

▼ DATA SCIENCE

by Dr Juan H Klopper

- Research Fellow
- School for Data Science and Computational Thinking
- Stellenbosch University



▼ INTRODUCTION

Welcome to this course resource, created for you by the School of Data Science and Computational thinking at Stellenbosch University.

There is little doubt that you have not come across the term **Data Science**. The ability to learn from data defines our scientific endeavour. **Data hides useful information**, whether that information be in spreadsheets, images, sounds, and even the written word. Modern tools are making extracting knowledge from data ever more powerful, and perhaps more importantly, easier.

Gaining access to the information in data has evolved. We currently have an abundance of riches, both in software and in hardware. Every day increases the amount of data that we generate and capture. We live in an age of information. Never before has it been so important to learn to tap into this information.

This course is all about tapping into that information using one of the most popular tools in existence. Python is an open source computer language. It is ideally suited to data science and its popularity has skyrocketed over the last decade. It is fast becoming one of the most important languages to speak, arguably starting to rival the importance of being able to communicate with the spoken and written word.

Great efforts are made all over the world to democratise access to this new language. Leading Universities have campus wide courses in data science using Python. Multiple online companies have sprung up teaching Data Science using Python. A quick YouTube search will show an endless supply of tutorials on data science and Python. Access to data science has truly been democratised.

This freedom of access to resources for what is a free and open language, is one of the cornerstones of the success of Python and Data Science as a whole. You will never get stuck on this new journey that you are about to embark upon. Help is available everywhere. Developing this new skill has never been easier.

In this course we will define Data Science, understand common terms and definitions used in Data Science, explore the tools of trade, and learn the most useful techniques in our pursuit to understand and learn from our data.

▼ DEFINING DATA SCIENCE

Data Science is about the exploration of data in order to learn from it, draw conclusion from it, predict with it, and use it to infer results to new situations, remaining cognisant of the uncertainty in our results.

As a science it brings together an array of approaches to understanding and using data. This includes data generation, data capture and storage, data verification and manipulation, statistics, and a variety of modern learning algorithms such as machine learning.

The common techniques and principles of this science makes it applicable to an enormous array of real-world situations. From earth and climate science to astronomy, cosmology, physics, biology, chemistry, engineering, healthcare, economics, politics, and so many more fields and subfields.

All of this made possible by today's modern computing infrastructure and access to the internet. This course is produced to run in the cloud. No need to even install anything on your local computer. These services for Data Science are available free of charge. Access to the internet is the only barrier to entry. With this statement I add my voice to the international plea to governments to see access to the internet as a basic human right. The more minds we add to our problems, the easier they will be to solve. With these solutions will come the solutions to poverty, hunger, climate change, illness, and so many more of our current and future challenges as a species on *a mote of dust suspended in a sunbeam*. *Carl Sagan*.

▼ THE TOOLS OF DATA SCIENCE

▼ COMPUTER LANGUAGES

Open source computer languages such as Python, R, and Julia are languages that allow us to interact with and instruct our computers to perform useful tasks. As with the spoken languages they have structure and rules.

The advantages of these languages for us as Data Scientists lie in their ease of access and use. As open source languages they are freely available, not controlled by some corporation. Instead they thrive on a community of developers and maintainers from all over the world that, in most cases, spend their free time to continuously update the tools. It is a tremendous achievement of humanity that such goodwill can lead to such success.

Newer versions of these languages are introduced at regular intervals, much like your cellphone operating system. The decision on what new changes are incorporated is decided by groups of developers. In the case of Python this used to be led by the creator of the language, Guido van Rossum, who perceived of the language all the way back in the 1980's. Julia is maintained by its developers out of the Massachusetts Institute of Technology.

All of these languages are expanded by packages, modules, and libraries. These are pieces of software written, once again mostly, by volunteers. They expand on the core language by providing convenient and useful additions to language. This course will teach you about the most common and powerful Python packages for Data Science.

Any computer language requires a piece of software in which we can type our code. The majority are integrated development environments (IDE). These environments are programs, much like Microsoft Word.

It has also become popular to use modern web browsers as coding environments. This allows us to write code in the cloud, without using the local resources of our machines. For the purposes of security and data safety, these browser environments can also be created on local machines as easy installs. For the purposes of this course, we will use Google Drive as there are no sensitive data that will be used.

▼ GOOGLE COLAB

Google Drive provides many tools such as Google Docs for writing reports and letters and Google Sheets for data capture as spreadsheets. It also provides a tool called **Colaboratory**.

Once you have created a free Google account, you will have access to Google email and a suite of other online apps. While you can use most browsers, Google Chrome will work best. Once signed into your account, you can type the following URL into the address bar of the browser:

<https://colab.research.google.com>

This will make Google Colab available as another Google online app in your account. A Google Colab file or document is termed a **notebook**. It is based on the Jupyter notebook (which you can install and use on your local machine).

