# Software Systems for Policymakers and Data Scientists

Harnessing Big Data for National Statistics and Evidence-Based Policymaking

18-19 August, 2025
Blantyre, MW

NSO │ UC Berkeley │ World Bank

# Overview

- Computing and software choices
- Statistical Computing/Data Science
  - Python
  - R
  - Stata
- Mechanics of running code in Python and R
- Further resources

**Goal is not for you to memorize these slides, but to understand the landscape of choices and have this document as a reference.**

Computing choices

- Open source versus closed source
- Paid vs free
- Cloud vs local

# Open source vs closed source

**Open Source:** The software's source code is publicly available, allowing anyone to view, modify, and distribute it. Often community-driven, it emphasizes transparency and collaboration.

**Closed Source:** The source code is kept secret and is proprietary to the company or developer. Users are given a license to use the software but cannot access or change its internal workings.

# Paid vs free

**Paid:** Users must pay a fee to use the software, either as a one-time purchase or a recurring subscription. Payment typically grants access to features, updates, and dedicated support.

**Free (Freeware/Gratis):** The software can be downloaded and used at no financial cost. However, it may have limitations, such as restricted features, ads, or a lack of professional support.

# Local vs cloud

**Local (On-Premise):** The software is installed and runs directly on a user's computer or a company's private servers. It does not require an internet connection for core functionality but may be less scalable and require manual updates.

**Cloud:** The software is hosted on remote servers and accessed via the internet. This model offers flexibility, as users can access their applications and data from any device with an internet connection.

# Statistical Computing/ Data Science

# Programming languages for statistical computing

Main choices:

- Python
- R
- Stata
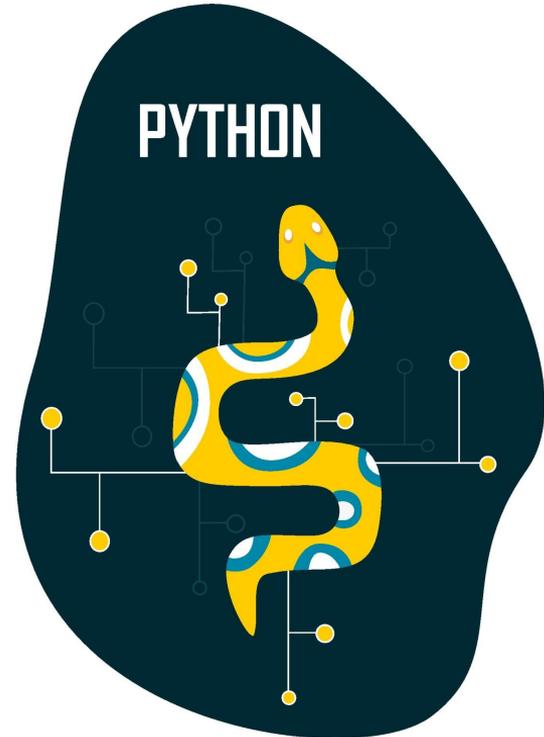
Other choices:

- SAS
- SPSS
- Julia
- JMP

# Python

**Versatility:** A general-purpose programming language. It is a great choice for end-to-end projects that require more than just statistical analysis, such as web scraping, automation, and machine learning.

**Key Libraries:** Statistical/numerical computing is powered by specialized libraries. **Pandas** and **NumPy** are foundational for data manipulation and numerical operations. **SciPy** and **Statsmodels** provide a comprehensive range of statistical tests and models. **Scikit-learn** is the standard for machine learning.

**Learning Curve:** The learning curve is moderate. Python's clean syntax is beginner-friendly, but mastering the full data science ecosystem requires learning multiple libraries.
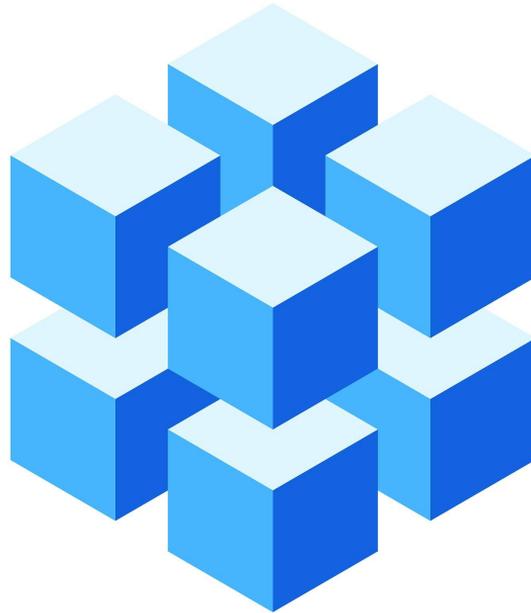
**Community & Use:** Has a massive, diverse user base in data science, software engineering, and academia. Widely used in the industry for production environments.

