**CHANDIGARH UNIVERSITY**
Discover. Learn. Empower.

# APEX INSTITUTE OF TECHNOLOGY
## DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

MACHINE LEARNING (21CSH-286)
**Faculty:** Prof. (Dr.) Vineet Mehan (E13038)

Lecture – 4
Data Pre-Processing

DISCOVER . **LEARN** . EMPOWER

1

---

## DBMS: Course Objectives

**COURSE OBJECTIVES**
The Course aims to:

1. Understand and apply various data handling and visualization techniques.
2. Understand about some basic learning algorithms and techniques and their applications, as well as general questions related to analysing and handling large data sets.
3. To develop skills of supervised and unsupervised learning techniques and implementation of these to solve real life problems.
4. To develop basic knowledge on the machine techniques to build an intellectual machine for making decisions behalf of humans.
5. To develop skills for selecting suitable model parameters and apply them for designing optimized machine learning applications.

2

---

## COURSE OUTCOMES

On completion of this course, the students shall be able to:-

| CO1 | Understand machine learning techniques and computing environment that are suitable for the applications under consideration. |
|-----|-----|

---

## Unit-1 Syllabus

| Unit-1 | Introduction to Machine Learning |
|--------|-----------------------------------|
| Introduction to Machine Learning | Definition of Machine Learning, Working principles of Machine Learning; Classification of Machine Learning algorithms: Supervised Learning, Unsupervised Learning, Reinforcement Learning, Semi-Supervised Learning; Applications of Machine Learning. |
| Data Pre-Processing and Feature Extraction | Data Sourcing and Cleaning, Handling Missing data, Encoding Categorical data, Feature Scaling, Handling Time Series data; Feature Selection techniques, Data Transformation, Normalization, Dimensionality reduction |
| Data Visualization | Data Frame Basics, Different types of analysis, Different types of plots, Plotting fundamentals using Matplotlib, Plotting Data Distributions using Seaborn. |

4

---

## SUGGESTIVE READINGS

• **TEXT BOOKS:**
• There is no single textbook covering the material presented in this course. Here is a list of books recommended for further reading in connection with the material presented:
• **T1:** Tom.M.Mitchell, "Machine Learning, McGraw Hill International Edition".
• **T2:** Ethern Alpaydin," Introduction to Machine Learning. Eastern Economy Edition, Prentice Hall of India, 2005".
• **T3:** Andreas C. Miller, Sarah Guido, Introduction to Machine Learning with Python, O'REILLY (2001).

• **REFERENCE BOOKS:**
• **R1** Sebastian Raschka, Vahid Mirjalili, Python Machine Learning, (2014)
• **R2** Richard O. Duda, Peter E. Hart, David G. Stork, "Pattern Classification, Wiley, 2nd Edition".
• **R3** Christopher Bishop, "Pattern Recognition and Machine Learning, illustrated Edition, Springer, 2006".

5

---

## Data Sourcing

• For data sourcing Panda is used.

• Panda is a python Library for analyzing data.

• **Name?**
• Panda = Panel Data + Python Data Analysis (Combination) gave the name.
• Panel data is a subset of longitudinal data where observations are for the same subjects each time.

By: Prof. (Dr.) Vineet Mehan

6

## Data Sourcing

- Use of Panda ?

- Pandas allow us to analyze big data and make conclusions based on statistical theories.

- Pandas can clean messy data sets, and make them readable and relevant.

- Pandas are used in Data Science.

## Data Sourcing

- Data Science: is a branch of computer science where we study how to store, use and analyze data for deriving information from it.

- How to install Pandas?
- 1.    Open cmd prompt
- 2.    Type
- >>> python –m pip install pandas

## Make a data Frame that tells the type of vehicles that passed a toll plaza.

- import pandas
- mydataset = { 'cars': ["Maruti", "Hundai", "Tata"],  'passings': [20, 12, 15]}
- myvar = pandas.DataFrame(mydataset)
- print(myvar)

```
>>> import pandas
>>> mydataset = { 'cars': ["Maruti", "Hundai", "Tata"],  'passings': [20, 12, 15
]}
>>> myvar = pandas.DataFrame(mydataset)
>>> print(myvar)
      cars  passings
0   Maruti        20
1   Hundai        12
2     Tata        15
>>>
```

