

Internship

→ dated 26.5.25

→ Data Preprocessing

Import necessary library

a) ~~is~~ import pandas as pd

import numpy as np

import seaborn as sns

import matplotlib.pyplot as plt

b) Reading the dataset

df = pd.read_csv("filename.csv")

head (printing first 5 rows)

df.head()

tail (last 5 rows)

df.tail()

c) Sanity Check of data

df.shape

df.info()

df.isnull().sum()

If there is a missing values more than 50% in a column, we need to do mean, median etc.,... instead of deleting the row or column. (But mostly we do this even if we have less than 1% of Null values in that column.)

```
df.isnull().sum()/df.shape[0]*100
```

#finding duplicates

```
df.duplicated().sum()
```

#identify garbage values.

for i in df.select-dtypes(include="object").columns:

```
print(df[i].value_counts())
```

```
print(" " * 10)
```

Since the garbage values like special characters is exist as string. so object type

It gives us the count of each unique value in that column. Even if it special charact it count that and give in output.

Exploratory data analysis

```
df.describe().T
```

→ get the description in numbers.

```
df.describe(include="object")
```

→ to get the description in string not in numbers

histogram to understand the distribution

we can understand how the data distrib

impact warnings
 warnings: filterwarnings("ignore") } don't give warnings.
 for i in df.select_dtypes(include="number").columns:

sns.histplot(data=df, x=i)
 plt.show()

Boxplot for outliers

for i in select_dtypes(include="number").columns:
 sns.boxplot(data=df, x=i)
 plt.show()

Note:

⇒ We use Histogram to understand distribution of data.

⇒ We use Boxplot to understand outliers of data.

Now we find any bivariate relationship

we can find it using scatterplot. (input to model) independent variable

for i in ['Year', 'Life expectancy']:

sns.scatterplot(data=df, x=i, y='Life expectancy')
 (plt.show())

Target variable

result we want

Note:

- ⇒ Positive Relationship → dots go up together.
- ⇒ Negative Relationship → dots go down together.
- ⇒ No Relationships → dots are scattered randomly.

Correlation with heatmap to interpret the relation and multicollinearity.

```
s = df.select_dtypes(include="number").corr()
```

```
plt.figureuse (figsize=(15,15))
```

```
sns.heatmap(s, annot=True)
```

Missing treatment → We don't do this for target variable.

Mean or Median

```
for i in ["BMI", ...]:
```

```
    df[i].fillna(df[i].median(), inplace=True)
```

```
df[i].fillna(df[i].mean(), inplace=True)
```

for i in [...]: categorical column
 df[i].fillna(df[i].mode()[0],
 (inplace=True)) that contains
 categorical data

KNNImputer

from sklearn.impute import KNNImputer

imputer = KNNImputer(n_neighbors=3)

for i in df.select_dtypes(include='number').columns:
 # looking at 3 closest rows that have similar values
 # uses their average to fill this column.

df[i] = imputer.fit_transform(df[i])

Outlier Treatment

We do this only for continuous numerical data column

=> Not for target variable, categorical and discrete variable

def whisker(col):
 # counted (x) measured (✓)

q1, q3 = np.percentile(col, [25, 75])

iqr = q3 - q1

lw = q1 - 1.5 * iqr

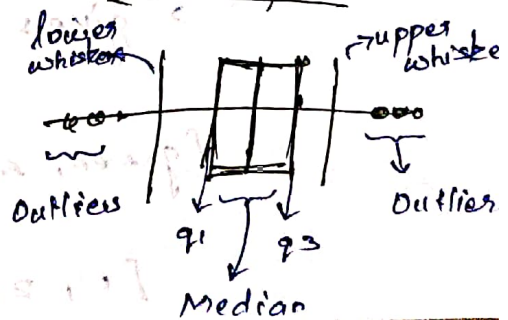
uw = q3 + 1.5 * iqr

Return lw, uw

lw -> Lower whisker

uw -> Upper whisker

Box Plot



for i in []:

lw, uw = whisker(df[i])

df[i] = np.where(df[i] < lw, lw, df[i])

df[i] = np.where(df[i] > uw, uw, df[i])

Duplicate & garbage value treatment

df.drop_duplicates()

For garbage value treatment

→ we have multiple ways (survey web)

Easy way is replacing it with median or

mode of that column.

Encoding of data (for categorical (text) data into numbers)

One hot encoding → create separate column for each categorical data [order does matter]

pd.get_dummies(data = df, columns = ["Country", "Status"], drop_first = True)

Label encoding

→ gives each category one number like [order matter]

[green, red, green]

↓

[1, 2, 1]

(like grade) → ordinal data

Normalization (Only for number type column)

Normalization

- Output Range: 0 to 1
- Formula: $\frac{x - \min}{\max - \min}$
- When to use
 - ↳ When you need fixed range.

x = current value from that column
min → min. value of that column
max → max. "]

Standardization

- Mean → 0, Std. Deviation → 1
- Formula → $\frac{x - \text{mean}}{\text{Std. Deviation}}$
- When to use
 - ↳ When data has outliers