# numpy, scipy

**NumPy (Numerical Python)** is the foundational library for scientific computing in Python. It provides a powerful **N-dimensional array object** and functions for working with these arrays. It's the backbone for most other data science libraries.

**SciPy (Scientific Python)** is a library of algorithms and functions built on top of NumPy. It offers a wide range of modules for scientific and technical computing, including optimization, linear algebra, signal processing, and statistics.

# pandas

A powerful and flexible library for **data manipulation and analysis**.

It introduces two key data structures: `DataFrame` for 2D tabular data (like a spreadsheet or SQL table) and `Series` for 1D data.

Pandas makes it easy to handle messy, real-world data by providing functions for tasks like cleaning, merging, and reshaping data.
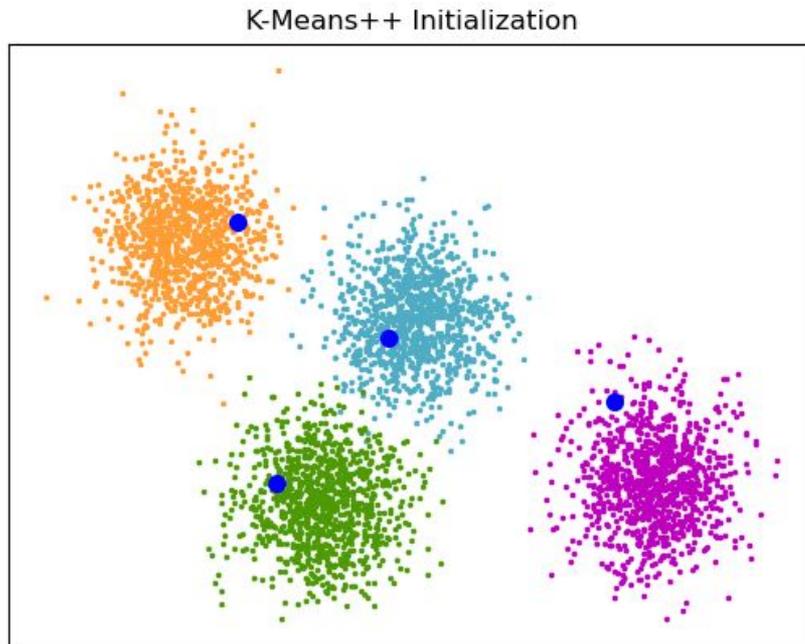
# scikit-learn

The de facto standard library for **classical machine learning** in Python.

It provides a consistent and simple interface to a vast collection of algorithms for classification, regression, clustering, dimensionality reduction, and model selection.

It's built on NumPy, SciPy, and matplotlib, making it a key component of the scientific Python stack.



K-Means++ Initialization

# tensorflow/pytorch

Tensorflow (Google) and Pytorch (Meta) are frameworks for machine learning and deep learning.

They are focused on building and training **neural networks** and are highly optimized for performance on various hardware (CPUs, GPUs, TPUs).

While a steeper learning curve than Scikit-learn, they offers the flexibility and power needed for cutting-edge research and large-scale applications.

# IPython

An interactive command shell that provides a significant upgrade to the standard Python interpreter.

It features powerful additions like **tab completion, rich media display (e.g., inline plots), and command history**, making it an essential tool for exploratory data analysis.

IPython is the kernel that powers the popular Jupyter Notebook.

# Jupyter Notebook

An interactive, web-based environment that allows users to create and share documents containing **live code, equations, visualizations, and narrative text**.

They are structured into **cells**, which can be either code or Markdown text. This enables a "computational narrative" where you can interleave explanations and analysis with the code and its output.

