1.  **Introducción al análisis de datos y data science**

2.  **Lenguaje de programación R**

3.  **ROOT Data Analysis Framework**

4.  **Manipulación y visualización de datos**

5.  **Modelamiento estadístico**

6.  **Machine Learning**

7.  **TMVA (Toolkit for Multivariate Data Analysis)**

**Libros:**

- Venables, Smith, et al. An Introduction to R. 2015

- Peng. R Programming for Data Science. 2015

- A ROOT Guide for Beginners, 2015.

- Hoecker et al. Toolkit for Multivariate Data Analysis with ROOT User´s Guide. 2007

- Box et al. Statistics for Experimenters. 2005

- Hastie et al. The Elements of Statistical Learning. 2008

- Witten, Frank & Hall. Data mining: practical machine learning tools and techniques. 2011.

**Internet:**

**Coursera:** Data Science at Johns Hopkins University:
**https://www.coursera.org/specializations/jhu-data-science**

**The Comprehensive R Archive Network**
**https://cran.r-project.org/**

**ROOT**
**https://root.cern.ch/**

**Stack Overflow**
**http://stackoverflow.com/**

**Kaggle**
**https://inclass.kaggle.com/**

Los archivos de la clase se encontrarán en Intranet del curso

En este curso se aplica la modalidad de nota única.

$$\text{Nota Final} = (1\ \text{Par} + 2\ \text{Ta} + 3.5\ \text{Ex} + 3.5\ \text{Pre}) / 10$$

Par=Participación en clase

Ta=Promedio de 3 tareas

Ex= Examen

Pre=Presentación del trabajo

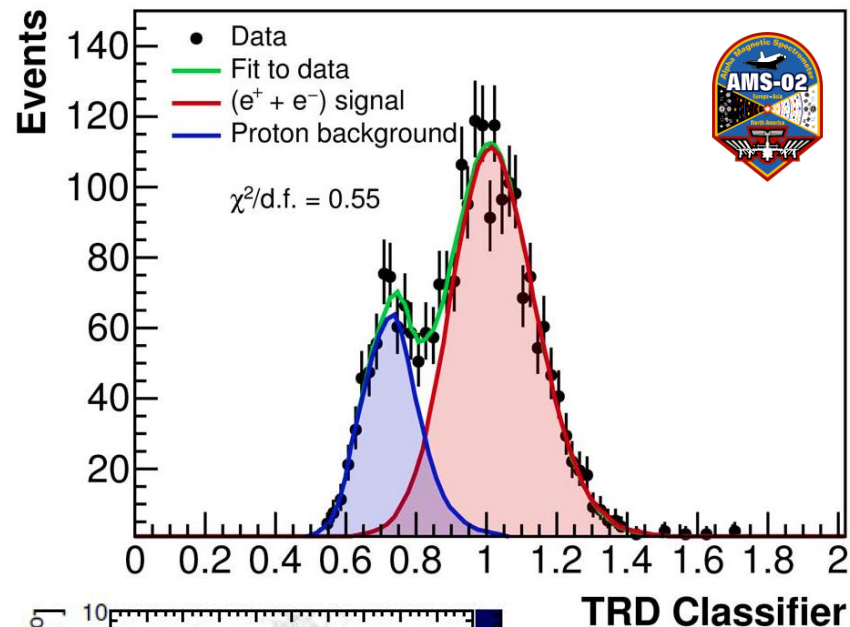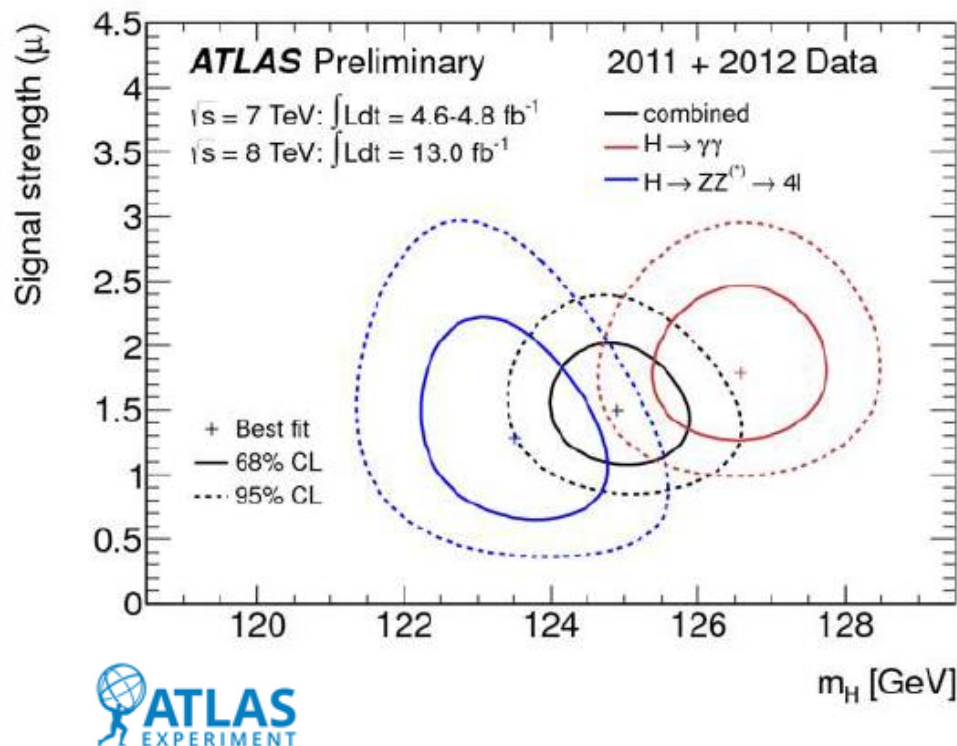**Tarea 1**:  8 abril
**Tarea 2**:  22 abril
**Tarea 3**:  13 mayo

