

Lecture #2 Data - preprocessing

1.0 Introduction

2.0 Data pre-processing

• 2.1 Data visualization and Exploration

2.2 Data cleaning

- Detect and removal of blanks
- Detect and removal of duplicates
- Detect and removal of non-numeric
- Detect and removal of Outliers

2.3 Feature selection

2.4 Data Smoothing

2.4.1 Moving average

2.4.2 Exponential smoothing

3.0 Data Splitting methods

4.0 Summary

Lecture 2 Feature Engineering / Data-preprocessing

Introduction

Data is used for ML modeling.
Raw data contains useful information and Undesired / Non-informative / Noise as well.

Before using Raw data for ML modeling, it must be preprocessed to:

- make it clean.
- Select the useful features and
- minimize the noise level.

Therefore this Chapter presents 3-parts.
Feature Engineering / Data Preprocessing

(1) Data Cleaning

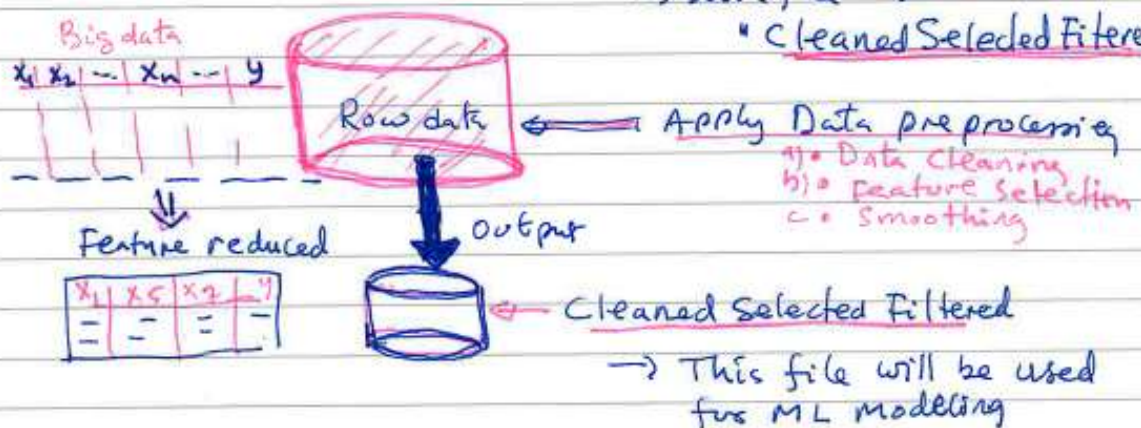
- Remove unrecorded / blanks
- Remove Non-numerics
- Remove duplicates
- Remove outliers
- Save file as "Cleaned.XLSX"

(2) Feature Selection

- Select appropriate / useful features. → Features reduced
- Save file as "Cleaned Selected"

(3) Apply smoothing filter

- Save file as "Cleaned Selected Filtered"



Lecture Data Preprocessing (Feature Engineering)

Before we start data processing,
let us see the difference between
Data and Information.

Data	Information
~ a collection of facts. eg. observations, figures,	~ A result of <u>processing</u> <u>analysing</u> and <u>interpreting</u> of <u>data</u>
~ data on it's own is <u>meaningless</u>	~ when data is analysed and interpreted, it becomes a <u>meaningful</u> information
~ Data <u>does not</u> depend on Information	~ Information <u>depends</u> on data
~ Data is not sufficient for making decision	~ But, you can make decision based on Information

Information is defined as knowledge gained
through study, research, or instruction

Essentially, Information is the result of
analysing and interpreting piece of data

Example 1 A set of data includes temperature readings in a location over several years.

→ Without any additional context, those temperatures have no meaning

→ However, when you analyze and organize that information, you could determine

→ Seasonal temperature pattern

or → even broader climate trends

⇒ Min / Max / average values...

→ To Summarize:

→ Only when the data is organized and compiled in a useful way, it can provide INFORMATION that is beneficial

Hopefully, the above example illustrates about Data and Information.

