

# Intro to R

Ezra Garcia | ASTRAI | W07

27.11.2024

# What is R?

- An open-source programming language and software environment.
- Developed by statisticians Ross Ihaka and Robert Gentleman in the early 1990s and has since become widely used in academia, research, and industry.
- Used for statistical computing, data analysis, and graphical representation.



# Common use cases

- **Data Analysis:** R allows users to manipulate, analyze, and interpret large datasets efficiently.
- **Statistical Modeling:** It provides a wide range of statistical techniques such as linear and nonlinear modeling, classification, clustering, and time-series analysis.
- **Data Visualization:** With powerful libraries like ggplot2, R can create high-quality plots, graphs, and charts.
- **Machine Learning:** R offers a range of machine learning algorithms for classification, regression, clustering, and dimensionality reduction.

# Field-specific applications

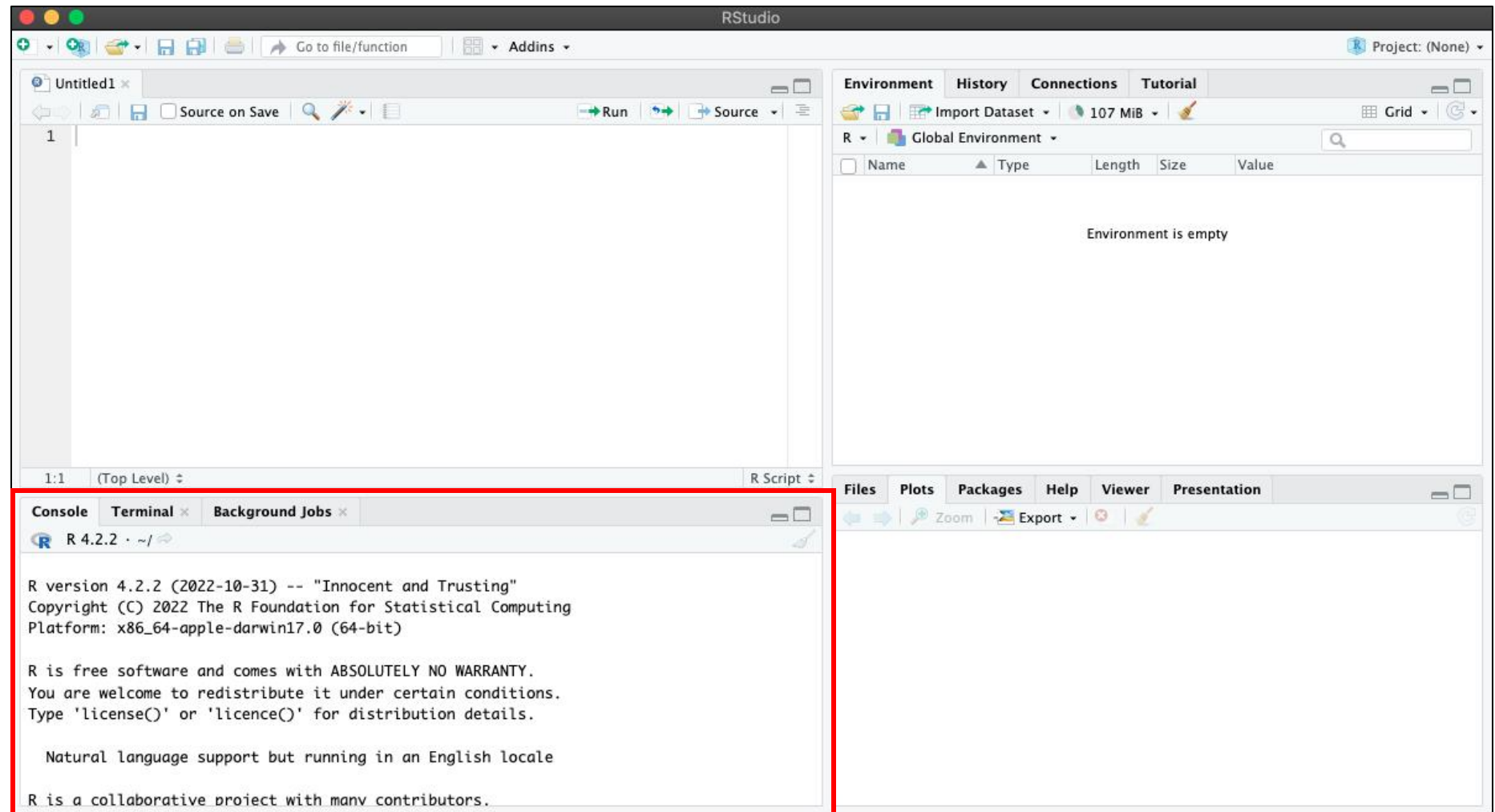
- **Bioinformatics:** R is commonly used in bioinformatics for analyzing biological data like gene expression datasets.
- **Finance:** Financial analysts use R for quantitative finance, portfolio management, and risk analysis.
- **Econometrics:** Economists use R for analyzing economic data, building econometric models, and conducting hypothesis testing.
- **Transportation:** R can be leveraged for traffic flow analysis, route optimization, demand forecasting, crash data analysis, spatial analysis, survey analysis, etc.

