# Data-science - R

#### **Duration: 3 Days**

#### INTRODUCTION

- Data science & its importance
- Key Elements of Data Science
- Introduction to ML
- Artificial Intelligence & Machine Learning Introduction
- Who uses AI?
- Al for Banking & Finance, Manufacturing, Healthcare, Retail and Supply Chain
- Supervised & Unsupervised Learning
- Regression & Classification Problems
- What makes a Machine Learning Expert?
- What to learn to become a Machine Learning Developer?
- Overview of Machine Learning Algorithms

#### R

- R basics (If-Else, Control Structures, Loops, Functions, Data Types)
- Data structures (Vector, Matrix, Dataframe, Lists)
- Indexing, Data Processing
- Mathematical computing basics
- Getting Started with Dplyr and tidyr
- Data Acquisition (Import & Export)
- Selection and Filtering
- Combining and Merging Data Frames

### **STATISTICS and EDA**

- Introduction to Visualization
- Visualization Importance
- Working with R visualization libraries like ggplot2
- Creating Line Plots, Bar Charts, Pie Charts, Histograms, Scatter Plots
- Understanding Box plots
- Understanding Probability Distributions, Violin Plots
- Correlations and Heatmaps
- Summary Statistics
- Central Tendency measures
- Measures of dispersion
- Normal Distributions and z-score
- Missing Value Imputation
- Outlier Detection and handling
- Advanced EDA techniques

- Machine Learning Algorithms Generic Concepts
  - $\circ$  Sample and Population
  - Bias-Variance Trade off
  - Overfitting and Underfitting
  - o Cross Validation
  - o Regularization techniques
  - o Hyperparameter tuning & grid search optimization
- Linear Regression
  - o Regression Problem Analysis
  - $\circ \quad \text{Mathematical modelling of Regression Model}$
  - o Gradient Descent Algorithm
  - o Use cases
  - o Regression Table
  - $\circ \quad \text{Model Specification} \\$
  - o L1 & L2 Regularization
  - o Building simple Univariate Linear Regression Model
  - o Multivariate Regression Mode
  - R2, p-value, RMSE and residual plots
- Logistic Regression
  - o Assumptions
  - $\circ \quad \text{Sigmoid function} \quad$
  - o ROC Curve
  - o Model Specification
  - o Confusion Matrix
  - o Accuracy, Recall, Precision and F1 Score
- Decision Trees
  - Forming a Decision Tree
  - o Components of Decision Tree
  - o Mathematics of Decision Tree
  - o Decision Tree Evaluation
- KNN
  - $\circ$  Components of Decision Tree
  - $\circ \quad \text{Mathematics of Decision Tree}$
  - $\circ \quad \text{Metrics for evaluation} \quad$
- Support Vector Machine
  - Concept and Working Principle
  - Mathematical Modelling
  - o Optimization Function Formation
  - The Kernel Method and Nonlinear Hyperplanes
  - Ensemble Models
    - Bagging
    - o Boosting
    - $\circ$  Stacking
    - Voting Classifier
    - o Random Forest
- Unsupervised Machine Learning algorithms
  - Clustering with K-meansClustering
  - Clustering with Hierarchical Clustering
  - Advanced clustering techniques and use cases
  - o Dimensionality Reduction
  - o PCA
- Text Mining and NLP
  - o Sentiment Analysis

- Topic Summarization
- o Topic Modelling
- $\circ$   $\,$  Bag of Words and Tf-IDF  $\,$
- Cosine Similarity of terms, documents concepts
- Text Cleaning and Preprocessing using Regex
- Tokenization, Stemming and Lemmatization

## CASE STUDY AND PROJECTS

Students would be given challenging real-life cases to solve – just to augment their learning skills