



Minia University

Faculty of Computers & information

Artificial Neural Networks and Deep Learning

Slides By:



T.A. Sarah Osama Talaat



E-mail: SarahOsama.fci@gmail.com

Slides were prepared based on set of references mentioned in the last slide



Labs, FCI, Mina University

Agenda

- ❑ Python Libraries
- ❑ Data Preprocessing.
 - ❑ Data cleaning
 - Dealing with missing data.
 - Dealing with categorical data.
 - ❑ Feature scaling
 - Rescaling
 - Mean normalization
 - Standardization
 - Scaling to unit length
 - ❑ Dimensionality reduction
 - Feature selection
 - Feature extraction
 - ❑ Partitioning a dataset in training and testing sets



Let's Start



Library	Uses
numpy	numpy is a library that contains a mathematical tools. Basically, we need this library to include any type of mathematical operations .
pandas	pandas considers as one of the best libraries that is used to import and manage the datasets .
matplotlib.pyplot	matplotlib.pyplot is a sub library which is actually used to plot a set of charts
xlsxwriter	xlsxwriter is a library that is used to write in Excel files.

sklearn.preprocessing

This a package/library provides several common utility functions and transformer classes to change **raw feature** vectors into a representation that is more suitable for the downstream estimators.

Uses

- Standardization, or mean removal and variance scaling
- Non-linear transformation
- **Normalization**
- Binarization
- **Encoding categorical features**
- **Imputation of missing values**
- Generating polynomial features
- Custom transformers

Class	Imputer	
Uses	The imputer class provides basic strategies for imputing missing values , either using the mean, the median or the most frequent value of the row/column .	
Instruction	from sklearn.preprocessing import Imputer	
Some of the Imputer parameters	<ul style="list-style-type: none"> missing_values strategy axis 	
The possible values for each parameter	Parameter	Values
	missing_values	integer or “NaN”, optional (default=“NaN”)
	strategy	string, mean, median or most_frequent, optional (default=“mean”)
	axis	integer, optional (default=0) The axis along which to impute. If axis=0, then impute along columns and if axis=1, then impute along rows.

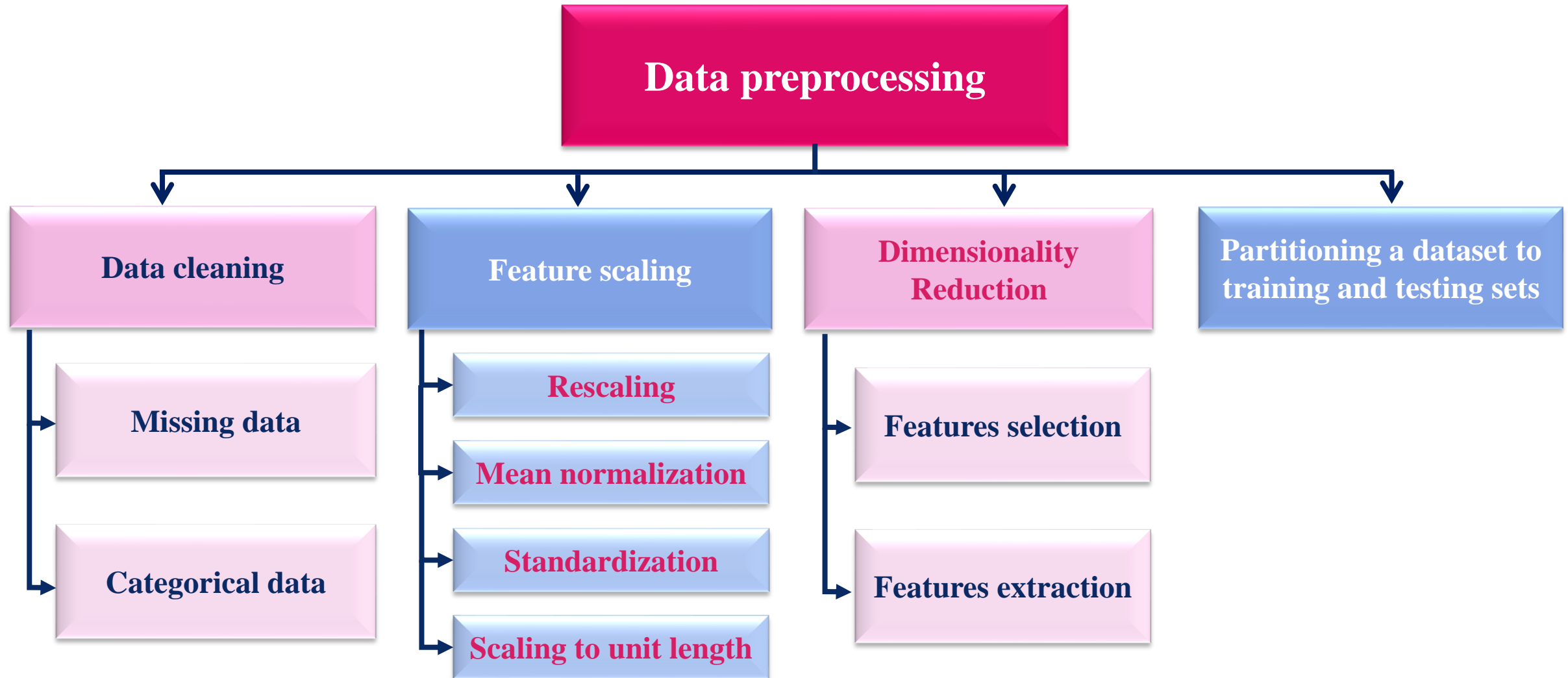
Why we need to the data preprocessing process?

- The main objectives of data preprocessing are to **manipulate** and **transform** raw data into **cleaned** and **scaled** format.
- In addition it is important to **compress** the data onto a **smaller dimensional** subspace while retaining most of the **relevant** information.



