

# Data Scientist

Roadmap

## Mathematics

- Linear Algebra
- Analytics Geometry
- Matrix
- Vector Calculus
- Optimization
- Regression
- Dimensionality Reduction
- Density Estimation
- Classification

## Probability

- Discrete Distribution
  - Binomial
  - Bernoulli
  - Geometric etc
- Continuous Distribution
  - Uniform
  - Exponential
  - Gamma
- Normal Distribution
- Introduction to Probability
- 1D Random Variable
- Function of One Random Variable
- Joint Probability Distribution

## Statistics

- Introduction to Statistics
- Data Description
- Random Samples
- Sampling Distribution
- Parameter Estimation
- Hypotheses Testing
- ANOVA
- Reliability Engineering
- Stochastic Process
- Computer Simulation
- Design of Experiments
- Simple Linear Regression
- Correlation
- Multiple Regression
- Nonparametric Statistics
  - Sign Test
  - The Wilcoxon Signed-Rank Test
  - The Wilcoxon Rank Sum Test
  - The Kruskal-Wallis Test
- Statistical Quality Control
- Basic of Graphs

## Programming

- | Python   | R   |
|--|---|
| <b>Python Basics</b> <ul style="list-style-type: none"><li>• List</li><li>• Set</li><li>• Tuples</li><li>• Dictionary</li><li>• Function, etc.</li></ul> | <b>R Basic</b> <ul style="list-style-type: none"><li>• Vector</li><li>• List</li><li>• Data Frame</li><li>• Matrix</li><li>• Array, etc</li></ul> |
| <b>NumPy</b>   | <b>dplyr</b>  |
| <b>Pandas</b>  | <b>ggplot2</b>  |
| <b>Matplotlib/Seaborn, etc.</b>  | <b>Tidyr</b>  |
|  | <b>Shiny, etc.</b>  |
| <b>DataBase</b>  | <b>Other</b>  |
| <b>SQL</b>   | <b>Data Structure</b> <ul style="list-style-type: none"><li>• Array, etc</li></ul>  |
| <b>MongoDB</b>   | <b>Web Scraping</b>   |
|  | <b>Linux</b>  |
|  | <b>Git</b>  |

## Machine Learning

- | Introduction  | Intermediate   |
|---|--|
| <ul style="list-style-type: none"><li>• How Model Works</li><li>• Basic Data Exploration</li><li>• First ML Model</li><li>• Model Validation</li><li>• Underfitting &amp; Overfitting</li><li>• Random Forests</li><li>• scikit-learn</li></ul> | <ul style="list-style-type: none"><li>• Handling Missing Values</li><li>• Handling Categorical Variables</li><li>• Pipelines</li><li>• Cross-Validation</li><li>• XGBoost</li><li>• Data Leakage</li></ul> |

## Deep Learning

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>• Artificial Neural Network</li><li>• Convolutional Neural Network</li><li>• Recurrent Neural Network</li><li>• Keras</li><li>• PyTorch</li><li>• TensorFlow</li></ul> | <ul style="list-style-type: none"><li>• A Single Neuron</li><li>• Deep Neural Network</li><li>• Stochastic Gradient Descent</li><li>• Overfitting and Underfitting</li><li>• Dropout Batch Normalization</li><li>• Binary Classification</li></ul> |
|--|--|

## Feature Engineering

- Baseline Model
- Feature Generation
- Categorical Encodings
- Feature Selection

## Natural language Processing

- Text Classification
- Word Vectors

## Data Visualization Tools

- Excel VBA
- BI (Business Intelligence)
  - Tableau
  - Power BI
  - Qlik View
  - Qlik Sense

## Deployment

- Microsoft Azure
- Heroku
- Google Cloud Platform
- Flask
- Django

## Other Points

- Domain Knowledge
- Communication Skill
- Reinforcement Learning
- Case Studies
  - Data Science at Netflix
  - Data Science at Flipkart
  - Project on Credit Card Fraud Detection
  - Project on Movie Recommendation , etc.

## Keep Practicing