

# Data Analyst Nanodegree Syllabus

*Discover Insights from Data*



## Before You Start

**Prerequisites:** Thank you for your interest in the Data Analyst Nanodegree! In order to succeed in this program, we recommend having experience programming in Python. If you've never programmed before, or want a refresher, there is an Introduction to Python Programming in the extracurricular section of the nanodegree program.

**Educational Objectives:** Learn to organize data, uncover patterns and insights, make predictions using machine learning, and clearly communicate critical findings.

**Length of Program\*:** 260 Hours

**Frequency of Classes:** Self-paced

**Textbooks required:** None

**Instructional Tools Available:** Video lectures, 1:1 appointments, forum support

\*This is a self-paced program and the length is an estimation of total hours the average student may take to complete all required coursework, including lecture and project time. Actual hours may vary.

## Intro Project: Analyze Bay Area Bike Share Data (10 hrs)

This project will introduce you to the key steps of the data analysis process. You'll do so by analyzing data from a bike share company found in the San Francisco Bay Area. You'll submit this project in your first 7 days, and by the end you'll be able to:

- Use basic Python code to clean a dataset for analysis
- Run code to create visualizations from the wrangled data
- Analyze trends shown in the visualizations and report your conclusions
- Determine if this program is a good fit for your time and talents

## Project: Compute Statistics from Card Draws (20 hrs)

In this project, you will demonstrate your knowledge of descriptive statistics by conducting an experiment dealing with drawing from a deck of playing cards and creating a write-up containing your findings. This project is self-graded.

### Supporting Lesson Content: Statistics

Lesson Title	Learning Outcomes
<b>INTRO TO RESEARCH METHODS</b>	→ Identify several statistical study methods and describe the positives and negatives of each
<b>VISUALIZING DATA</b>	→ Create and interpret histograms, bar charts, and frequency plots
<b>CENTRAL TENDENCY</b>	→ Compute and interpret the 3 measures of center for distributions: the mean, median, and mode
<b>VARIABILITY</b>	→ Quantify the spread of data using the range and standard deviation → Identify outliers in data sets using the interquartile range
<b>STANDARDIZING</b>	→ Convert distributions into the standard normal distribution using the Z-score → Compute proportions using standardized distributions
<b>NORMAL DISTRIBUTION</b>	→ Use normal distributions to compute probabilities → Use the Z-table to look up the proportions of observations above, below, or in between values
<b>SAMPLING DISTRIBUTIONS</b>	→ Apply the concepts of probability and normalization to sample data sets

## Project: Investigate a Dataset (30 hrs)

In this project, you'll choose one of Udacity's curated datasets and investigate it using NumPy and pandas. You'll complete the entire data analysis process, starting by posing a question and finishing by sharing your findings.

### Supporting Lesson Content: Introduction to Data Analysis

Lesson Title	Learning Outcomes
<b>DATA ANALYSIS PROCESS</b>	<ul style="list-style-type: none"><li>→ Identify the key steps in the data analysis process</li><li>→ Complete an analysis of Udacity student data using pure Python, with minimal reliance on additional libraries</li></ul>
<b>NUMPY AND PANDAS FOR 1D DATA</b>	<ul style="list-style-type: none"><li>→ Use NumPy arrays, pandas series, and vectorized operations to ease the data analysis process</li></ul>
<b>NUMPY AND PANDAS FOR 2D DATA</b>	<ul style="list-style-type: none"><li>→ Use two-dimensional NumPy arrays and pandas DataFrames</li><li>→ Understand how to group data and to combine data from multiple files</li></ul>

## Project: Wrangle OpenStreetMap Data (60 hrs)

In this project, you'll use data munging techniques, such as assessing the quality of the data for validity, accuracy, completeness, consistency and uniformity, to clean the OpenStreetMap data for a part of the world that you care about.

### Supporting Lesson Content: Data Wrangling with SQL

Lesson Title	Learning Outcomes
<b>DATA EXTRACTION FUNDAMENTALS</b>	<ul style="list-style-type: none"><li>→ Properly assess the quality of a dataset</li><li>→ Understand how to parse CSV files and XLS with XLRD</li><li>→ Use JSON and Web APIs</li></ul>
<b>DATA IN MORE COMPLEX FORMATS</b>	<ul style="list-style-type: none"><li>→ Understand XML design principles</li><li>→ Parse XML &amp; HTML</li><li>→ Scrape websites for relevant data</li></ul>
<b>DATA QUALITY</b>	<ul style="list-style-type: none"><li>→ Understand common sources for dirty data</li><li>→ Measure the quality of a dataset &amp; apply a blueprint for cleaning</li><li>→ Properly audit validity, accuracy, completeness, consistency, and uniformity of a dataset</li></ul>
<b>ANALYZING DATA</b>	<ul style="list-style-type: none"><li>→ Identify common examples of the aggregation framework</li></ul>

	<ul style="list-style-type: none"> <li>→ Use aggregation pipeline operators \$match, \$project, \$unwind, \$group</li> </ul>
<b>SQL FOR DATA ANALYSIS</b>	<ul style="list-style-type: none"> <li>→ Understand how data is structured in SQL</li> <li>→ Run queries to summarize data</li> <li>→ Use joins to combine information across tables</li> <li>→ Create tables and import data from csv</li> </ul>
<b>CASE STUDY: OPENSTREETMAP DATA</b>	<ul style="list-style-type: none"> <li>→ Use iterative parsing for large datafiles</li> <li>→ Understand XML elements in OpenStreetMap</li> </ul>

## Project: Explore and Summarize Data (50 hrs)

In this project, you'll use R and apply exploratory data analysis techniques to explore a selected data set for distributions, outliers, and anomalies.

