

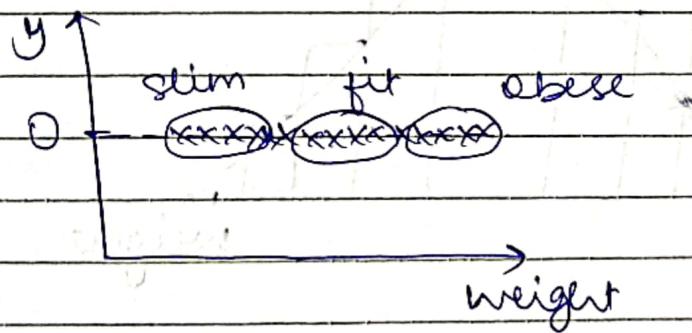
## # univariate, Bivariate & Multivariate analysis

### ↳ Univariate

- we only take one feature in account

#data

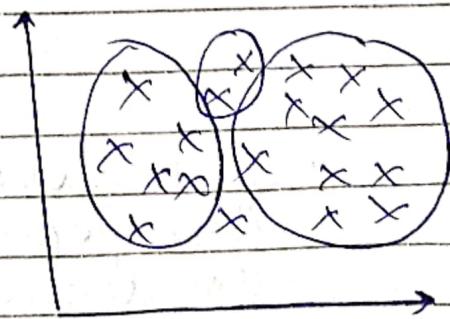
height	weight	o/p
180	90	Obese
160	50	slim
170	78	fit
190	90	fit
195	80	slim



- Overlapping of data pts.

### ↳ Bivariate analysis

#### ↳ weight & ht.



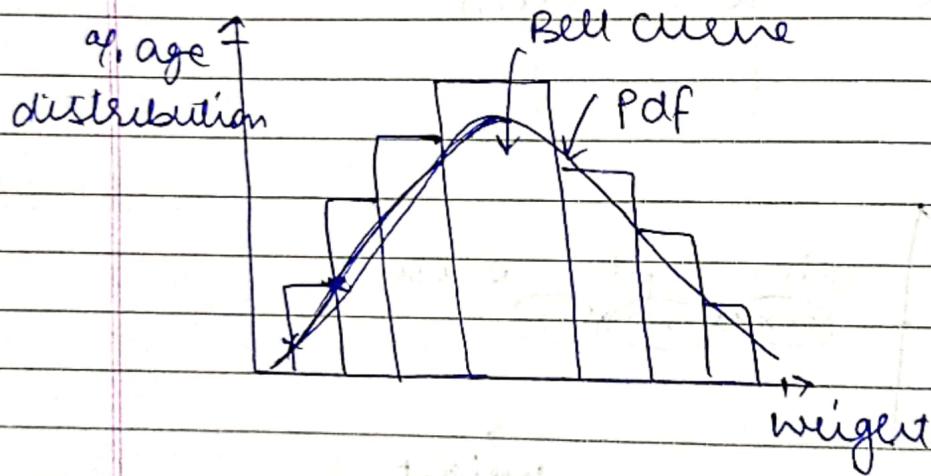
## ↳ Multivariate Analysis

- ↳ when we have multiple features
- ↳ we use pair plot or scatter plots.
- ↳ to plot various graphs

pair plots } Age vs Age      Age vs wt      Age vs ht  
               weight vs Age      weight vs wt      weight vs ht  
               wt vs Age      wt vs wt      wt vs ht

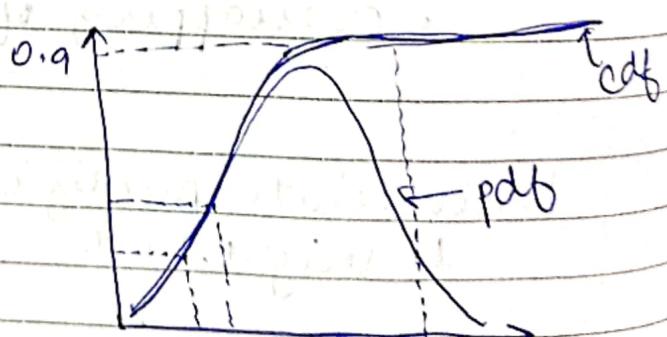
## # Histogram

### # Probability Density F<sup>n</sup>.



### # CDF

cumulative density f<sup>n</sup>



• It shows that 90% population is having weight less than 130kg and 10% has greater weight than 130kg.

