

Big Data

| | MON | TUE | WED | THU | FRI | ENVIRONMENT |
|--|---|--|---|--|--|---|
| Week 1 Python | <ul style="list-style-type: none"> Interpreter vs Compiler REPL Identifiers & Keywords Simple/Compound statements Types & Values | <ul style="list-style-type: none"> Namespaces String Literals Operators Lambdas Functions Exceptions | <ul style="list-style-type: none"> Lists Sets Tuples Dictionaries Classes & Objects Variable Scope/lifetime Simple/Compound Statements | <ul style="list-style-type: none"> File Handling: open()/close() Logging Collections Datetime RegEx | <ul style="list-style-type: none"> PyPI & pip Modules Unit Testing Pylint | <ul style="list-style-type: none"> VS Code |
| Week 2 Database-MongoDB | <p>Written evaluation</p> <p>Trainer interviews</p> <p>Quality Control Audit</p> <ul style="list-style-type: none"> Intro to RDBMS DML, DDL, DQL Basic SQL queries Aggregate functions Multiplicity | <ul style="list-style-type: none"> NoSQL MongoDB PyMongo Driver Connection String | <ul style="list-style-type: none"> Mongo Collections Mongo Documents Insert Find | <ul style="list-style-type: none"> Query Sort Limit | <ul style="list-style-type: none"> Update Delete Indexing | <ul style="list-style-type: none"> MongoDB |
| Week 3 UNIX/Hadoop Fundamentals | <p>Written evaluation</p> <p>Trainer interviews</p> <p>Quality Control Audit</p> <ul style="list-style-type: none"> Intro to Open Source Software Linux commands brainstorm Root [/] vs Home [~] Commands: mkdir, rm, cp, mv, cd, ls, cat, grep, echo | <ul style="list-style-type: none"> Commands: df, fdisk, sfdisk, cfdisk, lsblk, blkid, mdadm File Editors - vim, nano Intro to SSH (credentials/private key) Intro to data evolution Intro to big data | <ul style="list-style-type: none"> Hadoop Ecosystem Introduction Intro to HDFS Evolution of Hadoop HDFS Commands | <ul style="list-style-type: none"> Introduction to MapReduce Mapper/Intermediate/Reducer phases Partitioners Combiners | <ul style="list-style-type: none"> YARN Overview InputFormats | |
| Week 4 AVRO/Sqoop/Hive | <p>Written Evaluation</p> <p>Trainer interviews</p> <p>Quality Control Audit</p> <ul style="list-style-type: none"> Introduction to AVRO File Format Schema declaration in AVRO | <ul style="list-style-type: none"> Primitive and complex types AVRO and MapReduce | <ul style="list-style-type: none"> Importing/exporting data using Sqoop Introduction to Sqoop- basic commands | <ul style="list-style-type: none"> Introduction to Hive Basic Hive Queries Hive Commands | <ul style="list-style-type: none"> Data Types Managed vs External tables Partitions and Buckets | |

| | MON | TUE | WED | THU | FRI | ENVIRONMENT |
|--|---|---|--|---|---|---|
| Week 5 Pig/Spark Fundamentals | <p>Written Evaluation</p> <p>Trainer Interviews</p> <p>Quality Control Audit</p> <ul style="list-style-type: none"> Introduction to Pig Latin Commands Datatypes | <ul style="list-style-type: none"> sorting/filtering Describe/Explain/Illustrate local mode vs MapReduce mode | <ul style="list-style-type: none"> Introduction to Spark Hadoop vs Spark Spark Setup | <ul style="list-style-type: none"> Introduction to RDDs Basic RDD operations Local vs Cluster mode Working with Key/Value pairs | <ul style="list-style-type: none"> Transformations Actions Shared variables pySpark | |
| Week 6 Spark Fundamentals | <p>Quality Control Audit</p> <p>Trainer interviews</p> <p>Written Evaluation</p> <ul style="list-style-type: none"> Accumulators | <ul style="list-style-type: none"> Creating Spark EMR cluster Spark Cluster mode Introduction to YARN | <ul style="list-style-type: none"> Spark Cluster Manager Running Spark job on EMR Driver class configuration Executors | <ul style="list-style-type: none"> Configure number of executors Spark cluster configuration Configure memory: Driver & executors | <ul style="list-style-type: none"> Spark caching Memory management | <ul style="list-style-type: none"> AWS EMR |
| Week 7 Spark SQL/DataFrames | <p>Quality Control Audit</p> <p>Trainer interviews</p> <p>Written Evaluation</p> <ul style="list-style-type: none"> Introduction to Spark SQL Introduction to DataSets | <ul style="list-style-type: none"> Introduction to DataFrames Entry point: SparkSession Creating DataFrames | <ul style="list-style-type: none"> Creating DataSets Working with RDDs Using DataFrame aggregate functions | <ul style="list-style-type: none"> Bucketing Sorting and Partitioning | <ul style="list-style-type: none"> Working with JSON Datasets Working with Parquet Files | |
| Week 8 Streaming/Kafka | <p>Quality Control Audit</p> <p>Trainer interviews</p> <p>Written Evaluation</p> <ul style="list-style-type: none"> Introduction to Streaming Introduction to Kafka Kafka Fundamentals Topics/Brokers/Consumer/Producer | <ul style="list-style-type: none"> Apache Kafka architecture Pub-Sub messaging Creating a Kafka topic Retrieve list of topics | <ul style="list-style-type: none"> Create producer and consumer Send messages from Producer Producer API | <ul style="list-style-type: none"> Introduction to Spark Streaming Spark engine Sending data stream from Kafka to Spark | <ul style="list-style-type: none"> Processing data stream using Spark streaming | <ul style="list-style-type: none"> Kafka |
| Week 9 | <p>Project 3</p> <p>Written evaluation</p> <p>Trainer interviews</p> <p>Quality Control Audit</p> | Project 3 | Project 3 | Project 3 | Project 3 | |
| Week 10 | <p>Project 3</p> <p>QC Audit - Cumulative</p> | Project 3 | Project 3 | Project 3 | Project 3 | |

| | MON | TUE | WED | THU | FRI | ENVIRONMENT |
|---------|-----------|-----------|-----------|-------------------------------|-----|-------------|
| Week 11 | Project 3 | Project 3 | Project 3 | Project 3 Project showcase | | |

| PROJECT | TECHNOLOGIES |
|---|--|
|  Project 3 | Spark, Spark SQL, Kafka, Spark Streaming |