### Supporting Lesson Content: Data Analysis with R

Lesson Title	Learning Outcomes
<b>WHAT IS EDA?</b>	<ul style="list-style-type: none"> <li>→ Define and identify the importance of exploratory data analysis (EDA)</li> </ul>
<b>R BASICS</b>	<ul style="list-style-type: none"> <li>→ Install RStudio and packages</li> <li>→ Write basic R scripts to inspect datasets</li> </ul>
<b>EXPLORE ONE VARIABLE</b>	<ul style="list-style-type: none"> <li>→ Quantify and visualize individual variables within a dataset</li> <li>→ Create histograms and boxplots</li> <li>→ Transform variables</li> <li>→ Examine and identify tradeoffs in visualizations</li> </ul>
<b>EXPLORE TWO VARIABLES</b>	<ul style="list-style-type: none"> <li>→ Properly apply relevant techniques for exploring the relationship between any two variables in a data set</li> <li>→ Create scatter plots</li> <li>→ Calculate correlations</li> <li>→ Investigate conditional means</li> </ul>
<b>EXPLORE MANY VARIABLES</b>	<ul style="list-style-type: none"> <li>→ Reshape data frames and use aesthetics like color and shape to uncover information</li> </ul>
<b>DIAMONDS AND PRICE PREDICTIONS</b>	<ul style="list-style-type: none"> <li>→ Use predictive modeling to determine a good price for a diamond</li> </ul>

## Project: Test a Perceptual Phenomenon (20 hrs)

In this project, you'll use descriptive statistics and a statistical test to analyze the Stroop effect, a classic result of experimental psychology. Communicate your understanding of the data and use statistical inference to draw a conclusion based on the results.

### Supporting Lesson Content: Inferential Statistics

Lesson Title	Learning Outcomes
<b>ESTIMATION</b>	→ Estimate population parameters from sample statistics using confidence intervals
<b>HYPOTHESIS TESTING</b>	→ Use critical values to make decisions on whether or not a treatment has changed the value of a population parameter
<b>T-TESTS</b>	→ Test the effect of a treatment or compare the difference in means for two groups when we have small sample sizes

## Project: Identify Fraud from Enron Email (50 hrs)

In this project, you'll play detective and put your machine learning skills to use by building an algorithm to identify Enron employees who may have committed fraud based on the public Enron financial and email dataset.

### Supporting Lesson Content: Introduction to Machine Learning

Lesson Title	Learning Outcomes
<b>SUPERVISED CLASSIFICATION</b>	<ul style="list-style-type: none"><li>→ Implement the Naive Bayes algorithm to classify text</li><li>→ Implement Support Vector Machines (SVMs) to generate new features independently on the fly</li><li>→ Implement decision trees as a launching point for more sophisticated methods like random forests and boosting</li></ul>
<b>DATASETS AND QUESTIONS</b>	→ Wrestle the Enron dataset into a machine-learning-ready format in preparation for detecting cases of fraud
<b>REGRESSIONS AND OUTLIERS</b>	→ Use regression algorithms to make predictions and identify and clean outliers from a dataset
<b>UNSUPERVISED LEARNING</b>	→ Use the k-means clustering algorithm for pattern-searching on unlabeled data

---

**FEATURES, FEATURES,  
FEATURES**

- Use feature creation to take your human intuition and change raw features into data a computer can use
- Use feature selection to identify the most important features of your data
- Implement principal component analysis (PCA) for a more sophisticated take on feature selection
- Use tools for parsing information from text-type data

---

**VALIDATION AND  
EVALUATION**

- Implement the train-test split and cross-validation to validate and understand machine learning results
  - Quantify machine learning results using precision, recall, and F1 score
- 

## Project: Make an Effective Visualization (20 hrs)

In this project, you'll create a data visualization, using Tableau, from a data set that tells a story or highlights trends or patterns in the data. Your work should be a reflection of the theory and practice of data visualization, harnessing visual encodings and design principles for effective communication.

### Supporting Lesson Content: Data Visualization with Tableau

Lesson Title	Learning Outcomes
<b>DATA VISUALIZATION FUNDAMENTALS</b>	<ul style="list-style-type: none"><li>→ Understand the importance of data visualization</li><li>→ Know how different data types are encoded in visualizations</li></ul>
<b>DESIGN PRINCIPLES</b>	<ul style="list-style-type: none"><li>→ Select the most effective chart or graph based on the data being displayed</li><li>→ Use color, shape, size, and other elements effectively</li></ul>
<b>CREATING VISUALIZATIONS WITH TABLEAU</b>	<ul style="list-style-type: none"><li>→ Become proficient in basic Tableau functionality, including charts, filters, hierarchies, etc.</li><li>→ Create calculated fields in Tableau</li></ul>
<b>TELLING STORIES WITH TABLEAU</b>	<ul style="list-style-type: none"><li>→ Create Tableau dashboards and stories to effectively communicate data</li></ul>

---