## # Types of Encoding

### 1. Nominal encoding

- ↳ categories do not have any order.

- ↳ One hot encoding

- ↳ One hot encoding with many categories

- ↳ mean encoding

### 2. Ordinal encoding

- ↳ categories can be rearranged on the basis of ranks.

- ↳ Label encoding

- ↳ Target guided ordinal encoding

value ↑ → label encoding ↑

ex:

BE → 1

MAS → 3

PHD → 4 → highest imp. high no.

STAT → 2

## # Handle missing values in categorical var.

### 1. Delete rows

2. Replace with the most frequent values  
 ↳ problem: imbalanced dataset within the category.

3. Apply classifier algom to predict.

4. ~~then~~ apply unsupervised ML

↳ group them into categories

- # Handle Categorical Features; Many categories (Count/Frequency Encoding)
  - If we have categorical variables containing many multiple labels or high cardinality, then by using One hot encoding, we will expand the feature space dramatically.
- Approach; replace each label of the categorical variable by the count, this is the amount of times each label appears in a dataset. Or the frequency, this is the %age of observations within that category.

#

```
mapped = df['X2'].value_counts().to_dict()
```

```
{'aa': 1659,  
'ae': 496,  
'ai': 415  
} : {  
    'frequency'  
    'category'
```

```
df['X2'] = df['X2'].map(mapped)
```

↑ maps all the ~~the~~ categories with their frequency

\* Advantages:

- Simple to implement
- Does not increase the feature dimensionality

## L. Disadvantages:

- If some labels have the same count, then they will be replaced with the same count and they will lose some valuable info.
- Odds somewhat arbitrary no., weights to the different labels, they may not be related to their predictive power.

## # ~~Effects of missing Data on classification~~

## # Missing values

### # Different types of missing values :

#### 1. Missing Completely At Random (MCAR)

- A variable is missing completely at random if the probability of missing is the same for all the observations.
- When data is MCAR, there is absolutely no relationship between the missing data.

- 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.

## 2. Missing Data Not at Random (MNAR)

(Systematic missing values)

- There is some relationship between the data missing and any other values, observed or missing within the dataset.

## 3. Missing at random (MAR)

- It means that the reason for missing values can be explained by variable on which you have complete information as there is some relationship between the missing data and other values/ data.
- In this case, the data is not missing for all the observations. It is missing only within sub-samples of the data and there is some pattern in the missing values.

### # Techniques for handling missing value

1. Mean / Median / Mode replacement
2. Random sample Imputation
3. Capturing nan values with a new feature
4. End of distribution imputation
5. Arbitrary imputation
6. Frequent categories Imputation

## 1. Mean / Median / Mode

- Mean / Median imputation has the assumption that the data is MCAR.
- We solve this by replacing the NaN with the most frequency occurrence of the variables.

#

df.isnull().mean()

↳ mean of null values.

```
def impute_mean(df, variable, median):
    df[variable + "_median"] = df[variable].fillna(median)
```

median = df['Age'].median()

print(df['Age'].std())

print(df['Age\_median'].std())

↳ 14.5264

13.0196

] There is no drastic difference.

# we can also observe the diff. by graph

fig = plt.figure()

ax = fig.add\_subplot(111)

df['Age'].plot(kind='kde', ax=ax)

df['Age\_median'].plot(kind='kde', ax=ax,

color='red')

lines, labels = ax.get\_legend\_handles\_labels()

ax.legend(lines, labels, loc='best')

↳ Advantages :

- Easy to implement (robust to outliers)
- Faster way to obtain the complete df

↳ Disadvantages :

- Change or distortion in the original variance.
- Impacts correlation

2. Random sample imputation

↳ It consists of taking random observation from the dataset and we use this observation ~~from the dataset~~ and ~~and are repeated~~ to replace the nan values.

