

# Machine Learning using Python

## Course Outline



## Overview of the Training Course

This course provides a broad introduction to data science, machine learning, data mining, and statistical pattern recognition. It covers supervised and unsupervised learning, natural language processing, computer vision problems, time series analysis. You'll learn the essential theories, the interpretations of what comes out of machine learning execution and gain the practical know-how needed to quickly and powerfully apply these techniques to new problems in your own organization.

You will spend 40% time on theory lectures and demos and rest on lab exercises. Lab exercises will be based on problems hosted on platforms such as kaggle.com, UCI dataset.

Python programming language is used for lab exercises and demonstrations.

**Duration:** 4 Days

### Prerequisite

- Strong programming knowledge in Python is required to perform the exercises during the class. If you want to brush up python knowledge, here is a quick [tutorial](#).
- Basic knowledge of statistics, linear algebra is required to understand and implement some of the ML algorithms. You may revise the textbooks or online materials on these subjects.
- I strongly recommend taking one of the following online courses prior to the class.
  - Machine Learning by Andrew Ng available in [coursera.org](https://www.coursera.org)
  - Statistical Learning by Trevor Hastie and Rob Tibshirani available in [online.stanford.edu](https://online.stanford.edu)
  - Analytics Edge available at [edx.org](https://edx.org)

## Day 1: Intro to Machine Learning and Data Preprocessing

Topic	Description
Introduction to Machine Learning	<ul style="list-style-type: none"><li>• What is machine learning, data science and its use cases</li><li>• Supervised and Unsupervised Learning</li><li>• Types of machine learning – choice of algorithm<ul style="list-style-type: none"><li>◦ Regression, Classification, Recommender</li><li>◦ Parametric vs Non-parametric</li></ul></li><li>• Machine learning libraries</li></ul>
Introduction to Numpy and Scipy	<ul style="list-style-type: none"><li>• Fast numeric computation using numpy arrays</li><li>• Linear Algebra – matrix multiplication, decompositions/factorization, determinants</li><li>• Common statistical functions using Scipy</li></ul>
Exploratory data analysis	<ul style="list-style-type: none"><li>• Descriptive statistics - mean, median, range, variance, standard deviation, variance, kurtosis, skewness, percentile</li><li>• Correlation</li><li>• Analysis of variance and F statistics</li><li>• Chi-square test to test goodness of fit</li><li>• Hypothesis tests</li></ul>
Introduction to Pandas	<ul style="list-style-type: none"><li>• Explore data in Dataframe - index, slice, filter, sort, join, aggregate and transform</li><li>• Handle Missing Data</li><li>• Descriptive statistics<ul style="list-style-type: none"><li>◦ Correlation and covariance</li><li>◦ Unique values, value count, distribution</li></ul></li></ul>
Lab Exercises	<ul style="list-style-type: none"><li>• Analyzing the grouplens.org data using Pandas</li></ul>
Data Visualization (interactive session)	<ul style="list-style-type: none"><li>• Introduction to matplotlib and Seaborn</li><li>• Types of chart, when to use which one</li><li>• Interpretation from each type of chart</li><li>• Subplots and overlay plots</li><li>• Colors, styles, labels, legends</li></ul>

Feature Engineering (interactive session)	<ul style="list-style-type: none"> <li>• Characteristics of good features</li> <li>• Handle missing values</li> <li>• Handle outliers</li> <li>• Handle categorical data</li> <li>• Feature scaling - standard scaler, normalizer (L1/L2), min-max scaler</li> <li>• Probability distribution of variables</li> <li>• Normality test for continuous variables</li> <li>• Power transformation</li> <li>• Polynomial transformation</li> <li>• Discretization</li> <li>• Sampling - down sampling, up sampling, stratified sampling</li> <li>• Feature selection based on significance</li> <li>• Multi-collinearity detection using correlation, variance inflation ratio etc.</li> <li>• Dimensionality reduction using PCA</li> <li>• Generate new features</li> <li>• Handle Time series data</li> <li>• Denormalization of structured data of RDBMS systems</li> </ul>
Apache Spark for ML	<ul style="list-style-type: none"> <li>• Features of apache Spark</li> <li>• Data structures in Spark</li> <li>• Dataframe DSL and Spark SQL</li> <li>• Spark integration with Pandas, Numpy, Scipy</li> </ul>
Lab Exercise	<ul style="list-style-type: none"> <li>• Explore ad-click prediction data from Outbrain</li> </ul>

### Learning Goals:

- Spot the use cases of machine learning
- Explore and transform data using numpy, pandas, Matplotlib
- Clean and prepare data for further processing

## Day 2: Regression, Model Diagnostics and Tuning

Topic	Description
Linear Regression	<p>Predict numeric outcomes with linear regression</p> <ul style="list-style-type: none"> <li>• Finding best fit lines with linear regression</li> <li>• Model evaluation metrics – R<sup>2</sup>, RMSE, MAE, F-stat</li> <li>• Shrinking coefficients to understand the data – Ridge regression, Lasso regression, Forward stage-wise regression, L1/L2 Penalty for regularization</li> <li>• Bias / variance trade off</li> <li>• Tree based regression (CART regression)</li> <li>• Understanding the cost function and minimize cost function using batch gradient descent algorithm</li> <li>• Variations of gradient algorithms - stochastic gradient descent, mini batch gradient descent for classification problems</li> </ul>