# What makes R special?

- **Open Source:** R is free to use.
- **Comprehensive Libraries:** R offers a multitude of packages (e.g., ggplot2, dplyr, caret, shiny) for various types of data analysis and visualization.
- **Cross-Platform:** R works across various platforms (Windows, Mac, Linux).
- **Extensive Community Support:** R has an active user community, meaning constant improvements and a wealth of resources.
- **High-Quality Graphics:** R's visualization capabilities allow for the creation of publication-ready graphics.
- **Flexibility with Data:** R can handle various types of data, including structured (CSV, Excel) and unstructured data (text, images).

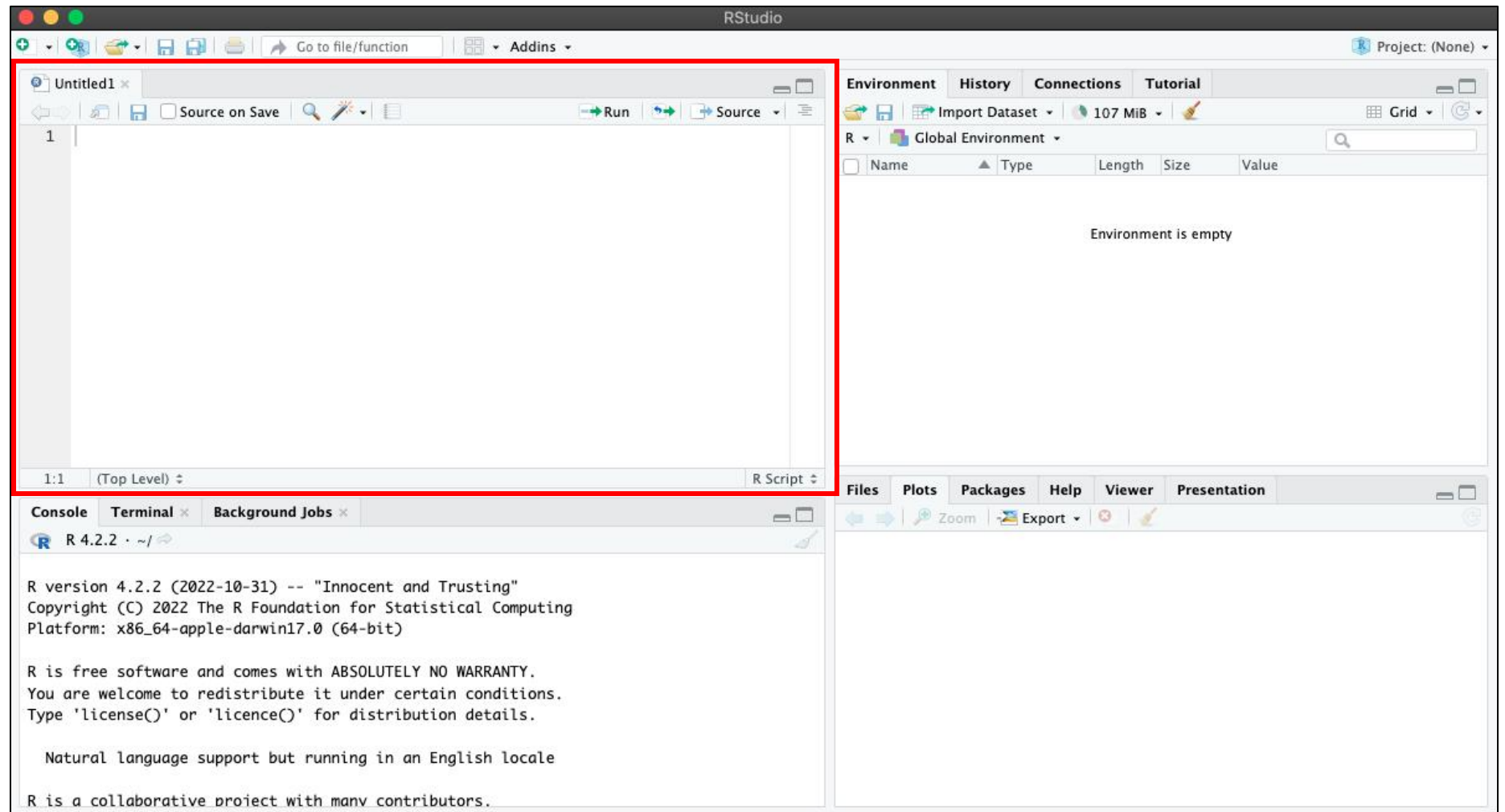
# The RStudio user-interface

Console window: where output and computations are displayed.

