



# University of Central Punjab

## Faculty of Information Technology

<b>Course Code</b>	CSDS4423
<b>Course Title</b>	Programming for Big Data
<b>Credit Hours</b>	3
<b>Prerequisites by Course(s) and Topics</b>	CSCP2023 & CSCP2021
<b>Assessment Instruments with Weights (homework, quizzes, midterms, final, programming assignments, lab work, etc.)</b>	Quizzes (10%) + Assignments + Tutorials (10%), Mid Term (20%), Final Term (45%), Project/Presentations (10%) and Class Participation (5%)
<b>Semester</b>	SPRING 2021
<b>Course Instructor</b>	Saeed Iqbal Khattak
<b>Course Coordinator</b>	Dr. Adnan N. Qureshi
<b>Lab Instructor</b>	
<b>Office Hours</b>	To be announced later.
<b>Plagiarism Policy</b>	If anybody found in act of plagiarism, he/she will be marked zero in all of his/her instruments of that category. Plagiarism offense in midterm and final term will result in (F) grade. Marks will be uploaded on portal and can be contested within a week or would be considered final.
<b>Current Catalog Description</b>	The course aims to introduce basic concepts that will help students to: <ol style="list-style-type: none"> <li>1. Understanding of Python</li> <li>2. Data manipulation using python</li> <li>3. Hands on practice with python libraries (Numpy, Pandas and Matplot).</li> <li>4. Data Visualization</li> <li>5. Introduction Big Data</li> <li>6. Introduction to Spark</li> <li>7. Introduction to RDDs</li> <li>8. Introduction to DataFrames</li> <li>9. Advanced Spark Topics</li> <li>10. Introduction to Spark MLlib</li> </ol>
<b>Textbook (or Laboratory Manual for Laboratory Courses)</b>	<ol style="list-style-type: none"> <li>1. Mining of Massive Datasets</li> <li>2. Data Analysis with open source tools</li> <li>3. Learning Apache Spark with Python</li> <li>4. Python for Data Analysis</li> <li>5. Python for Programmers</li> <li>6. Big Data, Mining, and Analytics</li> </ol>
<b>Reference Material</b>	<a href="https://saeediqbal.github.com">https://saeediqbal.github.com</a>
<b>Course Goals</b>	This course is for students who have some programming and database experience. The objective of this course is to give students some experience in data analysis and developing applications that utilize the vast amount of data that is available to general public to create programs that provides information used in improving the standard of application performance. Discovering how the efficiency of applications can be improved by understanding the data.

Learning Outcomes	<div>1. Learn basic concepts required for big data applications</div> <div>2. Design and develop algorithms to collect and present data into an information.</div> <div>3. Develop algorithms to display real-time content using Distributed Programming.</div> <div>4. Implement interactive real-time distributed framework.</div> <div>5. Configured advanced real-time applications frameworks.</div> <div>6. Learn, implement and integration of web services with big data applications.</div>											
Topics Covered in the Course, with Number of Lectures on Each Topic (assume 15-week instruction and 1.5 hour lectures)	Attached											
Programming Assignments Done in the Course	Yes											
Exam:	<div>Mid : 1.5Hrs</div> <div>Final : 3Hrs</div>											
Class Time Spent on (in credit hours)	<table><tr><th>Theory</th><th>Problem Analysis</th><th>Solution Design</th><th>Social and Ethical Issues</th></tr><tr><td>0.5</td><td>1</td><td>1</td><td>0.5</td></tr></table>				Theory	Problem Analysis	Solution Design	Social and Ethical Issues	0.5	1	1	0.5
Theory	Problem Analysis	Solution Design	Social and Ethical Issues									
0.5	1	1	0.5									
Oral and Written Communications												