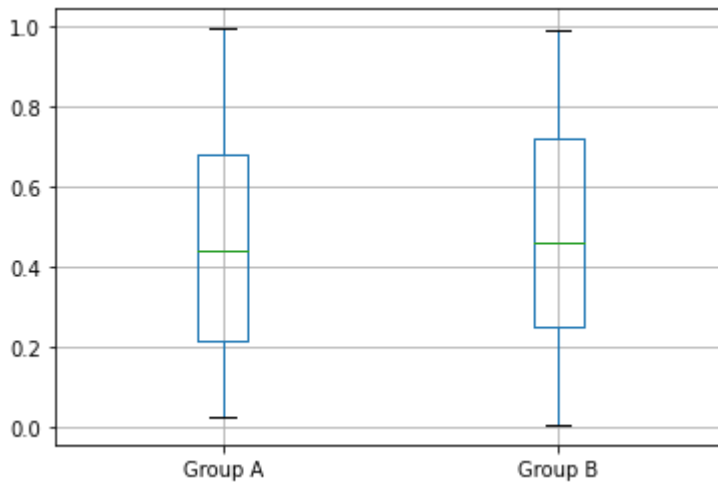
▼ COLAB NOTEBOOK

This document was created as a Colab notebook. It allows for writing of formatted words, sentences, paragraphs, and sections as you can see here. We can even import and show images (as seen by our logo at the top of this notebook). It also allows for the typing and execution of Python code as can be seen below. Don't be concerned about the code. We will learn all about it in upcoming notebooks.

```
1 %%config InlineBackend.figure_format = "retina"
2
3 import numpy as np
4 import pandas as pd
5
6 df = pd.DataFrame(
7     np.random.rand(100, 2),
8     columns=['Group A', 'Group B']
9 )
10
11 df.describe()
```

Group A Group B

```
1 df.boxplot();
```



Both *normal written language* and code are entered into **cells**. A cell can be a text cell or a code cell. In the two code cells above, we generated some random variables and first calculated summary statistics for the two variables and then generated a box-and-whisker plot of the data.

When hovering in the middle of the page, between two cells (or at the bottom of the last cell), two buttons pop up. They are marked + Code and + Text. They generate new cell right below the buttons.

The text cells have a tool bar, much like a simple program. The buttons in the toolbar allow for formatting of text and insertion of images and links.

Text cells can format text using Markdown language. Markdown language uses special characters to format text. For instance adding an underscore before and after a word or words, prints them in *italics*. Double underscores print in **bold**.

Text cells can also print LaTeX. LaTeX is a language for mathematical typesetting. In (1) below, we see the equation for the normal distribution given a mean, μ , and a standard deviation, σ .

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad (1)$$

The top of the notebook also has a toolbar with the usual File, Edit, View, and other buttons. We will become familiar with these throughout the course.

Since this is a Google app, there is also some familiar buttons on the top right. These allow us to share a notebook with our coworkers and collaborators. We can also leave comments, just as

with a Google Doc.

At the top of the left pane (which is collapsable), are four icons for a table of content, search, code snippets, and a folder explorer. The first is a very convenient *table of contents*. When we use titles and subtitles, these appear in the table of contents.

The beauty of the notebook format as coding environment, is that it generates a document that serves as a very useful report. We can include *normal text*, images, links, formatted input, and our code. It makes for expressive research documents that can be used for collaboration, sharing, and presentation.

Below, is a brief example of a research project. The data comes from a real-life survey run by an organisation called Kaggle. You can sign up for a free account. Kaggle even provides a very similar notebook environment to Google Colab.

The survey is run every year and is completed by Data Scientists from across the Globe. It gives a unique insight into who Data Scientists are, what tools they use, and even includes their salaries.

This example is for demonstration purposes only. As mentioned above, do not be concerned about the code. Appreciate the results and the use of a notebook as research document, combining text, code, and results to gain insight.

▼ DATA SCIENCE SURVEY EXAMPLE

▼ PACKAGES USED IN THIS NOTEBOOK

Packages are the extension to the language developed by enthusiasts, academics, and professionals. Below, we import these for use in the example project. You will notice pound symbols (hashtags) followed by some green text. This symbol in Python denotes a code comment. Python ignores everything in a line after a pound symbol. It is good practice to leave code comments, both to your future self, and to others who might read your code.

```
1 # Statistics module of the scipy package
2 from scipy import stats
```

```
1 # Plotting package import
2 import plotly.graph_objects as go
3 import plotly.express as px
```

```

4 import plotly.figure_factory as ff
5 import plotly.io as pio
6 pio.templates.default = 'plotly_white'

1 # Format tables printed to the screen (don't put this on the same line as the c
2 %load_ext google.colab.data_table

1 # For connecting to Google Drive
2 from google.colab import drive