Example 2 Table 1 shows data X, Y

X	Y
0	0
1	3
2	
3	4
3	4
4	xx7
5	2 56
6	2975
7	8
9	10
11	12
	15
18	20
18	20

→ to analyse the nature of the data
→ visualization allows about the data trend.

← As we can see data contains:

- unrecorded / blanks
- duplicates
- Non-numeric
- outliers
- and noises

→ These makes the quality of data 'poor'

← Blank (unrecorded)

← duplicated

← Object

← Not numeric object

← outlier

Blank unrecorded

← duplicated

"Poor" data quality, results in "poor" information.

Data must therefore be cleaned before analysing and interpreting

→ After cleaning the above dataset, the clean data will be

X	Y
0	0
1	3
3	4
7	8
9	10
11	12
18	20

→ However, if we have big-dataset, the data cleaning process will not be performed manually as shown above.

→ It will be performed by data analysis and processing tools, which are ~~performed~~ by the Computers.
performed

Lesson 1 — Data must be cleaned before analysing and interpreting

Therefore the first part of the lesson is Data Pre-processing / Feature Engineering.

This is the first step of Machine Learning workflow as shown below.

2.0 Data - Pre processing

Before ML modelling data must be pre processed (1) to make it clean
(2) to select the right features

Feature Engineering :

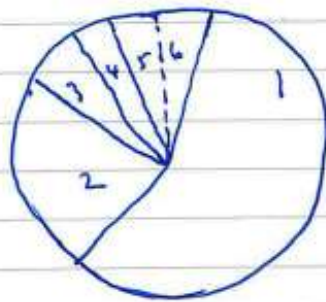
is the process of selecting, manipulating, and transforming raw data into features that can be used in supervised learning

Feature is any measurable input that can be used in a predictive model.

Feature Engineering, in a simple term, is the act of converting observations into desired features using statistical or ML approach.

→ The goal is to improve performance

Data Scientists spend most of their time on data, which is important to make models accurate. Example, Time spent



- 1 - Cleaning and organization
- 2 - collecting data set
- 3 - Mining data
- 4 - Refining algorithm
- 5 - Others
- 6 - Training Sets

The following present the typical data Preprocessing process. Here we will look at

- (1) Data cleaning
- (2) Feature Selection
- (3) Data Smoothing

Step by step guide: of (1) Data cleaning

Step 1 Import data important libraries
Panda and matplotlib

Step 2. Import data
`df = pd.read_excel('Filename.xlsx')`

Step 3 Information about the data .info()

- For modeling data must be numeric, not object
- Data must have the same format = float

in `Print(df.info())`

out:

	Header	Non-null	dtype
0	x ₁	1500 Non-null	float
1	x ₂	1490 Non-null	int64
2	x ₃	1300 Non-null	object
3	y	1495 nonnull	object

Except for x₁, the rest int64 and object should be converted to numeric

Step 4 Convert Object and int64 to numeric

(a) Convert int64 to float

Syntax:

`df['x2'] = df['x2'].astype(float)`

(b) Each objects convert to numeric

- objects are due to

- Blank
- Strings
- String with numerics

Syntax

```
df['Column name'] = pd.to_numeric(df['Column name'],  
                                  errors='coerce')
```

- Coerce will convert the instances (blank / string) to Nan. (Not a Numeric)
- Later the Nan will be removed

Example

```
df['x3'] = pd.to_numeric(df['x3'], errors='coerce')
```

```
df['y'] = pd.to_numeric(df['y'], errors='coerce')
```

- Check if int64 and objects are converted in df.info()

out	Header	Non-Null	dtype
0	x1	1500 Non-Null	float
1	x2	1490 Non-Null	float
2	x3	1300 Non-Null	float
3	x4	1495 Non-Null	float

- As we can see both int64 and object are converted to float

- If you check the dataset →

— all the data ^{int} are float 3.0

→ 3 → 3.0

→ All object / blank → Nan

- Step 5: Display statistical values/analysis
df.describe()

- Step 6: Blank (Nan) data Removal

(6a) df.dropna(axis=0, inplace=True)

Remove row do and return

- (6b) Check if blanks are removed

df.info()

Step 7 Duplicates data must be REMOVED

vs	den	neu	up

(6a) First, check if duplicates exist
— use `duplicated()`

* calculate duplicates

`dups = df.duplicated()`

* Report if exist. True = detected! False = not
`Print(df.dups.any())`

* List all duplicate row
`Print(df[dups])`

Out

	True	vs	den	neu	up
1103		1.05	2.38	0.33	258
1104		1.051	2.37	0.34	2.61

List-duplicate

(6b) How to DELETE row that contain Duplicate data?

