

late
2017

Machine Learning

Predictive Modelling Cases

- Movie prediction → Sales prediction
- Viewing traffic of website → predicting Stock Market

Predictive Modelling

- ↳ making use of past data
- ↳ predict the future using past data

- * Recommending movie → Horror Movie ↑↑↑ likes much
P.M. task.

- * Sales deduction →
 - Analysing past Data
 - No future outcome or forecastNot a Predictive Modelling task

Diagnostic Analysis

- * traffic website : Not a P.M.
- * Stock Market : P.M. task

types of Predictive Modelling

House Price Prediction

- Size
- Location
- No. of rooms
- Price

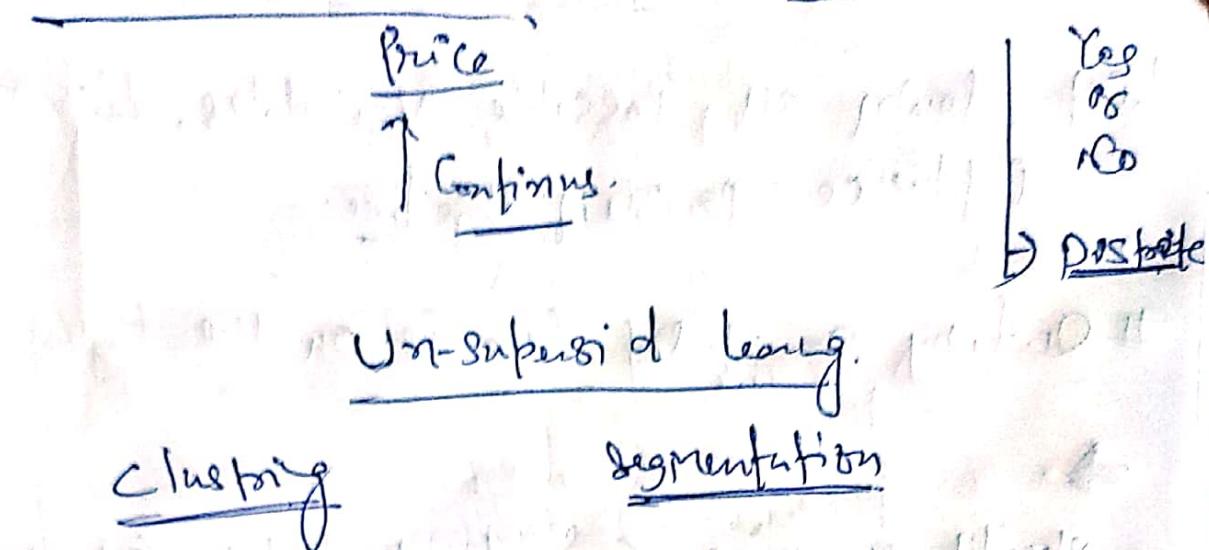
Supervised

Google News Clustering

- No data

Unsupervised

House price prediction



Stages of Predictive Modelling

- Problem Def
- Hypothesis Generation
- Data Exploration / Collection
- Data Exploration and Transformation
- Predictive Modelling
- Model Deployment / Implementation

Problem Statement-

- * Identify the right problem statement, ideally formulate the problem mathematically
 - Bad Problem Statement.

↳ ScreenShop.

Problem Definition

Bad Problem Statement: Want to improve the profitability of credit card customers

Want to increase the APR of credit cards

Want to deploy different APR for different segment of customers

Want to identify the customers segments having lower default rate

Want to have different APR and other benefits for different customers segments (On expected default rate) to maximize profit

00:55

01:29



Hypothesis Generation

List down all possible variables, which might influence problem objective.

Quality of Model depends on the Hypothesis



Should Hypothesis Gen be done before (or) after looking at the data?

Hypothesis gen. to be done before looking at the data.

Why?

→ Let you think of all the factors which might affect the problem without being biased.

