Government of Pakistan

# National Vocational and Technical Training Commission

## Prime Minister Hunarmand Pakistan Program,
"Skills for All"



## Course Contents/ Lesson Plan
**Course Title:** Big Data Analytics
**Duration:** 6 Months

| Trainer Name | |
|---|---|
| Course Title | **Big Data Analytics** |
| Objective of Course | **Employable skills and hands on practice for Web Development, Graphic Designing and Mobile App Development**<br><br>The main goal of this course is to help students learn, understand, and practice big data analytics and machine learning approaches, which include the study of modern computing big data technologies and scaling up machine learning techniques focusing on industry applications. Mainly the course objectives are: conceptualization and summarization of big data and machine learning, trivial data versus big data, big data computing technologies, machine learning techniques, and scaling up machine learning approaches. |
| **Learning Outcome of the Course** | The students learning outcomes are designed to specify what the students will be able to perform after completion of the course:<br><ul><li>Ability to identify the characteristics of datasets and compare the trivial data and big data for various applications.</li><li>Ability to select and implement machine learning techniques and computing environment that are suitable for the applications under consideration.</li><li>Ability to solve problems associated with batch learning and online learning, and the big data characteristics such as high dimensionality, dynamically growing data and in particular scalability issues.</li><li>Ability to understand and apply scaling up machine learning techniques and associated computing techniques and technologies.</li><li>Ability to recognize and implement various ways of selecting suitable model parameters for different machine learning techniques.</li><li>Ability to integrate machine learning libraries and mathematical and statistical tools with modern technologies like Apache Spark.</li></ul> |
| **Course Execution Plan** | Total Duration of Course: **6 Months (26 Weeks)**<br>Class Hours: **5 Hours per day**<br>Theory:  **20% Practical: 80%**<br>Weekly Hours: **25 Hours Per week**<br>Total Contact Hours: **650 Hours** |

Plot no. 38, Kirthar Road, H-9 Islamabad
051-9044250

| | |
|---|---|
| **Companies Offering Jobs in the respective trade** | 1. Upwork<br>2. Freelancer<br>3. Fiverr<br>4. Government Institutes<br>5. Software Houses<br>6. Companies all over the world are offering its jobs as they want to know the trends of market |
| **Job Opportunities** | Upskilling in Big Data and Analytics field is a smart career decision. According to Allied Market Research, the globalmarket of only Hadoop/Spark will reach $84.6 Billion by 2021 and there is a shortage of 1.4-1.9 million Hadoop/Spark data analysts in the U.S. alone. Here is selection of specialist opportunities in your area:<br>• Big Data Architect (Average Salary: 124000$ / Annum)<br>• Big Data Engineer (Average Salary: 117000$ / Annum)<br>• Big Data Developer (Average Salary: 88500$ / Annum) |
| **No of Students** | 25 |
| **Learning Place** | Classroom / Lab |
| **Instructional Resources** | **Development Platform:**<br>• https://github.com/ ,<br>• https://spark.apache.org/,<br>• https://www.edureka.co/apache-spark-scala-certification-training,<br>• https://www.youtube.com/watch?v=iP1wOSsKjW8&list=PLS1Qul Wo1RIahlYDqHWZb81qsKgEvPiHn,<br>• https://stackoverflow.com/<br><br>**Learning Material:**<br>• https://spark.apache.org/docs/latest/api/python/index.htmlhttps ://www.youtube.com/watch?v=9mELEARcxJo&list=PL9ooVrP1hQ OGyFc60sExNX1qBWJyV5IMb<br>• https://www.youtube.com/watch?v=Uct_EbThV1E&list=PLZ7s-Z1aAtmIbaEj_PtUqkqdmI1k7libK<br>• https://www.edureka.co/apache-spark-scala-certification-training<br>• https://www.youtube.com/watch?v=wjfeGxqAQOY&list=PLrjkTql 3jnm-CLxHftqLgkrZbM8fUt0vn<br>• https://www.youtube.com/watch?v=iP1wOSsKjW8&list=PLS1Qul Wo1RIahlYDqHWZb81qsKgEvPiHn |

