



COURSE HANDOUT

Data
Scientist
Program

Course Title:	Data Science
Course Code:	
Trainer:	_____
Commencement Date:	_____
End Date:	_____
Total Session Planned:	_____
Weekly Frequency:	Thrice a week
Moderator Name:	_____
CoE Head:	_____
Faculty Consultation Day	For every five session, there would be a consultation session

SR.NO	MODULE	TOPIC DETAILS	TIME
01	Introduction to Data Science & Machine Learning (ML) & Deep Learning(DL) & Artificial Intelligence(AI)	Introduction to Data Science & Machine Learning (ML) Introduction to Data Science, ML, DL & AI - why is it so important? Applications of Data science across industries <ul style="list-style-type: none"> • Business problems – Analytics scenarios • Analytics Industry in India, Job Market & Top Skills • Data science – CRIS DM Approach and DIPP framework • Data Scientist Toolbox, Tool of choice-Python: what & why? Data Scientist - Tasks and Capabilities	Daily 3/4 Hours
02	SQL	Introduction to SQL & Data Warehouse Concepts <ul style="list-style-type: none"> • Introduction to Data Warehouse • Dimensions & Facts • Normalization & Schemas • Modelling • ETL 	Daily 3/4 Hours
03	R Programming	Introduction to R Programming <ul style="list-style-type: none"> • Data Structures • Data Visualization • Statistics for Data Science -1 • Statistics for Data Science -2 • Regression Analysis • Classification • Clustering • Association Practice Assignment-1 	Daily 3/4 Hours

04
**Python
programming**
Introduction to Python

- Installation of Python framework and packages: Anaconda and pip
- Writing/Running python programs using Spyder, Command Prompt
- Working with Jupyter Notebooks
- Creating Python variables: Numeric, string and logical operations
- Basic Data containers: Lists, Dictionaries, Tuples & sets

Practice Assignment-2
Operations & Functions in Python

- Writing for loops in Python
- List & Dictionary Comprehension
- While loops and conditional blocks
- List/Dictionary comprehensions with loops
- Writing your own functions in Python
- Writing your own classes and functions as classobjects
- Practice assignment – 2A**
- Numerical Summary of Data**
- Summarizing numeric data and categorical data in pandas
- Group wise summary of mixed data
- Practice assignment – 2B**
- Data Visualization using Python**
- Need for visual summary
- Introduction to Seaborn
- Visual summary of different data combinations
- Practice assignment – 2C**
- Data Handling using NumPy and Pandas**
- Introduction to NumPy arrays, functions & properties
- Introduction to pandas, Data frame functions and properties

**Daily
3/4
Hours**

		<ul style="list-style-type: none">• Reading and writing external data Manipulating Data Columns Practice assignment – 2D Regular expressions Introduction <ul style="list-style-type: none">• Regular expression – Data Preparation	
05	Statistics and Linear Algebra	Basics of Statistics <ul style="list-style-type: none">• Introduction to Univariate Statistics, Shape• Central Tendency and variability Outliers Correlation	Daily 3/4 Hours
		Linear Algebra <ul style="list-style-type: none">• Introduction to Linear Algebra• Mathematics for Machine Learning Vectors and Matrices Matrices Operations Applications to Data Problems	
06	Machine Learning Basics	Basics of Machine Learning <ul style="list-style-type: none">• Business Problems to Data Problems Broad Categories of Business Problems Supervised and Unsupervised Machine Learning Algorithm Drivers of ML algorithms Cost Functions Brief introduction to Gradient Descent Importance of Model Validation Methods of Model Validation <ul style="list-style-type: none">• Introduction to Cross Validation and Average Error	Daily 3/4 Hours

07

Machine Learning – Algorithms (Supervised Learning)

Generalized Linear Models (Linear/Lasso/Ridge/Logistic)

Linear Regression

Limitation of simple linear models and need of regularization

Ridge and Lasso Regression (L1 & L2 Penalties)

Introduction to Classification with Logistic Regression

Methods of threshold determination

Performance measures for classification score models

Case Study 1 – Linear Regression, Ridge, Lasso and Logistic Regression

Practice assignment – 3

Decision Trees & Random Forests

Introduction to decision trees

- Tuning tree size with cross validation

Introduction to bagging algorithm

Random Forests

Grid search and randomized grid search

Extra Trees (Extremely Randomized Trees)

Case Study 2 – DT and RF

Practice assignment – 4

Boosting Machines in Python

Concept of weak learners

- Introduction to boosting algorithms

Adaptive Boosting

- Extreme Gradient Boosting (XGBoost)