— use `drop_duplicates()`

* check size of data before
duplicate removal

`Print(df.shape)`

* Delete duplicate rows

`df.drop_duplicates(inplace=True)`

→ check size after
removal

`Print(df.shape)`

Output

(1111, 4)

→ Before removal

(1103, 4)

→ After removal

Step 8 Outliers removal

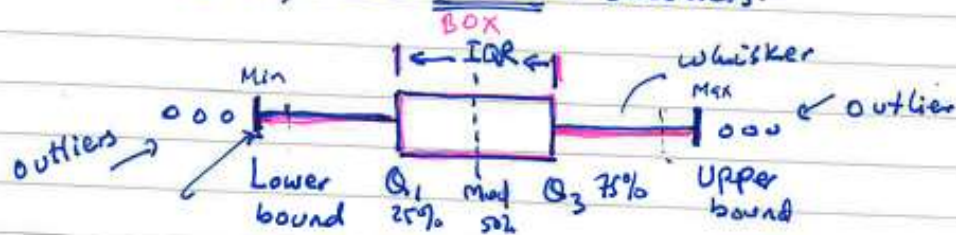


— Outliers are data that differ significantly from the rest of the dataset.

→ These dataset may be plotted as individual points beyond the ~~whisker~~ whiskers on the box-plot

→ Box-plot display variation in samples of the population

1st — We need to identify outliers
2nd, we remove outliers.



$$IQR = Q_3 - Q_1$$

To remove outlier → below lower bound

$$\text{Lower bound} = Q_1 - 1.5 IQR$$

~~upper~~

To remove above upper bound

$$\text{Upper bound} = Q_3 + 1.5 IQR$$

(1) Detection

IQR

$$Q_1 = np.percentile(df, 25)$$

$$Q_3 = np.percentile(df, 75)$$

$$\text{or } Q_1, Q_3 = np.percentile(df, [25, 75])$$

$$IQR = Q_3 - Q_1$$

$$LB = Q_1 - 1.5 \cdot IQR$$

$$UB = Q_3 + 1.5 \cdot IQR$$

Print('old shape: ', df.shape)

using target column.

upper bound

upper = np.where(df['Vp'] >= (Q3 + 1.5 * IQR))

lower bound

lower = np.where(df['Vp'] <= (Q1 - 1.5 * IQR))

(2) Removing outlier → use drop()

df.drop(upper[0], inplace = True)

df.drop(lower[0], inplace = True)

print('New shape : ' + df.shape)

[dataframe.drop(row_index, inplace = True)
↑

* The row from the dataset given
in the row index will be removed

* Inplace = True → used to instruct
Python to make required
change in the original dataset.

Ex. df.drop(list[0], inplace = True)
↑

- row index can be ONLY one
value

- OR List of values

OR

Outlier detection

Outliers = [x for x in df['Vp'] if x < lower or
x > upper]

print('Identified outlier: %d' % len(outliers))

Outlier Removal

df = df[(df['Vp'] > lower) & (df['Vp'] < upper)]

print('Non-outlier: %d' % len(df))

(2.3) Feature selection

After the dataset has been cleaned, the final data-preprocessing step is to select the right features that are useful to a model in order to predict the target variable.

→ Feature Selection is primarily focused on removing non-informative or redundant predictors (input features) from the dataset. So that the reduced number of input variables will be used for modeling.

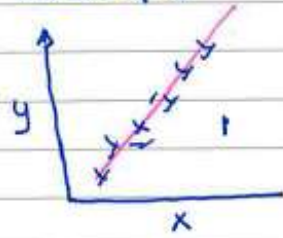
There are many methods. Among others, the most common technique is a correlation coefficient method, which will be used in this course.