| Scheduled Week | Module Title | Learning Units | Remarks |
|---|---|---|---|
| **Week 1** | ➢ Introduction | • **Motivational Lecture**<br>• **Course Introduction**<br>• **Success stories**<br>• **Job market**<br>• **Course Applications**<br>• **Institute/work ethics**<br>• Discussion on Python and its market position.<br>• Motivation regarding learning aspects of this course<br>• Setting up environment for Python.<br>• Installation of Anaconda<br>• What is Big Data?<br>• Characteristics of Big Data<br>• The Impact of Big Data<br>• Big Data - Beyond the Hype, Big Data Examples, Sources of Big Data<br>• Big Data Adoption, The Big Data and Data Science<br>• The Big Data Platform, Big Data and Data Science. Skills for Data Scientists | |
| **Week 2** | **Module -1**<br><br>**Chapter 1.1-** | **Overview of DBMS**<br>• Components of DBMS<br>• Database Architecture<br>• Types of Database Model<br>• ER Model: Basic Concepts<br>• ER Model: Creating ER Diagram<br>• The Extended ER Model<br>• Codd's 12 rule of RDBMS<br>• Basic Concepts of RDBMS<br>• Types of Database key<br>• Introduction to Normalization<br><br>**Basic SQL**<br><br>• SQL Introduction<br>• Create query<br>• Alter query<br>• Truncate, Drop and Rename query<br>• All DML command<br>• All TCL Command<br>• All DCL Command<br>• WHERE clause<br>• SELECT query<br>• LIKE clause | |

| | | | |
|---|---|---|---|
| | | • ORDER BY clause<br>• Group BY clause<br>• Having clause<br>• DISTINCT keyword<br>• AND & OR operator<br>• DIVISION operator<br><br>**Advanced SQL**<br><br>• SQL Constraints<br>• SQL function<br>• SQL Join<br>• SQL Alias<br>• SQL SET operation<br>• SQL Sequences<br>• SQL Views | |
| **Week 3** | **Chapter 1.2-** | • Types of IDE(s) and IDE that will be used in the duration of this course. e.g. Spyder, Jupyteretc<br>• Hello World Program "Print Command"<br>• Keyword Types<br>• Expressions and Variables<br>• Input Method<br>• Conditions and Branching<br>• Loops | |
| **Week 4** | **Chapter 2.1** | • String Operations<br>• Lists and Tuples<br>• Sets<br>• Dictionaries<br>• Reading and Writing files<br>• Functions<br>• Objects and Classes | |
| **Week 5** | **Chapter 2.2** | • Introduction with Numpy<br>• Numpy one dimensional Array<br>• Numpytwo-dimensional Array<br>• Numpy Array Operations | |
| **Week 6** | **Chapter 3.1** | • Descriptive Statistics<br>• Data Manipulation<br>• Data Wrangling | |
| **Week 7** | **Chapter 3.2** | • Working with Pandas<br>• Descriptive Statistics with Pandas<br>• Group by with Python | |

| | | | |
|---|---|---|---|
| | | • Data Manipulation with Pandas | |
| **Week 8** | **Chapter 4** | • Data Wrangling with Pandas<br>• Discussion regarding exam | |
| **Week 9** | **Chapter 5.1** | • Introduction to Matplotlib<br>• Basic Plotting with Matplotlib<br>• Line Plots<br>• Area Plots<br>• Histograms | |
| **Week 10** | **Chapter 5.2** | • Bar Charts<br>• Pie Charts<br>• Box Plots<br>• Scatter Plots<br>• Word Cloud | |
| **Week 11** | **Chapter 6.1** | • What is Spark and what is its purpose?<br>• Components of the Spark unified stack<br>• Resilient Distributed Dataset (RDD)<br>• Scala and Python overview | |
| **Week 12** | **Chapter 6.2** | • Understand how to create parallelized collections and external datasets<br>• Work with Resilient Distributed Dataset (RDD) operations<br>• Utilize shared variables and key-value pairs | |
| **Week 13** | **Chapter 6.3** | • Describe and run some Spark examples<br>• Pass functions to Spark<br>• Create and run a Spark standalone application | |
| **Week 14** | **Chapter 6.4** | • Understand and use the various Spark libraries | |
| **Week 15** | **Mid-Term Assignment** | | |
| **Week 16** | **Chapter 7**<br>Apache Shark Next-Generation Big Data Framework | • Apache Spark Next-Generation Big Data Framework<br>• History of Spark<br>• Why we should prefer spark?<br>• Introduction to Apache Spark<br>• Components of Spark<br>• Application of In-memory Processing<br>• Hadoop Ecosystem vs Spark<br>• Advantages of Spark<br>• Spark Architecture | |