↳ It assumes that the data that are missing are completely at random (MCAR)

#

```
def impute_nan(df, variable):
    median = df[variable].median()
    df[variable + '_median'] =
        df[variable].fillna(median)
    df[variable + '_random'] = df[variable]
    random_sample =
        df[variable].dropna().sample(
            n=df[variable].isnull().sum(),
            random_state=0)
    random_sample.index =
        df[df[variable].isnull()].index
```

`df['cc[ df['variable'].isnull(), variable + '_random' ] = random.sample`

# create plot to see the change in distribution

#### b) Advantages :

- Easy to implement.
- <sup>less</sup> distortion in variance.

#### b) Disadvantages :

- <sup>in</sup> Every situation randomness won't work.

#### 3. Capturing NAN values with a new feature

# It works well ~~to~~ if the data are not missing completely at random (MNAR).

#

`df['Age_NAN'] = np.where(df['Age'].isnull(), 10)`

# now we can handle missing values

#### b) Advantages :

- Easy to implement.
- Captures the importance of missing values.

with mean, median,  
random  
imputation

#### b) Disadvantages :

- Creates additional features.

#### 4. End of distribution Imputation

↳ Here we take the data points from the end of distribution and replace the nan values.

tail of distribution.

↳ Here missing values are not at random (MNAR)

#

`df['Age'].hist(bins=50)`

`extreme=df['Age'].mean() + 3*df['Age'].std()`

↳ mean after 3<sup>rd</sup> s.d.

`sns.boxplot('Age', data=df)`

`median = df['Age'].median()`

`def impute_nan(df, variable, median, extreme):  
 df[variable + '_end_dist'] = df[variable].`

`fillna(extreme)`

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

`impute_nan(df, 'Age', median, extreme)`

↳ Advantages:

- Easy to implement.
- captures the importance of missingness if there is one.

↳

### 4. Disadvantages:

- Distorts the original distribution of the variable.
- If missingness is not important, it may mask the predictive power of the original variable by distorting its distribution.
- If the no. of NA is big, it will mask true outliers in the distribution.
- If the no. 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

- It consists of replacing the nan values by an arbitrary value.

```
def impute_nan(df, variable):
```

```
    df[variable + "_zero"] = df[variable].fillna(0)
```

```
    df[variable + "_hundred"] = df[variable].fillna(100)
```

### 6. Advantages:

- Easy to implement.
- Captures the importance of missingness if there is one.

### b) Disadvantages:

- Distorts the original distribution of the variable.
- If missingness is not important, it may mask the predictive power of the original variable by distorting its distribution.
- Need to decide which value to use.

### # Handling categorical missing values

#### 6. Frequent category Imputation

```
def impute_na(df, variable):
    most_freq_cate = df[variable].value_counts()
    index[0]
```

```
df[variable].fillna(most_freq_cate,
                    inplace=True)
```

OR

```
df[variable].fillna(df[variable].mode()[0],
                    inplace=True)
```

### b) Advantages :

- Easy to implement.
- Faster way to implement.

### b) Disadvantages:

- Since we are using the most frequent

labels, it may use them in over represented way, if there are many nans.

- It displays the relation of the most frequent label.

## 7. Adding a variable to capture nan

```
df['BsmtQual_var'] = np.where(df['BsmtQual']
                               .isnull(), 1, 0)
```

# adds new col BsmtQual\_var with val 0

# where there is nan value.

Then we can impute the nan with mode().

```
df['BsmtQual'].fillna(df['BsmtQual'].mode()[0], inplace=True)
```

## Suppose if you have more frequent categories, we just replace NAN with a new category

```
def imputenan(df, variable):
    df[variable + "newvar"] = np.where(
        df[variable].isnull(), "Missing",
        df[variable])
```

## # Handling categorical Features

## 1. One hot encoding

```
pd.get_dummies(df['Sex'],  
               drop-first=True)
```

- If ~~more~~ many categories are there then we can take top 10 or 15 most frequent categories and one hot encode them.

```
list_10 = [df['x1'].value_counts().  
           sort_values(ascending=False).  
           index]
```

for category in list\_10:

```
df[category] = np.where(df['x1'] ==  
                       category, 1, 0)
```

## 2. Ordinal number encoding

```
dictionary = { 'Monday' : 1,  
              'Tuesday' : 2,  
              'Wednesday' : 3,  
              :  
              'Sunday' : 7}
```

$df['week\_ordinal'] = df['weekday'].map(dictionary)$

### 3. Count of Frequency Encoding

- we map the categories with their frequency.