```

▼ CONNECTING TO GOOGLE DRIVE

The data is in the Google Drive in which this notebook was created. It was saved in a subfolder. We need to log into our account again to generate a security key which we copy and paste into the cell below. Simply follow the on-screen instructions and allow the required permissions. There is also some code to navigate to the subfolder.

```

1 # Connect to your Google Drive
2 drive.mount('/gdrive', force_remount=True)

Mounted at /gdrive

1 '/gdrive/My Drive/Stellenbosch University/School for Data Science and Computati
   /gdrive/My Drive/Stellenbosch University/School for Data Science and Computat

```

▼ IMPORTING THE DATA

The data is saved in a comma separated values (CSV) spreadsheet file.

```
1 df = pd.read_csv('kaggle_survey_2020_responses.csv', low_memory=False)
```

The data is displayed as a spreadsheet to make sure that it imported properly.

```

1 # Drop first row containing metadata
2 df.drop(0, axis=0, inplace=True)
3 # First five rows of data
4 df[:5]

```

	Time from Start to Finish (seconds)	Q1	Q2	Q3	Q4	Q5	Q6	Q7_Part_1	Q7_Part_2
1	1838	35-39	Man	Colombia	Doctoral degree	Student	5-10 years	Python	R
2	289287	30-34	Man	United States of America	Master's degree	Data Engineer	5-10 years	Python	R
3	860	35-39	Man	Argentina	Bachelor's degree	Software Engineer	10-20 years	NaN	NaN
4	507	30-34	Man	United States of America	Master's degree	Data Scientist	5-10 years	Python	NaN
		30-			Master's	Software	3-5		

The number of observations and variables are returned to give us an indication of the dimensions of our data set.

```
1 # The shape of the dataframe object
2 df.shape

(20036, 355)
```

We note that there are 200036 observations (respondents to the survey), each completing data for 355 variables.

▼ RESPONDENT QUALIFICATIONS

Calculating the frequency of different (highest) qualifications.

```
1 # Frequency of sample space elements for Q4 (highest qualification)
2 df.Q4.value_counts()

Master's degree      7859
Bachelor's degree   6978
Doctoral degree     2302
Some college/university study without earning a bachelor's degree 1092
Professional degree    699
I prefer not to answer 399
No formal education past high school 240
Name: Q4, dtype: int64
```

