

Self-Created Course Plan: How to Become a Data Analyst / Data Scientist

Bernardo Di Chiara

Scope

This is a self-created plan of on-line courses which provides a time-efficient, convenient and affordable way to achieve the necessary competences to become a data analyst / data scientist. The courses are from well-known Universities and are all contained in Coursera platform ¹. The plan makes use of single courses and of specializations and is the result of heavy modifications to different proposals by Coursera, based on feedback got through job ads, job interviews and other sources.

Courses' Structure

The courses are instructor-led, have a perfect balance of theory and practice and their structure ensures that the student has mastered the contents after all the sections have been completed.

Each course contains at least video lessons with transcriptions and presentation material, weekly exams and a final peer-reviewed assignment. For almost all the courses an additional text book is available as well as practical exercises (like R labs). The best courses have also additional practice quiz and in-video quiz, in addition to the graded ones.

Plan's Structure

The courses have been organized in five modules.

The "Statistics with R" specialization by Duke University together with selected courses of the "Data Science" specialization by Johns Hopkins University allow acquiring the necessary competences to start to work as a data analyst/data scientist: statistics, R programming, data analysis techniques. Two courses complete this first set of ten on-line courses: an SQL course and a course about Hadoop. The courses are organized in 2 modules. The second module contains also some lessons and exercises from the statistical courses of the Data Science specialization which cover arguments not contained in the corresponding courses from Duke University. The statistical courses from Duke University are preferred since they are of better quality and didactically are definitely superior ². The second module contains also the only course which is not an on-line Coursera course. It is a classroom course about Big Data that, while it is not essential to proceed with the study plan, it adds a lot of value since it allows putting different tools and methods into context.

The third module allows learning how to visualize data and developing prediction models by using a second language: Python. Also, more advanced machine learning techniques (including vector machines, neural networks) and text mining are included. This is achieved with the specialization "Applied Data Science with Python" by the University of Michigan, which is preceded by a very basic Python course.

The fourth module is based on the specialization "Excel to MySQL: Analytic Techniques for Business" from Duke University, which teaches about business metrics and about using fancy data visualization tools (Tableau), among other things. Also, a course about more advanced statistics has been included.

Costs

Access to the video lessons is free (but for some courses is limited to the first week only). In order to access to the graded exams, the final peer-reviewed project and to obtain the certificate, it is necessary to enroll to the premium version. The cost is about 40€ per month for specializations. For single courses the allowed calendar time to complete the course with no extra fee might be slightly longer.

¹ The only exception is the Big Data course of the second module which is a classroom course in Helsinki.

² The first course of the first module is suggested for those who do not have a good mathematical background. Otherwise, it can be skipped. Similarly, the first course of the third module is suggested for those who do not have previous programming experience.

Basic: basic statistics, data analysis techniques and R (6,5 months)

- Data Science Math Skills (Duke) (4 weeks) COMPLETED Sep. 2017
 - Probability Theory
- Introduction to Probability and Data (Statistics with R, Duke) (5 weeks) COMPLETED Oct. 2017
 - Exploratory Data Analysis, Data Preparation, Data Visualization, RStudio
- Inferential Statistics (Statistics with R, Duke) (5 weeks) COMPLETED Nov. 2017
 - Hypothesis Tests, Confidence Intervals, ANOVA, Chi-square test, Bootstrapping, ...
- The Data Scientist's Toolbox (Data Science, Johns Hopkins) (4 weeks) COMPLETED Dec. 2017
 - Overview of Data Analysis Techniques, GitHub
- R Programming (Data Science, Johns Hopkins) (4 weeks) COMPLETED Jan. 2018
 - Creating R functions, using loop functions, debugging, profiling, ...
- Getting and Cleaning Data (Data Science, Johns Hopkins) (4 weeks) COMPLETED Feb. 2018
 - R Interfaces, Data Cleansing

Intermediate: machine learning, SQL and Big Data (4,5 months)