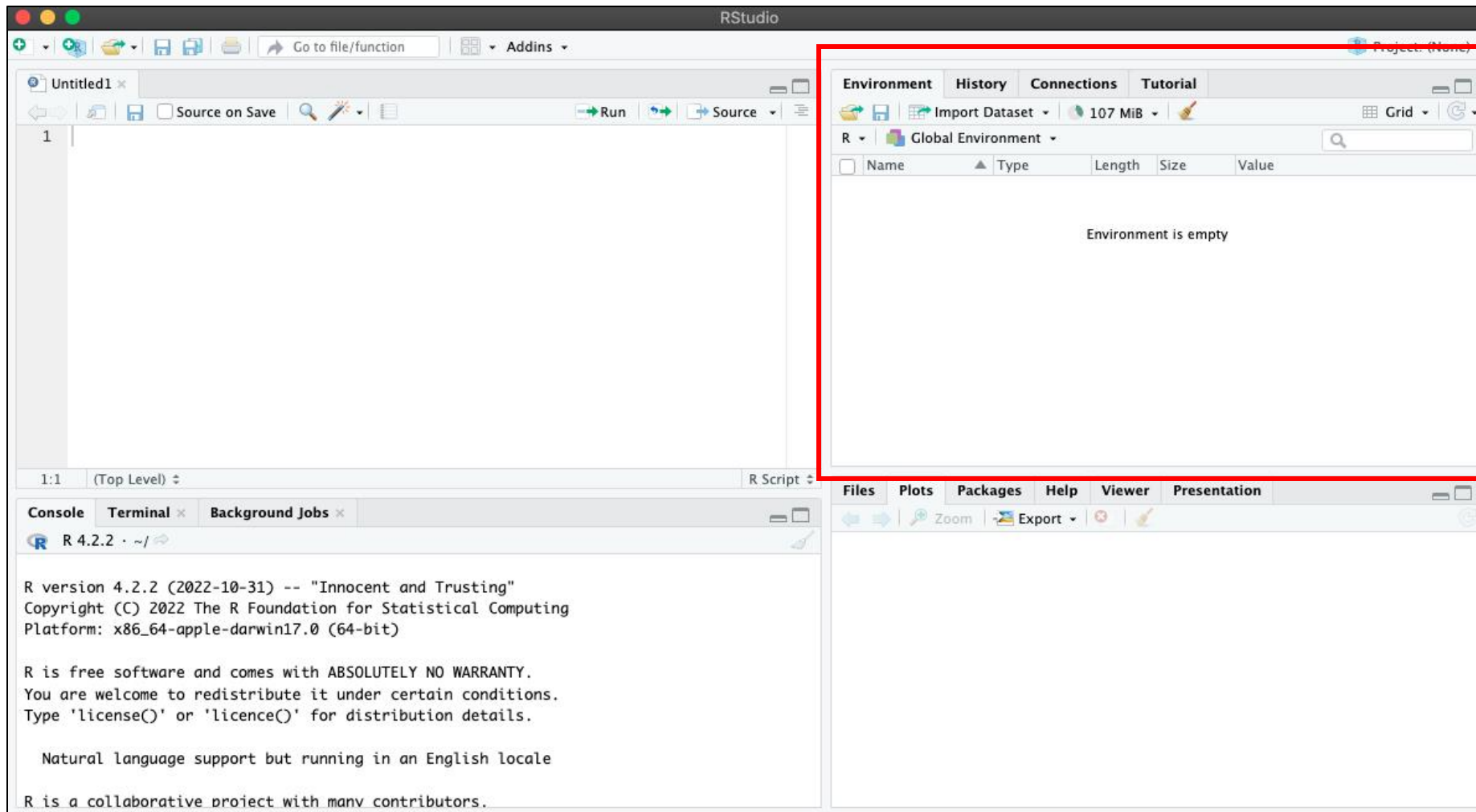


# The RStudio user-interface

Source (code-editing) window:  
used for editing a script.