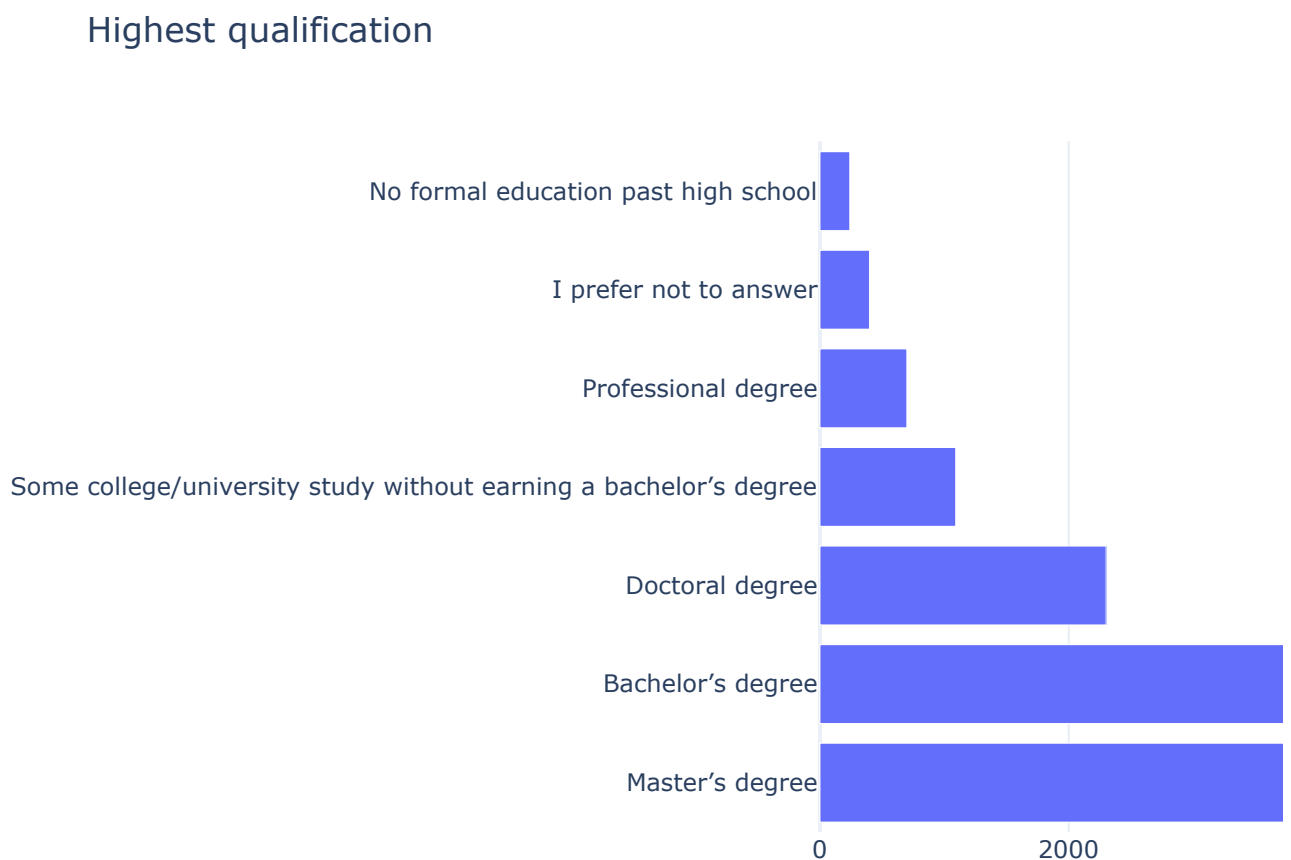

The majority of candidates have Master's and Bachelor's degrees. More than 2000 have a PhD.

A bar plot provides a visual analysis of the data. With the plotly packages, we can hover over elements in the plot for more information. We can also pan, zoom, and save the image to disk.

```

1 go.Figure(
2     go.Bar(
3         y=df.Q4.value_counts().index.tolist(),
4         x=df.Q4.value_counts().values.tolist(),
5         orientation='h'
6     )
7 ).update_layout(
8     title='Highest qualification',
9     xaxis={'title': 'Count'}
10 )

```



▼ PYTHON AND R

Next up, we investigate the number of Data Scientist that use Python and R (and the previous leader, C) in their daily work. We also include Structured Query Language used in databases.

```
1 go.Figure(
```

```
2     go.Bar(
3         x=[ 'Total' ],
4         y=[df.shape[0]],
5         name='Total'
6     )
7 ).add_trace(
8     go.Bar(
9         x=[ 'Python' ],
10        y=df.Q7_Part_1.value_counts().values.tolist(),
11        name='Python'
12    )
13 ).add_trace(
14     go.Bar(
15        x=[ 'R' ],
16        y=df.Q7_Part_2.value_counts().values.tolist(),
17        name='R'
18    )
19 ).add_trace(
20     go.Bar(
21        x=[ 'C' ],
22        y=df.Q7_Part_4.value_counts().values.tolist(),
23        name='R'
24    )
25 ).add_trace(
26     go.Bar(
27        x=[ 'SQL' ],
28        y=df.Q7_Part_3.value_counts().values.tolist(),
29        name='SQL'
30    )
31 ).update_layout(title='Commonly used Data Science languages',
32                 xaxis={'title': 'Count'},
33                 yaxis={'title': 'Language'})
```

Commonly used Data Science languages

30.

We note that Python is the leading language used in Data Science.

▼ CODING ENVIRONMENTS

We investigate the commonly used coding environments and compare Jupyter notebooks (on which Google Colab is based), as well as RStudio, PyCharm, and Microsoft Visual Studio Cloud. Free versions of these environments are either in the cloud or can be installed on your local machine.