Pearson's Correlation Coefficient, $r \Rightarrow$
→ is the most common way of measuring a linear correlation b/w x, y

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$$

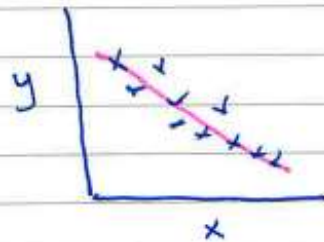
- It's value is between -1 and 1 that measures the STRENGTH and DIRECTION between two variables

Example



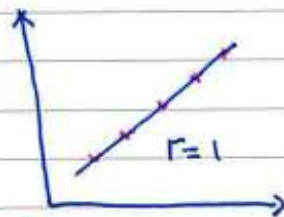
Positive Correlation

→ as x increase, y also increase
 r is between 0 and 1

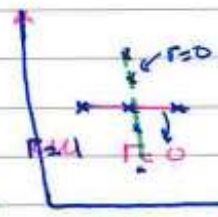


Negative Correlation

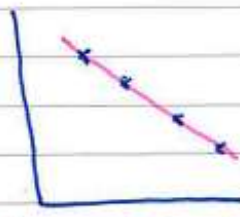
→ as x increase, y decrease
 r is between 0 and -1



Perfect positive
 Correlation

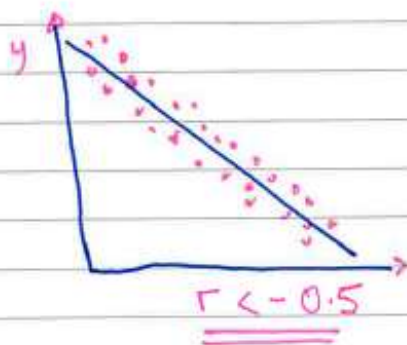


No correlation
 $r = 0$

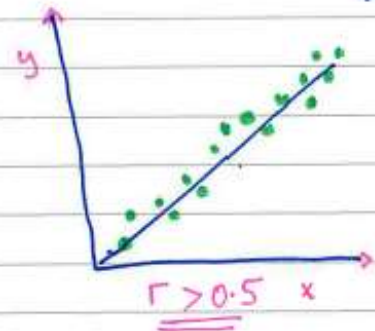


Perfect negative
 correlation

Example when r is greater than 0.5 or less than -0.5, the points are closer to the line of best fit

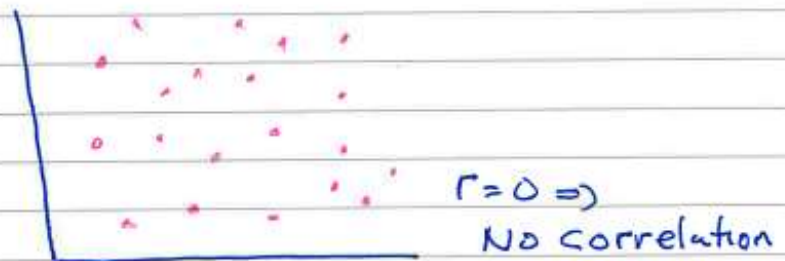


$r < -0.5$



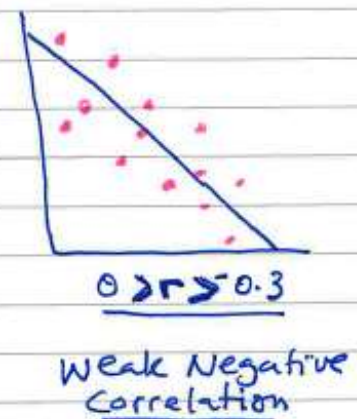
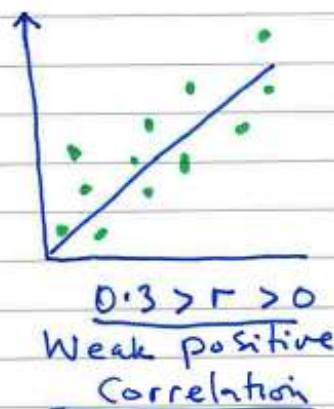
$r > 0.5$