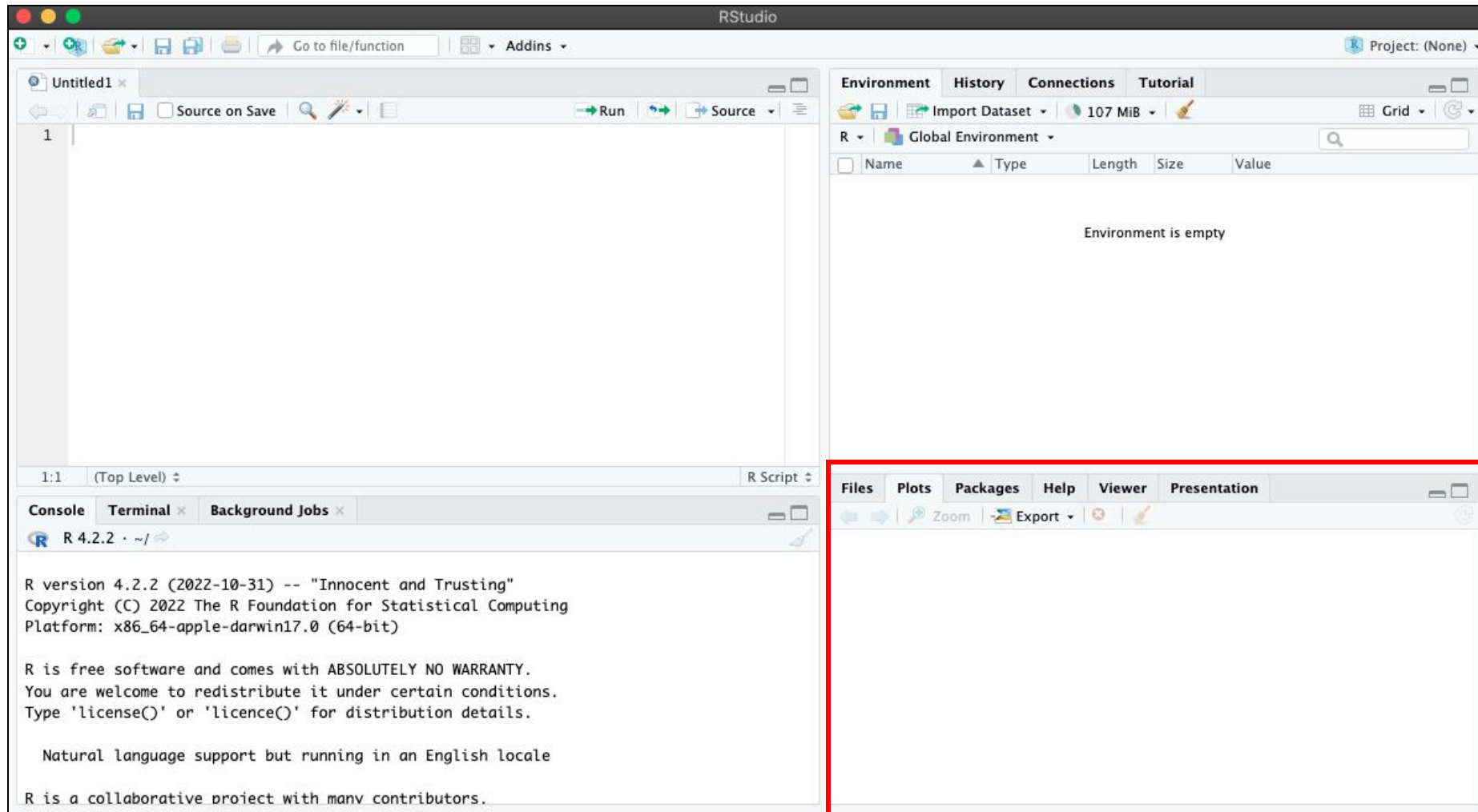
# The RStudio user-interface



Workspace and History window: where objects available for computations are displayed.