Jupyter Notebooks are an indispensable tool for data science because they provide a highly **reproducible and shareable workflow** for exploratory data analysis, data cleaning, and model prototyping.

Also runs Julia and R code.

# Anaconda

A free and open-source **distribution** of Python and R for scientific computing.

It simplifies the setup process by coming pre-packaged with over 7,500 data science packages, including all the libraries mentioned above.

Its key feature is **Conda**, a powerful package and environment manager that helps you create isolated, conflict-free environments for different projects.
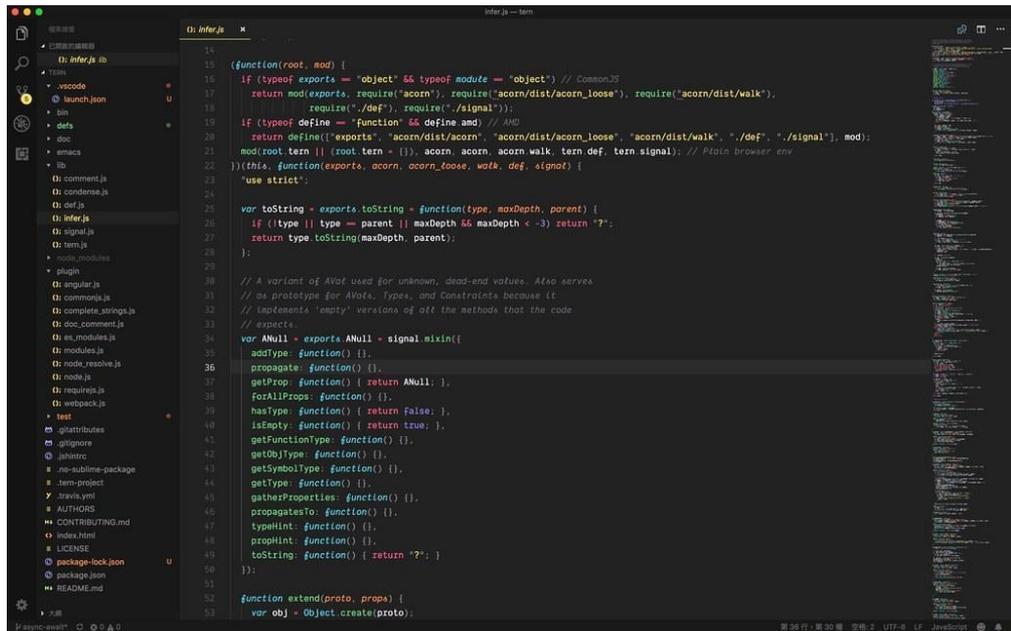
# VSCode

A free, lightweight, and highly customizable **code editor** from Microsoft.

With its extensive marketplace of extensions, it can be transformed into a powerful integrated development environment (IDE) for Python, data science, and many other programming languages.

It provides excellent support for debugging, Git integration, working with Jupyter Notebooks, and LLMs (like ChatGPT).

# R

**Purpose-Built:** Created by statisticians for statistical computing and graphics. It is unparalleled for its depth and breadth of statistical packages.

**Key Packages:** The **tidyverse** is a popular collection of packages (**dplyr**, **ggplot2**, **tidyr**) that provide a consistent and intuitive approach to data manipulation and visualization. The **CRAN** (Comprehensive R Archive Network) repository hosts over 18,000 user-contributed packages for virtually any statistical need.

**Learning Curve:** Can have a steeper initial learning curve due to its unique syntax, but it's exceptionally efficient for statistical tasks once mastered.

**Community & Use:** The primary language for statisticians and academic researchers. It's often favored for its state-of-the-art statistical methods and publication-quality graphics.

# tidyverse

A collection of packages in R designed with a **consistent philosophy and grammar** to make data manipulation and visualization more intuitive.

It operates on the principle of **"tidy data,"** where each variable is a column, each observation is a row, and each type of observational unit is a table.