| | | | |
|---|---|---|---|
| | | • Spark Cluster in Real World<br>• Demo: Running a Scala Programs in Spark Shell<br>• Demo: Setting Up Execution Environment in IDE<br>• Demo: Spark Web UI<br>• Key Takeaways<br>• Knowledge Check<br>• Practice Project: Apache Spark Next-Generation Big Data Framework | |
| **Week 17** | **Chapter 8**<br><br>Spark Core Processing RDD | • Introduction to Spark RDD<br>• RDD in Spark<br>• Creating Spark RDD<br>• Pair RDD<br>• RDD Operations<br>• Demo: Spark Transformation Detailed Exploration Using Scala Examples<br>• Demo: Spark Action Detailed Exploration Using Scala<br>• Caching and Persistence<br>• Storage Levels<br>• Lineage and DAG<br>• Need for DAG<br>• Debugging in Spark<br>• Partitioning in Spark<br>• Scheduling in Spark<br>• Shuffling in Spark<br>• Sort Shuffle<br>• Aggregating Data with Paired RDD<br>• Demo: Spark Application with Data Written Back to HDFS and Spark UI<br>• Demo: Changing Spark Application Parameters<br>• Demo: Handling Different File Formats<br>• Demo: Spark RDD with Real-world Application<br>• Demo: Optimizing Spark Jobs<br>• Key Takeaways<br>• Knowledge Check<br>• Practice Project: Spark Core Processing RDD | |
| **Week 18** | **Chapter 9**<br><br>Spark SQL Processing DataFrames | • Spark SQL Processing DataFrames<br>• Spark SQL Introduction<br>• Spark SQL Architecture<br>• Dataframes<br>• Demo: Handling Various Data Formats | |

| | | | |
|---|---|---|---|
| | | • Demo: Implement Various Dataframe Operations<br>• Demo: UDF and UDAF<br>• Interoperating With RDDs<br>• Demo: Process Dataframe Using SQL Query<br>• RDD vs Dataframe vs Dataset<br>• Practice Project: Processing Dataframes<br>• Key Takeaways<br>• Knowledge Check<br>• Practice Project: Spark SQL - Processing Dataframes | |
| **Week 19** | **Chapter 10.1**<br><br>**Part 1**<br><br>Spark Mlib Modelling BigData with Spark | ● Spark Mlib Modeling Big Data With Spark<br>● Role of Data Scientist and Data Analyst in Big Data<br>● Analytics in Spark<br>● Machine Learning<br>● Supervised Learning<br>● Demo: Classification of Linear SVM<br>● Demo: Linear Regression With Real World Case Studies<br>● Unsupervised Learning Demo: Unsupervised Clustering K-means | |
| **Week 20** | **Chapter 10.2**<br><br>**Part 2**<br><br>Spark Mlib Modelling BigData with Spark | ● Reinforcement Learning<br>● Semi-supervised Learning<br>● Overview of Mlib<br>● Mlib Pipelines<br>● Key Takeaways<br>● Knowledge Check<br>● Practice Project: Spark Mlib - Modelling Big data With Spark | |
| **Week 21** | Employable Project/Assignment (6 weeks i.e. 21-26) in addition of regular classes.<br>**OR**<br>On job training ( 2 weeks) | ● Guidelines to the Trainees for selection of students employable project like final year project (FYP)<br>● Assign Independent project to each Trainee<br>● A project based on trainee's aptitude and acquired skills.<br>● Designed by keeping in view the emerging trends in the local market as well as across the globe.<br>● The project idea may be based on Entrepreneur.<br>● Leading to the successful employment.<br>● The duration of the project will be 6 | |

|  |  | weeks | |
|---|---|---|---|
|  |  | ● Ideas may be generated via different sites such as: https://1000projects.org/ https://nevonprojects.com/ https://www.freestudentprojects.com/ https://technofizi.net/best-computer-science-and-engineering-cse-project-topics-ideas-for-students/ <br>• Final viva/assessment will be conducted on project assignments. <br>• At the end of session the project will be presented in skills competition <br>• The skill competition will be conducted on zonal, regional and National level. <br>• The project will be presented in front of Industrialists for commercialization <br>• The best business idea will be placed in NAVTTC business incubation center for commercialization. <br>-------------------------------------------------------- <br>**OR** <br>**On job training for 2 weeks:** <br>• Aims to provide 2 weeks industrial training to the Trainees as part of overall training program <br>• Ideal for the manufacturing trades <br>• As an alternate to the projects that involve expensive equipment <br>• Focuses on increasing Trainee's motivation, productivity, efficiency and quick learning approach. | |
| **Week 22** | **Chapter 11.1** <br><br>**Part 1** <br><br>Stream Processing Frameworks and Spark Streaming | ● Streaming Overview <br>● Real-time Processing of Big Data <br>● Data Processing Architectures <br>● Demo: Real-time Data Processing <br>● Spark Streaming <br>● Demo: Writing Spark Streaming Application <br>● Introduction to DStreams <br>● Transformations on DStreams <br>● Design Patterns for Using Foreachrdd <br>● State Operations <br>● Windowing Operations <br>● Join Operations Stream-dataset Join <br>● Demo: Windowing of Real-time Data | |