**Examen**: 27 mayo

**Presentación Trabajo**: 1 julio

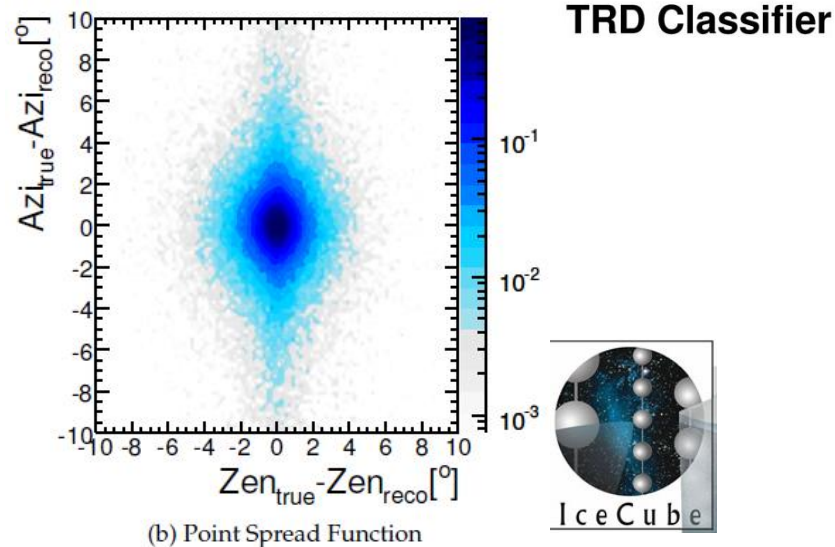**Asesorías**: Oficina 305, tercer piso Edificio Física

# 1. Introducción al análisis de datos y data science

ROOT does what physicists do:

It makes plots.



(b) Point Spread Function

**Josh Cogan**, PhD from Stanford Physics
8.8k Views · Upvoted by Yan Ren, PhD in Physics

ROOT is unavoidable in experimental particle physics. I encourage you to round out your skills by doing as much as possible in Python. C++ itself is already a dinosaur in the software engineering world. And don't listen to your professors/post-docs if they say otherwise, when was the last time they applied to an engineering job? Unfortunately ROOT can't even claim to be a dinosaur; its some awkward half-bird/half-reptile species that has only eked out an existence in an extremely small niche of academia.

**Jay Wacker**, Researcher in particle physics.
5.6k Views · Upvoted by Hongwan Liu, Graduate student in Theoretical Cosmology, Andy Buckley

Almost no one likes ROOT. The only advantage it has going for it is that there is a huge code base built up and that starting over from scratch is nearly impossible. The good news is that there are things like PyROOT that is a Python module that allows you to interact with ROOT. ROOT was a big step forward from FORTRAN based PAW which lingered until the late 90s. I think no one would think about using ROOT if they were starting again, but it is impossible to change out in either experiment at the LHC because they are running and no other experiment is big enough to devote resources to creating a new

**Kevin Sapp**, Aspirateur
2.6k Views · Upvoted by Hor

Con: ROOT breeds poor programmers within the physics community. Josh C. calls ROOT a "non-transferrable skill", which is only part of the bigger problem: many, if not most, new students of particle physics learn C/C++, often as a first programming language, by mimicking ROOT scripts or compiled code containing ROOT libraries. Unfortunately, in many ways this is the worst system to learn programming from; CINT has very lax rules about syntax, the memory management is done automatically but half-heartedly, and

link

ET LUX IN TENEBRIS LUCET · MCMXVII

**PUCP**

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

## MATH & STATISTICS

☆ Machine learning
☆ Statistical modeling
☆ Experiment design
☆ Bayesian inference
☆ Supervised learning: decision trees, random forests, logistic regression
☆ Unsupervised learning: clustering, dimensionality reduction
☆ Optimization: gradient descent and variants

