# Big Data Analytics Course Outline

## Session1:  Introduction to Big Data, Hadoop and MapReduce

- Which data is called as Big Data, Applications of Big Data
- Introduction to AWS Academy, connection & basic operations
- Traditional Data warehousing & Big Data
- Introduction to Distributed Computing Environment, Hadoop & Its Ecosystem
- HDFS Architectures, HDFS Commands
- MapReduce

## Session 2: Data Ingestion with Hive Sqoop

- Intro to Hive
- Hive Architecture
- HQL
- Bucketing and partitioning in Hive
- Importing/Exporting data to/from Hive using Sqoop
- File formats -ORC, Parquet, Avro
- Flume* (Optional)

## Session 3: Introduction to NoSQL, MongoDB & PyMongo

- Introduction to NoSQL
- Cap Theorem & NoSQL Database types
- MongoDB & Its Features
- MongoDB & Its Features
- MongoDB collections, documents and operations
- Intro to PyMongo, Install PyMongo & Python Driver, connect to MongoDB
- Perform basic Create, Retrieve, Update and Delete (CRUD) operations using PyMongo
- Sharding* (optional)

*Optional/not part of end semester exam

## Session 4: Apache Spark & Spark SQL

- o  Introduction to Spark
- o  Spark Architecture
- o  PySpark and Data Bricks
- o  Introduction to Spark SQL
- o  Spark SQL as an ETL tool
- o  Spark SQL Performance Tuning

## Session 5: PySpark ML

- o  Intro to Spark ML
- o  Spark ML Pipeline - Transformers, Estimators
- o  Spark ML Component Flow
- o  Spark ML Data Types
- o  Spark ML Algorithms
- o  Building Pipeline
- o  Model Persistence

## Session 6: Stream Processing: Spark-Streaming and Kafka*

- o  What is Streaming Data
- o  Intro to Spark Streaming & Its working
- o  Spark Streaming + Kafka Example
- o  Intro to Kafka & Its working
- o  Kafka Commands

## Session 7: Case Study*

*Optional/not part of end semester exam