| | | | |
|---|---|---|---|
| | | Processing <br> ● Streaming Sources Demo: Processing Twitter Streaming Data <br> ● Structured Spark Streaming <br> ● Use Case Banking Transactions <br> ● Structured Streaming Architecture Model and Its Components <br> ● Output Sinks | |
| **Week 23** | **Chapter 11.2** <br><br> **Part 2** <br><br> Stream Processing Frameworks and Spark Streaming | ● Structured Streaming APIs <br> ● Constructing Columns in Structured Streaming <br> ● Windowed Operations on Event-time <br> ● Use Cases <br> ● Demo: Streaming Pipeline <br> ● Practice Project: Spark Streaming <br> ● Key Takeaways <br> ● Knowledge Check <br> ● Practice Project: Stream Processing Frameworks and Spark Streaming | |
| **Week 24** | **Chapter 12.1** <br><br> **Part 1** <br><br> Spark GraphX | ● Spark GraphX <br> ● Introduction to Graph <br> ● GraphX in Spark <br> ● GraphX Operators <br> ● Join Operators <br> ● GraphX Parallel System <br> ● Algorithms in Spark | |
| **Week 25** | **Chapter 12.2** <br><br> **Part 2** <br><br> Spark GraphX | ● Pregel API <br> ● Use Case of GraphX <br> ● Demo: GraphX Vertex Predicate <br> ● Demo: Page Rank Algorithm <br> ● Key Takeaways <br> ● Knowledge Check <br> ● Practice Project: Spark GraphX Project Assistance <br> ● **Final Project Assessment** | |
| **Week 26** | Entrepreneurship and Final Assessment in project | ● Job Market Searching <br> ● Self-employment <br> ● Freelancing sites <br> ● Introduction <br> ● Fundamentals of Business Development <br> ● Entrepreneurship <br> ● Startup Funding <br> ● Business Incubation and Acceleration <br> ● Business Value Statement <br> ● Business Model Canvas <br> ● Sales and Marketing Strategies <br> ● How to Reach Customers and Engage CxOs | |

| | | • Stakeholders Power Grid<br>• RACI Model, SWOT Analysis, PEST Analysis<br>• SMART Objectives<br>• OKRs<br>• Cost Management (OPEX, CAPEX, ROCE etc.)<br>• Final Assessment | |

**List of Machinery / Equipment**

| Sr. No | Name of item as per curriculum | Quantity physically available at the training location |
|---|---|---|
| 1 | Computers Minimum Corei5<br><br>• LCD Display 17" with built in speakers | 25 |
| 2 | DSL Internet Connection (Minimum 1 MB) | Available on every PC |
| 3 | **Accessories/Devices**<br><br>• Connectors<br>• Multimedia<br>• Printer (NW printer)<br>• Audio/visual aid<br>• White Board<br>• Pin Board<br>• Flip Chart Board<br>• Hard copy of Training Material<br>• Mobile Phones | 25 each |
| 4 | **Wires, data cables, power plugs, power supply** | For every PC |
| 5 | **UPS** | Available |
| 6 | **Generator / Solar Backup** | Available |
| 7 | **Air Conditioner (2 Tons)** | Available |

Plot no. 38, Kirthar Road, H-9 Islamabad
051-9044250

1. **Software List**

| Sr. No | Software Name |
|--------|---------------|
| 1. | MS Office 2016 (Installed on each PC) |
| 2. | Operating System (Windows, Linux or other Operating Systems) |
| 3. | Programming Languages including NetBeans, Android studio (Licensed |
| 4. | Web Servers including IIS, Apache (Licensed software installed on each PC) |
| 5. | Databases including MySQL, ERWIN (Licensed software installed on each PC) |
| 6. | FTP Client including FileZilla, File Manager (Licensed software installed on each PC) |
| 7. | Web hosting manager/control panel |
| 8. | Web browser including Internet Explorer, Google Chrome, Mozilla Firefox, Netscape, Opera (installed on each PC) |
| 9. | Firewall (each PC) |
| 10. | Security scanning tools including Antivirus (each PC)<br>Networking |
| 11. | Required Software's:<br>• Anaconda Jupyter<br>• MySQL Database<br>• MS Office<br>• MS Visio<br>• MySQL |

2. **Minimum Qualification of Teachers / Instructor**

   The qualification of teachers / instructor of this course should be minimum **of bachelors in Computer science with minimum 3 years of development experience** in relevant trade.
   - Bachelors of Computers Science / Networks (Hons)

3. **Supportive Notes**

**Teaching Learning Material**

| Books Name | Author |
|------------|--------|
| Python Crash Course | Eric Matthes |

| | |
|---|---|
| Big Data Analysis with Python | Ankit Shukla,Ivan Marin and Sarang VK |
| Big Data Course ( Edureka Online Course) | |