Example when $r=0$, the line of best fit is not helpful in describing the relationship between the variables



Example - Weak correlation

- When r is between 0 and 0.3 or between 0 and -0.3, the points are far from the best fit line



<u>Summary</u>	<u>Pearson Correlation Coeff r</u>	<u>Strength</u>	<u>direction</u>
Greater than 0.5		Strong	Positive
between 0.3 and 0.5		Moderate	Positive
between 0 and 0.3		Weak	Positive
0		None	None
between 0 and -0.3		Weak	Negative
between -0.3 and -0.5		Moderate	Negative
less than -0.5		Strong	Negative

Example Feature Selection with Sklearn

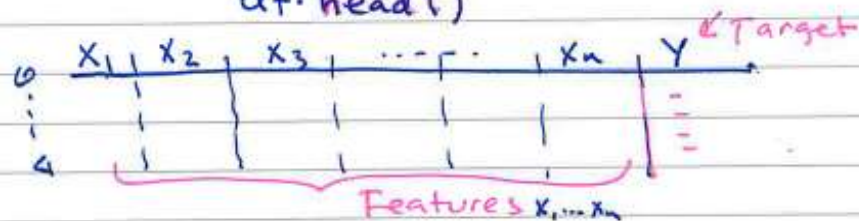
Step 1 Import important libraries.

```
Import Pandas as pd  
Import Matplotlib.pyplot as plt  
Import Seaborn as sns
```

Step 2. Import data file and display(.head())

```
df = pd.read_excel('cleaned.xlsx')
```

```
df.head()
```



Step 3 Calculate correlation between the Target and Feature variables

```
# Using Pearson Correlation, generate HEAT-MAP
```

```
plt.figure(figsize=(12, 10))
```

```
Cor = df.corr()
```

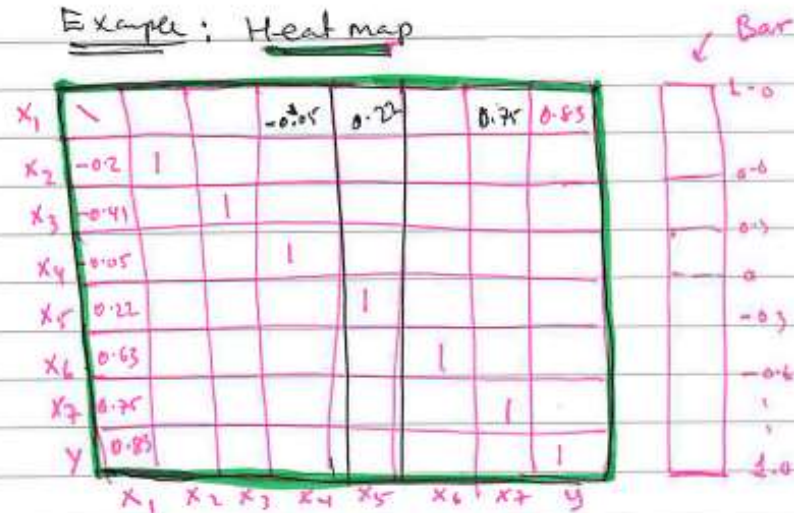
```
sns.heatmap(Cor, annot=True,  
             cmap=plt.cm.Reds)
```

```
plt.show()
```

annot=True → If true, write the "r" values in each cell

cmap → The mapping from data value to color space

Example: Heat map



Step 4: Apply criteria to select the appropriate features that would be correlated with target variable (y)

Correlation with output/Target variable

$$\text{Cor_target} = \text{abs}(\text{cor}['Y'])$$

$$\text{Relevant_feature} = \text{Cor_target}[\text{Cor_target} > 0.5]$$

Relevant_features

Outputs: $X_1 = 0.83$
 $X_6 = 0.63$
 $X_7 = 0.75$
 $Y = 1$

Step 5 Save Selected / cleaned features

Select Column

$\text{SelectedColumns} = ['X_1', 'X_6', 'X_7', 'Y']$

Selecting final Column

$\text{Data_For ML} = \text{df}[\text{SelectedColumns}]$

$\text{Data_For ML. head()}$

	X1	X6	X7	Y
0				
1				
2				
3				
4				

Step 6. Saving the Selected Column for/to
be used ML. modeling

DataToBeSaved = Data_forML

Export the Data Frame to an Excel file

DataToBeSaved.to_excel('FinalCleanedSelected.xlsx')

(2.4) Data Smoothing / Filtering

Sometimes data may contain spikes, which is due to / could be due to Noises.

