# Advanced Supervised and Unsupervised Learning Algorithms
By Dr. Vishwanath Rao

Python for Data Science

Python Data Science-Centric Libraries
SciPy
NumPy
pandas
Scikit-learn
Matplotlib
Seaborn
Python Dev Tools and REPLs
IPython
Jupyter Notebooks
Anaconda
Summary

Data Visualization in Python

Why Do I Need Data Visualization?
Data Visualization in Python
Getting Started with matplotlib
A Basic Plot
Scatter Plots
Figures
Saving Figures to a File
Seaborn
Getting Started with seaborn
Histograms and KDE
Plotting Bivariate Distributions
Scatter Plots in seaborn
Pair plots in seaborn
Heatmaps
A Seaborn Scatterplot with Varying Point Sizes and Hues
Summary

Introduction to NumPy

What is NumPy?
The First Take on NumPy Arrays
The ndarray Data Structure
Understanding Axes
Indexing Elements in a NumPy Array
Re-Shaping
Commonly Used Array Metrics
Commonly Used Aggregate Functions
Sorting Arrays

Overview of the scikit-learn Library

Developer Support
scikit-learn Estimators, Models, and Predictors
Annotated Example of the LinearRegression Estimator
Annotated Example of the Support Vector Classification Estimator
Data Splitting into Training and Test Datasets
Data Splitting in scikit-learn
Cross-Validation Technique
Summary

## Classification Algorithms (Supervised Machine Learning)

Classification (Supervised ML) Use Cases
Classifying with k-Nearest Neighbors
k-Nearest Neighbors Algorithm Visually
Decision Trees
Decision Tree Terminology
Decision Tree Classification in the Context of Information Theory
Using Decision Trees
Properties of the Decision Tree Algorithm
The Simplified Decision Tree Algorithm
Random Forest
Properties of the Random Forest Algorithm
Support Vector Machines (SVMs)
SVM Classification Visually
Properties of SVMs
Dealing with Non-Linear Class Boundaries
Logistic Regression (Logit)
The Sigmoid Function
Logistic Regression Classification Example
Logistic Regression's Problem Domain
Naive Bayes Classifier (SL)
Naive Bayesian Probabilistic Model in a Nutshell
Bayes Formula
Document Classification with Naive Bayes
Summary

## Deep Study on Supervised learning

Linear Models
Linear and Quadratic Discriminant Analysis
Kernel ridge regression
Support Vector Machines
Stochastic Gradient Descent
Nearest Neighbors
Gaussian Processes
Cross decomposition
Naive Bayes
Decision Trees
Ensemble methods
Multiclass and multioutput algorithms
Feature selection

Semi-supervised learning
Isotonic regression
Probability calibration
Neural network models (supervised)


Unsupervised Machine Learning Algorithms

PCA
PCA and Data Variance
PCA Properties
Importance of Feature Scaling Visually
Unsupervised Learning Type: Clustering
Clustering vs Classification
Clustering Examples
k-means Clustering
k-means Clustering in a Nutshell
k-means Characteristics
Global vs Local Minimum Explained
Summary


Deep study on Unsupervised learning

Gaussian mixture models
Manifold learning
Clustering
Biclustering
Decomposing signals in components (matrix factorization problems)
Covariance estimation
Novelty and Outlier Detection
Density Estimation
Neural network models (unsupervised)