## PROGRAMMING & DATABASE

☆ Computer science fundamentals
☆ Scripting language e.g. Python
☆ Statistical computing package e.g. R
☆ Databases SQL and NoSQL
☆ Relational algebra
☆ Parallel databases and parallel query processing
☆ MapReduce concepts
☆ Hadoop and Hive/Pig
☆ Custom reducers
☆ Experience with xaaS like AWS

## DOMAIN KNOWLEDGE & SOFT SKILLS

☆ Passionate about the business
☆ Curious about data
☆ Influence without authority
☆ Hacker mindset
☆ Problem solver
☆ Strategic, proactive, creative, innovative and collaborative

## COMMUNICATION & VISUALIZATION

☆ Able to engage with senior management
☆ Story telling skills
☆ Translate data-driven insights into decisions and actions
☆ Visual art design
☆ R packages like ggplot or lattice
☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

**MarketingDistillery.com** is a group of practitioners in the area of e-commerce marketing. Our fields of expertise include: marketing strategy and optimization; customer tracking and on-site analytics; predictive analytics and econometrics; data warehousing and big data systems; marketing channel insights in Paid Search, SEO, Social, CRM and brand.

*Marketing* DISTILLERY

---

https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century

MENU

**Harvard Business Review**

DATA

# Data Scientist: The Sexiest Job of the 21st Century

by **Thomas H. Davenport** and **D.J. Patil**

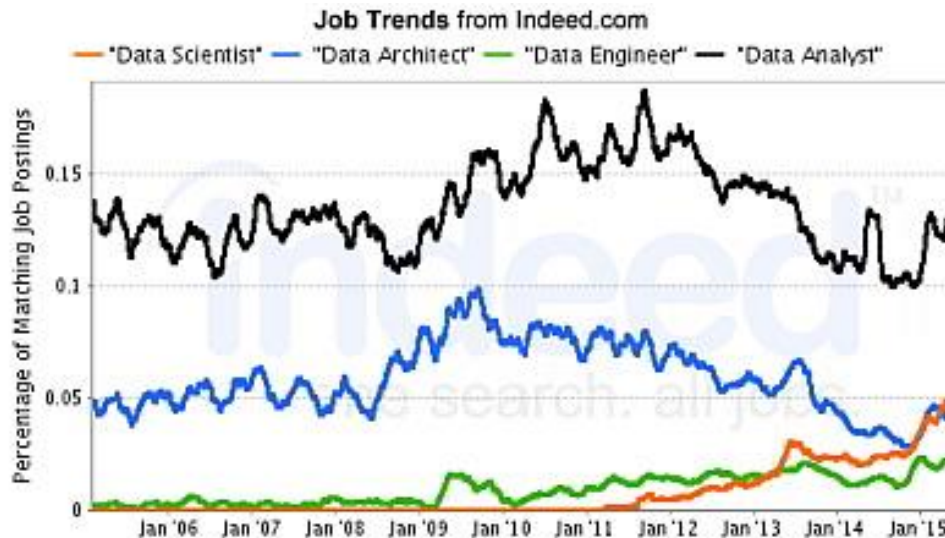FROM THE OCTOBER 2012 ISSUE

SUMMARY · SAVE · SHARE · COMMENT · TEXT SIZE · PRINT · $8.95 BUY COPIES

W hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join.

[link](#)

Hacking Skills · Math & Statistics Knowledge

Machine Learning

Data Science

Danger Zone! · Traditional Research

**Substantive Expertise**

**Harvard Business Review**

GETTING CONTROL OF **BIG DATA**

Job Trends from Indeed.com
— "Data Scientist"   — "Data Architect"   — "Data Engineer"   — "Data Analyst"

**Big Data, Big Paycheck**

Median salary for analytics professionals and those specifically within data science, by level of experience.

| | | |
|---|---|---|
| Up to 3 years | Analytics professionals | $65,000 |
| | Data scientists | $80,000 |
| 4 to 8 years | | $85,000 |
| | | $120,000 |
| 9+ years | | $115,000 |
| | | $150,000 |

Note: Data do not include managers   Source: Burtch Works   The Wall Street Journal

De un mail de Hans Beck (The ALICE juniors) sobre "Qualities for next spokesperson": …

8. Support careers within and outside of academia

The job market is tough and quite some ALICE people experience unemployment or uncertain futures. How can we make sure the hard work everybody puts in their project is being rewarded? …. A career outside academia is a viable option; but for many of us it is uncharted territory. Can we prepare the next generation better to be fit also for industry? We miss a regular space to share experiences in the job search and highlight opportunities…