Data Preprocessing

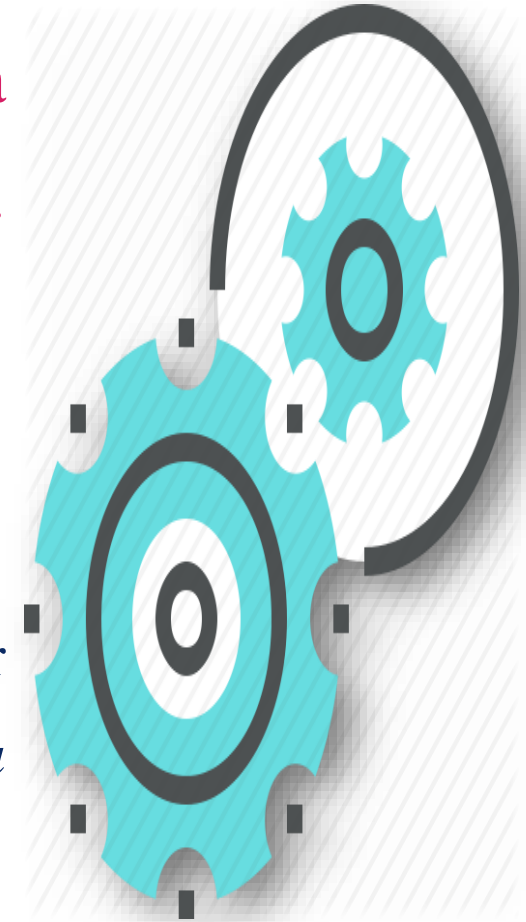
Types of data preprocessing?

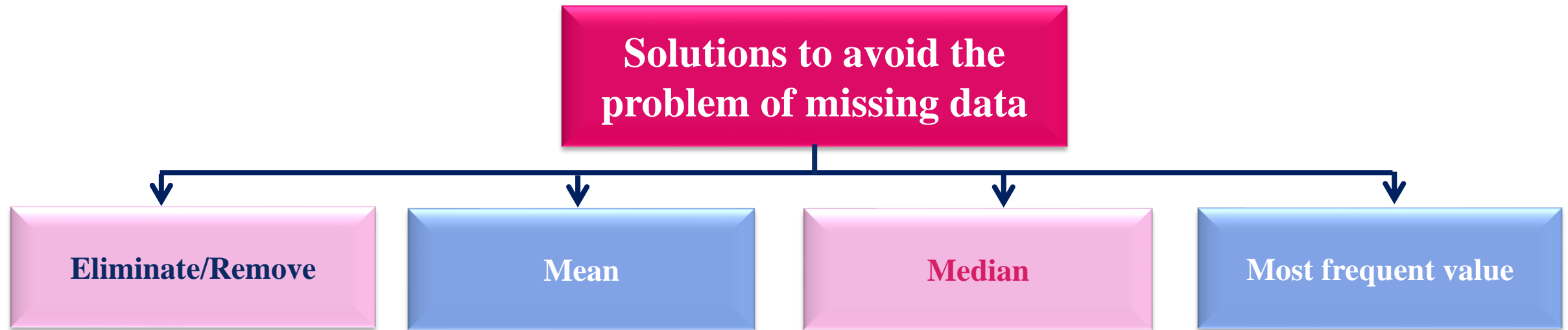


❑ In real-world applications are **familiar** that the collected data samples contain one or more missing values for various reasons.

These reasons include:

- There could have been an error in the data collection process,
- certain measurements are not applicable and
- particular fields could have been simply left blank in a survey, for instance. We typically see *missing values as the blank spaces in our data table or as placeholder strings such as **NaN** (Not A Number).*





□ There are several available solutions to avoid the problem of missing data :

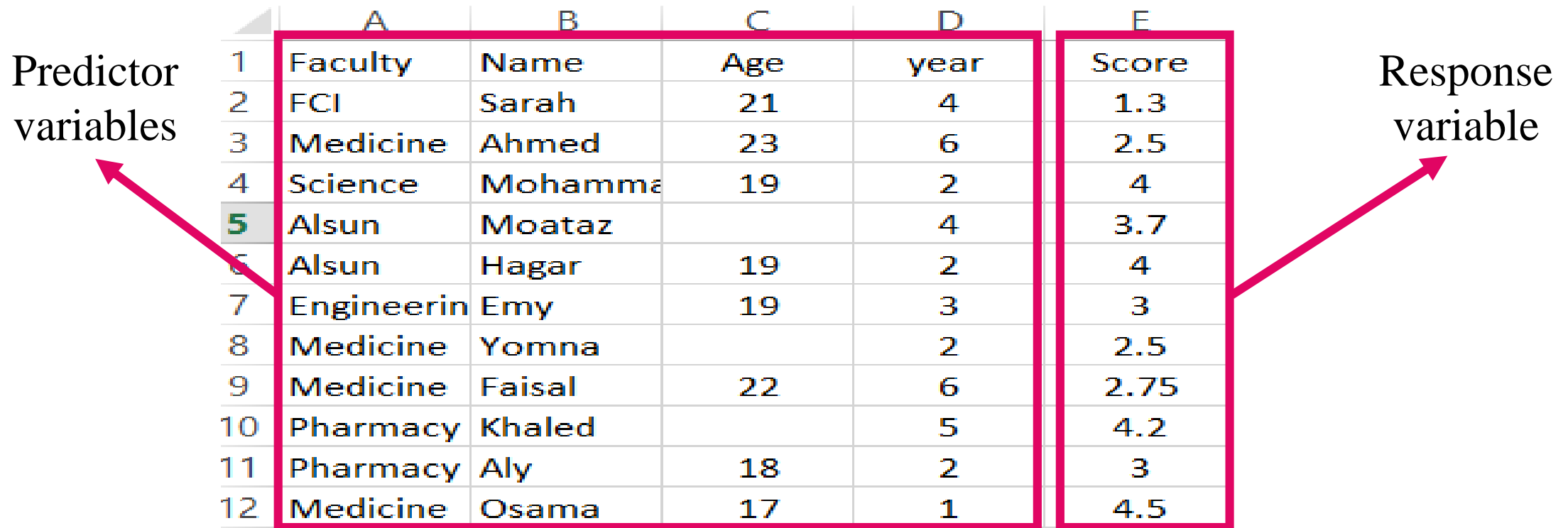
- **Eliminate/Remove:** you can simply eliminate/remove the corresponding features (columns) or samples (rows) that contain missing data from the used dataset (BUT it's quite dangerous, because we might lose too much valuable data).
- **Mean/Median/Most frequent value:** you can take the mean/median/most frequent value from the column/row that contains missing data.
 - **Mean** is used to replace the missing values using the mean/average along the axis (i.e. column or row).
 - **Median** is used to replace missing values using the median along the axis.
 - **Most frequent value** is used to replace missing values using the most frequent value along the axis.

Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- In this dataset, we have 5 columns which includes 4 features (independent variables) and one response (dependent variable).



	A	B	C	D	E
1	Faculty	Name	Age	year	Score
2	FCI	Sarah	21	4	1.3
3	Medicine	Ahmed	23	6	2.5
4	Science	Mohammed	19	2	4
5	Alsun	Moataz		4	3.7
6	Alsun	Hagar	19	2	4
7	Engineering	Emy	19	3	3
8	Medicine	Yomna		2	2.5
9	Medicine	Faisal	22	6	2.75
10	Pharmacy	Khaled		5	4.2
11	Pharmacy	Aly	18	2	3
12	Medicine	Osama	17	1	4.5

Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- In this dataset, we have 5 columns which includes 4 features (independent variables) and one response (dependent variable).

	A	B	C	D	E
1	Faculty	Name	Age	year	Score
2	FCI	Sarah	21	4	1.3
3	Medicine	Ahmed	23	6	2.5
4	Science	Mohammed	19	2	4
5	Alsun	Moataz		4	3.7
6	Alsun	Hagar	19	2	4
7	Engineering	Emy	19	3	3
8	Medicine	Yomna		2	2.5
9	Medicine	Faisal	22	6	2.75
10	Pharmacy	Khaled		5	4.2
11	Pharmacy	Aly	18	2	3
12	Medicine	Osama	17	1	4.5

Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- In this dataset, we have 5 columns which includes 4 features (independent variables) and one response (dependent variable).

	A	B	C	D	E
1	Faculty	Name	Age	year	Score
2	FCI	Sarah	21	4	1.3
3	Medicine	Ahmed	23	6	2.5
4	Science	Mohamma	19	2	4
5	Alsun	Moataz		4	3.7
6	Alsun	Hagar	19	2	4
7	Engineerin	Emy	19	3	3
8	Medicine	Yomna		2	2.5
9	Medicine	Faisal	22	6	2.75
10	Pharmacy	Khaled		5	4.2
11	Pharmacy	Aly	18	2	3
12	Medicine	Osama	17	1	4.5

Solution by using python:

- From “**scikit-learn**” package we use “**Imputer**” class for data missing treatment. It is included in “**preprocessing**” module.

Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- We can solve the missing data problem in Python by using the following instructions;

```
27 import numpy as np
28 import pandas as pd
29 import matplotlib.pyplot as plt
30
31 # Importing the dataset
32 dataset = pd.read_csv('Dataset_1.csv')
33
34 # Create a matrix to store the predictor (independent) variables
35 X = dataset.iloc[:, :-1].values
36
37 # Create a victore to store the response (dependent) variable
38 y = dataset.iloc[:, 3].values
39 # Dealing with missing data by the eliminating method
40 # Remove the rows that have missing cells
41 dataset.dropna()
42 # Remove the columns that have missing cells
43 dataset.dropna(axis=1)
44 # Dealing with missing data by using mean, median or most frequent
45
```

Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- We can solve the missing data problem in Python by using the following instructions;

```
23 # Data Preprocessing
24 # Importing the Libraries
25 import numpy as np
26 import pandas as pd
27 import matplotlib.pyplot as plt
28
29 # Importing the dataset
30 dataset = pd.read_csv('Dataset_1.csv')
31 # Create a matrix to store the predictor (independent) variables
32 X = dataset.iloc[:, :-1].values
33 # Create a victore to store the response (dependent) variable
34 y = dataset.iloc[:, 3].values
35 # Dealing with missing data by using mean, median or most frequent
36 from sklearn.preprocessing import Imputer
37
38 # strategy can be "mean", "median" or "most_frequent"
39 #If "mean", then replace missing values using the mean along the axis.
40 #If "median", then replace missing values using the median along the axis.
41 #If "most_frequent", then replace missing using the most frequent value along
42
43 imputer = Imputer(missing_values = 'NaN', strategy = 'mean', axis = 0)
44 imputer = imputer.fit(X[:, 2:3])
45 X[:, 2:3] = imputer.transform(X[:, 2:3])
```


Data Preprocessing

Back to data cleaning: missing data

□ Example(2.1):

- We can solve the missing data problem in Python by using the following instructions;

```
23# Data Preprocessing
24# Importing the Libraries
25import numpy as np
26import pandas as pd
27import matplotlib.pyplot as plt
28
29# Importing the dataset
30dataset = pd.read_csv('Dataset_1.csv')
31# Create a matrix to store the predictor (independent) variables
32X = dataset.iloc[:, :-1].values
33# Create a victore to store the response (dependent) variable
34y = dataset.iloc[:, 3].values
35# Dealing with missing data by using mean, median or most frequent
36from sklearn.preprocessing import Imputer
37
38# strategy can be "mean", "median" or "most_frequent"
39#If "mean", then replace missing values using the mean along the axis.
40#If "median", then replace missing values using the median along the axis.
41#If "most_frequent", then replace missing using the most frequent value along
42
43imputer = Imputer(missing_values = 'NaN', strategy = 'mean', axis = 0)
44imputer = imputer.fit(X[:, 2:3])
45X[:, 2:3] = imputer.transform(X[:, 2:3])
```

	A	B	C	D	E
1	Faculty	Name	Age	year	Score
2	FCI	Sarah	21	4	1.3
3	Medicine	Ahmed	23	6	2.5
4	Science	Mohammad	19	2	4
5	Alsun	Moataz	19.75	4	3.7
6	Alsun	Hagar	19	2	4
7	Engineering	Emy	19	3	3
8	Medicine	Yomna	19.75	2	2.5
9	Medicine	Faisal	22	6	2.75
10	Pharmacy	Khaled	19.75	5	4.2
11	Pharmacy	Aly	18	2	3
12	Medicine	Osama	17	1	4.5

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- In this dataset, we have three independent variables and one dependent variable.
- In this dataset all input variables consider as categorical features.