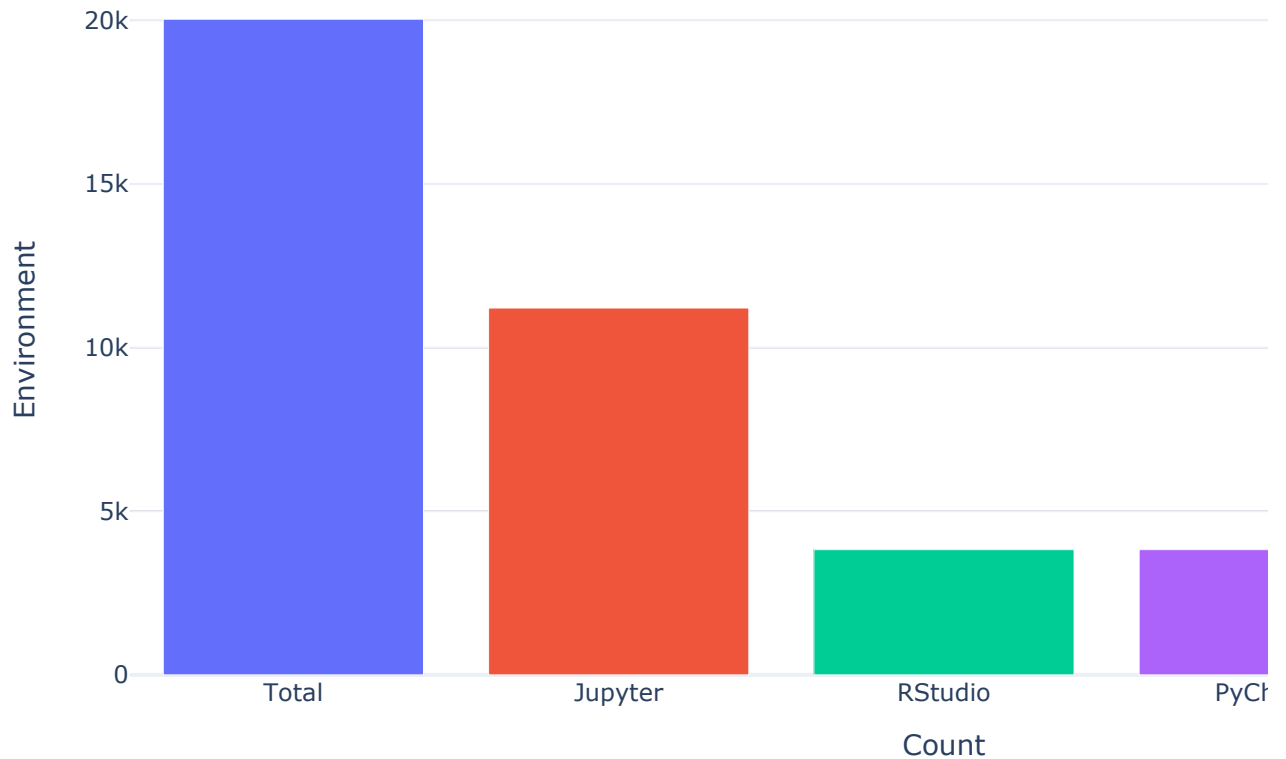
31.

```

1 go.Figure(
2     go.Bar(
3         x=[ 'Total' ],
4         y=[df.shape[0]],
5         name='Total'
6     )
7 ).add_trace(
8     go.Bar(
9         x=[ 'Jupyter' ],
10        y=df.Q9_Part_1.value_counts().values.tolist(),
11        name='Jupyter'
12    )
13 ).add_trace(
14    go.Bar(
15        x=[ 'RStudio' ],
16        y=df.Q9_Part_2.value_counts().values.tolist(),
17        name='R'
18    )
19 ).add_trace(
20    go.Bar(
21        x=[ 'PyCharm' ],
22        y=df.Q7_Part_5.value_counts().values.tolist(),
23        name='R'
24    )
25 ).add_trace(
26    go.Bar(
27        x=[ 'VS Code' ],
28        y=df.Q7_Part_4.value_counts().values.tolist(),
29        name='R'
30    )
31 ).update_layout(title='Commonly used Data Science coding environments',
32                 xaxis={'title': 'Count'},
33                 yaxis={'title': 'Environment'})

```

Commonly used Data Science coding environments



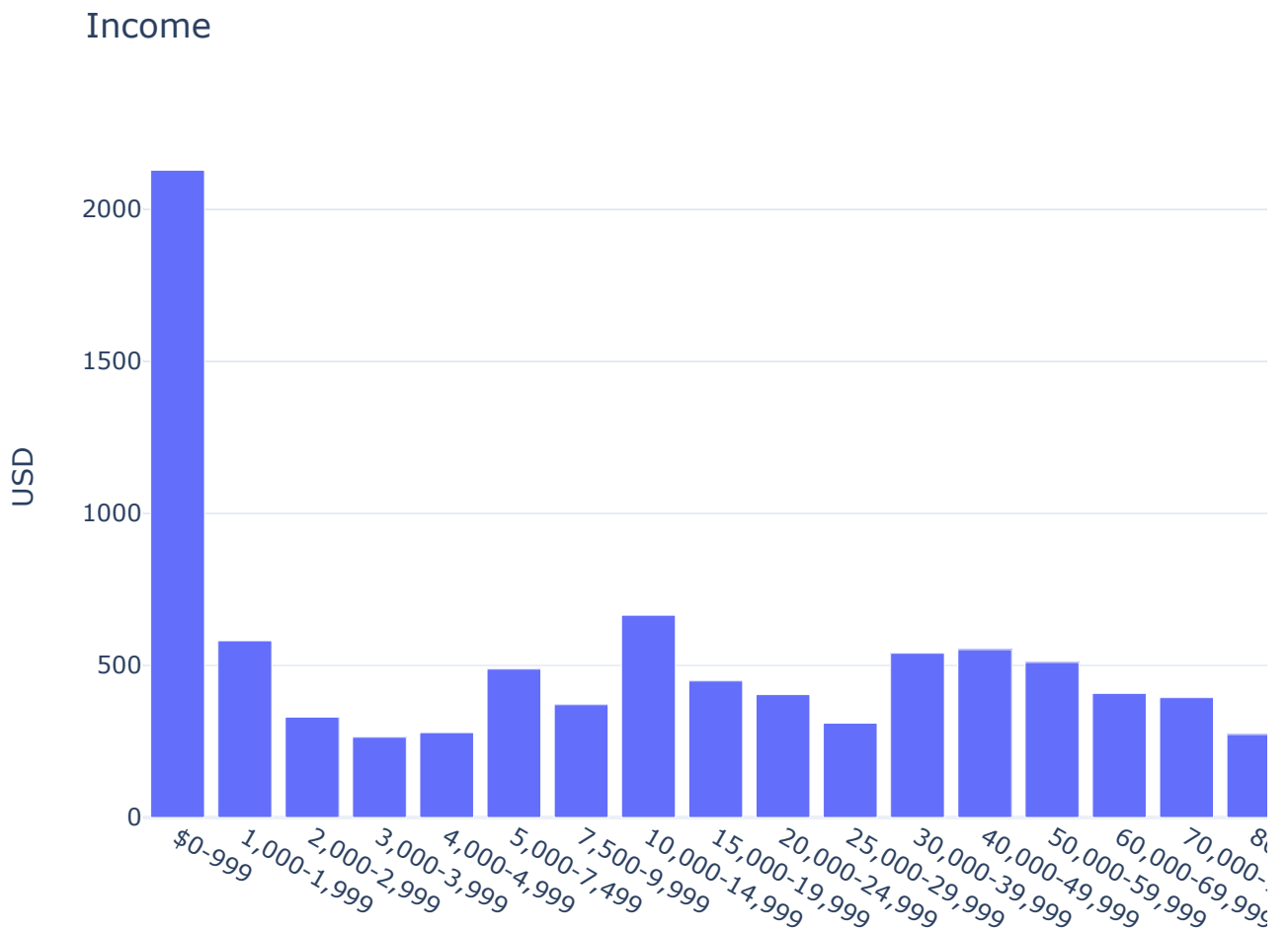
Jupyter notebooks are very popular. They provide an excellent coding environment for the analysis of data.