@PUCP : maestría en Informática con mención en Ciencias de la Computación
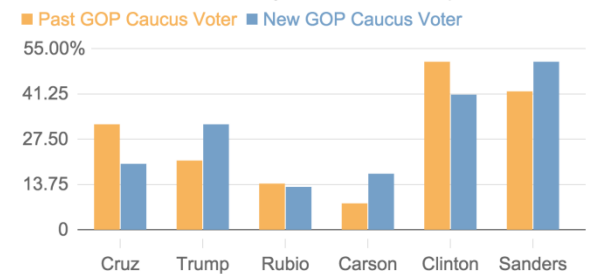
| Paso | Ejemplo |
|---|---|
| **Definir pregunta** | Correlación temporal entre flujo de fotones y neutrinos astrofísicos |
| **Definir juego de datos ideal** | Datos de Fermi y IceCube |
| **Obtener datos accesibles** | Fermi públicos formato FITS, IceCube (en txt) no todo |
| **Limpiar datos** | Eliminar períodos off-line, calibración, primer fondo |
| **Análisis exploratorio** | Graficar número de eventos versus tiempo en intervalos de minutos |
| **Modelamiento estadístico** | Modelar señal y fondo, método likelihood, cortes con BDT |
| **Graficar resultados** | Flujo versus tiempo |
| **Interpretar resultados** | Datos compatibles con ruido de fondo (sin correlación) |
| **Crear código reproducible** | Macro/programa con ROOT |
| **Distribuir resultados** | Publicación |

- Descriptivo

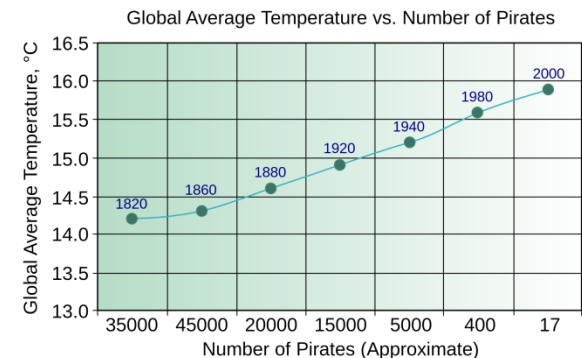- Exploratorio

- Inferencial

- Predictivo

- Causal

- Mecanista



Peru - 2014



Iowa Candidate Preferences, by Past Caucus Participation

"Prediction is not inference"
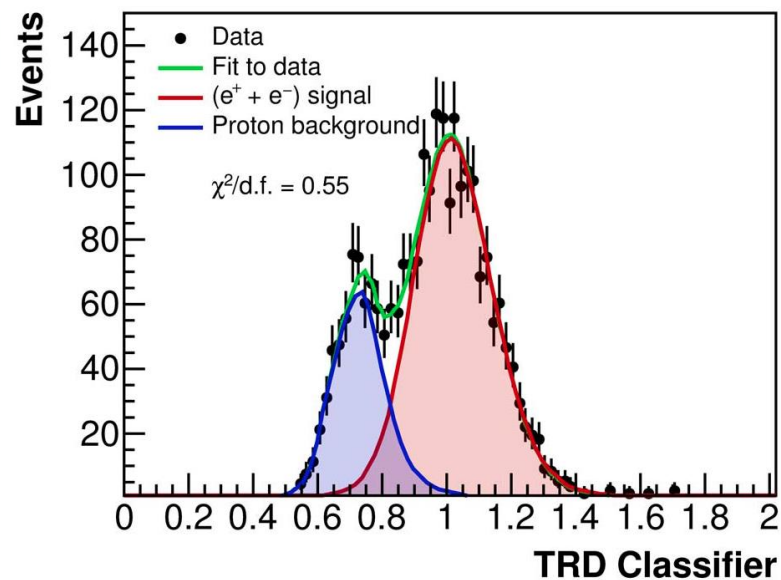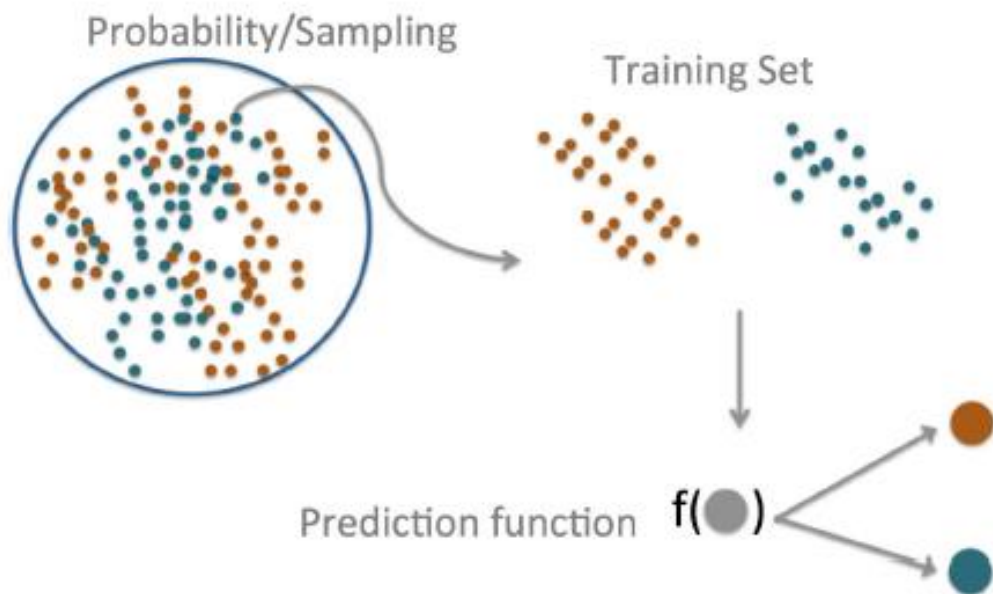