Case Study 3 – Boosting Machines

Practice assignment – 5

K Nearest Neighbors

Introduction to idea of observation-based learning •

Distances and Similarities

Daily
3/4
Hours

		<p>• Nearest Neighbors (KNN) for classification and Regression)</p> <p>Case Study 4 - KNN</p> <p>Practice assignment – 6</p> <p>Support Vector Machines</p> <ul style="list-style-type: none"> • Introduction to SVM for classification <p>Case Study 5 - SVM</p> <p>Practice assignment – 7</p> <p>Neural Networks</p> <ul style="list-style-type: none"> • Introduction to Neural Networks <p>Single layer neural network</p> <p>Multiple layer Neural network</p> <p>Back propagation Algorithm</p> <ul style="list-style-type: none"> • Neural Networks implementation in Python <p>Case study 6 - NN</p>	
08	<p>Machine Learning – Algorithms (Unsupervised Learning)</p>	<p>Dimensionality Reduction</p> <p>Need for dimensionality reduction</p> <ul style="list-style-type: none"> • Introduction to Principal Component Analysis (PCA) • Difference between PCAs and Latent Factors <p>Introduction to Factor Analysis</p> <p>Case study 7 – PCA</p> <p>Case Study 8 - FA</p> <p>Segmentation in Python</p> <p>Patterns in the data in absence of a target</p> <p>Segmentation with Hierarchical Clustering and Kmeans</p> <p>Measure of goodness of clusters</p> <p>Limitations of K-means</p> <p>Introduction to density-based clustering (DBSCAN)</p> <p>Case study 9 – K-Means</p> <p>Case study 10 - DBSCAN</p>	<p>Daily 3/4 Hours</p>

09	Web scraping & API	<p>Data collection with web scraping & APIs</p> <p>Gathering text data using web scraping with urllib</p> <p>Processing raw web data</p> <ul style="list-style-type: none"> Interacting with Google search using urllib with custom user agent <p>Collecting twitter data with Twitter API</p> <p>Case study 11 – web scrapping</p> <p>Case study 12 – API to extract Data</p>	Daily 3/4 Hours
10	Natural Language Processing	<p>Natural Language Processing (Text Mining)</p> <ul style="list-style-type: none"> Quick Recap of string data functions and Introduction to Text Mining <p>Feature Engineering for text Data</p> <ul style="list-style-type: none"> Feature creation with TFIDF for text data <p>Case Study 13 – Text Data to model data</p> <p>Sentiment Analysis (NLP Supervised Learning) - Naïve Bayes/RF</p> <p>Introduction to Naive Bayes</p> <p>Case Study 14 – Naïve Bayes Classifier using Text Data (SPAM/Not SPAM)</p> <p>Introduction to Topic Modeling</p> <ul style="list-style-type: none"> Topic to word matrix and Document to topic matrix <p>Case Study 16 – LDA</p>	Daily 3/4 Hours
11	Ensemble Methods	<p>Ensemble Methods & Bokeh</p> <ul style="list-style-type: none"> Making use of multiple ML models taken together Simple Majority vote and weighted majority vote <p>Blending & Stacking</p> <p>Case study 17 – ensemble method</p>	Daily 3/4 Hours

12	introduction Big Data Analytics	Big Data Analytics <ul style="list-style-type: none"> • Big Data Hadoop Architecture, MapReduce • Apache Spark, PySpark, MLLib and Spark Tools PySpark Integration with Jupyter Notebook	Daily 3/4 Hours
13	Version Control & Data Product	Version control with Git & Interactive Data Product – prototyping solutions as Data Product Need and Importance of Version Control Setting up git and github accounts on local machine Creating and uploading GitHub Repos <ul style="list-style-type: none"> • Push and pull requests with GitHub App Merging and forking projects Pipeline and Pickle <ul style="list-style-type: none"> • Examples of static and interactive data products 	Daily 3/4 Hours
14	AI & Deep Learning	Deep Learning & Artificial Intelligence Installation of Tensor flow <ul style="list-style-type: none"> • Basics of Tensor flow with real-time project Image classification with Tensor flow real-time project Speech to Text with Tensor flow real-time project OCR with OpenCv with real-time project <ul style="list-style-type: none"> • Object detection with Tensor flow real-time project • Deep Learning Concepts like Neural Networks, AI 	Daily 3/4 Hours

		<p>image captioning with Tensor flow real-time project Deep Learning: Searching for Images Searching for images: A case study in deep learning</p> <ul style="list-style-type: none"> • Learning very non-linear features with neural networks <p>Application of deep learning to computer vision Deep learning performance Demo of deep learning model on Image Net data Deep learning ML block diagram</p> <ul style="list-style-type: none"> • Deploying Tensorflow deep learning models in production <p>DNN/CNN/RNN Building A chat bot with NLP Case Study 18: Ecommerce product recommendation Case Study 19 :Chat Bot Implementation in Tensor Flow Case Study 20 :Chat Bot Implementation in Pytorch</p>	
15	Business Analytics with Adv.Excel	<p>Introduction to Business Analytics Introduction</p> <ul style="list-style-type: none"> • Introduction to business analytics • Formatting conditional formatting and logical functions <p>Analyzing data with pivot tables Dashboarding Business analytics with Adv.excel Data analysis using statistics</p>	Daily 3/4 Hours
16	Model Deployments Production in Cloud	<ul style="list-style-type: none"> • Model Building and Deployment in AZURE /AWS/GCP <p>Deep Learning Model Building and Deployment in AZURE/AWS/GCP</p>	Daily 3/4 Hours

17	Advanced Predictive Analysis	Introduction to Advanced Predictive Analysis Data Understanding Data Preparation Data Transformation Item sets and Associations	Daily 3/4 Hours
18	Story Telling Using BI	Story Telling Using Business Intelligence (BI) Tool – Power BI/Tableau Introduction to BI and BI Tool Exploratory Data Analysis (EDA) using BI Tool <ul style="list-style-type: none"> • Creating Dashboard using predictive model results Case Study 21 – EDA on Sales Data 	Daily 3/4 Hours

Other related topics

	<ul style="list-style-type: none"> • • • • 	
Interview Preparation		Interview Questions