~~map~~  
~~country = df['Country'].value\_counts().to\_dict()~~

~~df['Country'] = df['Country'].map(  
 country\_map)~~

#### Advantages:

- Easy to implement.
- Not increasing feature space.

#### Disadvantages:

- It will ~~not~~ provide same weight if the frequencies of categories are same.

### 4. Target Guided Ordinal Encoding

- Ordering the labels according to the target.
- replace the labels by the joint probability of being 1 or 0.

~~labels~~  
~~ordinal = df.groupby(['~~Sex~~'])['Survived'].mean().reset\_index()~~

ordinal\_dict = {k: i for i, k in enumerate(ordinal\_labels)}

ordinal\_dict

{'T': 0, 'M': 1, 'A': 2, ... 'D': 8}

↳ max value

`df['cabin_ordinal'] = df['cabin'].map(`  
   `Ordinal-dice)`

### 5. Mean Encoding

- We replace the categories by replacing it with its mean.

`mean_ordinal = df.groupby(df['cabin'])`  
   `['Survived'].mean().to_dict()`

`mean_ordinal`

`{'A': 0.466, }`

`'B': 0.744, }`

`:`

`'M': 0.299`

`'T': 0.0}`

`df['mean_encode'] = df['cabin'].map(`  
   `mean_ordinal)`

### 6 Advantages :

- Captures the information within the labels. Therefore rendering more predictive features.
- Creates a monotonic relationship between the variable and the target.

### 6 Disadvantages :

- Prone to overfitting.

## 6. Probability Ratio Encoding

Probability Ratio =  $\frac{P(\text{success})}{P(\text{failure})}$

`prob_df = df.groupby(['cabin'])['survived'].mean()`

`prob_df = pd.DataFrame(prob_df)`

~~prob\_df = prob\_df.dropna()~~

`prob_df['Died'] = 1 - prob_df['survived']`

`prob_df['Prob ratio'] = prob_df['survived'] / prob_df['Died']`

`prob_map = prob_df['Prob ratio'].to_dict()`

`df['cabin-encoded'] = df['cabin'].map(prob_map)`

### # Types of Transformation

#### 1. # Normalization And Standardization

##### ↳ Standardization:

- we try to bring all the variables and features to a similar scale.

- standardization means centering the variable at zero

$$Z = \frac{(x - \bar{x})}{\sigma_x} \quad | \quad \frac{x - x_{\text{mean}}}{\text{std}}$$

#

from sklearn.preprocessing import  
StandardScaler

scaler = StandardScaler()

df\_scaled = scaler.fit\_transform(df)

pd.DataFrame(df\_scaled)

## 2. b) Min Max scaling

- widely used in D.L (CNN)

- ~~It~~ Scales the values bet 0 & 1.

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

from sklearn.preprocessing import  
MinMaxScaler

min\_max = MinMaxScaler()

df\_minmax = pd.DataFrame(min\_max.

fit\_transform(df), columns=

df.columns)

df\_minmax

## 3. b) RobustScaler

- It is used to scale the features w median and quantiles

- Scaling ~~using~~ 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 bet. the 75<sup>th</sup> quantile and 25<sup>th</sup> quantile.

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

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

#

from sklearn.preprocessing import  
RobustScaler

scaler = RobustScaler()

df\_robust\_scaled = pd.DataFrame(

scaler.fit\_transform(df),

columns=df.columns)

df\_robust\_scaled.head()

#### 4. Gaussian Transformation

• some ML algos like linear and logistic assume that the features are normally distributed.

- If we want to check gaussian or normal distributed we use Q-Q plot.

```
import scipy.stats as stat
import pylab
```

```
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()
```

- In the plot (1,2,2) if all the pts lie in a straight line then we can say our data is normally distributed.

b) Logarithmic transformation  
 $df['Age\_log'] = np.log(df['Age'])$   
 $plot\_data(df, 'Age')$

- If there is error in log fn then use

$$df['Age\_log'] = np.log(df['Age'] + 1)$$

b) Reciprocal transformation  
 $df['Age\_Reciprocal'] = 1/df['Age']$   
 $plot\_data(df, 'Age\_Reciprocal')$