Lima, PE
TIEMPO METEOROLÓGICO
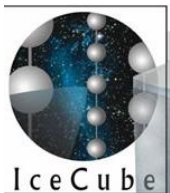


Global Average Temperature vs. Number of Pirates

"Correlation does not imply causation"

$$F_g = G \frac{m_1 m_2}{r^2}$$

Categorías tomadas de J. Leek (JHU)

Raw data per event ~1 Mb, with rate of $6 \times 10^8$ events/s

IceCube data volumes are 1TB per day raw data and 135GB per day of filtered data.

AMS-02 data sample: raw frames, reconstructed and simulated data ~150 TB/year.

Human genome data size ~770 Mb

As of June 2015, English wikipedia size ~10 TB uncompressed.

By 2016, global IP traffic will reach 1.1 zettabytes ($10^{21}$b) per year

# shell / terminal

**Command Line interface**:
- Explorar directorios
- Crear y editar archivos, directorios, programas
- Ejecutar programas

- Git Bash Cygwin: Windows
- Terminal: Linux/Mac

Para Windows: Bitwise SSH client
Notepad++

**Path**:
Directorio superior: root = "/"
Home = "~"

**Comandos**: command -flags arguments
ej. ls –lhrt ~
    chmod –R 777 ./

Básicos: pwd, clear, ls, cd, mkdir, touch, cp, rm, mv, date, echo



Linux Bash Shell Cheatsheet        .bashrc

Para compartir, almacenar y trabajar en grupo en código cambiante:
**(Sub)Version Control Software**: **Repositorio central**

Crear cuenta (free) en **GitHub**

https://github.com/



Otros sistemas:

CVS: Concurrent Versions System
(GNU license)

SVN: Subversion
(Apache license)

versus

Otro lugar donde guardar código libre:



ejemplo

https://sourceforge.net/

**PUCP**

Sistema que guarda cambios hechos archivos en el tiempo y permite acceder a versiones anteriores

Git Cheatsheet

Instalar Git (Bash o GUI):
guardar en **repositorio local**

Configuración:

$ git config --global user.name jlbazo
$ git config --global user.email jbazo@pucp.edu.pe
$ git config --list
$ exit

Usar mismo username y correo de GitHub

GitHub: web-based hosting service
         **repositorio remoto**

**PUCP**

Permite subir/bajar **repositorios locales** hacia/de **remotos** (web)
Homepage de repositorios públicos: compartir, exposición

**¿Cómo crear un repositorio online?**

Profile page -> click create a new repo

Escribir nombre y descripción
Seleccionar "Public"
Marcar "Initialize …"

### Create a new repository
A repository contains all the files for your project, including the revision history.

Owner          Repository name
🖥 jlbazo ▾  /  FIS725  ✓

Great repository names are short and memorable. Need inspiration? How about **sturdy-computing-machine**.

**Description** (optional)

Files for Data Analysis classes

⦿ 📖 **Public**
   Anyone can see this repository. You choose who can commit.

◯ 🔒 **Private**
   You choose who can see and commit to this repository.

☑ **Initialize this repository with a README**
   This will let you immediately clone the repository to your computer. Skip this step if you're importing an existing repository.

Add .gitignore: **None** ▾     Add a license: **None** ▾   ⓘ

**Create repository**

# Git: copia local

**¿Cómo copiar repositorio creado en disco local?**

Usando **Git Bash**:

$ mkdir FIS725          ← Crear directorio donde copiar repositorio
$ cd FIS725

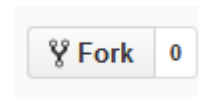$ git init          ← inicializar repositorio local

$ git remote add origin https://github.com/jlbazo/FIS725.git          ← dirección externa

$ git pull origin master

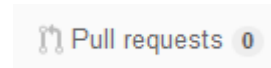También se pueden copiar repositorios de otros usuarios (**fork**)
(hacerlo primero en GitHub)

