

BOOTCAMP TRAINING IN PROGRAMMING, STATISTICS AND DATA SCIENCE

JUNE 2021

1 Introduction

This is a training designed for Accounting Guatemala that covers topics in computer programming, math and statistics, data science and machine learning. It includes the explanation of the basic theory necessary to implement simple machine learning models, as well as instruction in how to handle Python programming in most common environments, that covers various popular modules and libraries in data science.

Course:	Data Science Bootcamp
Main topics:	Mathematics for data science Descriptive and inference statistics Python programming Basic machine learning models
Time:	Four weeks, 5 days per week, 2 hour each session for a total of 40 hours of session work.
Instructor:	Alan Reyes-Figueroa, Ph.D. Professor of Mathematics and data science at Universidad del Valle, Data scientist at Bam.

2 Skills to Develop

- Learn what data science is, the various activities of a data scientist's job, and methodology to work as a data scientist.
- Develop hands-on skills using the theory and tools, languages, and Python libraries used by professional data scientists.
- Import and clean data sets, analyze and visualise data, and build and evaluate machine learning models and pipelines using Python.
- Apply various data science skills, techniques, and tools to complete a project and publish a report.

3 Syllabus

1. Installation and set up of Python working environments for data science: Anaconda Python, Jupyter notebooks and Jupyter-lab, VSCode and Jupyter built in.
2. Libraries and Python modules for data science: Numpy, Pandas, Matplotlib, Seaborn, Plotly, Statsmodels, Scikit-learn. SQL basic queries and modules for working with SQL in Python.
3. Math and statistical concepts: Vectors and matrix calculus and linear algebra, plotting functions, basic probability, distributions, descriptive statistics, inference statistics, hypothesis testing.
4. Python programming: Variables and functions, conditionals, cycles *for*, *while* and exceptions. Basic structures: lists, tuples, dictionaries, strings. Vectors and matrices.
5. Data science: Read and modify datasets with Pandas. Deal with missing information. Data exploration: histograms, covariance analysis, visualisation and plots. Principal components analysis.

6. Machine learning models: Scikit-learn, K-means and other clustering algorithms, K-nn nearest neighbours, logistic regression, Support vector machines. Evaluation metrics and cross-validation. Linear regression with Statsmodels. Other models such as trees and random forests.

4 References

- P. Bruce, H. Bruce (2020). *Practical Statistics for Data Scientists*. O'Reilly.
- A. Martelli, A. Ravenscroft, D. Ascher (2005). *Python CookBook*. O'Reilly.
- M. Harrison, T. Petrou (2020). *Pandas CookBook*. Packt.
- C. Bishop (2000). *Pattern Recognition and Machine Learning*. Springer
- T. Hastie, R. Tibshirani, J. Friedman (2013). *The Elements of Statistical Learning*. Springer.

5 Proposed Calendar

	Content		Content		Content		Content		Content
1	Installation and Setup Anaconda Jupyter Manage environments VSCode + Jupyter	2	Python crash course I Variables, data types Conditionals, for, range While, break, exception Functions	3	Python crash course II List, comprehension Dictionaries, tuples Strings Files I/O	4	Numpy Vectors and matrices Tensors Operations, reshape Basic plots	5	Visualisation Matplotlib Seaborn Pandas built-in Plotly
6	Pandas I Dataframes Access cells Select subdataframes Combine and merge	7	Pandas II Filter commands Handle missing data Dummy variables Handle time and dates	8	SQL Basic queries Aggregate functions Group by, order by Join and merge tables	9	SQL in Python SQLite MySQL	10	Probability Probability basics Bayes' law Discrete distributions Continuous distrib.
11	Descriptive statistics Histograms Mean, median, modes Variance Covariance, correlation	12	Inference statistics Confidence intervals Bootstrap, sampling Visualisation tools QQ-plots, densities	13	Hypothesis testing Hypothesis tests: to test normality compare distributions to compare samples	14	Data exploration Distributions Covariance analysis Stats and summaries Visualisation	15	Principal components PCA Dimension reduction Data exploration Biplots
16	Clustering K-means t-SNE Other methods	17	Regression I Linear regression Residual analysis Categorical variables Predictions	18	Regression II Hypothesis testing What if assumptions do not hold? Ridge and LASSO	19	Classification I K-nn Naive Bayes Logistic regression Train and test	20	Classification II SVM Decision trees Random forests Cross-validation