

B.M.S COLLEGE OF ENGINEERING BENGALURU
Autonomous Institute, Affiliated to VTU



A Technical Seminar Report
“Using Python for research”

Submitted in partial fulfillment for the award of degree of

Bachelor of Engineering
in
Computer Science and Engineering

Submitted by:
Nithin B S
(1BM20CS100)

Work carried out at



Internal Guide

Prof. M Lakshmi Neelima
Assistant Professor

Department of Computer Science and Engineering
B.M.S College of Engineering
Bull Temple Road, Basavanagudi, Bangalore 560 019
2022-2023

B.M.S COLLEGE OF ENGINEERING
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



DECLARATION

I, NITHIN B S (1BM20CS100) student of 6th Semester, B.E, Department of Computer Science and Engineering, B.M.S College of Engineering, Bangalore, hereby declare that, this technical seminar entitled “**Using Python for research**” has been carried out under the guidance of Prof. M Lakshmi Neelima, Assistant Professor, Department of CSE, BMS College of Engineering, Bangalore during the academic semester March - July 2023. I also declare that to the best of our knowledge and belief, the technical seminar report is not from part of any other report by any other students.

Signature of the Candidate

NITHIN B S (1BM20CS100)

BMS COLLEGE OF ENGINEERING
DEPARTMENT OF COMPUTER SCIENCE AND
ENGINEERING



CERTIFICATE

This is to certify that the Technical Seminar titled “**Using Python for research**” has been carried out by NITHIN B S (1BM20CS100) during the academic year 2022-2023.

Signature of the Guide

Signature of the Head of the
Department

Signature of Examiners with date

1. Internal Examiner

2. External Examiner

Abstract

Course Title: HarvardX: Using Python for Research

MOOC Platform: edX

Course Duration: 12 weeks

Completion Date: 01-06-2021

The course "HarvardX: Using Python for Research" is a comprehensive online program designed to equip learners with the necessary skills and knowledge to leverage Python programming for research purposes. Python has emerged as a popular language in the research community due to its versatility, ease of use, and rich ecosystem of libraries. This course offers a hands-on learning experience, enabling participants to explore Python's capabilities for data manipulation, analysis, visualization, statistical modelling, and machine learning.

The course begins by introducing the basics of Python 3, covering essential concepts such as syntax, data types, variables, control structures, and functions. Participants gain a solid foundation in Python programming, enabling them to effectively write code and understand the core principles of the language.

Throughout the course, learners delve into key Python libraries widely used in research, including NumPy and Pandas. They discover how to manipulate and analyze data efficiently, perform complex operations, and explore datasets. The curriculum emphasizes the practical application of these libraries, enabling participants to gain insights and draw meaningful conclusions from real-world datasets.

Furthermore, the course delves into the realm of machine learning, providing an overview of concepts and algorithms. Participants learn how to apply supervised and unsupervised learning techniques using the sci-kit-learn library. They gain hands-on experience in tasks such as classification, regression, and clustering, enabling them to apply machine learning to their research projects.

The course incorporates case studies from diverse research domains, allowing participants to see Python in action in real-world scenarios. These case studies provide valuable insights into applying Python for research purposes, addressing challenges, and discovering best practices. By the end of the course, learners complete a final project that integrates the knowledge and skills acquired throughout the program. They apply Python to solve a research problem, analyze a dataset, or build a predictive model. The final project showcases their ability to utilize Python effectively in research settings.

Chapter 1

Introduction To The Course:

1.1 Overview

The course "HarvardX: Using Python for Research" is an immersive online program designed to equip participants with the necessary skills to leverage Python programming for research purposes. Python has gained popularity in the research community due to its versatility, simplicity, and robust set of libraries. This course provides a comprehensive exploration of Python's capabilities for data manipulation, analysis, visualization, statistical modelling, and machine learning.

1.2. Objectives

The number one goal of the course is to equip students with the knowledge and realistic abilities required to recognize, implement, and evaluate machine-learning algorithms.

By the end of the path, members will be able to:

- Develop proficiency in Python programming for research applications.
- Gain practical experience in data manipulation, analysis, and visualization using Python libraries.
- Learn statistical modelling techniques and apply them to research datasets.
- Understand the basics of machine learning and apply machine learning algorithms using Python.
- Complete a final project demonstrating the application of Python in a research context.