→ Stop time waste in analyzing all available data

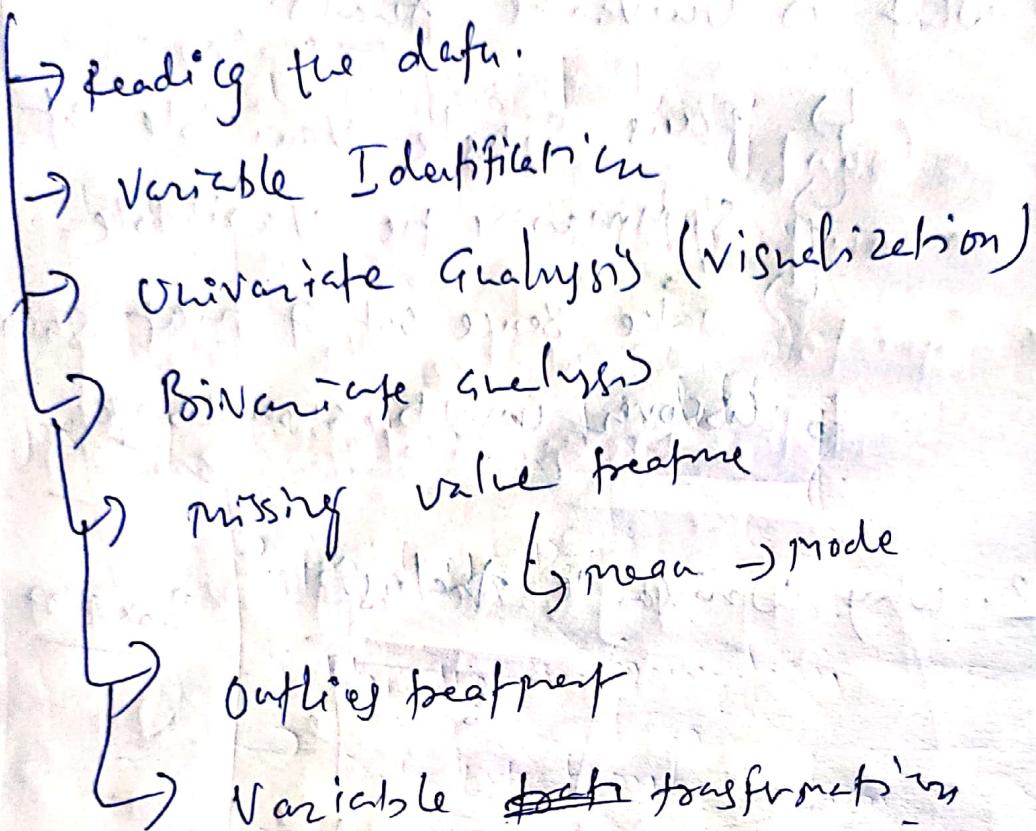
Data Extraction / Collection

Extract / Collect data from diff sources and combine those for exploration and Model building.

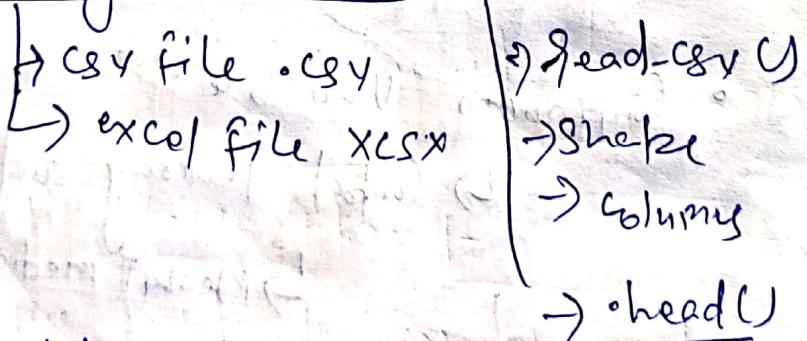


Data Exploration

- Making sense out of data.
- Good Analyst → knows/hws her data



reading the data



Variable Identification

import pandas as pd

```
file = pd.read_csv('example')
```

```
file.shape
```

```
file.head()
```

→ file types



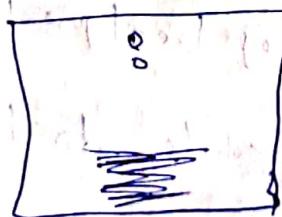
Univariate Analysis

1. What is Univariate Analys?

- ↳ finds single Variable at time
- ↳ Summarize the Variable
- ↳ Make sense out of that summary to discover insights, anomalies etc.

2. Why Univariate Analysis?

Box plot (SS)



Univariate helps to detect anomalies in data

• Continuous Variable

↳ Central tendency and dispersion

↳ mean / median / mode / SD

↳ Distribution of variable

↳ Symmetric / right skewed / left skewed

presence of missing value → knee.

↳ outliers.

method (Continuous Variable Methods)

↳ numerical method

- ↳ mean, med, mode, SD, and missing value

↳ graphical method

- ↳ Distribution of Variable presence of outliers.

df.describe() (pandas Method)

↳ Cont means SD min etc

↳ Histogram → Distribution of Boxplot → outliers Left. Variable

Pd.DataFrame.hist()

Example :-

df.describe() :- Only Continuous Variable.

