

**Common Pool of Generic Electives (GE) Courses
Offered by Department of Computer Sciences
Category-IV**

GENERIC ELECTIVES (GE-2a): Data Analysis and Visualization

Credit distribution, Eligibility and Pre-requisites of the Course

Course title & Code	Credits	Credit distribution of the course			Eligibility criteria	Pre- requisite of the course
		Lecture	Tutorial	Practical/ Practice		
GE2a Data Analysis and Visualization using Python	4	3	0	1	Class XII pass with Mathematics	knowledge of Python

Learning Objectives

This course is designed to introduce the students to real-world data analysis problems, their analysis and interpretation of results in the field of exploratory data science using Python.

Learning outcomes

On successful completion of the course, students will be able to:

- Apply descriptive statistics to obtain a deterministic view of data
- Apply basic and advanced level statistical function on data
- Perform data handling using Numpy arrays
- Do data cleaning and transformation before extracting useful information
- Visualize data for ease of understanding the revealed information

SYLLABUS OF GE-2a

UNIT – I & II (09 Hours)

Introduction to basic statistics and analysis: Fundamentals of Data Analysis, Statistical foundations for Data Analysis, Types of data, Descriptive Statistics, Python Libraries: NumPy, Pandas, Matplotlib

Array manipulation using NumPy: NumPy array: Creating NumPy arrays, various data types of NumPy arrays

UNIT – I & II (09 Hours)

Introduction to basic statistics and analysis: contd..

Correlation and covariance, Linear Regression, Statistical Hypothesis Generation and Testing

Unit 2 Array manipulation using Numpy: contd..

Indexing and slicing, swapping axes, transposing arrays, data processing using Numpy arrays

UNIT – III (15 Hours)

Data Manipulation using Pandas: Data Structures in Pandas: Series, Data Frame, Index objects, loading data into Panda's data frame, Working with Data Frames: Arithmetics, Statistics, Binning, Indexing, Reindexing, Filtering, Handling missing data, Hierarchical indexing, Data wrangling: Data cleaning, transforming, merging and reshaping

UNIT – IV (12 Hours)

Plotting and Visualization: Using Matplotlib to plot data: figures, subplots, markings, color and line styles, labels and legends, Plotting functions in Pandas: Lines, bar, Scatter plots, histograms, stacked bars, Heatmap

Practical component (if any) – 30 Hours

Use data set of your choice from Open Data Portal ([https:// data.gov.in/](https://data.gov.in/), UCI repository) or load from scikit, seaborn library for the following exercises to practice the concepts learnt.

1. Load a Pandas data frame with a selected dataset. Identify and count the missing values in a data frame. Clean the data after removing noise as follows
 - a. Drop duplicate rows.
 - b. Detect the outliers and remove the rows having outliers
 - c. Identify the most correlated positively correlated attributes and negatively correlated attributes
2. Import iris data using sklearn library or (Download IRIS data from: <https://archive.ics.uci.edu/ml/datasets/iris> or import it from sklearn.datasets)
 - a. Compute mean, mode, median, standard deviation, confidence interval and standard error for each feature
 - b. Compute correlation coefficients between each pair of features and plot heatmap
 - c. Find covariance between length of sepal and petal
 - d. Build contingency table for class feature
3. Load Titanic data from sklearn library , plot the following with proper legend and axis labels:
 - a. Plot bar chart to show the frequency of survivors and non-survivors for male and female passengers separately
 - b. Draw a scatter plot for any two selected features
 - c. Compare density distribution for features age and passenger fare

- d. Use a pair plot to show pairwise bivariate distribution
4. Using Titanic dataset, do the following
 - a. Find total number of passengers with age less than 30
 - b. Find total fare paid by passengers of first class
 - c. Compare number of survivors of each passenger class

Project students are encouraged to work on a good dataset in consultation with their faculty and apply the concepts learned in the course.

Essential/recommended readings

1. McKinney W. *Python for Data Analysis: Data Wrangling with Pandas, NumPy and IPython*. 2nd edition, O'Reilly Media, 2018.
2. Molin S. *Hands-On Data Analysis with Pandas*, Packt Publishing, 2019.
3. Gupta S.C., Kapoor V.K., *Fundamentals of Mathematical Statistics*, Sultan Chand & Sons, 2020.

Suggestive readings

- (i) Chen D. Y, *Pandas for Everyone: Python Data Analysis*, Pearson, 2018.
- (ii) Miller J.D. *Statistics for Data Science*, Packt Publishing, 2017.

GENERIC ELECTIVES (GE-2b): Data Analysis and Visualization using Spreadsheet

Credit distribution, Eligibility and Pre-requisites of the Course

Course title & Code	Credits	Credit distribution of the course			Eligibility criteria	Pre-requisite of the course
		Lecture	Tutorial	Practical/ Practice		
GE2b Data Analysis and Visualization using Spreadsheet	4	3	0	1	Class XII pass with Mathematics	Nil

Learning Objectives