## 15.1.1    Introduction to R

Jeremy is ecstatic that he has been given the opportunity to lead the data analytics team. He's confident that his 10 years working with the manufacturing and research team has provided him with sufficient expertise in the subject matter. However, he is much less confident about his statistics background and his programming ability. He took some stats and programming courses in college, but it has been a long time since he had to think about either subject.

But Jeremy has never been one to walk away from a challenge. He knows that if he fully commits to learning stats and R programming, he'll feel comfortable in no time! With his previous experience in programming, learning R seems to be a good starting point to prepare for his new role.

R is a programming language that has a variety of uses in data science. R has solidified itself in academia and industry as one of the go-to programming languages for statistical modelling and hypothesis testing. In recent years, R developers have extended R's capabilities to generate machine learning algorithms and other advanced models to ensure that R can be used in every stage of data analytics.

## Benefits of R

R is a versatile and extensible programming language with many benefits. One of the benefits of using an interpreted programming language such as R (or Python) is that the analysis scripts are written in plaintext. The versions of plaintext files are easy to control using tools such as Git, which

means that **R analysis scripts** (or **RScripts**) are highly reproducible and easy to share with peers and collaborators.

Another benefit to using R is that the R programming language was specifically designed for data analysis. This means that the process of loading in a dataset, visualizing the data, and performing statistical tests is straightforward and easy to interpret. In fact, many of the statistical tests in Python have been directly ported from R due to how well they were implemented. In addition to the native statistical functions, there are many other useful data transformation and modelling libraries, such as the tidyverse package, that simplify the process of ETL and visualizations.

## Drawbacks of R

Still, R is not perfect. The biggest drawback is its licensing. R and most of R's libraries are licensed as General Public License, version 2 (GPL 2). This means that if you program or model anything using R, GPL forces your application, program, or script to be open source.

In many personal and academic uses, this is not a problem because you're either (a) not trying to monetize your program or (b) going to publish your analysis and findings. However, if you are working for a company with intellectual property, or proprietary data and programs, this can be an issue. Therefore, many companies use R for internal analysis and regulatory testing, but use Python for any application or script that contains proprietary information.

Despite the licensing drawback, R is still a highly valuable programming language for data analysis and is used by data professionals at all levels across many fields.
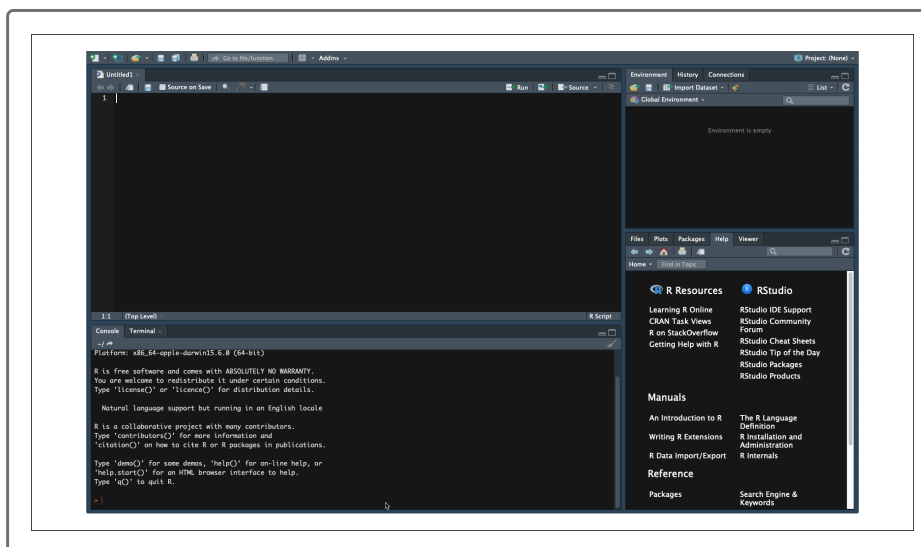
## RStudio Integrated Development Environment

Just as Jupyter Notebooks are an integrated development environment (IDE) used to help design and test Python scripts, RStudio is an IDE used to help design and test RScripts. RStudio provides users a graphical user interface (GUI) with multiple dynamic windows and perpetual access to their RScript and R console.
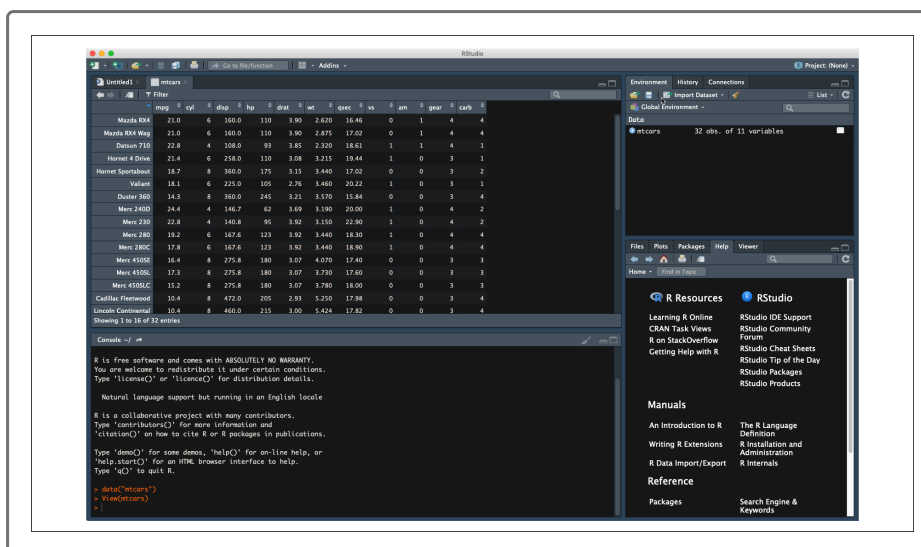
Similar to Jupyter Notebooks, RStudio enables users to test their analysis scripts line by line while allowing users to view different environment variables and outputs. This means that for each line of code written and

executed, users can verify the results and troubleshoot any problems quickly and easily.

The following image shows what an empty session in RStudio will look like:



And here's what an active RStudio looks like. Note the DataFrame in the top left pane and the mtcars object in the top right pane:



Now that you know what R and RStudio do, let's install them on your machine.