- The commonly used approach to minimize the spikes is by applying a filter that smooths the signal.

→ The two types of smoothing functions are

- (a) Moving average
- (b) Exponential smoothing.

Moving average is the mean of a certain window data set.

$$y_{mva} = \frac{1}{m} \sum_{i=1}^m y_i, \quad m \text{ is window size}$$

1	y_1	}
2	y_2	
3	y_3	
4	y_4	
5	y_5	
6	y_6	}
7	y_7	
8	y_8	}
9	y_9	
10	y_{10}	}
11	y_{11}	
12	y_{12}	}
13	y_{13}	
14	y_{14}	}
15	y_{15}	

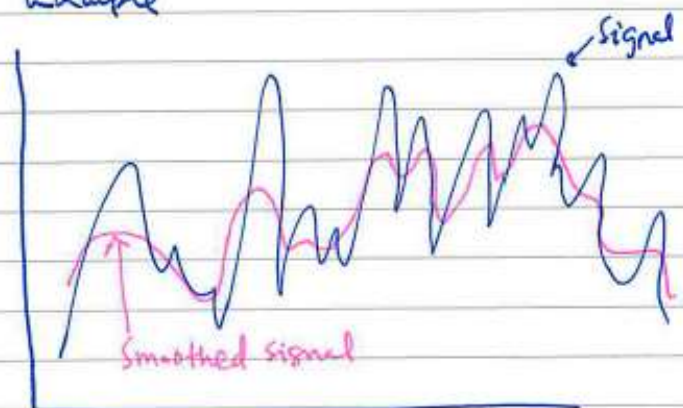
$$m=3 \quad y_{mva}^1 = \frac{1}{3} (y_1 + y_2 + y_3)$$

$$y_{mva}^2 = \frac{1}{3} (y_2 + y_3 + y_4)$$

$$y_{mva}^3 = \frac{1}{3} (y_3 + y_4 + y_5)$$

⋮

Example



Example: Here we will apply MOVING average Smoothing filter on the final selected features

→ The file name was "FinalCleanedSelected.xlsx"

Step 1 Import important libraries

```
import pandas as pd
import matplotlib.pyplot as plt
```

Step 2: Import data file

```
df = pd.read_excel('FinalCleanedSelected.xlsx')
```

Step 3 Apply MOVING AVERAGE filter

```
DataFilter = df.rolling(window=5).mean()
```

Step 4. Save the filtered data as new filename

```
DataFilter.to_excel('FinalCleanedSelectedFiltered.xlsx')
```

The filtered file is saved as

FinalCleanedSelectedFiltered.xlsx

Note that For the ML-modeling, you will use either:

(a) FinalCleanedSelected.xlsx → Unfiltered

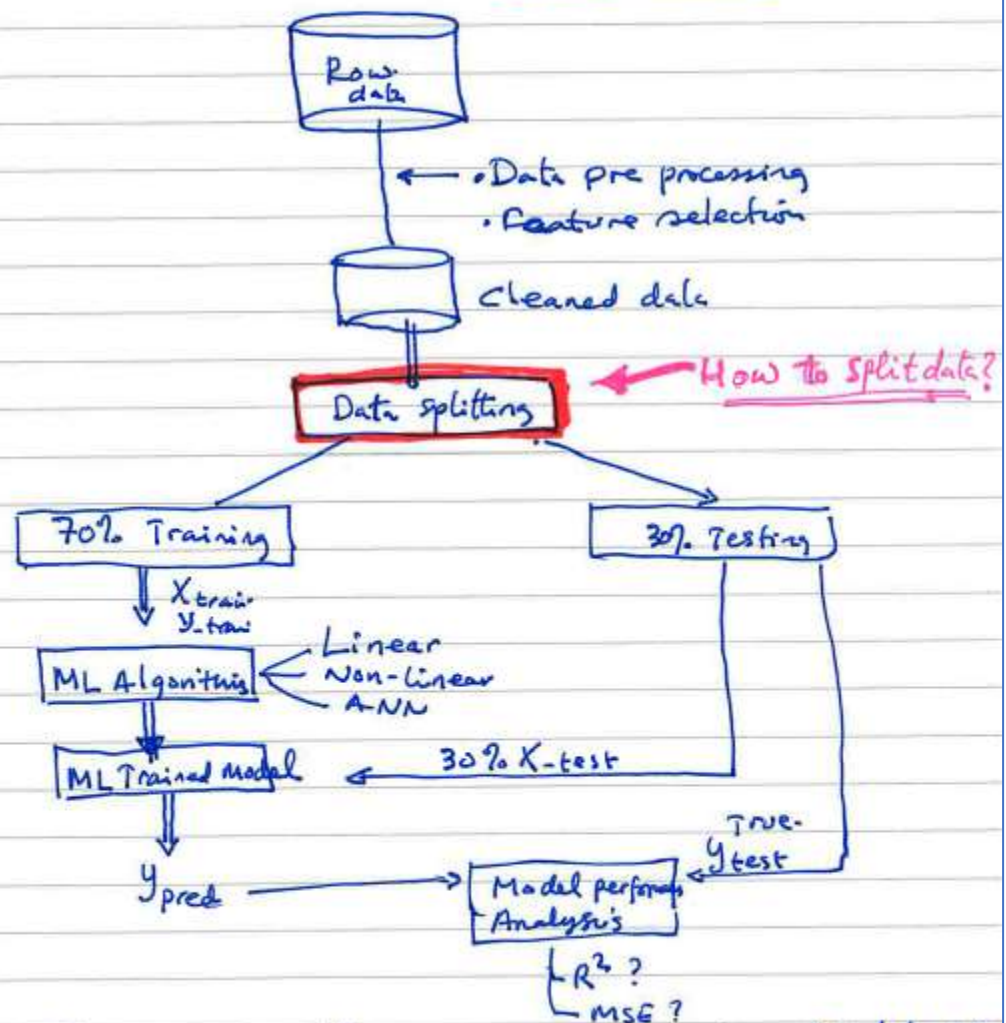
(b) FinalCleanedSelectedFiltered.xlsx → filtered

3.0 Data Splitting

For machine learning modelling, the cleaned and selected features will be split into training and testing

→ Training dataset (70%) will be used for modeling

→ The remaining dataset (30%) will be used for testing the model



In the next section, we will learn How to split dataset?

Assume that the final cleaned (and selected) feature data set has been saved as "cleaned.xlsx"