1.3 Topics Covered/Learned

Throughout the route, contributors will delve into the following key subjects:

Week 0: Introductions and Self-Assessment

- Introduction to the course structure, objectives, and learning resources.
- Self-assessment to gauge the participants' prior knowledge and experience with Python and research concepts.
- Interactions with fellow learners and instructors to establish a supportive learning community.

Week 1: Basics of Python 3

- Introduction to Python 3 and its syntax, data types, variables, and control structures.
- Understanding functions and how to define and use them in Python.
- Hands-on exercises to practice Python programming fundamentals and strengthen coding skills.

Week 2: Python Libraries and Concepts Used in Research

- Exploring essential Python libraries commonly used in research, such as NumPy and Pandas.
- Understanding the functionalities of these libraries for efficient data manipulation, analysis, and exploration.
- Introduction to Jupyter Notebooks as a powerful tool for interactive coding and data analysis.

Week 3: Case Studies Part 1

- Engaging with real-world case studies that showcase the application of Python in research across different domains.
- Analyzing and manipulating datasets using Python libraries like Pandas and NumPy.
- Implementing data visualization techniques using Matplotlib and Seaborn to gain insights from the data.

Week 4: Case Studies Part 2

- Continuation of case studies, exploring advanced data analysis and visualization techniques.
- Applying statistical concepts and methods to perform hypothesis testing and draw meaningful conclusions.
- Working with larger datasets and addressing challenges encountered during the analysis process.

Week 5: Statistical Learning

- Introduction to statistical learning and its applications in research.
- Understanding the principles of supervised and unsupervised learning algorithms.
- Implementing machine learning techniques using the scikit-learn library for tasks like classification, regression, and clustering.

Chapter 2

METHODOLOGY/TECHNIQUES OR ALGORITHM:

In the course "HarvardX: Using Python for Research," participants are exposed to various methodologies, techniques, and algorithms commonly used in research and data analysis. Here are some of the key methodologies and techniques covered:

1. Data Manipulation and Analysis: -

- Utilizing NumPy arrays for efficient data manipulation and mathematical operations. -
- Applying Pandas library for handling datasets, filtering data, and performing aggregations.
- Employing data transformation techniques such as reshaping, merging, and grouping.

2. Data Visualization: -

- Creating visualizations using Matplotlib and Seaborn libraries.
- Techniques for customizing plots, adding labels, legends, and titles.
- Exploring different plot types, such as line plots, scatter plots, bar plots, and histograms.

3. Statistical Analysis: -

- Understanding statistical concepts such as probability distributions, hypothesis testing, and confidence intervals.
- Implementing statistical tests using libraries like SciPy and statsmodels.
- Exploring descriptive statistics and statistical measures to analyze data.

4. Machine Learning: -

- Introduction to supervised learning algorithms like linear regression, logistic regression, and decision trees.
- Understanding unsupervised learning techniques like clustering and dimensionality reduction.
- Evaluation of machine learning models using metrics like accuracy, precision, recall, and F1 score.

5. Research Project Development: -

- Developing a research project from conceptualization to implementation.
- Applying data manipulation, analysis, visualization, and statistical modeling techniques to solve research problems. Integrating machine learning algorithms to gain insights from research data.

Chapter 3

DESCRIPTION OF TOOLS SELECTED:

In the course "HarvardX: Using Python for Research," several tools are selected to facilitate data manipulation, analysis, visualization, and machine learning tasks. Here are some of the key tools covered in the course:

1. Python:

Python is the primary programming language used throughout the course. It is renowned for its simplicity, versatility, and extensive range of libraries that make it a powerful tool for research and data analysis.

2. Jupyter Notebooks:

Jupyter Notebooks are interactive computing environments that enable participants to write and execute Python code in a browser-based interface. Jupyter Notebooks facilitate the creation of reproducible research workflows by combining code, visualizations, and documentation in a single document.

3. NumPy:

NumPy is a fundamental library for numerical computing in Python. It provides support for multidimensional arrays, efficient mathematical operations, and a wide range of numerical algorithms. NumPy is widely used for data manipulation, transformation, and mathematical computations in research applications.

4. Pandas:

Pandas is a popular Python library for data manipulation and analysis. It offers high-performance data structures, such as DataFrame, and functions for reading, cleaning, filtering, aggregating, and transforming data. Pandas simplifies data preprocessing and exploratory data analysis tasks.

5. Matplotlib:

Matplotlib is a plotting library that enables the creation of a wide variety of static, animated, and interactive visualizations in Python. It provides flexible and customizable plotting functions for line plots, scatter plots, bar plots, histograms, and more. Matplotlib is extensively used for data visualization in research.

