0,
 'one stop': 1, ...})

np.nan: 1

* Duration

"2h 50m"

Q18

> $df[{}^{\circ} \text{hour}^{\prime}] = df[{}^{\circ} \text{duration}^{\prime}]$.str.split(' - ')
 .str[0].
 .str.split('h')
 .str[0]

* Same with

min

1 | $\xleftarrow{\cdot \text{str}[1]}$
1 | $\xrightarrow{\cdot \text{str}[0]}$

* $df[{}^{\circ} \text{at-}f \rightarrow {}^{\circ} \text{ - } {}^{\prime}] = 68$

* $df[{}^{\circ} \text{loc} \rightarrow {}^{\circ}] \Rightarrow$ only needed count
 waste

Fixing non value

* $df[{}^{\circ} \rightarrow {}^{\prime}] = df[{}^{\circ} \rightarrow {}^{\prime}]$.replace(np.nan, 1)

(
 .replace= → , value= →)
 [- , -]

Categorical values

↳ OneHotEncoder encoding

> from sklearn.preprocessing import
OneHotEncoder

> encoder = OneHotEncoder()

> encoded_df = pd.DataFrame(
encoder.fit_transform(df[[" ", " "]]).toarray(),
columns=encoder.get_feature_names_out()))

>
final_df = pd.concat([encoded_df,

axis=1)

~~Playstore~~

Playstore

(df[[" "]].str.isnumeric(); sum())

123 ✓
3 Max

> df[~

not

> df = df.drop(df.index[-1])

record

for spec char in char_to_replace
 for col in col_to_clean
 $\text{df}[\text{cal}] = \text{df}[\text{cal}].\text{str.replace}(\text{spec}, '')$

> January 1, 2018 \rightarrow 2018-01-01

~~$\text{df}[\text{cal}] = \text{pd.to_datetime(df[cal])}$~~

* 2018-01-01



$\text{df['day']} = \text{df['last update']}.dt.day$

• month

• year

> $\text{df}[\text{cal}].value_counts()$

in numerical values

$\text{df}[\text{cal}].value_counts(\text{normalize=True}) * 100$

in percentage

subplot title

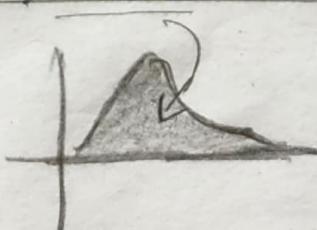
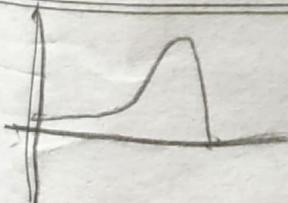
plt: ~~subplots~~ (\rightarrow)

Numerical

Kernel density estimator

> sns.kdeplot($x = \text{df}[\text{cal}]$) shape=Toone

shade=Toone, color="#-")



Categorical

→ sns.countplot(x=df[—], palette='Set2')

pie

> df[—].value_counts().plot.pie(y=df[—],
figsize=(15,6))

Top 10 app Categories

→ pd.DataFrame(df['Category'].value_counts())

> c.rename(columns={'Category': 'Count'}, inplace=True)



bar plot

> sns.barplot(data=df[:10], x=df.index[:10],
y="Count")

by target no of installation Category

- > `dfw = pd.DataFrame(df.groupby('category')[['install']].sum())`
- > `dfw.sort_values(by='installs', ascending=False, inplace=True)`
- > `dfw.reset_index(inplace=True)`

```
sns.barplot( data = dow,  
            x = dow[ "category" ][ :10 ],  
            y = dow[ "installs" ][ :10 ] )
```

<u>Index</u>	<u>Value</u>
0	"Hi")
1	"bye")
2	"20"

* ~~Ch~~ Revision

- ⑥ > from operator. import concat
from functools import reduce
> flat list = reduce (concat, nested-list)

$$\textcircled{2} \quad dt = \{e^A : 20 - \}$$

unpacking operator

`d3 = {** d1, ** d2}`

\downarrow

`{'A': 20, 'X1': 100, ... }`