	A	B	C	D
1	Country	Color	Size	Prise
2	Egypt	red	S	150.75
3	USA	black	M	200
4	KSA	blue	L	170
5	Canda	white	S	300
6	Egypt	green	M	120
7	China	green	XI	200.5
8	Germany	blue	2XL	499.99

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- At the first, we deal with the ordinal feature (feature number 3 at column 2).

```
1 # Data Preprocessing
2
3 # Importing the Libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Dataset_2.csv')
10 X = dataset.iloc[:, :-1].values
11 Y = dataset.iloc[:, 3].values
12
13 # Dealing with categorical data
14 # Mapping ordinal feature/independent variable
15
16 # Define a new dictionary to store the sizes and its values
17 size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
18 dataset['Size'] = dataset['Size'].map(size_mapping)
19 X[:,2] = dataset.iloc[:,2].values
20
```

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- At the first, we deal with the ordinal feature (feature number 3 at column 2).

```
1 # Data Preprocessing
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Dataset_2.csv')
10 X = dataset.iloc[:, :-1].values
11 Y = dataset.iloc[:, 3].values
12
13 # Dealing with categorical data
14 # Mapping ordinal feature/independent variable
15
16 # Define a new dictionary to store the sizes and its values
17 size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
18 dataset['Size'] = dataset['Size'].map(size_mapping)
19 X[:,2] = dataset.iloc[:,2].values
20
```

	Country	Color	Size	Prise
0	Egypt	red	1	150.75
1	USA	black	2	200.00
2	KSA	blue	3	170.00
3	Canda	white	1	300.00
4	Egypt	green	2	120.00
5	China	green	4	200.50
6	Germany	blue	5	499.99

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- If we want to transform the integer values back to the original string representation at a later stage, we can simply define a reverse-mapping dictionary.

```
1 # Data Preprocessing
2 |
3 # Importing the Libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7 # Importing the dataset
8 dataset = pd.read_csv('Dataset_2.csv')
9 X = dataset.iloc[:, :-1].values
10 Y = dataset.iloc[:, 3].values
11 # Dealing with categorical data
12 # Mapping ordinal feature/independent variable
13 # Define a new dictionary to store the sizes and its values
14 size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
15 dataset['Size'] = dataset['Size'].map(size_mapping)
16 X[:,2] = dataset.iloc[:,2].values
17 # If we want to transform the integer values back to the original
18 # string representation at a later stage, we can simply define
19 # a reverse-mapping dictionary
20 int_size_mapping = {v: k for k, v in size_mapping.items()}
21 dataset['Size'] = dataset['Size'].map(int_size_mapping)
22 X[:,2] = dataset.iloc[:,2].values
```

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- If we want to transform the integer values back to the original string representation at a later stage, we can simply define a reverse-mapping dictionary.

```
1# Data Preprocessing
2|
3# Importing the Libraries
4import numpy as np
5import matplotlib.pyplot as plt
6import pandas as pd
7# Importing the dataset
8dataset = pd.read_csv('Dataset_2.csv')
9X = dataset.iloc[:, :-1].values
10Y = dataset.iloc[:, 3].values
11# Dealing with categorical data
12# Mapping ordinal feature/independent variable
13# Define a new dictionary to store the sizes and its values
14size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
15dataset['Size'] = dataset['Size'].map(size_mapping)
16X[:,2] = dataset.iloc[:,2].values
17# If we want to transform the integer values back to the original
18# string representation at a later stage, we can simply define
19# a reverse-mapping dictionary
20int_size_mapping = {v: k for k, v in size_mapping.items()}
21dataset['Size'] = dataset['Size'].map(int_size_mapping)
22X[:,2] = dataset.iloc[:,2].values
```

	Country	Color	Size	Prise
0	Egypt	red	S	150.75
1	USA	black	M	200.00
2	KSA	blue	L	170.00
3	Canda	white	S	300.00
4	Egypt	green	M	120.00
5	China	green	XL	200.50
6	Germany	blue	2XL	499.99

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- Now, we deal with the nominal features (feature number 1 and 2 at column 0 and 1).

```
1 # Data Preprocessing
2
3 # Importing the Libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7 # Importing the dataset
8 dataset = pd.read_csv('Dataset_2.csv')
9 X = dataset.iloc[:, :-1].values
10 Y = dataset.iloc[:, 3].values
11 # Dealing with categorical data
12
13 # Encoding nominal features/independent variable
14 from sklearn.preprocessing import LabelEncoder
15 Country_labelencoder = LabelEncoder()
16 X[:, 0] = Country_labelencoder.fit_transform(X[:, 0])
17 Color_labelencoder = LabelEncoder()
18 X[:, 1] = Color_labelencoder.fit_transform(X[:, 1])
19
```


Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

- Now, we deal with the nominal features (feature number 1 and 2 at column 0 and 1).

```
In [38]: # Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# Dealing with categorical data
# Encoding nominal features/independent variable
from sklearn.preprocessing import LabelEncoder
Country_labelencoder = LabelEncoder()
X[:, 0] = Country_labelencoder.fit_transform(X[:, 0])
Color_labelencoder = LabelEncoder()
X[:, 1] = Color_labelencoder.fit_transform(X[:, 1])
```

```
Out[38]: X
array([[2, 3, 1],
       [5, 0, 2],
       [4, 1, 3],
       [0, 4, 1],
       [2, 2, 2],
       [1, 2, 4],
       [3, 1, 5]], dtype=object)
```


Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

```
1 # Data Preprocessing
2
3 # Importing the Libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7 # Importing the dataset
8 dataset = pd.read_csv('Dataset_2.csv')
9 X = dataset.iloc[:, :-1].values
10 Y = dataset.iloc[:, 3].values
11 # Dealing with categorical data
12 # Mapping ordinal feature/independent variable
13 # Define a new dictionary to store the sizes and its values
14 size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
15 dataset['Size'] = dataset['Size'].map(size_mapping)
16 X[:,2] = dataset.iloc[:,2].values
17
18 # Encoding nominal features/independent variable
19 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
20 Country_labelencoder = LabelEncoder()
21 X[:, 0] = Country_labelencoder.fit_transform(X[:, 0])
22 Color_labelencoder = LabelEncoder()
23 X[:, 1] = Color_labelencoder.fit_transform(X[:, 1])
24
25 Country_color_onehotencoder = OneHotEncoder(categorical_features = [0,1])
26 X = Country_color_onehotencoder.fit_transform(X).toarray()
```