→ Import the file as

`df = pd.read_excel('cleaned.xlsx')`

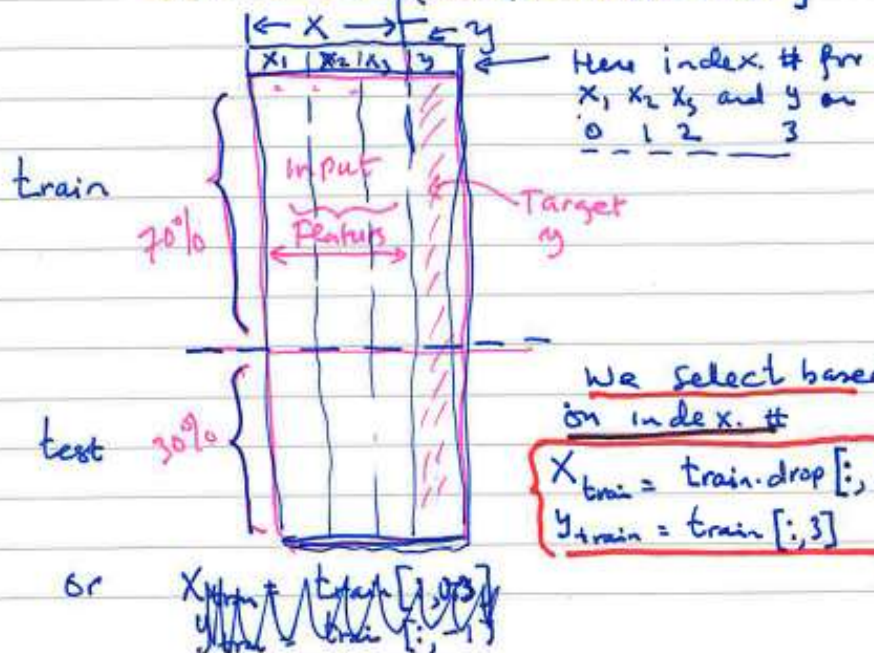
→ Here, we will see different methods of splitting training ($X_{\text{train}}, y_{\text{train}}$) and testing ($X_{\text{test}}, y_{\text{test}}$) datasets.

→ By convention we use capital letter X → for input features, and small letter y → for target output

Method #1

`train = df[: (int(len(df)*0.7))]`

`test = df[(int(len(df)*0.7)):]`



Method #2

← input

x_1	x_2	x_3	y
-------	-------	-------	-----

index # for x_1, x_2, x_3 and y are
index ←

0	1	2	3
---	---	---	---

$X_{\text{train}} = \text{train} ['x_1', 'x_2', 'x_3'].asfloat()$

$y_{\text{train}} = \text{train} ['y'].asfloat()$

$X_{\text{test}} = \text{test} ['x_1', 'x_2', 'x_3'].asfloat()$

$y_{\text{test}} = \text{test} ['y'].asfloat()$

Method #3

$X_{\text{train}} = \text{train}[:, 0:3] \rightarrow$ excluding index #3,
0, 1, 2 \rightarrow for x_1, x_2, x_3 respectively
 $y_{\text{train}} = \text{train}[:, 3]$ \swarrow use index #3, which is y
 $X_{\text{test}} = \text{test}[:, 0:3]$
 $y_{\text{test}} = \text{test}[:, 3]$

Method #4

$X_{\text{train}} = \text{train}.iloc[:, 0:3].value$

$y_{\text{train}} = \text{train}.iloc[:, 3].value$

$X_{\text{test}} = \text{test}.iloc[:, 0:3].value$

$y_{\text{test}} = \text{test}.iloc[:, 3].value$

Method #5 → train_test_split method

This method is based on built-in sklearn library / function

» For this the function should be first imported, which is train_test_split

- Import Library

```
from sklearn.model_selection import  
train_test_split
```

- Import data

```
df = pd.read_excel('cleaned.xlsx')
```

Step 1

- Separate X and y features / target

```
X = df[['x1', 'x2', 'x3']] ← double  
y = df['y'] P
```

Step 2

- Split data for training and split

```
X_train, X_test, y_train, y_test = train_test_split(  
X, y, test_size=0.3,  
random_states=1)
```

After this Modeling follows, Example
What is the Purpose of Random States?

"Random State" is a parameter in train_test_split that CONTROLS the random number generator used to shuffle the data before splitting it.

> With random_state = None (default), we get different train and test across different executions value

> With random_state = 0, we get the same train and test ^{sets} across different executions

→ The most popular are 0, or 42

> With random_state = 42, we get the same train and test sets across the different execution, but in this time the train and test sets are different from the previous case with random_state = 0

The train and test sets directly affect the model's performance score

Because we set different train and test sets with different integer values for random_state in the train-test-split() function, the value of the random state hyperparameter indirectly affects the model performance score.

4.0 Summary

Data pre processing is the first stage of machine learning modeling.

The main focus on data preprocessing that we learned are

- Data Cleaning
- Feature selection
- Data Smoothing / optional.

Once the data is selected, the next stage is Modeling. This is performed by splitting the selected dataset into Training (for modeling) and testing (for prediction).

For data splitting, chapter 3 of this note presents different methods.