## Import pandas as pd and use pd

```
>>> import pandas as pd
>>> data = {
  "calories": [420, 380, 390],
  "duration": [50, 40, 45]
}
>>> #load data into a DataFrame object:
df = pd.DataFrame(data)
>>> print(df)
   calories  duration
0       420        50
1       380        40
2       390        45
>>>
```

## Read data from a CSV File

```
>>> import pandas as pd
>>> df = pd.read_csv('D:\\3_Educational\\1_EngineeringSubjects\\Python Programmi
ng\\Lecture 24 Panda\\data.csv')
>>> print(df.to_string())
    Duration  Pulse  Maxpulse  Calories
0         60    110       130     409.1
1         60    117       145     479.0
2         60    103       135     340.0
3         45    109       175     282.4
4         45    117       127     406.0
5         60    102       127     300.0
6         60    110       136     374.0
7         45    104       134     253.3
8         30    109       133     195.1
9         60     98       124     269.0
10        60    103       147     329.3
11        60    100       120     250.7
12        60    106       128     345.3
13        60    104       132     379.3
14        60     98       123     275.0
15        60     98       120     215.2
16        60    100       120     300.0
17        45     90       112       NaN
18        60    103       123     323.0
19        45     97       125     243.0
20        60    108       131     364.2
21        45    100       119     282.0
22        60    130       101     300.0
23        45    105       132     246.0
24        60    102       126     334.5
25        60    100       120     250.0
26        60     92       118     241.0
27        60    103       132       NaN
28        60    100       132     280.0
```

## Reading CSV but print without converting to string

```
>>> import pandas as pd
>>> df = pd.read_csv('D:\\3_Educational\\1_EngineeringSubjects\\Python Programmi
ng\\Lecture 24 Panda\\data.csv')
>>> print(df)
     Duration  Pulse  Maxpulse  Calories
0          60    110       130     409.1
1          60    117       145     479.0
2          60    103       135     340.0
3          45    109       175     282.4
4          45    117       148     406.0
..        ...    ...       ...       ...
164        60    105       140     290.8
165        60    110       145     300.0
166        60    115       145     310.2
167        75    120       150     320.4
168        75    125       150     330.4

[169 rows x 4 columns]
```

## Checking the pandas version

```
>>> import pandas as pd
>>> print(pd.__version__)
1.3.5
>>>
```

## Pandas Data Frames

• A Pandas DataFrame is a 2 dimensional data structure, like a 2 dimensional array, or a table with rows and columns.

• Create a simple Panda Data Frame

```
>>> import pandas as pd
>>> data={
...    "calories":[100,200,300],
...    "duration":[10,20,30]
... }
>>> df=pd.DataFrame(data)
>>> print(df)
   calories   duration
0      100        10
1      200        20
2      300        30
>>> _
```

## Load the CSV file into data Frame

```
>>> import pandas as pd
>>> df=pd.read_csv('D:\\3_Educational\\1_EngineeringSubjects\\Python for Data Sc
ience\\1_Panda\\1_data.csv')
>>> print(df)
     Duration  Pulse  Maxpulse  Calories
0          60    110       130     409.1
1          60    117       145     479.0
2          60    103       135     340.0
3          45    109       175     282.4
4          45    117       148     406.0
..        ...    ...       ...       ...
164        60    105       140     290.0
165        60    110       145     300.0
166        60    115       145     310.2
167        75    120       150     320.4
168        75    125       150     330.4

[169 rows x 4 columns]
>>>
```

## Data Cleaning

• Data cleaning means fixing bad data in your data set.

• Bad data could be:
  • Empty cells
  • Data in wrong format
  • Wrong data
  • Duplicates

The data set contains some empty cells ("Date" in row 22, and "Calories" in row 18 and 28).

The data set contains wrong format ("Date" in row 26).

The data set contains wrong data ("Duration" in row 7).



The data set contains duplicates (row 11 and 12).

# 1. Remove Rows

- One way to deal with empty cells is to remove rows that contain empty cells.

- This is usually OK, since data sets can be very big, and removing a few rows will not have a big impact on the result.

- See Row 17 and 27 (removed)

Pandas *dropna()* method allows the user to analyze and drop Rows/Columns with Null values

By default, the dropna() method returns a new DataFrame, and will not change the original.



By default, the dropna() method returns a new DataFrame, and will not change the original.

If you want to change the original DataFrame, use the inplace = True argument.



See Row 17 replaced with 130

The fillna() method allows us to replace empty cells with a value.

It will Replace NULL values with the number 130.