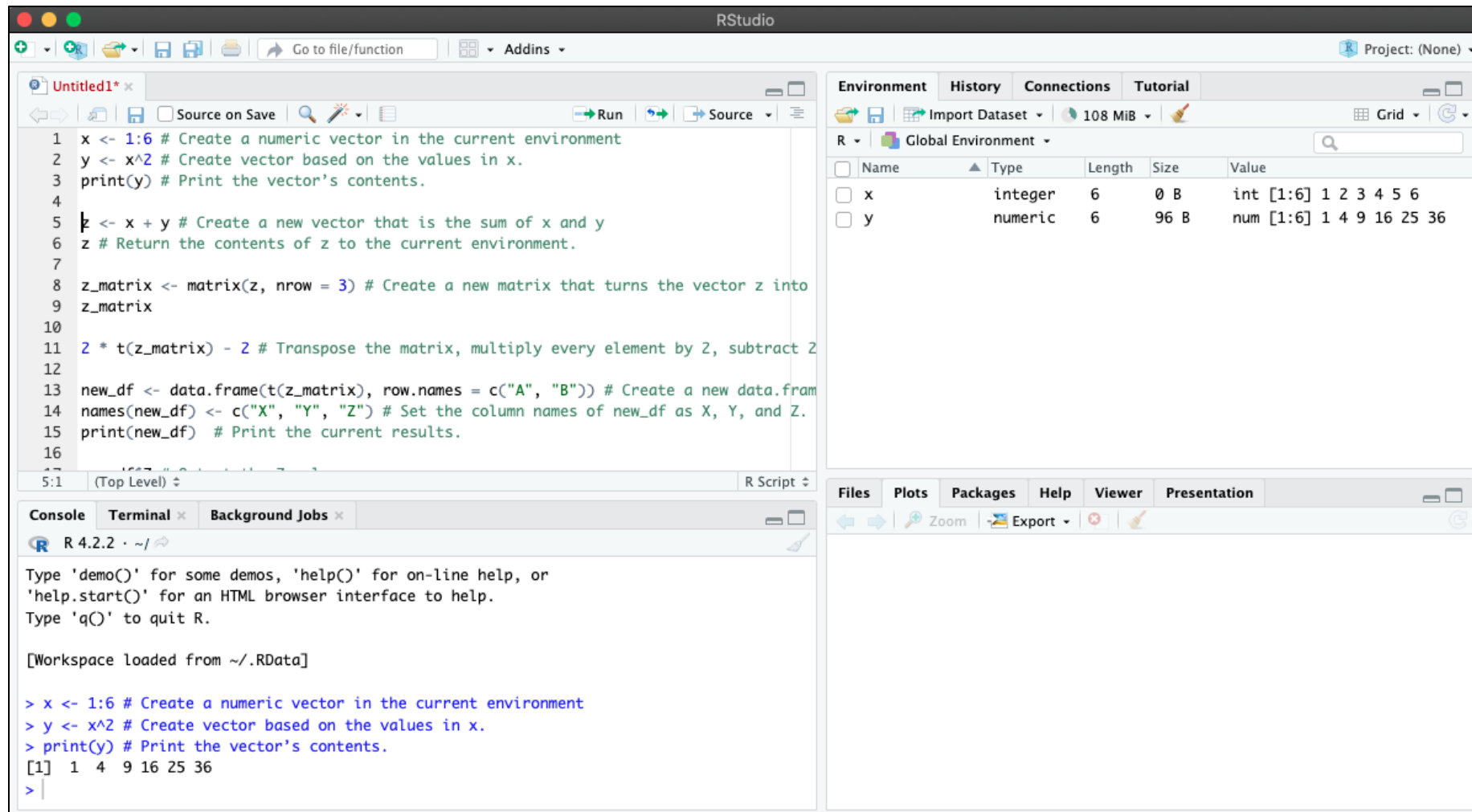


# The RStudio user-interface



The Plots and Files window: multi-use panes.