Week	Topics Covered	Instruments
1	<ul style="list-style-type: none"> <li>• Discussion on Python and its market position.</li> <li>• Motivation regarding learning aspects of this course</li> <li>• Setting up environment for Python.</li> <li>• Installation of Anaconda</li> <li>• What is Data?</li> <li>• What is Big Data?</li> <li>• Characteristics of Big Data</li> <li>• What are the Vs of Big Data?</li> <li>• The Impact of Big Data</li> <li>• Big Data - Beyond the Hype, Big Data Examples, Sources of Big Data</li> <li>• Big Data Adoption, The Big Data and Data Science</li> <li>• The Big Data Platform, Big Data and Data Science. Skills for Data Scientists</li> </ul>	
2	<ul style="list-style-type: none"> <li>• Machine-Generated Data: People Generated Data , Organization generated data</li> <li>• Characteristics of Big Data types , volume , velocity , variety, veracity ,value</li> <li>• Building a Big Data Strategy , Component of big data</li> <li>• Types of IDE(s) and WIDE that will be used in the duration of this course. e.g. Spyder, Jupyter etc</li> <li>• Hello World Program "Print Command"</li> <li>• Keyword Types</li> <li>• Expressions and Variables</li> </ul>	Assignment 1
3	<ul style="list-style-type: none"> <li>• Input Method</li> <li>• Conditions and Branching</li> <li>• Loops</li> <li>• String Operations</li> <li>• Lists and Tuples</li> <li>• Sets</li> <li>• Dictionaries</li> </ul>	Quiz 1
4	<ul style="list-style-type: none"> <li>• Data Analysis Process</li> <li>• Steps of processes: acquiring, exploring, pre-processing, Analyzing, communicating and turning into action</li> <li>• What is a Distributed File System?</li> <li>• Scale-able Computing over the Internet, Programming Models for Big Data</li> <li>• Reading and Writing files</li> <li>• Functions</li> <li>• Objects and Classes</li> </ul>	
5	<ul style="list-style-type: none"> <li>• Working with Pandas</li> <li>• Descriptive Statistics with Pandas</li> <li>• Group by with Python</li> <li>• Data Manipulation with Pandas</li> </ul>	Assignment 2
6	<ul style="list-style-type: none"> <li>• Data Wrangling with Pandas</li> <li>• Data Manipulation with Pandas</li> </ul>	Quiz 2
7	<ul style="list-style-type: none"> <li>• Introduction with Numpy</li> <li>• Numpy one dimensional Array</li> <li>• Numpy two dimensional Array</li> <li>• Numpy Array Operations</li> </ul>	Quiz 3

8	<ul style="list-style-type: none"> <li>• Introduction to Matplotlib</li> <li>• Basic Plotting with Matplotlib</li> <li>• Line Plots</li> <li>• Area Plots</li> <li>• Histograms</li> <li>• Bar Charts</li> <li>• Pie Charts</li> <li>• Box Plots</li> <li>• Scatter Plots</li> <li>• Word Cloud</li> </ul>	
	<b>Revision</b>	
	<b>MID TERM</b>	
9	<ul style="list-style-type: none"> <li>• Introduction to Hadoop</li> <li>• Hadoop: Why, Where and Who?</li> <li>• The Hadoop Ecosystem: Welcome to the zoo!</li> <li>• The Hadoop Distributed File System: A Storage System for Big Data</li> <li>• YARN: A Resource Manager for Hadoop</li> </ul>	Review
10	<ul style="list-style-type: none"> <li>• The Hadoop Distributed File System: A Storage System for Big Data</li> <li>• YARN: A Resource Manager for Hadoop</li> <li>• What is Spark and what is its purpose?</li> </ul>	Assignment 3
11	<ul style="list-style-type: none"> <li>• Components of the Spark unified stack</li> <li>• Resilient Distributed Dataset (RDD)</li> <li>• What is Pig and Hive</li> <li>• Architecture of Pig and Hive</li> </ul>	Quiz 4
12	<ul style="list-style-type: none"> <li>• Understand how to create parallelized collections and external datasets</li> <li>• Work with Resilient Distributed Dataset (RDD) operations</li> <li>• Utilize shared variables and key-value pairs</li> </ul>	Assignment 4
13	<ul style="list-style-type: none"> <li>• Describe and run some Spark examples</li> <li>• Pass functions to Spark</li> <li>• Create and run a Spark standalone application</li> </ul>	Quiz 5
14	<ul style="list-style-type: none"> <li>• Introduction to Apache Kafka</li> <li>• Components of Apache Kafka</li> <li>• Internal Architecture of Apache Kafka</li> </ul>	Assignment 5 and Quiz 6
15	<ul style="list-style-type: none"> <li>• Introduction to Apache Zookeeper</li> <li>• Components of Apache Zookeeper</li> <li>• Internal Architecture of Apache Zookeeper</li> </ul>	Quiz 7

## Final Term Project

1. Students will create a final project on a topic of their choice using the technologies and techniques covered in this course.
2. Student will create a short (less than 10 minute) presentation explaining how to use their project and where requirement were met.
3. Final project presentations will take place during the week of finals. Exact date TBA.