## Slide 25

```
C:\Users\user>python
Python 3.7.4 (tags/v3.7.4:e09359112e, Jul  8 2019, 20:34:20) [MSC v.1916 64 bit
(AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import pandas as pd
>>> df = pd.read_csv("D:\\3_Educational\\1_EngineeringSubjects\\Python for Data
>>> df["Calories"].fillna(130, inplace= True)
>>> print(df.to_string())
```

**Values are replaced at position 17, 27, 91, 118, and 141 in the Calories column only.**

---

## Slide 26

### 5. Replace Using Mean, Median, or Mode

- A common way to replace empty cells, is to calculate the mean, median or mode value of the column.

- Mean → Average

- Median → Center value

- Mode → Most common occurring value

By: Prof. (Dr.) Vineet Mohan

---

## Slide 27

```
>>> import pandas as pd
>>> df = pd.read_csv("D:\\3_Educational\\1_EngineeringSubjects\\Python for Data
Science\\1_Panda\\d_data.csv")
>>> x = df["Calories"].mean()
>>> df["Calories"].fillna(x, inplace = True)
>>> print(df.to_string())
```

**Empty Values are replaced with mean at position 17, 27, 91, 118, and 141 in the Calories column only.**

**Mean here is 375.790244**

---

## Slide 28

```
>>> import pandas as pd
>>> df = pd.read_csv("D:\\3_Educational\\1_EngineeringSubjects\\Python for Data
Science\\1_Panda\\d_data.csv")
>>> x = df["Calories"].median()
>>> df["Calories"].fillna(x, inplace = True)
>>> print(df.to_string())
```

**Empty Values are replaced with median at position 17, 27, 91, 118, and 141 in the Calories column only.**

**Median here is 318.6**

---

## Slide 29

```
>>> import pandas as pd
>>> df = pd.read_csv("D:\\3_Educational\\1_EngineeringSubjects\\Python for Data
Science\\1_Panda\\d_data.csv")
>>> x = df["Calories"].mode()[0]
>>> df["Calories"].fillna(x, inplace = True)
>>> print(df.to_string())
```

**Empty Values are replaced with mode at position 17, 27, 91, 118, and 141 in the Calories column only.**

**Mode here is 300.0**

---

## Slide 30

### Wrong Data

- "Wrong data" does not have to be "empty cells" or "wrong format", it can just be wrong, like if someone registered "199" instead of "1.99".

- Sometimes you can spot wrong data by looking at the data set, because you have an expectation of what it should be.

- If you take a look at our data set, you can see that in row 7, the duration is 450, but for all the other rows the duration is between 30 and 60.

By: Prof. (Dr.) Vineet Mohan

One way to fix wrong values is to replace them with something else.

In our example, it is most likely a typo, and the value should be "45" instead of "450", and we could just insert "45" in row 7:

## For Larger Data

- For small data sets you might be able to replace the wrong data one by one, but not for big data sets.

- To replace wrong data for larger data sets you can create some rules, e.g. set some boundaries for legal values, and replace any values that are outside of the boundaries.

By: Prof. (Dr.) Vineet Mehan



## Removing Rows

- Another way of handling wrong data is to remove the rows that contains wrong data.

- This way you do not have to find out what to replace them with, and there is a good chance you do not need them to do your analyses.

- Value at position no 7 is removed

By: Prof. (Dr.) Vineet Mehan

## Duplicate Data

- Duplicate rows are rows that have been registered more than one time.

- By taking a look at our test data set, we can assume that row 11 and 12 are duplicates.

- To discover duplicates, we can use the duplicated() method.

- The duplicated() method returns a Boolean values for each row.

By: Prof. (Dr.) Vineet Mehan                                        37



Above program Returns True for every row that is a duplicate, otherwise False

38



The duplicate row (row no 12) is now removed

39

## Summary

- Methods of Sourcing Data

- Methods of Cleaning Data

40

## Task

- Applying various methods that are used for sourcing the data by taking a suitable arrays\datasets etc. (BT-Level3)

- Design a model that is used to clean Empty cells, Data in wrong format, Wrong data, and Duplicates. (BT-Level6)

By: Prof. (Dr.) Vineet Mehan                                        41

## REFERENCES

- https://www.javatpoint.com/machine-learning

- https://www.tutorialspoint.com/machine_learning/index.htm

- https://www.w3schools.com/python/

42

# THANK YOU

For queries
Email: vineet.e13038@cumail.in