# The RStudio user-interface

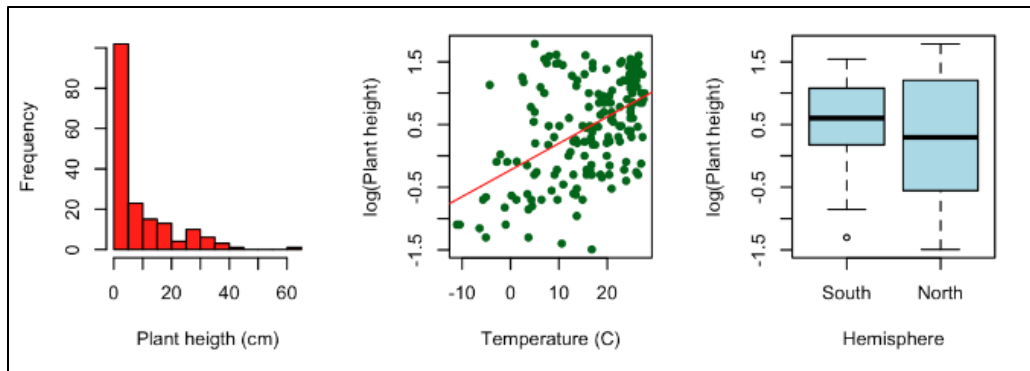


# Basic commands

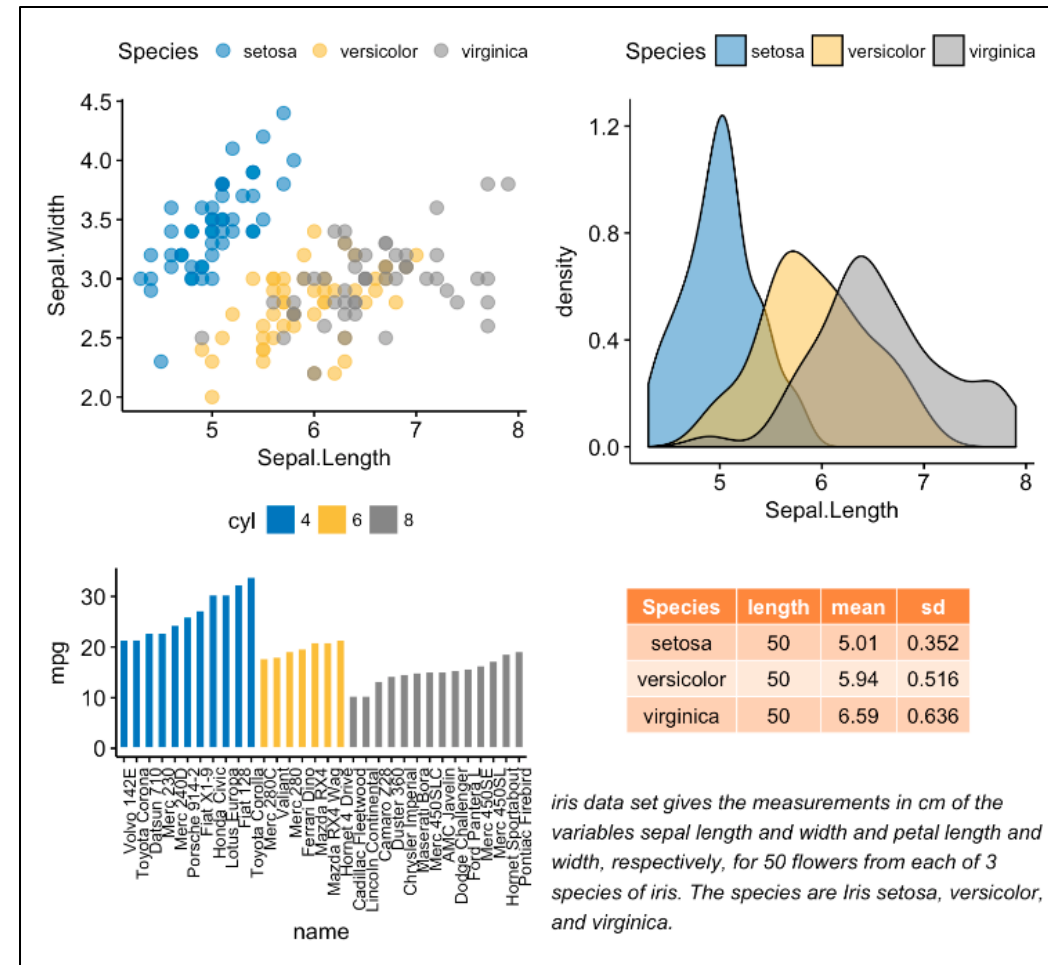
Assigning values	Creating vectors	Matrices
<ul style="list-style-type: none"><li>• <code>x &lt;- 5</code></li><li>• <code>y = 10</code></li></ul>	<ul style="list-style-type: none"><li>• <code>v &lt;- c(1, 2, 3, 4)</code></li><li>• <code>v[2]</code></li><li>• <code>v[v &gt; 3]</code></li></ul>	<ul style="list-style-type: none"><li>• <code>m &lt;- matrix(1:9, nrow=3)</code></li><li>• <code>m[1, ]</code></li><li>• <code>m[, 2]</code></li></ul>
Basic statistical functions	Data Frames	Installing/loading packages
<ul style="list-style-type: none"><li>• <code>mean(v)</code></li><li>• <code>median(v)</code></li><li>• <code>sd(v)</code></li><li>• <code>summary(v)</code></li></ul>	<ul style="list-style-type: none"><li>• <code>df &lt;- data.frame(Name=c('Alice', 'Bob'), Age=c(25, 30))</code></li><li>• <code>df\$Name</code></li><li>• <code>df[1, ]</code></li></ul>	<ul style="list-style-type: none"><li>• <code>install.packages("ggplot2")</code></li><li>• <code>library(ggplot2)</code></li></ul>

# Commands for visualization

- `plot(v)` # Plot a vector
- `hist(v)` # Histogram
- `boxplot(v)` # Boxplot



<https://environmentalcomputing.net/graphics/basic-plotting/>



<https://www.sthda.com/english/articles/24-ggpubr-publication-ready-plots/81-ggplot2-easy-way-to-mix-multiple-graphs-on-the-same-page/>

# Packages

Data Manipulation	Data Visualization	Statistical Analysis
<ul style="list-style-type: none"><li>• <b>dplyr</b>: For data manipulation (e.g., filtering, summarizing).</li><li>• <b>tidyr</b>: For reshaping and tidying data.</li><li>• <b>data.table</b>: For fast and efficient manipulation of large datasets.</li></ul>	<ul style="list-style-type: none"><li>• <b>ggplot2</b>: For creating sophisticated visualizations.</li><li>• <b>plotly</b>: For interactive plots.</li><li>• <b>lattice</b>: For multivariate data visualization.</li></ul>	<ul style="list-style-type: none"><li>• <b>stats</b>: Built-in package for basic statistical methods.</li><li>• <b>car</b>: For regression analysis and diagnostics.</li><li>• <b>psych</b>: For psychological and multivariate statistics.</li></ul>
Machine Learning	Time Series Analysis	Other
<ul style="list-style-type: none"><li>• <b>caret</b>: Comprehensive suite for machine learning workflows.</li><li>• <b>randomForest</b>: For random forest models.</li><li>• <b>xgboost</b>: For gradient boosting models.</li></ul>	<ul style="list-style-type: none"><li>• <b>forecast</b>: For forecasting models like ARIMA.</li><li>• <b>zoo</b>: For managing and analyzing time series data.</li></ul>	<ul style="list-style-type: none"><li>• <b>shiny</b>: For building interactive web applications.</li><li>• <b>tm</b>: For text mining and natural language processing.</li><li>• <b>sf</b>: For spatial data handling.</li></ul>