Key packages include `dplyr` for data manipulation, `tidyr` for tidying data, and `ggplot2` for plotting. The consistent syntax and "piping" operator (`%>%`) streamline a user's workflow.
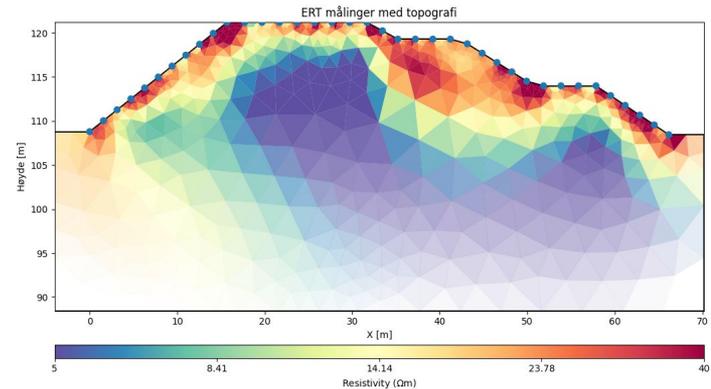
# ggplot2

A powerful and widely used R package for creating high-quality **data visualizations**. It is a core component of the Tidyverse.

It is based on the **"Grammar of Graphics,"** a conceptual framework that allows you to build complex plots by combining simple components like data, geometric objects (points, lines), and aesthetic mappings (color, size, shape).

Unlike other plotting libraries, `ggplot2` enables you to create sophisticated and aesthetically pleasing graphics with minimal code, making it a go-to for exploratory data analysis and publication-ready plots.

# glm

A fundamental function in the base R installation for fitting **generalized linear models**.

GLMs are a flexible extension of ordinary linear regression that allows for **response variables with different error distributions**, such as binomial (for logistic regression), Poisson (for count data), and gamma.

The `glm()` function provides a single, consistent interface for fitting a variety of models, making it a cornerstone for a wide range of statistical analyses beyond standard linear models.

# Machine learning libraries in R

**caret (Classification and Regression Training):**
Provides a unified interface for training and evaluating various machine learning models, simplifying the workflow for tasks like classification and regression.

**randomForest:**
Implements the Random Forest algorithm, an ensemble method effective for both classification and regression.

**xgboost:**
A highly efficient and popular library for gradient boosting, known for its performance in competitive machine learning scenarios.
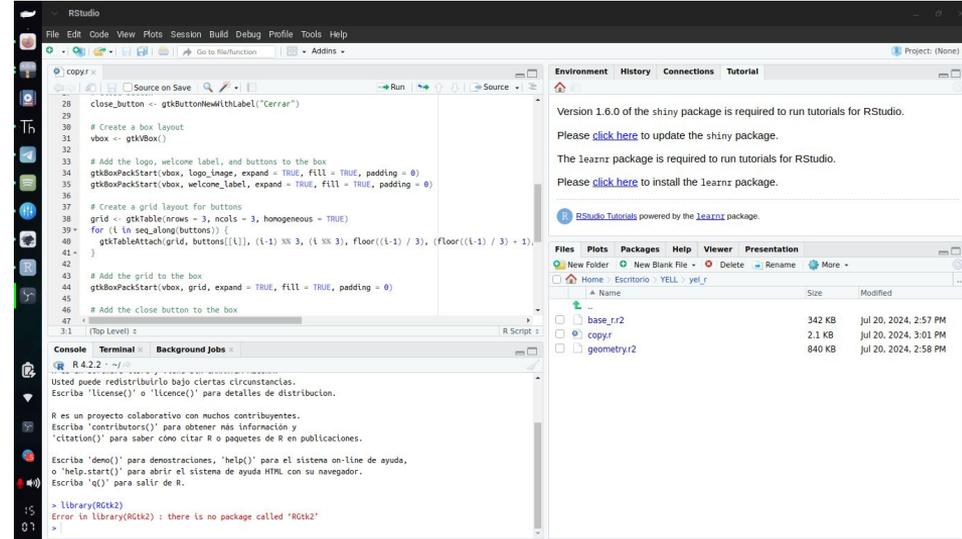
**e1071:**
Contains implementations of various algorithms, including Support Vector Machines (SVMs) and Naive Bayes classifiers.

# RStudio

A popular and comprehensive **integrated development environment (IDE)** specifically for R.

It simplifies the process of coding in R by providing a structured layout with a console, a script editor, a history pane, a file browser, and a "plots" pane for viewing visualizations.