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

```
In [78]: # Data Preprocessing
# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
# Importing the dataset
dataset = pd.read_csv('Dataset_2.csv')
X = dataset.iloc[:, :-1].values
Y = dataset.iloc[:, 3].values
# Dealing with categorical data
# Mapping ordinal feature/independent variable
# Define a new dictionary to store the sizes and its values
size_mapping = {'S':1, 'M':2, 'L':3, 'XL':4, '2XL':5}
dataset['Size'] = dataset['Size'].map(size_mapping)
X[:,2] = dataset.iloc[:,2].values
# Encoding nominal features/independent variable
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
Country_labelencoder = LabelEncoder()
X[:, 0] = Country_labelencoder.fit_transform(X[:, 0])
Color_labelencoder = LabelEncoder()
X[:, 1] = Color_labelencoder.fit_transform(X[:, 1])

Country_color_onehotencoder = OneHotEncoder(categorical_features = [0,1])
X = Country_color_onehotencoder.fit_transform(X).toarray()
X
```

```
Out[78]: array([[ 0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  1.],
 [ 0.,  0.,  0.,  0.,  0.,  1.,  1.,  0.,  0.,  0.,  0.,  2.],
 [ 0.,  0.,  0.,  0.,  1.,  0.,  0.,  1.,  0.,  0.,  0.,  3.],
 [ 1.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  1.],
 [ 0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  2.],
 [ 0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  0.,  4.],
 [ 0.,  0.,  0.,  1.,  0.,  0.,  0.,  1.,  0.,  0.,  0.,  5.]])
```

Data Preprocessing

Data cleaning: categorical data

□ Example(2.2):

```
In [78]: # Data Preprocessing
# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
# Importing the dataset
```

```
Out[78]: array([[ 0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  1.],
 [ 0.,  0.,  0.,  0.,  0.,  1.,  1.,  0.,  0.,  0.,  0.,  2.],
 [ 0.,  0.,  0.,  0.,  1.,  0.,  0.,  1.,  0.,  0.,  0.,  3.],
 [ 1.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  1.,  1.],
 [ 0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  0.,  2.],
 [ 0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  0.,  4.],
 [ 0.,  0.,  0.,  1.,  0.,  0.,  0.,  1.,  0.,  0.,  0.,  5.]])
```

Country feature after
dummy encoding

Color feature
after dummy
encoding

```
[ 0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  0.,  2.],
[ 0.,  1.,  0.,  0.,  0.,  0.,  0.,  0.,  1.,  0.,  0.,  4.],
[ 0.,  0.,  0.,  1.,  0.,  0.,  0.,  1.,  0.,  0.,  0.,  5.]])
```

Data Preprocessing

Feature scaling: rescaling

□ Example(2.3):

- In the following dataset, we have three independent variables and one dependent variable.

	A	B	C	D
1	High temp	Solar radiation	Wind speed	Rain
2	13	192	35	18.00
3	39	178	56	12.00
4	9	199	74	23.00
5	14	100	57	12.00
6	1	187	54	25.00
7	-4	139	57	14.00
8	46	156	58	13.00
9	6	195	51	24.00
10	-4	166	38	10.00
11	11	163	68	12.00
12	29	178	42	22.00

□ Example(2.3):

```
1 # Data Preprocessing
2 # Importing the libraries
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 # Importing the dataset
7 dataset = pd.read_csv('Dataset_3.csv')
8 X = dataset.iloc[:, :-1].values
9 Y = dataset.iloc[:, 3].values
10 # Feature scaling: rescaling
11 # The min-max scaling procedure is implemented in scikit-learn and can be used
12 # as follows:
13 from sklearn.preprocessing import MinMaxScaler
14 # Here we scale all input features
15 minMaxScaled_rescaling = MinMaxScaler()
16 X_scaled = minMaxScaled_rescaling .fit_transform(X)
```

Data Preprocessing

Feature scaling: rescaling

□ Example(2.3):

```
C:\Users\sarah\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning: Data with input dtype int64 was converted to float64 by MinMaxScaler.
```

```
warnings.warn(msg, DataConversionWarning)
```

```
array([[ 0.34      ,  0.92929293,  0.        ],
       [ 0.86      ,  0.78787879,  0.53846154],
       [ 0.26      ,  1.        ,  1.        ],
       [ 0.36      ,  0.        ,  0.56410256],
       [ 0.1       ,  0.87878788,  0.48717949],
       [ 0.        ,  0.39393939,  0.56410256],
       [ 1.        ,  0.56565657,  0.58974359],
       [ 0.2       ,  0.95959596,  0.41025641],
       [ 0.        ,  0.66666667,  0.07692308],
       [ 0.3       ,  0.63636364,  0.84615385],
       [ 0.66      ,  0.78787879,  0.17948718]])

[ 0.        ,  0.39393939,  0.56410256],
[ 1.        ,  0.56565657,  0.58974359],
[ 0.2       ,  0.95959596,  0.41025641],
[ 0.        ,  0.66666667,  0.07692308],
[ 0.3       ,  0.63636364,  0.84615385],
[ 0.66      ,  0.78787879,  0.17948718]])
```

□ Example(2.3):

```
1 # Data Preprocessing
2 # Importing the Libraries
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 # Importing the dataset
7 dataset = pd.read_csv('Dataset_3.csv')
8 X = dataset.iloc[:, :-1].values
9 Y = dataset.iloc[:, 3].values
10
11 # Feature scaling mean normalization
12 from sklearn.preprocessing import Normalizer
13 # Here we scale all input features
14 normalization = Normalizer()
15 X_scaled = normalization.fit_transform(X)
```