⑂ Fork   0
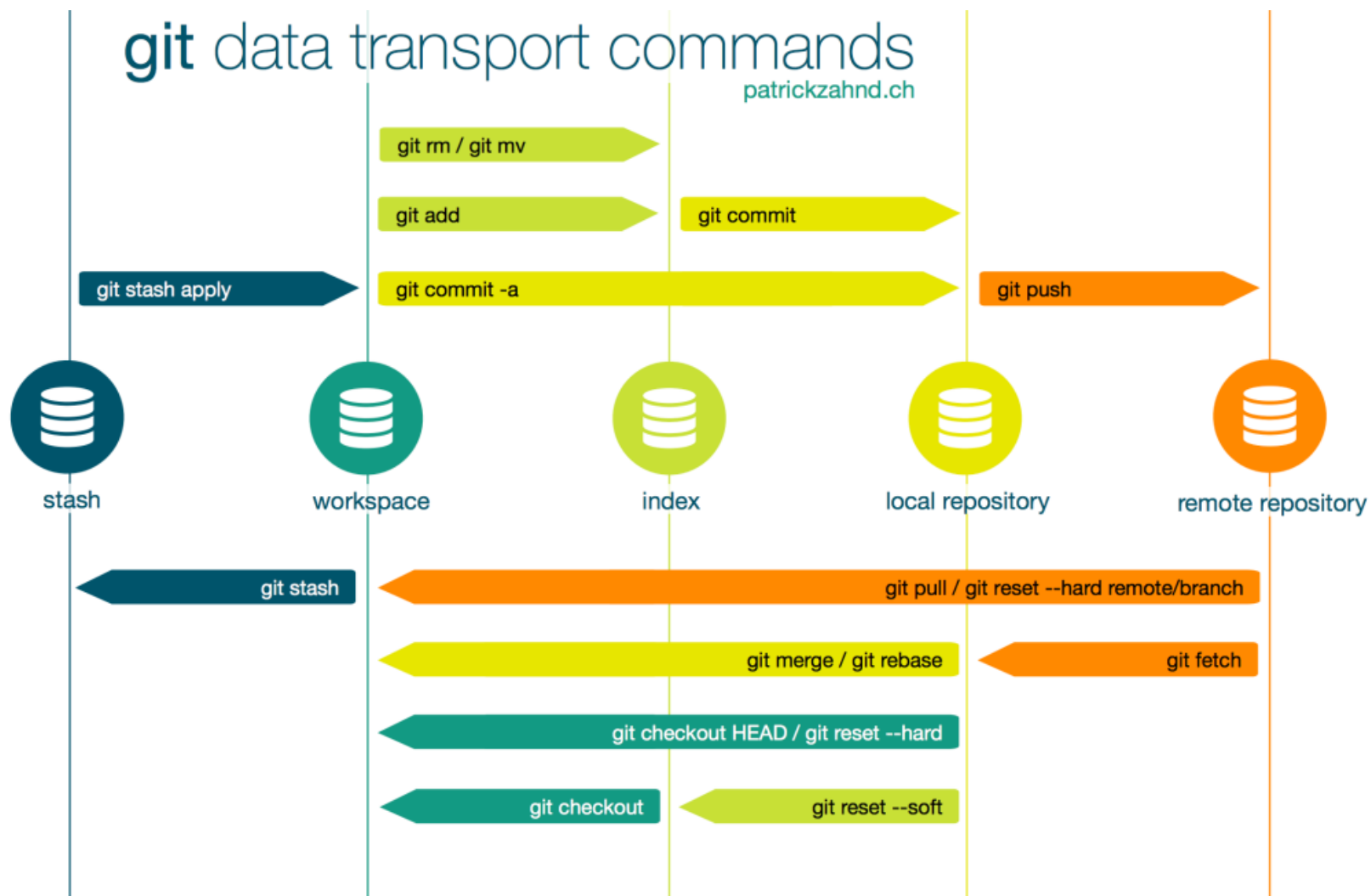
$ git clone https://github.com/jlbazo/forked_repo.git          ← Copia repositorio

Para incluir cambios de forked repo pedir "Pull request" en GitHub

⑂ Pull requests 0

git data transport commands
patrickzahnd.ch

# Git: comandos

$ git add -A ⟵ añadir (y actualizar) nuevos archivos al index, antes de *commit*

$ git commit –m "message" ⟵ Guardar cambios en repositorio local

$ git push ⟵ Guardar cambios en repositorio remoto (GitHub)

$ git status

```
$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
nothing to commit, working directory clean
```

[Tutorial](#)

Crear nueva versión (branch) del proyecto para editar sin interferir con otros:

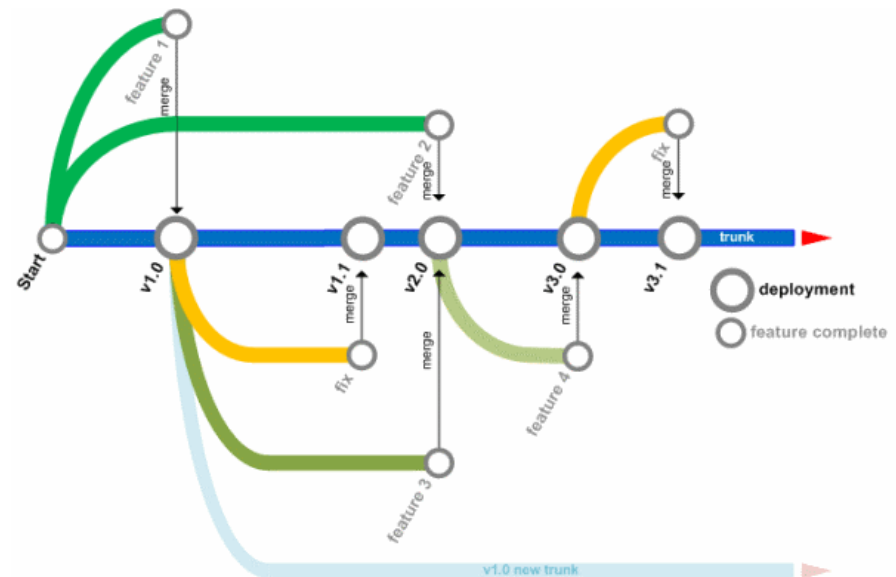$ git checkout –b *branchname*      ⟵      crear nueva rama

$ git branch      ⟵      ver rama actual

$ git checkout master      ⟵      Regresar a rama principal



$ git pull origin *branchname*

$ git push origin *branchname*

**PUCP**

**Instalar** R desde  CRAN     Comprehensive R Archive Network

Instalar R Studio       **R** Studio          mejor interfaz gráfica

Instalar **paquetes adicionales** (más de 8000) cuando sean necesarios:

> available.packages()

> install.packages("ggplot2")

> library(ggplot2)        ⟵          Cargar paquete

> search()        ⟵     Ver paquetes cargados

> find.package("ggplot2")        ⟵        comprobar si paquete está instalado

Instalando paquetes
 externos a CRAN        ⟶
(ej. BioConductor)

> source("http://bioconductor.org/biocLite.R")

> biocLite("rhdf5")

**Para Windows**

Instalar Rtools          Necesarias para los paquetes de R

Seleccionar versión de R correspondiente

Dejar al instalador seleccionar el *path* (marcar casilla correspondiente)

> install.packages("devtools")

> library(devtools)

> find_rtools()          ⟵          Si regresa TRUE, instalación exitosa

"Try first thyself, and after call in God; for to the worker God himself lends aid."
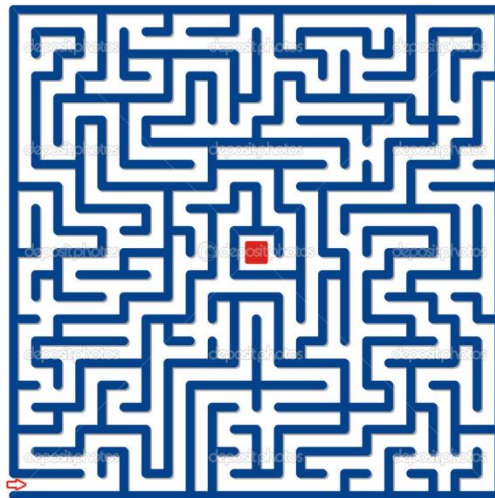
Εὐριπίδης

Siempre buscar en inglés para obtener más resultados

No hay necesidad de reinventar la rueda. Ya existe una solución para tu problema.

1. Buscar en internet

2. Leer manuales y archivos de ayuda

3. Preguntar a colegas expertos

4. Preguntar en un foro (antes agotar otros medios)

Escribir un título descriptivo, adjuntar código reproducible y datos del sistema usado.

Tener paciencia e insistir.