▼ INCOME

The survey asked respondents to select a salary bracket. We visualise the frequency of all of the brackets.

```
1 # Saving an ordered list of the brackets as dictionary key-value pairs
2 income_brackets = {
3     '$0-999': 2128,
4     '1,000-1,999': 581,
5     '2,000-2,999': 330,
6     '3,000-3,999': 264,
7     '4,000-4,999': 279,
8     '5,000-7,499': 488,
9     '7,500-9,999': 371,
10    '10,000-14,999': 665,
11    '15,000-19,999': 449,
12    '20,000-24,999': 404,
13    '25,000-29,999': 310,
14    '30,000-39,999': 540,
15    '40,000-49,999': 552,
```

```
16 '50,000-59,999':510,  
17 '60,000-69,999':408,  
18 '70,000-79,999':394,  
19 '80,000-89,999':273,  
20 '90,000-99,999':280,  
21 '100,000-124,999':573,  
22 '125,000-149,999':315,  
23 '150,000-199,999':347,  
24 '200,000-249,999':115,  
25 '250,000-299,999':48,  
26 '300,000-500,000':55,  
27 '> $500,000':50  
28 }  
29  
30 go.Figure(  
31     go.Bar(  
32         x=list(income_brackets.keys()),  
33         y=list(income_brackets.values())  
34     )  
35 ).update_layout(  
36     title='Income',  
37     yaxis={'title':'USD'}  
38 )
```



There is a definite bump at the 100000 USD mark. There are some well paid Data Scientist.

▼ INCOME DIFFERENCES FOR WOMEN

The underrepresentation of groups other than men is a continued problem in Data Science. We use the 100000 USD mark to distinguish high income individuals from the rest.

```
1 # All high income brackets (above $100000)
2 high_income = ['100,000-124,999',
3               '125,000-149,999',
4               '150,000-199,999',
5               '200,000-249,999',
6               '250,000-299,999',
7               '300,000-500,000',
8               '> $500,000']
```

```
1 # Add a new column for income above and below $100000
2 df['HighIncome'] = np.where(df.Q24.isin(high_income), 'Yes', 'No')
```

For the sake of simplicity and to answer our specific research question regarding women, we only extract those that identified as *Man* or *Woman*.

```
1 # Select only binary sex participants
2 binary_gender = df.loc[(df.Q2 == 'Man') | (df.Q2 == 'Woman')]
```

A contingency table of observed values shows the scale of the problem.

```
1 # Contingency table of observed values
2 pd.crosstab(
3     binary_gender.Q2,
4     binary_gender.HighIncome,
5     margins=True
6 )
```

1 to 3 of 3 entries

Q2	No	Yes	All
Man	14494	1295	15789
Woman	3706	172	3878
All	18200	1467	19667

Show per page

```
1 # Relative frequency observed value contingency table
2 pd.crosstab(
3     binary_gender.Q2,
4     binary_gender.HighIncome,
5     margins=True,
6     normalize=True
```

7)

1 to 3 of 3 entries 

Q2	No	Yes	All
Man	0.73697055982102	0.06584634158743072	0.8028169014084507
Woman	0.1884374841104388	0.00874561448111049	0.19718309859154928
All	0.9254080439314588	0.07459195606854122	1.0

Show per page

When selecting only those who identified as either man or women, only 19.7% were woman. Of these, 4.4% were in the high income brakcet. For men, the figure as almost double at 8.2%.