- Linear Regression and Modelling (Statistics with R, Duke) (4 weeks) COMPLETED Mar. 2018
 - Simple and Multiple Linear Regression, Logistic Regression
- More about Regression Models: studying Logistic Regression from the 3 videos in the Regression Models course of the Data Science specialization and from chapter 8.4 of OpenIntro Statistics book DONE
- More about EDA: studying Hierarchical Clustering, K-Means Clustering and Dimension Reduction (SVD and PCA) from the Exploratory Data Analysis course of the Data Science specialization (8 videos and 4 R labs) DONE
- Practical Machine Learning (Data Science, Johns Hopkins) (4 weeks) COMPLETED Apr. 2018
 - Clustering, Principal Component Analysis, Singular Value Decomposition, Decision Trees, Bagging, Random Forests, Boosting, Forecasting, ...
- SQL for Data Science (University of California, Davis) (4 weeks) COMPLETED May 2018
 - Data modelling, ER diagrams, retrieving data from multiple tables (subqueries, joins, unions), filtering, sorting, calculations, aggregations, manipulating text and dates, case statements, creating tables, SQLite
- Big Data (classroom course with instructors by TalentGate/Haaga-Helia University of Applied Sciences) (3 days) COMPLETED May 2018
 - Handling data volume, velocity and variety, supervised and unsupervised machine learning hands-on with KNIME
- Hadoop Platform and Application Framework (University of California, San Diego) (5 weeks) COMPLETED Jun. 2018
 - HDFS, MapReduce, Spark, HiveQL, HBase, Pig, ...

Advanced: more machine learning and Python (5,5 months)

- Getting Started with Python (Programming for Everybody, Univ. of Michigan) (7 weeks > 1 week) COMPLETED Sep. 2018
 - Conditional statements, iterations, functions
- Python Data Structures (Programming for Everybody, Univ. of Michigan) (7 weeks > 1 week) ONGOING
 - Files, lists, dictionaries, tuples
- Introduction to Data Science in Python (Applied Data Science with Python, University of Michigan) (4 weeks) high priority!
- Applied Plotting, Charting & Data Representation in Python (Applied Data Science with Python, University of Michigan) (4 weeks) high priority!
- Applied Machine Learning in Python (Applied Data Science with Python, University of Michigan) (4 weeks) high priority!
- Applied Text Mining in Python (Applied Data Science with Python, University of Michigan) (4 weeks)
- Applied Social Network Analysis in Python (Applied Data Science with Python, University of Michigan) (4 weeks)

Pro: business metrics and fancy visualization tools, more advanced statistics (6 months)

- Bayesian Statistics: From Concept to Data Analysis (University of California, Santa Cruz) (4 weeks)
- Business Metrics for Data-Driven Companies (Excel to MySQL: Analytic Techniques for Business, Duke) (4 weeks) important!
- Mastering Data Analysis in Excel (Excel to MySQL: Analytic Techniques for Business, Duke) (6 weeks) important!
- Data Visualization and Communication with Tableau (Excel to MySQL: Analytic Techniques for Business, Duke) (5 weeks) important!
- Managing Big Data with MySQL (Excel to MySQL: Analytic Techniques for Business, Duke) (5 weeks)

Plan's Flexibility

Note that the courses could be arranged in different order, provided that the courses within a specialization are taken in order and that statistics courses are completed before the first machine learning course. Also, the Python specialization assumes some programming and statistical knowledge and therefore it cannot be taken if the first two modules are not completed ³. The Hadoop course requires some Python knowledge and therefore postponing it after the first Python course of the third module could be an alternative.

The order described above is optimal for learning since it allows going back to a major topic and refresh it after a while. On the other hand, sticking to a certain specialization till all the wished courses are done and then moving to another one would be more efficient cost wise. There are many sessions per year for each course (monthly, bi-weekly or weekly, depending on the course).

The duration of the courses indicated here is the duration suggested by Coursera, which does not require a full-time commitment. Courses can be run slower or quicker but the peer-reviewed final project shall be submitted according to fixed schedules.

³ The (11 week) Machine Learning course from Stanford University could be an alternative to the first 3 courses of the specialization from the University of Michigan for the third module but it might make harder to follow the other two Python courses in the module.