link

# To R or not to R?

| | SAS | R | SPSS |
|---|---|---|---|
| **CURRENT VERSION** | 9.4<br>JULY, 2013 | Spring Dance, R-3.1.0<br>APRIL, 2014 | 22.0<br>AUGUST, 2013 |
| **HISTORY** | **Creator**: Jim Goodnight and Jim Barr, North Carolina State University<br><br>**Year Released**: mass distributed since 1972<br><br>**Must Knows**:<br><br>• SAS started because of a need for a computerized statistics program to analyze vast amounts of agricultural data<br>• The SAS institute was founded in 1976 and currently has 13,733 employees<br>• In 2013, SAS invested 25% of revenue in R&D | **Creator**: Ross Ihaka and Robert Gentleman, University of Auckland, New Zealand and the R foundation<br><br>**Year Released**: 1995<br><br>**Must Knows**:<br><br>• R is an implementation of the S programming language created at Bell Labs<br>• The design and evolution of R is controlled by the R-core group and R foundation<br>• The source code for the R software environment is written primarily in C, Fortran and R. | **Creator**: Norman H. Nie, Dale H. Bent, and Hadlai "Tex" Hull<br><br>**Year Released**: 1968<br><br>**Must Knows**:<br><br>• In 1976 SPSS jeopardized the University of Chicago's status as a tax-exempt organization<br>• SPSS was acquired by IBM in 2009 for US$1.2 billion<br>• In 1993 SPSS was taken public on the NASDAQ exchange |
| **PURPOSE AND USABILITY** | • SAS accumulated since the 1970s a large amount of high-quality production code for multiple purposes<br>• SAS has a strong leading position in the commercial analytics space. Code legacy plays an important role here<br>• SAS has strong data handling capabilities. Furthermore, it releases its software updates in a controlled environment, which make them well tested. Nevertheless, SAS is an expensive solution. | • R has been used in academics and research for a long time. Today, its finding its way into commercial applications as well. See R as the open-source counterpart of SAS.<br>• R has advanced graphical capabilities thanks to for example packages like ggplot2, googleVis and rCharts.<br>• Due to its open-source nature, R has a large and supportive community. The latest techniques are developed and released quickly. | • SPSS is a great tool for non-statisticians since it has a user-friendly Interface and easy-to-use drop down menus.<br>• Just like SAS, SPSS has a rather hefty price tag.<br>• SPSS has applications in many fields, but mainly plays a leading role in social sciences. |

**COMPANIES USING IT**

HP · Scotiabank · AIG · STAPLES · Facebook · FDA · Google · The New York Times · ABN·AMRO · Canon · BT · CREDIT SUISSE

**EASE OF LEARNING**

Although it's not like learning Microsoft Word, getting a basic understanding of how to work with SAS **shouldn't take you to long**. However, to become really good you will need to work through a lot of specifics.

There are **many official and unofficial tutorials** available, and official certifications can be obtained via SAS training institutes.
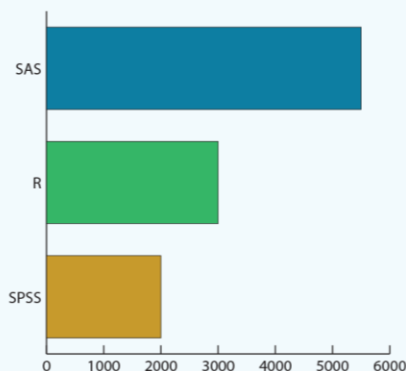
R has a reputation for being hard to learn. Instead of setting up a complete analysis at once, R users need to **learn how to analyze data interactively**. For most data analysts, this is a mind shift they first need to undergo.

The open-source community of R is rapidly lowering this learning curve by creating **high-quality introductory tutorials** and interactive coding tutorials.
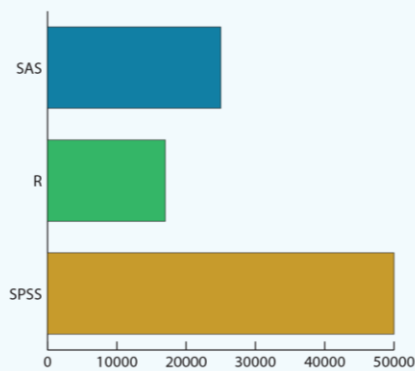
SPSS is by far the **easiest to learn** among the 3 languages listed here. So if you only open a statistical program twice a month SPSS is the way to go.

One of the biggest advantages in terms of learning is its **similarities with Excel**, something most of us are familiar with.

**MARKETABILITY**



Number of analytics jobs on Indeed.com 2/2014

Use of analytic software in academia 05/2013. Based on number of google scholar hits
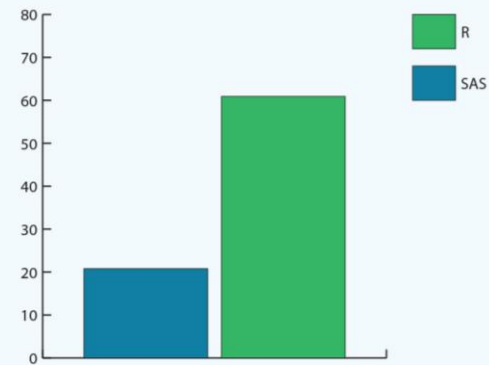
**POPULARITY**

kaggle
50% of Kaggle winners use R

stackoverflow
The number of R related posts on Stack Overflow is more than 7-fold the number of posts on SAS

Percentage of "What programming/statistics languages you used for an analytics / data mining / data science work in 2013?" (KD nuggets)