A χ^2 test for independence can investigate the dependence between sex and income.

```
1 stats.chi2_contingency([[14494, 1295], [3706, 172]])

(63.444674249142516,
 1.6493472192256754e-15,
 1,
 array([[14611.26760563, 1177.73239437],
        [ 3588.73239437, 289.26760563]]))
```

There was a significant association between gender and high income status, χ^2 value 63, 4, p value < 0.01 .

▼ UNITED STATES OF AMERICA VS SOUTH AFRICA

We investigate of there is a dependence between country (United States of America and South Africa) and income level (again using the 100000 USD cut-off).

```
1 usa_sa_df = df.loc[df.Q3.isin(['South Africa', 'United States of America'])]

1 pd.crosstab(
2     usa_sa_df.Q3,
3     usa_sa_df.HighIncome,
4     margins=True
5 )
```

1 to 3 of 3 entries 

Q3	No	Yes	All
South Africa	133	8	141
United States of America	1406	831	2237
All	1539	839	2378

Show per page

```

1 stats.chi2_contingency(
2     pd.crosstab(
3         usa_sa_df.Q3,
4         usa_sa_df.HighIncome
5     )
6 )

(56.17463007230668,
 6.631168127655343e-14,
 1,
 array([[ 91.25273339,  49.74726661],
        [1447.74726661,  789.25273339]]))

```

A total of 2237 respondents where from the USA compared to only 141 from South Africa. Thirty-seven percent of the American respondents where in the high income bracket compared to only 5.7% of the South African respondents, χ^2 value 56.2, p value < 0.01 .

▼ PREDICTING INCOME BRACKET USING MACHINE LEARNING

We generate a random forest machine learning model and use age, sex, country, highest qualification, and whether the respondent uses Python as feature variable in a supervised learning application, predicting income level (high or not).

```

1 # Selecting the feature and target variables
2 variables = [
3     'Q1',
4     'Q2',
5     'Q3',
6     'Q4',
7     'Q7_Part_1',
8     'HighIncome'
9 ]
10
11 rf_data = df[variables]

```

We will use the TensorFlow DecisionForest module from Google to create the model.

```

1 # Install this package if required
2 !pip install tensorflow_decision_forests

```

```

Collecting tensorflow_decision_forests
  Downloading https://files.pythonhosted.org/packages/3a/34/10f20fe95d9882b82/
  |████████████████████████████████████████████████████████████████████████████████| 6.2MB 6.5MB/s
Requirement already satisfied: tensorflow~=2.5 in /usr/local/lib/python3.7/dist-packages (2.5.0)
Requirement already satisfied: absl-py in /usr/local/lib/python3.7/dist-packages (0.10.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (1.19.5)
Requirement already satisfied: wheel in /usr/local/lib/python3.7/dist-packages (0.34.2)
Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (1.1.5)
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (1.16.0)
Requirement already satisfied: grpcio~=1.34.0 in /usr/local/lib/python3.7/dist-packages (1.34.0)

```



```

Requirement already satisfied: flatbuffers~=1.12.0 in /usr/local/lib/python3.
Requirement already satisfied: gast==0.4.0 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: termcolor~=1.1.0 in /usr/local/lib/python3.7/d
Requirement already satisfied: keras-nightly~=2.5.0.dev in /usr/local/lib/pyt
Requirement already satisfied: astunparse~=1.6.3 in /usr/local/lib/python3.7/
Requirement already satisfied: keras-preprocessing~=1.1.2 in /usr/local/lib/p
Requirement already satisfied: tensorboard~=2.5 in /usr/local/lib/python3.7/d
Requirement already satisfied: h5py~=3.1.0 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: opt-einsum~=3.3.0 in /usr/local/lib/python3.7/
Requirement already satisfied: wrapt~=1.12.1 in /usr/local/lib/python3.7/dist
Requirement already satisfied: tensorflow-estimator<2.6.0,>=2.5.0rc0 in /usr/
Requirement already satisfied: google-pasta~=0.2 in /usr/local/lib/python3.7/
Requirement already satisfied: typing-extensions~=3.7.4 in /usr/local/lib/pyt
Requirement already satisfied: protobuf>=3.9.2 in /usr/local/lib/python3.7/di
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/pytho
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.7/di
Requirement already satisfied: tensorboard-plugin-wit>=1.6.0 in /usr/local/li
Requirement already satisfied: tensorboard-data-server<0.7.0,>=0.6.0 in /usr/
Requirement already satisfied: werkzeug>=0.11.15 in /usr/local/lib/python3.7/
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.
Requirement already satisfied: setuptools>=41.0.0 in /usr/local/lib/python3.7
Requirement already satisfied: google-auth-oauthlib<0.5,>=0.4.1 in /usr/local
Requirement already satisfied: google-auth<2,>=1.6.3 in /usr/local/lib/python
Requirement already satisfied: cached-property; python_version < "3.8" in /us
Requirement already satisfied: importlib-metadata; python_version < "3.8" in
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /us
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/
Requirement already satisfied: requests-oauthlib>=0.7.0 in /usr/local/lib/pyt
Requirement already satisfied: cachetools<5.0,>=2.0.0 in /usr/local/lib/pytho
Requirement already satisfied: rsa<5,>=3.1.4; python_version >= "3.6" in /usr
Requirement already satisfied: pyasn1-modules>=0.2.1 in /usr/local/lib/python
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-pac
Requirement already satisfied: oauthlib>=3.0.0 in /usr/local/lib/python3.7/di
Requirement already satisfied: pyasn1>=0.1.3 in /usr/local/lib/python3.7/dist
Installing collected packages: tensorflow-decision-forests
Successfully installed tensorflow-decision-forests-0.1.7