Plotting Histogram

df['Age'].plot.hist() : Distribution of Variable

Box Plot

df['Age'].plot.box() : outliers

One Variable

↳ Categorical Value:

↳ sex
↳ class
↳ cancer

- Count
- Count %

• Tabular Method Frequency Table

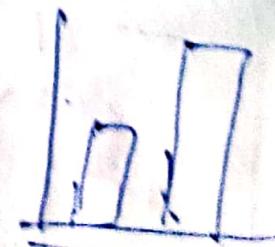
• Graphical Method Bar Plot

• df['sex'].Value_Count()

Male 577 }
Female 319 } Count

Male 0.64 } Count %
Female 0.35 } =

Bar graph



df['Sex'].Value_Count() / len(df['Sex'])

df['Sex'].Value_Count().plot_bar()

(df['Sex'].Value_Count() / len(df['Sex']))
plot_bar()

Bivariate Analysis

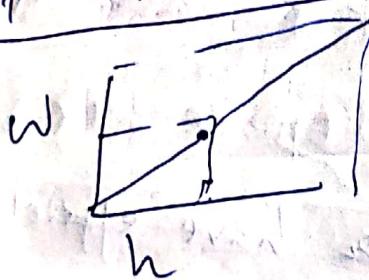
1. What is Bivariate Analysis?

- When two variables are studied together for their empirical relationships
- When you want to see whether two variables are associated with each other.

Ex: $\begin{cases} \uparrow \text{height} \\ \uparrow \text{weight} \end{cases}$

2. Why Bivariate

↳ It helps in prediction



3. Type of Bivariate Analysis

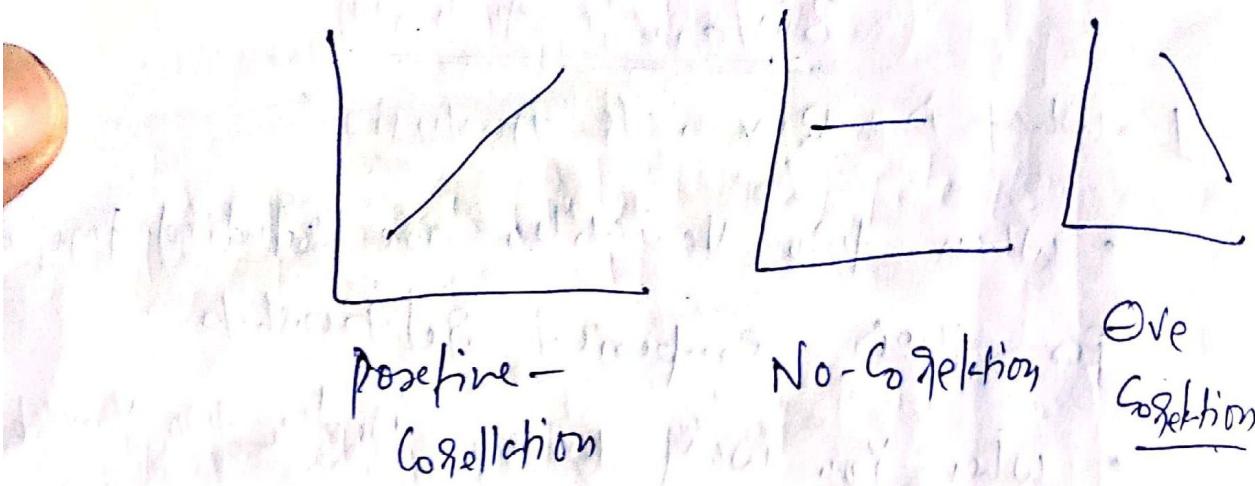
↳ Age - fare price (Scatter Plot)

↳ Gender - Mean Age (Bar Plot)

↳ Gender - survived (Two-way table)

Analysis test

$$r = \frac{\text{Cov}(x, y)}{\sqrt{s_x^2 s_y^2}}$$



* Greater Gender Mean

Analysis for $t = \frac{\bar{Y}_1 - \bar{Y}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

$t < 0.5$

* Does gender have effect on few several stuffs.

