

Comprehensive Machine Learning for CyberSecurity and GEN AI (5 Days)

By Dr. Vishwanath Rao

Overview

Comprehensive Machine Learning (ML) with Python training course builds on our

Comprehensive Data Science with Python class and teaches attendees how to write machine learning applications in Python.

Prerequisites

Basic knowledge on python

Objectives

- Understand machine learning as a useful tool for predictive models
- Know when to reach for machine learning as a tool
- Implement data preprocessing for an ML workflow
- Understand the difference between supervised and unsupervised tasks
- Implement several classification algorithms
- Evaluate model performance using a variety of metrics
- Compare models across a workflow
- Implement regression algorithm variations
- Understand clustering approaches to data
- Interpret labels generated from clustering
- Transform unstructured text data into structured data
- Understand text-specific data preparation
- Visualize frequency data from text sources
- Perform topic modeling on a collection of documents
- Use labeled text to perform document classification

Outline

- Introduction
- Review of Core Python Concepts
 - o Anaconda Computing Environment
 - o Importing and manipulating Data with Pandas
 - o Exploratory Data Analysis with Pandas and Seaborn
 - o NumPy ndarrays versus Pandas Dataframes

- An Overview of Machine Learning
 - o Machine Learning Theory
 - o Data pre-processing
 - Missing Data
 - Dummy Coding
 - Standardization
 - Data Validation Strategies
 - o Supervised Versus Unsupervised Learning

- Modeling for explanation (descriptive models)
 - o Understanding the linear model
 - o Describing model fit
 - o Adding complexity to the model
 - o Explaining the relationship between model inputs and the outcome
 - o Making predictions from the model

- Supervised Learning: Regression
 - o Linear Regression
 - o Penalized Linear Regression
 - o Stochastic Gradient Descent
 - o Decision Tree Regressor
 - o Random Forest Regression
 - o Gradient Boosting Regressor
 - o Scoring New Data Sets
 - o Cross Validation
 - o Variance-Bias Tradeoff
 - o Feature Importance

- Supervised Learning: Classification
 - o Logistic Regression
 - o LASSO
 - o Support Vector Machine
 - o Random Forest
 - o Ensemble Methods
 - o Feature Importance
 - o Scoring New Data Sets
 - o Cross Validation

- Unsupervised Learning: Clustering
 - o Preparing Data for Ingestion

- o K-Means Clustering
- o Visualizing Clusters
- o Comparison of Clustering Methods
- o Agglomerative Clustering and DBSCAN
- o Evaluating Cluster Performance with Silhouette Scores
- o Scaling
- o Mean Shift, Affinity Propagation and Birch
- o Scaling Clustering with mini-batch approaches

- Clustering for Treatment Effect Heterogeneity
 - o Understand average versus conditional treatment effects
 - o Estimating conditional average treatment effects for a sample
 - o Summarizing and Interpreting

- Data Munging and Machine Learning Via H2O
 - o Intro to H2O
 - o Launching the cluster, checking status
 - o Data Import, manipulation in H2O
 - o Fitting models in H2O
 - o Generalized Linear Models
 - o naïve bayes
 - o Random forest
 - o Gradient boosting machine (GBM)
 - o Ensemble model building
 - o automl
 - o data preparation
 - o leaderboards
 - o Methods for explaining modeling output
- Introduction to Natural Language Processing (NLP)
 - o Transforming Raw Text Data into a Corpus of Documents
 - o Identifying Methods for Representing Text Data
 - o Transformations of Text Data
 - o Summarizing a Corpus into a TF—IDF Matrix
 - o Visualizing Word Frequencies

- NLP Normalization, Parts-of-speech and Topic Modeling
 - o Installing And Accessing Sample Text Corpora
 - o Tokenizing Text
 - o Cleaning/Processing Tokens
 - o Segmentation
 - o Tagging And Categorizing Tokens
 - o Stopwords

- o Vectorization Schemes for Representing Text
- o Parts-of-speech (POS) Tagging
- o Sentiment Analysis
- o Topic Modeling with Latent Semantic Analysis

- NLP and Machine Learning
 - o Unsupervised Machine Learning and Text Data
 - o Topic Modeling via Clustering
 - o Supervised Machine Learning Applications in NLP
- Conclusion