□ Example(2.3):

```
In [8]: runfile('C:/Users/sarah/OneDrive/Machine Learning and Pattern Recognition/Part
1 - Data Preprocessing/Feature_scaling_rescaling.py', wdir='C:/Users/sarah/OneDrive/
Machine Learning and Pattern Recognition/Part 1 - Data Preprocessing')
C:\Users\sarah\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475:
DataConversionWarning: Data with input dtype int64 was converted to float64 by
MinMaxScaler.
    warnings.warn(msg, DataConversionWarning)
```

```
In [9]: X_scaled
```

```
Out[9]:
```

```
array([[ 0.06646335,  0.98161254,  0.17893979],
       [ 0.20458141,  0.93373054,  0.29375792],
       [ 0.04235212,  0.93645244,  0.34822855],
       [ 0.12073902,  0.86242154,  0.49158028],
       [ 0.0051376 ,  0.96073194,  0.27743061],
       [-0.02661585,  0.92490076,  0.37927585],
       [ 0.26639934,  0.90344123,  0.33589482],
       [ 0.02975479,  0.96703066,  0.25291571],
       [-0.02348233,  0.97451677,  0.22308215],
       [ 0.06216178,  0.92112449,  0.38427279],
       [ 0.1566104 ,  0.96126381,  0.22681506]])
```


Data Preprocessing

Feature scaling: standardization

□ Example(2.3):

```
1 # Data Preprocessing
2 # Importing the Libraries
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 # Importing the dataset
7 dataset = pd.read_csv('Dataset_3.csv')
8 X = dataset.iloc[:, :-1].values
9 Y = dataset.iloc[:, 3].values
10
11 # Feature scaling standardization
12 from sklearn.preprocessing import StandardScaler
13 # Here we scale all input features
14 standardization = StandardScaler()
15 X_scaled = standardization.fit_transform(X)
```

Data Preprocessing

Feature scaling: standardization

□ Example(2.3):

```
In [24]: runfile('C:/Users/sarah/OneDrive/Machine Learning and Pattern
Recognition/Part 1 - Data Preprocessing/Feature_scaling_rescaling.py',
wdir='C:/Users/sarah/OneDrive/Machine Learning and Pattern Recognition/
Part 1 - Data Preprocessing')
C:\Users\sarah\Anaconda3\lib\site-packages\sklearn\utils\validation.py:
475: DataConversionWarning: Data with input dtype int64 was converted
to float64 by MinMaxScaler.
      warnings.warn(msg, DataConversionWarning)
```

```
In [25]: X_scaled
```

```
Out[25]:
```

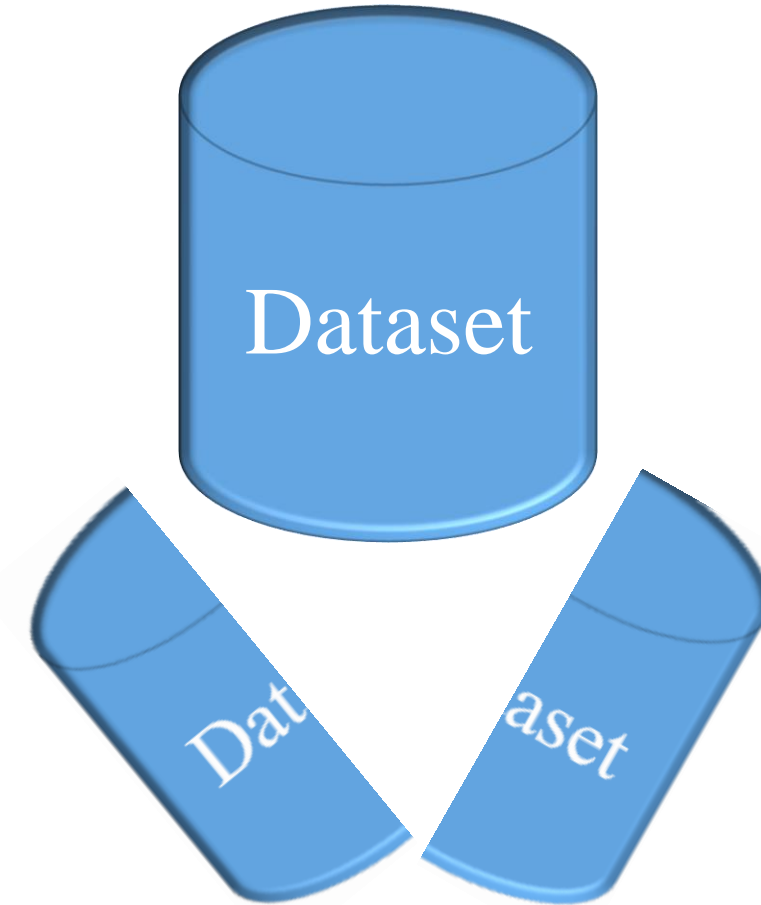
```
array([[ 0.06646335,  0.98161254,  0.17893979],
       [ 0.20458141,  0.93373054,  0.29375792],
       [ 0.04235212,  0.93645244,  0.34822855],
       [ 0.12073902,  0.86242154,  0.49158028],
       [ 0.0051376 ,  0.96073194,  0.27743061],
       [-0.02661585,  0.92490076,  0.37927585],
       [ 0.26639934,  0.90344123,  0.33589482],
       [ 0.02975479,  0.96703066,  0.25291571],
       [-0.02348233,  0.97451677,  0.22308215],
       [ 0.06216178,  0.92112449,  0.38427279],
       [ 0.1566104 ,  0.96126381,  0.22681506]])
```

□ Dimensionality Reduction:

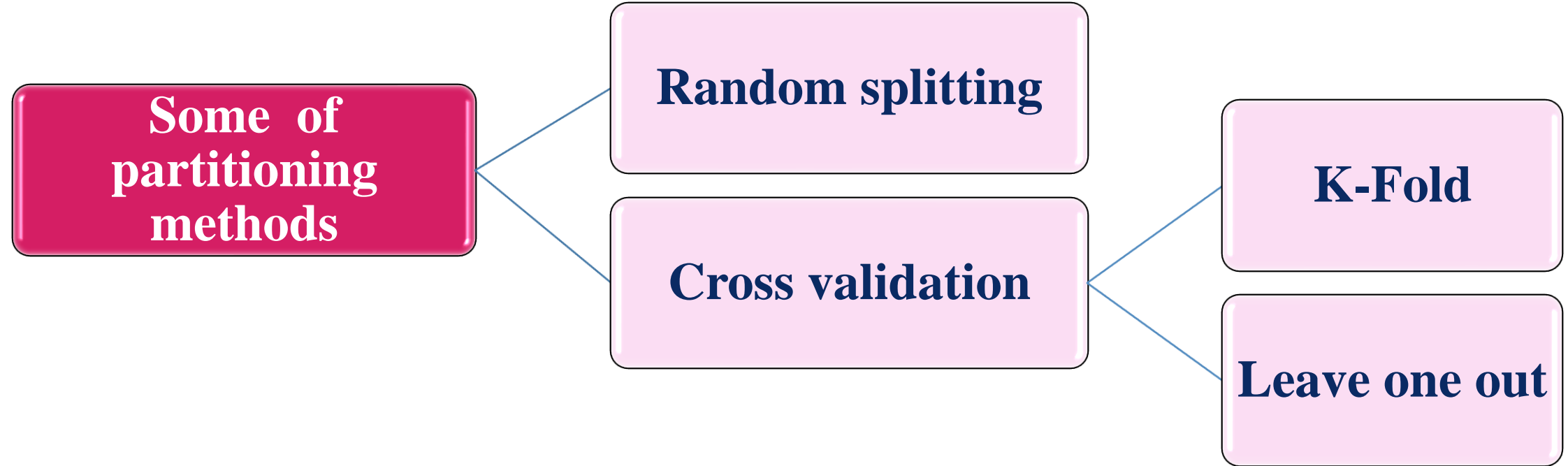
- In machine learning and statistics, **dimensionality reduction** or dimension reduction is the process of **reducing** the number of **random variables** under consideration.
- There are two main categories of dimensionality reduction techniques: **feature selection** and **feature extraction**.
 - **Feature Selection:** we **select** a subset of the original feature set.
 - **Feature Extraction:** we **derive** information from the feature set to construct a new feature subspace

□ Partitioning a dataset in training and testing sets:

- Partitioning is used to determine whether our machine learning algorithm **not only performs well** on the training set but also **generalizes** well to new data,
- We want to **randomly divide** the dataset into a separate **training and testing sets**.
- We use the **training** set to **train** and **optimize** our machine learning model, while we keep the **testing** set until the very end to **evaluate** the final model.



□ Partitioning methods:



□ Random splitting

- Split arrays or matrices into random train and test subsets
- A convenient way to **randomly** partition the dataset into a separate **testing and training** dataset is to use the “**train_test_split**” function from **scikit-learn's “cross_validation”** submodule.
- Another way is to use the to use the “**train_test_split**” function from **scikit-learn's “model_selection”** submodule



□ Example(2.3):

- Back to the dataset that is used in example (2.3) to apply the data partitioning methods on it

```
1 # Data Preprocessing
2
3 # Importing the Libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Dataset_3.csv')
10 X = dataset.iloc[:, :-1].values
11 Y = dataset.iloc[:, 3].values
12
13 # Feature scaling: standardization
14 from sklearn.preprocessing import StandardScaler
15 # Here we scale all input features
16 standardization = StandardScaler()
17 X_scaled = standardization.fit_transform(X)
18
19 # Partitioning a dataset in training and testing sets
20 # Splitting the dataset into the Training set and Test set
21 from sklearn.model_selection import train_test_split
22 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size = 0.2, random_state = 0)
```


□ Example(2.3):

- Back to the dataset that is used in example (2.3) to apply the data partitioning methods on it

Name	Type	Size	Value
X	int64	(11, 3)	array([[13, 192, 35],
X_scaled	float64	(11, 3)	array([[-0.0970969, 1.084888381, -1.65043736],
X_test	int64	(3, 3)	array([[1, 187, 54],
X_train	int64	(8, 3)	array([[29, 178, 42],
Y	float64	(11,)	array([18., 12., 23., ..., 10., 12., 22.])
dataset	DataFrame	(11, 4)	Column names: High temp, Solar radiation, Wind speed , Rain
y_test	float64	(3,)	array([25., 12., 23.])
y_train	float64	(8,)	array([22., 13., 12., 24., 10., 12., 18., 14.])

Profiler

Variable explorer

File explorer

Help

Data Preprocessing

Partitioning a dataset in training and testing sets

Artificial Neural Networks
&
Deep Learning

□ Example(2.3):

- Back to the dataset that is used in example (2.3) to apply the data partitioning methods on it

The image displays four NumPy array viewer windows arranged in a 2x2 grid, showing the partitioning of a dataset into training and testing sets. Each window has a title bar, a table of data, and buttons for Format, Resize, Background color, OK, and Cancel.

X_test - NumPy array

	0	1	2
0	1	187	54
1	11	163	68
2	9	199	74

X_train - NumPy array

	0	1	2
0	29	178	42
1	46	156	58
2	39	178	56
3	6	195	51
4	-4	166	38
5	14	100	57
6	13	192	35
7	-4	139	57