# How can R be used in transportation research?

## Strengths that other languages don't offer?

- R excels in statistical computing, making it ideal for analyzing transportation-related data where statistical modeling and hypothesis testing are required.
- R provides high-quality, publication-ready graphics. Packages like `ggplot2`, `shiny`, and `plotly` allow detailed visualizations of transportation networks and traffic patterns.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$
$$s^2 = \frac{\sum_{i=1}^{n_1} (x_i - \bar{x}_1)^2 + \sum_{j=1}^{n_2} (x_j - \bar{x}_2)^2}{n_1 + n_2 - 2}$$

```
> t.test(data$X, data$Y) # T-test to compare means of X and Y

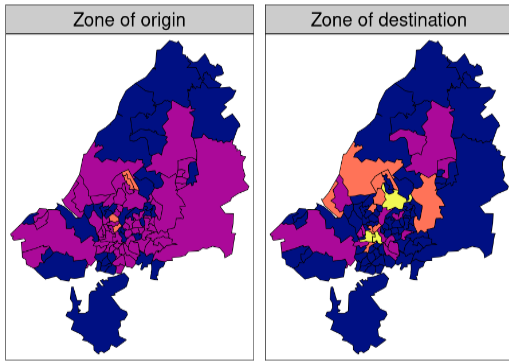
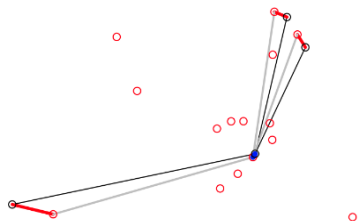
Welch Two Sample t-test

data: data$X and data$Y
t = -22.966, df = 158.15, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -52.04668 -43.80368
sample estimates:
mean of x mean of y
 50.32515  98.25033
```

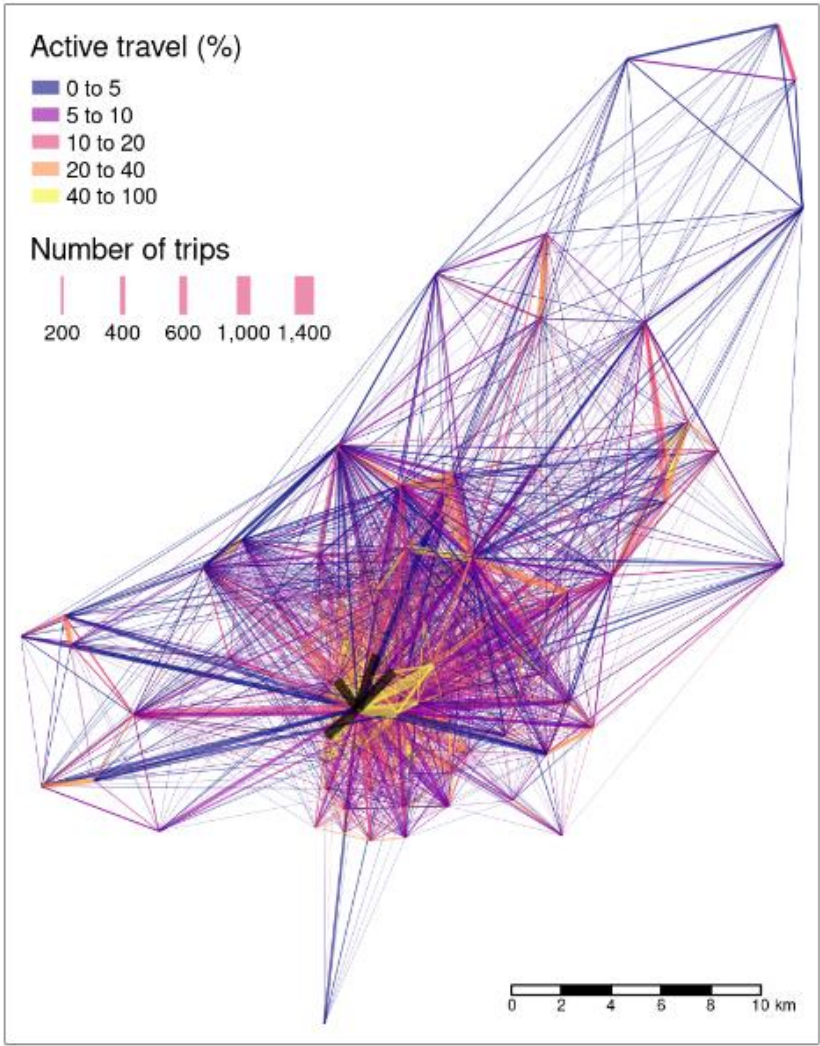
Feature	R	Python
<b>Statistical Modeling</b>	R is unparalleled in statistical methods, making it the go-to for rigorous analysis.	Python relies on libraries like statsmodels and scipy, which may not be as robust.
<b>Geospatial Analysis</b>	Specialized tools like sf and tmap offer deeper support for geospatial data.	Libraries like GeoPandas and folium are powerful but less mature than R's.
<b>Data Visualization</b>	ggplot2 and shiny allow for high-quality, interactive visualizations.	Python's matplotlib and seaborn are strong, but less intuitive for beginners.
<b>Ease of Use for Statistics</b>	R is designed with statistical applications in mind, making it more intuitive.	Python requires more effort to set up equivalent workflows for statistics.
<b>General-Purpose Programming</b>	Less versatile for general programming outside data science.	A general-purpose language, Python is better for tasks beyond analysis.
<b>Community</b>	Strong in academia and research (especially social sciences and transportation).	Dominant in industry, with extensive support for machine learning and automation.