```

```

1 # Importing the required packages
2 import tensorflow_decision_forests as tfdf
3 import tensorflow as tf

1 classes = rf_data.HighIncome.unique().tolist() # A list of the two income classes
2 classes

    ['No', 'Yes']

1 rf_data.HighIncome = rf_data.HighIncome.map(classes.index)

```

The data set is split into a training set, which the decision forest model will use to *learn* from, and a test set, against which we can evaluate the trained machine learning model.

```
1 def split(ds, r=0.3):
```

```

1 test_ind = np.random.rand(len(ds)) < r
2
3
4 return ds[~test_ind], ds[test_ind]

1 rf_data_train, rf_data_test = split(rf_data)

```

There is class imbalance with 92% of respondents not in the high income bracket. Simply predicting *NOT HIGH INCOME BRACKET* for all cases, makes you 93% correct all the time. Our model will have to beat this metric.

```

1 rf_data_train.HighIncome.value_counts(normalize=True)

0    0.925234
1    0.074766
Name: HighIncome, dtype: float64

```

The decision forest class is instantiated and the model compiles, prior to passing the training data to the model.

```

1 rf_model = tfdf.keras.RandomForestModel()

1 rf_model.compile(
2     metrics=['accuracy']
3 )

1 rf_data_train = tfdf.keras.pd_dataframe_to_tf_dataset(
2     rf_data_train,
3     label='HighIncome'
4 )
5 rf_data_test = tfdf.keras.pd_dataframe_to_tf_dataset(
6     rf_data_test,
7     label='HighIncome'
8 )

1 rf_model.fit(
2     x=rf_data_train
3 )

219/219 [=====] - 5s 1ms/step
<tensorflow.python.keras.callbacks.History at 0x7fbc2022a5d0>

```

The model is now evaluated on data it has never seen (the test data).

```

1 evaluation = rf_model.evaluate(
2     rf_data_test,
3     return_dict=True # Returning a dictionary of metrics
4 )

```

```
95/95 [=====] - 1s 5ms/step - loss: 0.0000e+00 - acc:
```

```
1 evaluation.keys() # Metric dictionary keys
```

```
dict_keys(['loss', 'accuracy'])
```

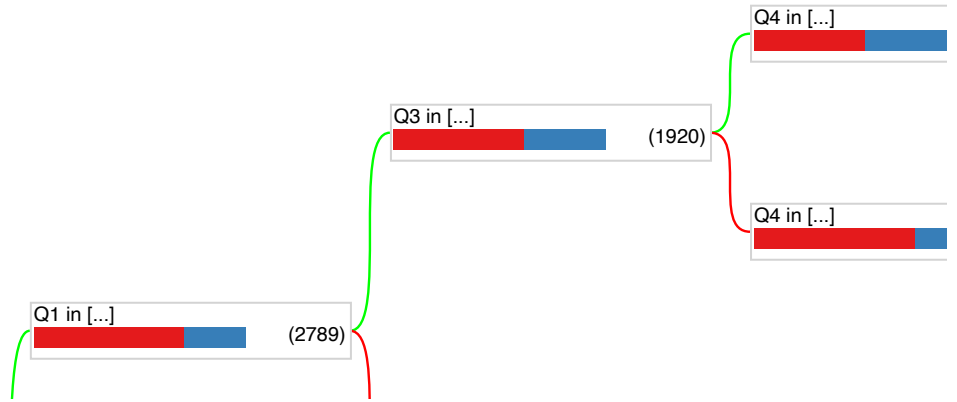
```
1 evaluation.values() # Metric values
```

```
dict_values([0.0, 0.9239148497581482])
```

We see an accuracy in excess of our baseline model.

Below, we plot the decision algorithm learned by this model.

```
1 tfdf.model_plotter.plot_model_in_colab(  
2     rf_model,  
3     tree_idx=0,  
4     max_depth=4  
5 )
```



```
1 rf_model.make_inspector().variable_importances()
```

```
{'NUM_AS_ROOT': [{"Q3" (4; #2), 300.0}]}
```

We note that the country of origin was the *most important variable* in the model.

CONCLUSION

This notebook introduced the topic, its tools, and a brief example of an Data Science project, incorporating a machine learning model.

At the end of this course you will be empowered to create all of this and much, much more.

1