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

Chi-Sq test

		survived	0	1	
		Sex	f	m	
		0	87	233	
		1	968	109	

Jupyter Notebook

```
→ import pandas as pd  
→ file = pd.read_csv('titanic.csv')  
→ file.head()  
→ file.dtypes  
# Conf - Continuous bivariate analysis  
# Age (Age) (Sex)  
→ file.plot.scatter('Age', 'Fare')  
→ df.corr() (or) file.corr()  
→ file['Age'].corr(file['Fare'])  
→ # Categorical - Continuous bivariate analysis  
→ file.groupby('sex')['Age'].mean()  
→ file.groupby('sex')['Age'].mean().plot.box  
→ from scipy.stats import ttest_ind  
males = file[file['sex'] == 'male']  
female = file[file['sex'] == 'female']  
ttest_ind(males['Age'], female['Age'], nan_policy
```

Categorical - Categorical: Bivariate Analysis

→ pd. crosstab(file['Sex'], file['survived'])

→ from scipy.stats import chi2_contingency
chi2_contingency(pd.crosstab(file['Sex'],
file['survived']))

Missing Value

→ Reasons of missing value:-

→ Non-response

→ Errors in Data Collection

→ Errors in reading Data.

→ Types of missing value

1. missing completely at random (MCAR)

2. missing at random (MAR)

3. missing not at random (MNAR)

MAR:
→ missing value have no relation to
1. the Variable in which missing value exists
2. other Variable in the dataset

MAR:
missing value have no relation to
missing value in which m.v exist
1. the Variable in which m.v exist
2. missing values have selection to
missing values other than the Variable
1. the Variable other than the Variable
in which missing value exists

MAR:
missing value have selection
1. the Variable in which
missing value exists

Identifying M.V
. isnull()
. describe()

Deal with missing value

1 Imputing. 2 Delecting
imputation
Conti \leftarrow mean
Regression Model
Cato \leftarrow mode
Classification Model

Deletion

- ↳ Row-wise Deletion
- ↳ Column-wise Deletion

Jupyter Notebook

- import pandas as pd
- file = pd.read_csv('data.csv')
- file.shape
- file.describe() → Conf

Pass	Sex	Age	Price	fare
10	10	9	8	7
- file.isnull().sum() → Identifying mpv
- # drop all row
- file.dropna().isnull().sum() → giving sum
- file.dropna(how='all').shape → remove all missing value
- # Drop columns with missing value
- file.dropna(axis=1).shape
- file.dropna(axis=1, how='all').shape

#How to fill Missing Value.

→ file.fillna(0) # Temporary

→ file.fillna(0, inplace=True) # Permanently

→ file['Age'].fillna(0) # Particular Column

→ # mean of column

→ file['Age'].fillna(file['Age'].mean())

Outliers treatment

→ Reasons for Outliers:

1. Data-Entry Errors

2. Measurement Errors

3. Processing Errors

4. Change in underlying population,

→ Types of Outliers

↳ Univariate
↳ Bivariate

→ Identifying Outliers

Graphical methods (Boxplot, Scatter Plot)

Treating outliers

1. Del. obs
2. Outliers and Bounding Box
3. Imputing outliers like miss. Value
4. treat term Separately

Jupyter notebook

Univariate outliers - Part 1

```
→ df['Age'].plot.box() # Univariate outliers
```

```
→ df.plot.scatter('Age', 'Fare') # Bivariate outliers
```

Removing outliers

```
df = df[df['Fare'] < 300]
```

Replacing outliers in age with mean age values

```
df.loc[df['Age'] > 65, 'Age'] = np.mean(df['Age'])
```

f formula for predicting outliers

$$x \approx Q1 - 1.5 * IQR \text{ or } x \approx Q3 + 1.5 * IQR$$

Variables

Transformations

What is Variable transformation?

- Variable transformation is a process by which we replace a variable with some function of that variable.

Example: replacing a variable x

with logarithm.

2. we choose the distribution (or) relationship of a variable with others.

Why?

Variable transformation is used to

change the scale of variable.

1. Change the scale of variable
if 10 Variable are measured in Km
and 1 in miles,

2. transforming Non-linear relationships into linear relationships.

