# Overview

Augustin Luna
09 January, 2016

Research Fellow
Department of Biostatistics and Computational Biology
Dana-Farber Cancer Institute

# Topics to be Covered

- R: Language Basics, Plotting, Getting Help

- Using the RStudio Editor

- Machine Learning Fundamentals

  - Dimension Reduction

  - Regression

  - Clustering

- Accessing Datasets: CellMiner (cell lines/drugs), CBioPortal (patient samples), Pathway Commons (pathways)

- Developing Web Applications

# What is R?

- Free, open source

- Started in 1993

- Geared towards scientific computing
  - Created by Ross Ihaka and Robert Gentleman (statisticians)

- Interpreted; similar to Python and MATLAB

# Why is R Popular?

- Free, open source
- Interactive data analysis
  - Script-driven rather than menu-driven helps reproducibility
- Flexible and powerful plotting support
- Excellent package management system
  - Large and growing collection of statistical analysis methods
  - Simple package installation; dependency management
  - R scripts usually portable to other platforms
  - Package repositories ensure functionality, documentation, and interoperability

# Extending R and Package Repositories

- Comprehensive R Archive Network (CRAN)
  - 5,800 R packages (as of June 2014)
  - Many packages call C, C++, Fortran, or Java code for speedups
- Bioconductor
  - 800+ R packages focused on bioinformatics
  - 50+ packages dedicated to pathway analysis
- Devtools
  - R package that allows package installation from code repositories

# RStudio

- https://www.rstudio.com/
- Available for Windows, OSX, and Linux
- Simplifies common tasks: plotting, package installation, accessing files, viewing variables, etc.

# Installing R and RStudio

- Install R

  - https://cran.rstudio.com/

- Install RStudio

  - https://www.rstudio.com/products/rstudio/download/

- RStudio does not come with R and R must be installed first

# RStudio Overview

# Table View of Variables

- Highlighted boxes open a table view of variable contents

# Change Current Directory

- Highlighted boxes open a table view of variable contents

# Making a New R Script

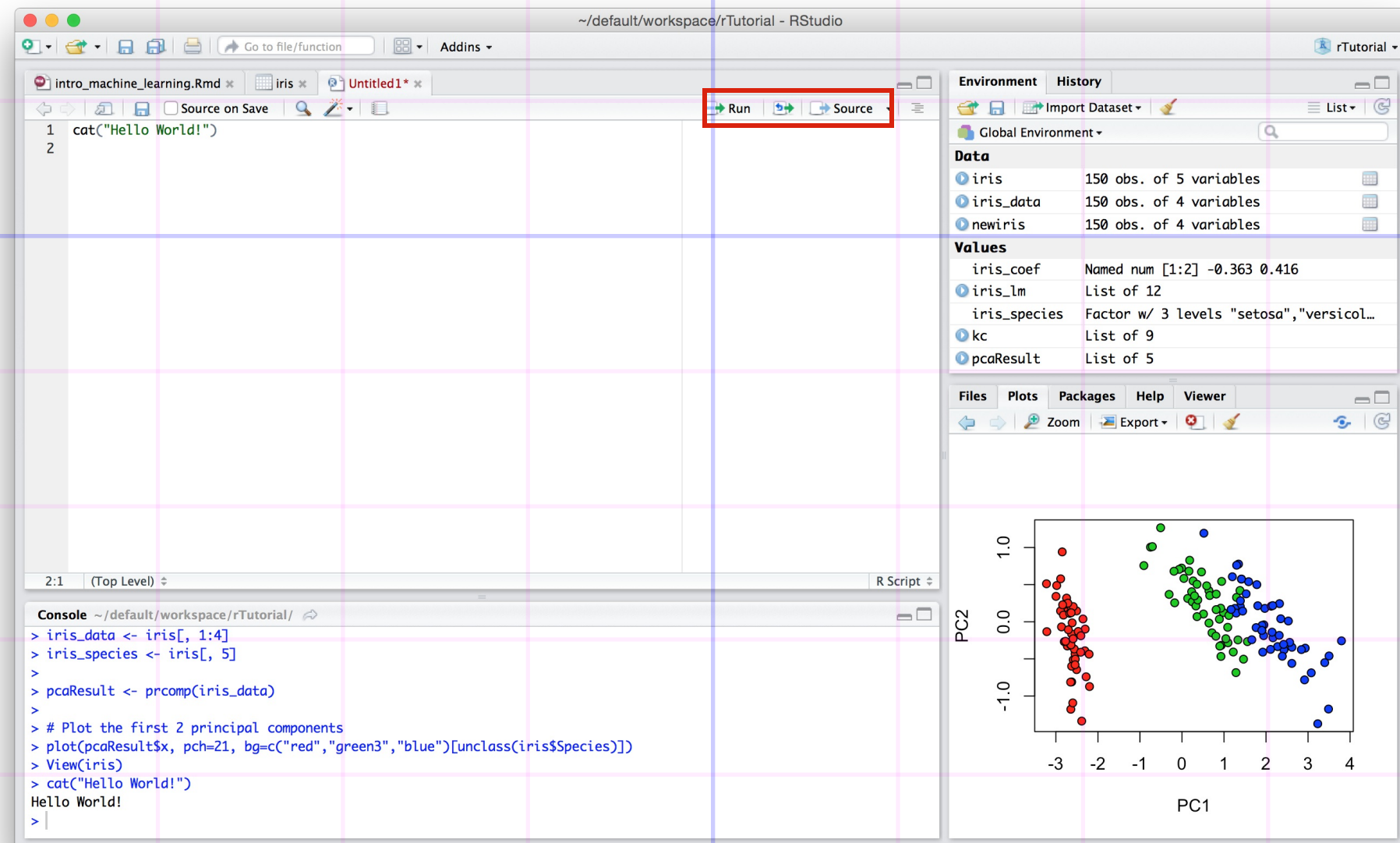# First Script: Hello World!

- cat() prints a simple message in the console

```
cat("Hello World!")
```

```
Hello World!
```

# Running Hello World Script

- "Run" button runs current line or selected lines

- "Source" button runs all lines in file

# Installing Packages

- CRAN packages can be installed using RStudio or `install.packages()`