RStudio is highly valued by R users for its powerful features like debugging tools, project management, and seamless integration with packages like the Tidyverse, making it the **standard environment for R programming**.

# Stata

**Proprietary & Specialized:** A commercial, user-friendly statistical software package particularly strong in econometrics, biostatistics, and panel data analysis.

**Strengths:** Offers a straightforward, command-driven syntax and a comprehensive set of built-in commands. It excels at data management and producing reproducible research through "do-files" that document all steps. It also includes an easy-to-use graphical user interface (GUI).

**Learning Curve:** Generally considered the easiest to learn for those with a statistics background who are not strong programmers, thanks to its intuitive commands and strong documentation.

**Community & Use:** Widely used in economics and other social sciences, as well as in public health research. Its key advantage is a reliable, all-in-one platform with consistent features.

# Example: Python vs R vs Stata

Let's look at how to generate synthetic data, and run multiple regression in Python, R, and Stata.

# Example: <u>Python</u> vs R vs Stata

```python
# Import necessary libraries
import numpy as np
import pandas as pd
from sklearn.linear_model import LinearRegression

# Set seed for reproducibility
np.random.seed(12345)

# Create a dictionary of variables
df = pd.DataFrame(data={f'x{i}': np.random.normal(size=200) for i in range(1, 5)})

# Generate the dependent variable y
df['y'] = 1.5 + 0.8 * df['x1'] - 0.7 * df['x2'] + 1.1 * df['x3'] - 1.7 * df['x4']
+ np.random.normal(size=200)
```

# Example: <u>Python</u> vs R vs Stata

```python
# Define the independent variables (X) and dependent variable (y)
X, y = df[['x1', 'x2', 'x3', 'x4']], df['y']

# Create a LinearRegression model object and fit the model to the data
model = LinearRegression()
model.fit(X, y)

# Print the model coefficients and intercept
print("Coefficients:", model.coef_)
print("Intercept:", model.intercept_)
```

# Example: Python vs R̲ vs Stata

```
set.seed(12345)

# Create a data frame with 200 observations
df <- data.frame(
  x1 = rnorm(200),
  x2 = rnorm(200),
  x3 = rnorm(200),
  x4 = rnorm(200)
)

# Generate the dependent variable y
df$y <- 1.5 + 0.8 * df$x1 - 0.7 * df$x2 + 1.1 * df$x3 - 1.7 * df$x4 + rnorm(200)

# Run the regression
model <- lm(y ~ x1 + x2 + x3 + x4, data = df)

# Print the summary of the regression results
summary(model)
```

# Example: Python vs R vs <u>Stata</u>

```stata
clear

set seed 12345

set obs 200

forv i=1/4 {

    gen x`i'=rnormal()

}

gen y=1.5+0.8*x1-0.7*x2+1.1*x3-1.7*x4+rnormal()


reg y x*
```

# Jupyter notebooks and RStudio

# Further resources

# Python

**Python for Data Analysis** by Wes McKinney: Written by the creator of the pandas library, this book is great for understanding how to work with large datasets.

**Python Data Science Handbook** by Jake VanderPlas: This comprehensive book provides a guide to essential Python data science tools like NumPy, pandas, and Matplotlib.

**Data Science from Scratch** by Joel Grus: This book covers the fundamental concepts of data science by having you implement the algorithms from scratch using Python.

**Learn Pandas** by Hernan Rojas: Comprehensive guide to how to use pandas.

# R

**R for Data Science** by Hadley Wickham and Garrett Grolemund: Written by key contributors to the R ecosystem, this book focuses on using the tidyverse for data analysis.

**An Introduction to Statistical Learning: With Applications in R** by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani: A classic resource for statistical learning with practical examples in R.

**swirl**: An R package that teaches you R from within the R console itself, offering an interactive, hands-on learning experience.

**RStudio Community**: A forum where you can ask questions and get help from other R users.

**The R Bloggers website**: A great aggregator of blogs from the R community.

# Stata

**Statalist**: The official forum for Stata users to engage in discussions about statistics and Stata.

**Stata Blog and Video Tutorials**: The official Stata website offers short, focused video tutorials and blog posts on specific features and techniques.

**Stata Documentation and Web Resources**: The official documentation is comprehensive, and there are web resources that provide step-by-step instructions.