### b) Square Root transformation

`df['Age_sqrt'] = df['Age'] ** (1/2)`  
`plot_data(df, 'Age_sqrt')`

### b) Exponential Transformation

`df['Age_exp'] = df['Age'] ** (1/1.2)`  
`plot_data(df, 'Age_exp')`

### b) Box Cox Transformation

~~`df['Age_boxcox'] = stats.boxcox(df['Age'])`~~

`df['Age_boxcox'], params = stats.boxcox(df['Age'])`  
`plot_data(df, 'Age_boxcox')`  
`print(params)`

## # Handling Imbalanced data

`!pip install imbalanced-learn`

### 1. Under Sampling

- reduce the points of the maximum labels.

- we should not use under sampling bcz. of less data.

- we should only use it when we have very low data set.

```
# from collections import Counter
from imblearn.under_sampling import
    NearMiss
```

nm = NearMiss(0.8)

x\_train\_nm, y\_train\_nm = nm.fit\_sample(x\_train, y\_train)

ytr\_nm

```
print("Class before fit {{".format(
    Counter(y_train)))}
```

```
print("Class after before fit {{".format(
    Counter(y_train_nm)))})
```

↳ {{0: 199017, 1: 3473}}  
 ↳ {{0: 433, 1: 347}}

formula:  $0.8 \times 433 = 347$

## 2. Over Sampling

#

```
from imblearn.over_sampling import
    RandomOverSampler
```

os = RandomOverSampler(0.5)

x\_train\_os, y\_train\_os = os.fit\_sample(
 x\_train, y\_train)

```
print("No class before fit {{".format(
    Counter(y_train)))})
```

```
print("Class after fit {{".format(
    Counter(y_train_os)))})
```

↳ { 0: 199017, 1: 347 }  
 ↳ { 0: 199017, 1: 98508 }

### 3. SMOTE Tomek

from imblearn.combine import SMOTE Tomek  
 OS = SMOTE Tomek(0.5)

X\_train\_st, y\_train\_st = OS.fit\_sample  
 (X\_train, y\_train)

print(\_\_\_\_)

print(\_\_\_\_)

↳ { 0: 199017, 1: 347 }  
 ↳ { 0: 199017, 1: 148640 }

### 4. Ensemble Techniques

from imblearn.ensemble import  
 EasyEnsembleClassifier

easy = EasyEnsembleClassifier()

easy.fit(X\_train, y\_train)

## # Outliers and Its Impacts

Q. Which ML algos are sensitive from outliers?

↳ Naive Bayes classifier — Not sensitive

↳ SVM — Not sensitive

✓ ↳ Linear Regression — Sensitive

✓ ↳ Logistic Regression — Sensitive

↳ Decision Trees — Not sensitive

- ↳ Ensemble (RF, XGBoost, GB) — Not sensitive
- ↳ KNN — <sup>Not</sup> sensitive
- ✓ ↳ KMeans — sensitive
- ✓ ↳ Hierarchical — sensitive
- ✓ ↳ PCA — sensitive
- ✓ ↳ Neural Networks — sensitive

### 1. Gaussian distributed Data

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

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

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

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

$$\text{upper\_bridge} = \text{df['Age'].quantile}(0.75) + (1.5 * \text{IQR})$$

we will consider the values after  
<sup>boundary</sup>~~lower\_bridge~~ & <sup>boundary</sup>~~upper\_bridge~~ as outliers.

## Q. Skewed data

$IQR = df['fare'].quantile(0.75) - df['fare'].quantile(0.25)$

lower-bridge : \_\_\_\_\_

upper-bridge : \_\_\_\_\_

## # FEATURE ENGINEERING STEP BY STEP

### I. EDA (Exploratory Data Analysis)

- ↳ Numerical features, Categorical features

- ↳ Histogram, pdf, etc.

- ↳ Missing values (visualise)

- ↳ Outliers (Box plot)

- ↳ Cleaning or not

### II. Handling the missing values

### III. Handling Imbalanced data

### IV. Treating the outliers

### V. scaling <sup>down</sup> the data

### VI. converting categorical features to numerical features

# Feature selection : selecting the best features

- ① Correlation
- ② K-Neighbour
- ③ Chi-square
- ④ Genetic algm
- ⑤ feature importance (extra tree classifier)