6. Seaborn:

Seaborn is a statistical data visualization library built on top of Matplotlib. It provides a higher-level interface for creating aesthetically pleasing and informative statistical graphics. Seaborn simplifies the creation of complex visualizations, including heatmaps, distribution plots, and categorical plots.

Chapter 4

DETAILED DESCRIPTION OF MODULES IMPLEMENTED

Module 1: Basics of Python 3

In this module, participants are introduced to the fundamentals of Python programming. They learn about Python syntax, variables, data types, control structures, and functions. The module also covers the basics of Jupyter Notebooks, an interactive environment commonly used in data analysis and research.

Module 2: Python Libraries and Concepts Used in Research

This module focuses on the essential Python libraries used in research, namely NumPy and Pandas. Participants learn how to manipulate and analyze data efficiently using NumPy arrays and Pandas DataFrames. They explore various data manipulation techniques, such as indexing, filtering, aggregating, and reshaping data.

Module 3: Case Studies Part 1

In this module, participants dive into real-world case studies from different research domains. They apply the concepts and techniques learned in the previous modules to analyze and gain insights from these datasets. The case studies cover data cleaning, exploratory data analysis, and visualization using tools like NumPy, Pandas, and Matplotlib.

Module 4: Case Studies Part 2

Building upon the previous module, participants continue with more advanced case studies. They learn additional data manipulation techniques, such as merging datasets and handling missing values. The module also introduces statistical analysis using libraries like SciPy and statsmodels to perform hypothesis testing and calculate confidence intervals.

Module 5: Statistical Learning

This module focuses on statistical learning techniques using Python. Participants learn about supervised and unsupervised learning algorithms, model evaluation, and overfitting. They explore concepts such as linear regression, logistic regression, decision trees, clustering, and dimensionality reduction. The scikit-learn library is utilized for implementing these machine learning algorithms.

Final Project and Course Wrap-Up

In the final module, participants work on a comprehensive final project that integrates the knowledge and skills acquired throughout the course. They apply Python for research purposes, combining data manipulation, analysis, visualization, and statistical modeling techniques to solve a specific research problem or analyze a dataset. The final project serves as a culmination of the course, allowing participants to showcase their proficiency in using Python for research.

Chapter 5

NEW LEARNINGS FROM THE COURSE

In the course "HarvardX: Using Python for Research" we can expect to gain several new learnings and insights. Here are some of the key takeaways from the course:

1. Proficiency in Python for Research:

Participants develop a strong foundation in Python programming specifically tailored for research purposes. They learn the syntax, data types, control structures, and functions necessary to manipulate and analyze data effectively.

2. Data Manipulation and Analysis Techniques:

Participants acquire practical skills in data manipulation and analysis using Python libraries like NumPy and Pandas. They learn techniques such as filtering, aggregating, reshaping, merging datasets, and handling missing values to prepare data for further analysis.

3. Data Visualization:

Participants learn how to create visually appealing and informative plots using libraries like Matplotlib and Seaborn. They explore different types of plots and customize them to effectively communicate research findings.

4. Statistical Analysis:

Participants gain an understanding of statistical concepts and learn how to apply statistical techniques using libraries like SciPy and statsmodels. They become familiar with hypothesis testing, confidence intervals, and descriptive statistics to derive meaningful insights from data.

5. Machine Learning Fundamentals:

Participants are introduced to the basics of machine learning and gain exposure to popular algorithms such as linear regression, logistic regression, decision trees, clustering, and dimensionality reduction. They learn how to evaluate and interpret machine learning models.

6. Practical Application of Python in Research:

Through real-world case studies and a final project, participants learn to apply Python programming and the techniques learned in the course to solve research problems. They develop the ability to integrate data manipulation, analysis, visualization, and statistical modeling to derive insights and make informed decisions.

7. Collaborative Learning and Community Interaction:

Participants have the opportunity to engage with fellow learners through discussion forums, promoting collaboration, knowledge sharing, and networking.

REFERENCES

- [1] McKinney, W. (2017). Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython. O'Reilly Media.
- [2] VanderPlas, J. (2016). Python Data Science Handbook: Essential Tools for Working with Data. O'Reilly Media.
- [3] Grus, J. (2015). Data Science from Scratch: First Principles with Python. O'Reilly Media.
- [4] Millman, K. J., & Aivazis, M. (2011). Python for Scientists and Engineers. Computing in Science & Engineering, 13(2), 9-12.
- [5] Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. Computing in Science & Engineering, 9(3), 90-95.