3. convert asymmetric distributions into skewed distributions

Common method of Variable Transformation

1. Logarithm - taking log of the variable

Reduces right skewness of the variable

2. Squared Root -

3. ~~Sq.~~ Cube Root

4. Binning - ~~Histogram~~

Used to convert continuous into categorical

Jupyter Notebook.

```
→ df['Age'].plot.hist()
```

```
→ np.log(df['Age'])
```

```
→ np.sqrt(df['Age']).plot.hist()
```

```
→ df['Age'].plot.hist()
```

```
→ bins=[0, 15, 80]
```

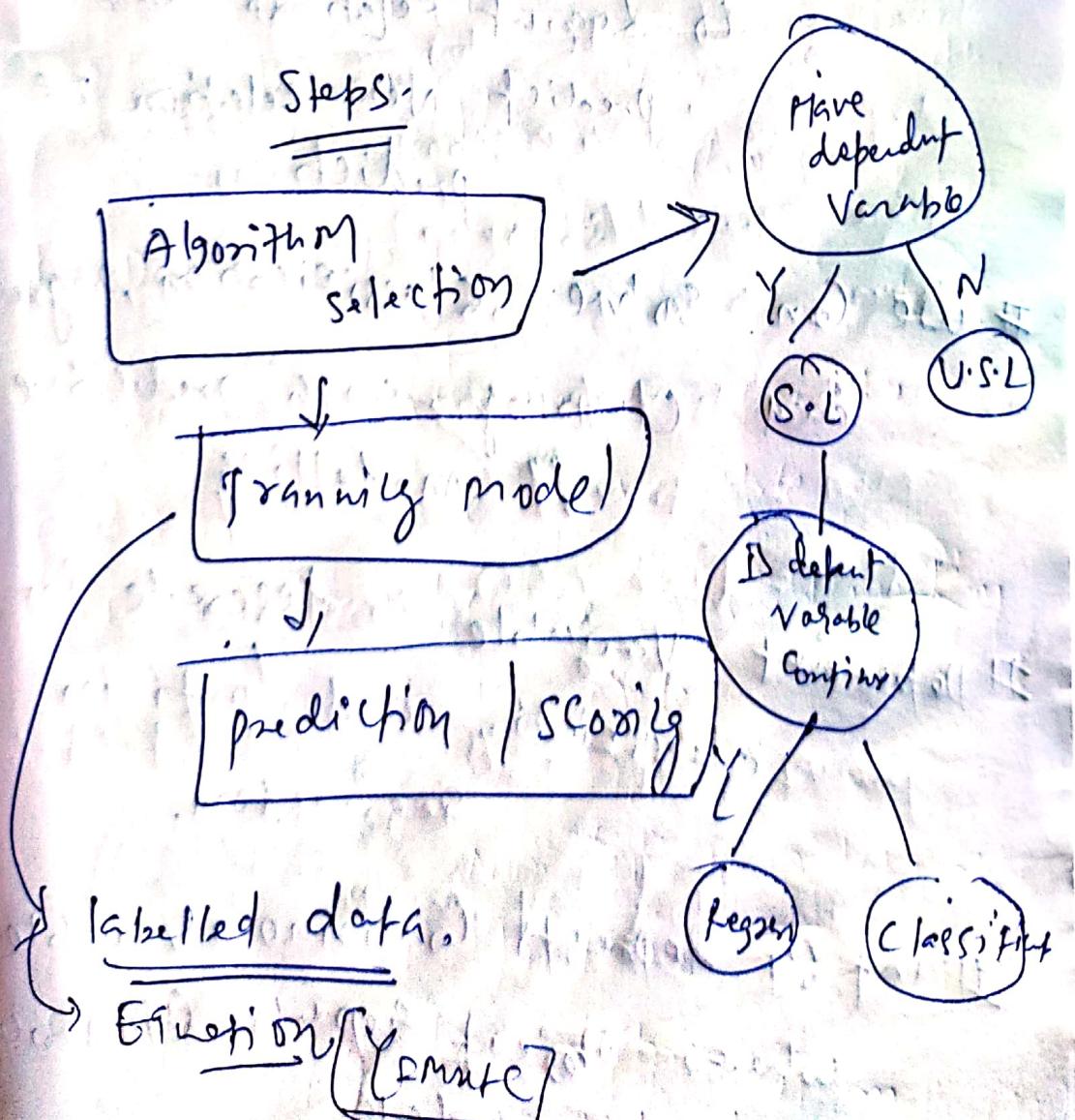
```
→ group=['child', 'Adult']
```

```
→ df['type']=pd.cut(df['Age'], bins=bins, labels=group)
```

`df.head()`
`df['type'].value_counts()`

Basics of Model'

→ Model Building is a process to create a mathematical model for estimating / predicting the future behavior based on past data



Dataset

Train

- past data (Known dependent Variable)

- Used to train model

• future data
(Unknown dependent variable)

- Used to

fitting score

Making Predictions

↳ learnt relations will predict New relation i.e prediction

we can solve Classification, Regression and Clustering with predictive Modelling.

Dependent Variable is discrete in classification and conti in regression.

int by represents Continuous Variable where Object represents Categorical data