# A transportation case study

A Bristol, UK case study aim at reducing motorized vehicle congestion and encouraging bicycle traffic.



o	d	all	bicycle	foot	car_driver	train
E02003043	E02003043	1493	66	1296	64	8
E02003047	E02003043	1300	287	751	148	8
E02003031	E02003043	1221	305	600	176	7
E02003037	E02003043	1186	88	908	110	3
E02003034	E02003043	1177	281	711	100	7



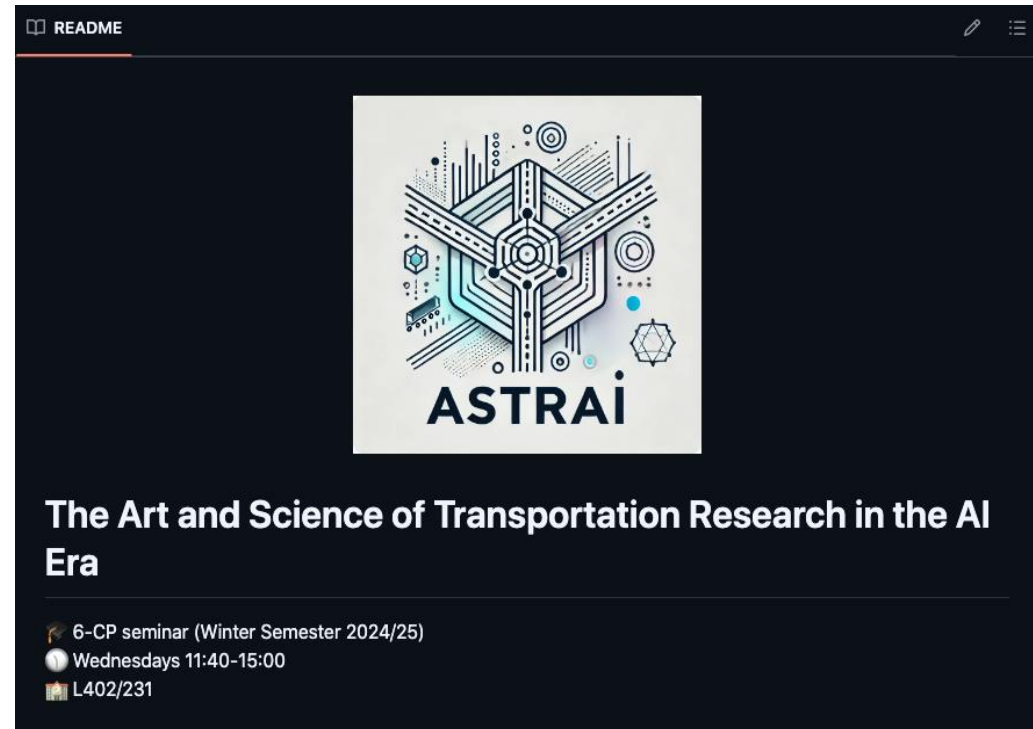
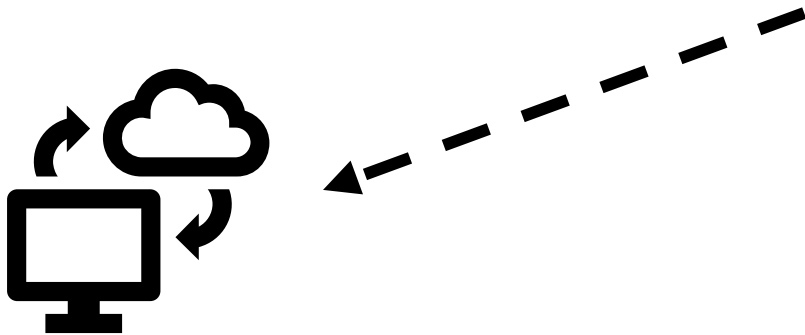


# Links to installation

- R: <https://cran.r-project.org/>
- RStudio: <https://posit.co/download/rstudio-desktop/>

# Learning by doing

- Download R file "syntax.R" from GitHub
- Download R file "practice.R" from GitHub
- Open files in RStudio
- Run scripts



# References

- <https://cran.r-project.org/doc/manuals/R-intro.pdf>
- <https://developer.r-project.org/equalAssign.html>
- <https://bookdown.org/robinlovelace/geocompr/transport.html>