Topic	Description
Lab Exercise	<ul style="list-style-type: none"> <li>• Forecast insurance premium charges based on customer profile</li> <li>• Forecasting home price based on data from Kaggle.com challenge</li> </ul>
Model Tuning and Deployment	<ul style="list-style-type: none"> <li>• Combining transformers and estimators in a pipeline</li> <li>• Using k-fold cross validation to assess model performance</li> <li>• Debugging algorithms with learning and validation curve</li> <li>• Diagnosis bias and variance problems with learning curves</li> <li>• Fine tuning hyper parameters using grid search</li> <li>• Production readiness of ML application</li> <li>• Exposing ML model as web service</li> </ul>
Lab Exercises	<ul style="list-style-type: none"> <li>• Tune models built in the previous exercises</li> </ul>

### Learning Goals:

- Learn to predict numeric data (for example, revenue, price etc.)
- Learn to evaluate a regression model
- Learn to find important predictors for regression model
- Learn to validate and tune a regression model for best performance

## Day 3: Classification and NLP

Topic	Description
Classification Algorithm	Predict categorical outcomes with classification <ul style="list-style-type: none"> <li>• Logistic Regression for binary and multi class classification</li> <li>• Classification using decision trees, random forest</li> <li>• Classifying with KNN</li> <li>• Model evaluation using accuracy score, precision, recall, F-score, ROC plot, AUC, confusion matrix</li> <li>• Tune the hyper parameters using cross validation and grid search</li> </ul>
Lab Exercise	<ul style="list-style-type: none"> <li>• Predicting customer credit status customer profile</li> <li>• Find fraudulent transactions in credit card transactions</li> <li>• Reveal local attitudes from personal ads from craigslist.org RSS feed</li> </ul>
Text Analytics	<ul style="list-style-type: none"> <li>• Text analysis steps for topic modeling</li> <li>• Term Frequency – Inverse Term Frequency</li> <li>• POS tagging, lemmatization, stemming</li> <li>• Collecting raw text</li> <li>• Representing text</li> <li>• Categorize Documents by type</li> <li>• Determining sentiments</li> <li>• Gaining insights</li> </ul>
Lab Exercise	<ul style="list-style-type: none"> <li>• Working with IMDB comments dataset to analyze sentiments</li> <li>• StumbleUpon Evergreen classification Challenge</li> <li>• BBC Datasets - content based classification problem</li> </ul>

Topic	Description
Boosting, Bagging and Ensemble	<ul style="list-style-type: none"> <li>• What is boosting, bagging and pasting?</li> <li>• Putting multiple models in action</li> <li>• Building a hybrid model and compare performance metrics</li> </ul>
Lab Exercise	<ul style="list-style-type: none"> <li>• Build a boosted model for a classifier</li> </ul>

### Learning Goals:

- Learn to predict categorical data types, for example – customer segmentation, detection of email spam, stock market gainer or loser, fraud prediction
- Learn which metrics to use under what scenarios to measure performance
- Learn to find important predictor
- Evaluate and tune classification models
- Learning common tasks for text analytics

## Day 4: Clustering and Neural Networks

Topic	Description
Clustering	<ul style="list-style-type: none"><li>• Working with unlabeled data – clustering analysis</li><li>• K-means clustering</li><li>• Hard (K-means) vs Soft clustering (Fuzzy c-means)</li><li>• Find optimal number of clusters</li><li>• Hierarchical Clustering</li></ul>
Lab Exercise	<ul style="list-style-type: none"><li>• Clustering users based social-media profile</li></ul>
Intro to Deep Learning	<ul style="list-style-type: none"><li>• What is Artificial Neural Network?</li><li>• What is deep learning?</li><li>• Introduction to Tensorflow and Keras</li><li>• Optimization to gradient descent algorithms<ul style="list-style-type: none"><li>◦ Batch vs mini-batch vs stochastic gradient descent</li><li>◦ Advanced initializers</li><li>◦ Momentum based methods</li><li>◦ Adaptive learning rate</li></ul></li><li>• Convolutional Neural Networks (CNN) and its use cases</li><li>• Recurrent Neural Networks (RNN) and its use cases</li></ul>
Lab Exercises	<ul style="list-style-type: none"><li>• TensorFlow basics</li><li>• Building regression model using TensorFlow</li><li>• Building multi class classifier using TensorFlow</li><li>• Classifier for image of hand-written digits (MNIST)</li><li>• Classifier for image of objects (CIFAR10)</li></ul>

### Learning Goals:

- Learn to apply clustering to find natural groupings within a dataset
- Solve computer vision related problems

## Appendix: Lab Setup

For lab exercises each student needs a laptop or a desktop with standard configuration 8 GB+ RAM and i5 CPU. Students will require administrative access to install software. For the lab exercise student will require to install the following software. Internet access is required for downloading test datasets, sample code, browse and documentations. Access to USB pen drive is required to copy materials from training instructor's machine.

- Anaconda distribution for python 3.6+ 64-bit version [full version](#)