y_test - NumPy array

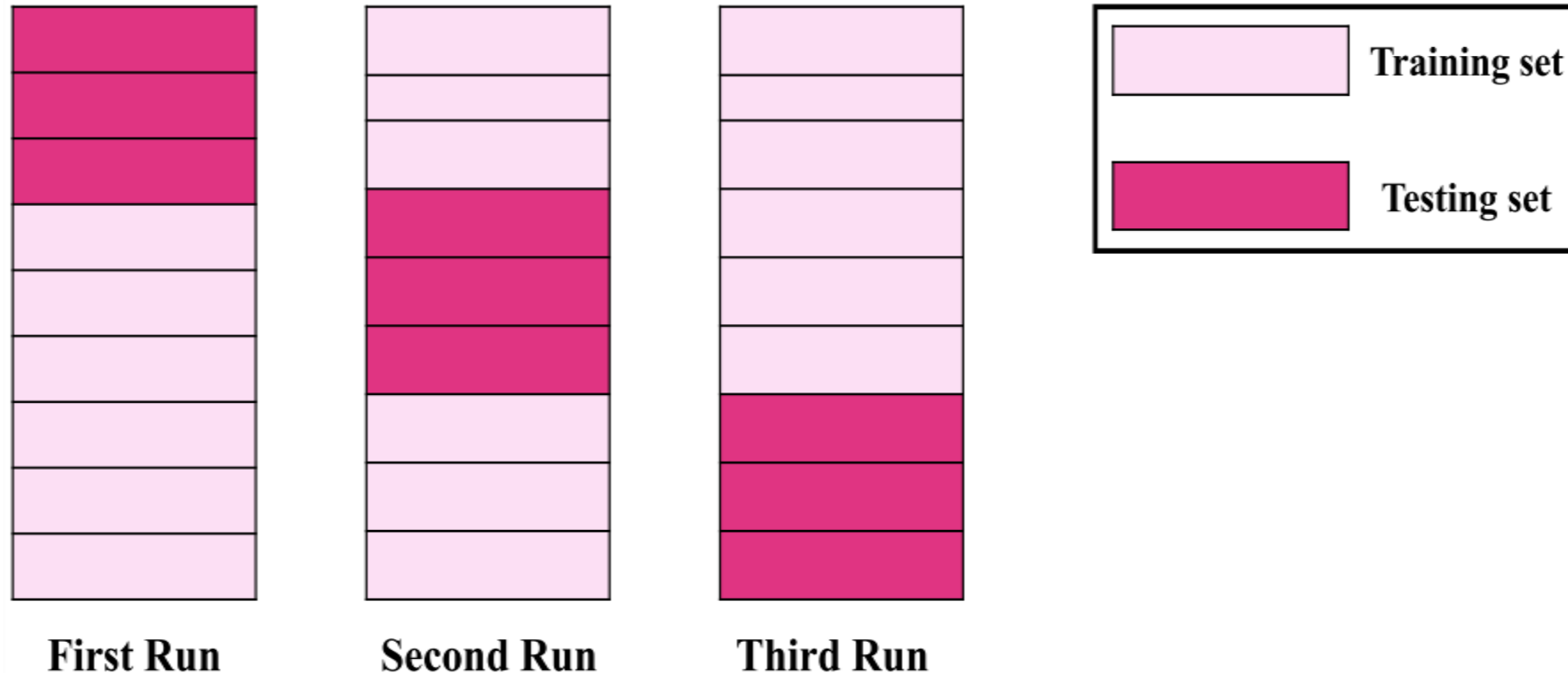
	0
0	25
1	12
2	23

y_train - NumPy array

	0
0	22
1	13
2	12
3	24
4	10
5	12
6	18
7	14

□ K-Fold cross validation:

- The following figure shows the k-fold cross validation schema when $k = 3$.



□ Example(2.3):

- Back to the dataset that is used in example (2.3) to apply the data partitioning methods on it

```
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 # Importing the dataset
7 dataset = pd.read_csv('Dataset_3.csv')
8 X = dataset.iloc[:, :-1].values
9 Y = dataset.iloc[:, 3].values
10 # Feature scaling: standardization
11 from sklearn.preprocessing import StandardScaler
12 standardization = StandardScaler()
13 X_scaled = standardization.fit_transform(X)
14 # Partitioning a dataset in training and testing sets
15 # Splitting the dataset into the Training set and Test set
16 from sklearn.model_selection import KFold
17 # shuffle is used to whether to shuffle the data before splitting into batches.
18 kfold = KFold(n_splits=3, shuffle=False, random_state=None)
19 # Returns the number of splitting iterations in the cross-validator
20 k = kfold.get_n_splits(X) # or # k = kfold.get_n_splits([X,Y,3])
21 # Generate indices to split data into training and test set.
22 indices = kfold.split(X)
23 i = 1
24 for train_index, test_index in kfold.split(X):
25     print("The fold number: ", i)
26     i += 1
27     print("TRAIN:", train_index, "TEST:", test_index)
28     X_train, X_test = X[train_index], X[test_index]
29     y_train, y_test = Y[train_index], Y[test_index]
30     print("X_train:", X[train_index], "\n Y_train:", Y[train_index])
31     print("X_test:", X[test_index], "\n Y_test:", Y[test_index])
```

□ Leave one out(LOO) cross validation:

- LOO is a simple method of cross-validation.
- Each learning set is created by taking **all the** samples **except one**, the test set being the sample **left out**.
- Thus, for samples, we have different training sets and different tests set.
- This cross-validation procedure does not waste much data as only one sample is removed from the training set



□ Example(2.3):

- Back to the dataset that is used in example (2.3) to apply the data partitioning methods on it

```
1 # Data Preprocessing
2 # Importing the Libraries
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import pandas as pd
6 # Importing the dataset
7 dataset = pd.read_csv('Dataset_3.csv')
8 X = dataset.iloc[:, :-1].values
9 Y = dataset.iloc[:, 3].values
10 # Feature scaling: standardization
11 from sklearn.preprocessing import StandardScaler
12 standardization = StandardScaler()
13 X_scaled = standardization.fit_transform(X)
14 # Partitioning a dataset in training and testing sets
15
16 from sklearn.model_selection import LeaveOneOut
17
18 LOO = LeaveOneOut()
19
20 for train, test in LOO.split(X):
21     print("%s %s" % (train, test))
```

□ Example(2.3):

- In [4]: `runfile('C:/Users/sarah/OneDrive/Machine Learning and Pattern Recognition/Part 1 - Data Preprocessing/Partitioning_dataset_L00.py', wdir='C:/Users/sarah/OneDrive/Machine Learning and Pattern Recognition/Part 1 - Data Preprocessing')`

1	2	3	4	5	6	7	8	9	10	[0]
0	2	3	4	5	6	7	8	9	10	[1]
0	1	3	4	5	6	7	8	9	10	[2]
0	1	2	4	5	6	7	8	9	10	[3]
0	1	2	3	5	6	7	8	9	10	[4]
0	1	2	3	4	6	7	8	9	10	[5]
0	1	2	3	4	5	7	8	9	10	[6]
0	1	2	3	4	5	6	8	9	10	[7]
0	1	2	3	4	5	6	7	9	10	[8]
0	1	2	3	4	5	6	7	8	10	[9]
0	1	2	3	4	5	6	7	8	9	[10]

Testing sets

Training set

C:\Users\sarah\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475:

DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.
warnings.warn(msg, DataConversionWarning)

Any Questions!?



Thank you