



# CURSOS DE VERANO 2014

**APROXIMACIÓN PRÁCTICA A LA CIENCIA DE DATOS Y BIG DATA: HERRAMIENTAS KNIME, R, HADOOP Y MAHOUT**

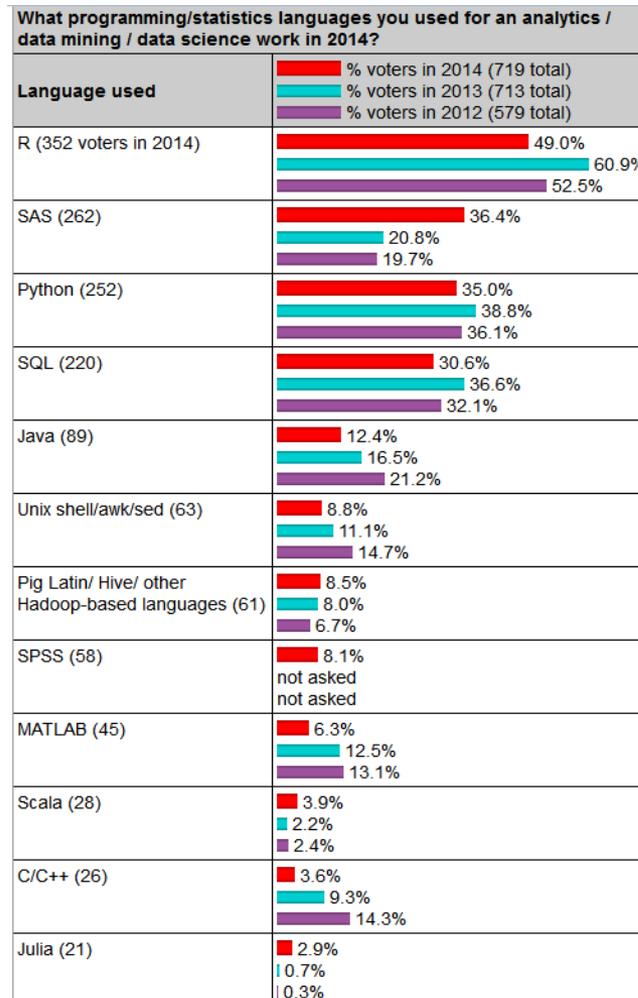
**Introducción a R**

**Francisco Charte Ojeda**

# Contenidos

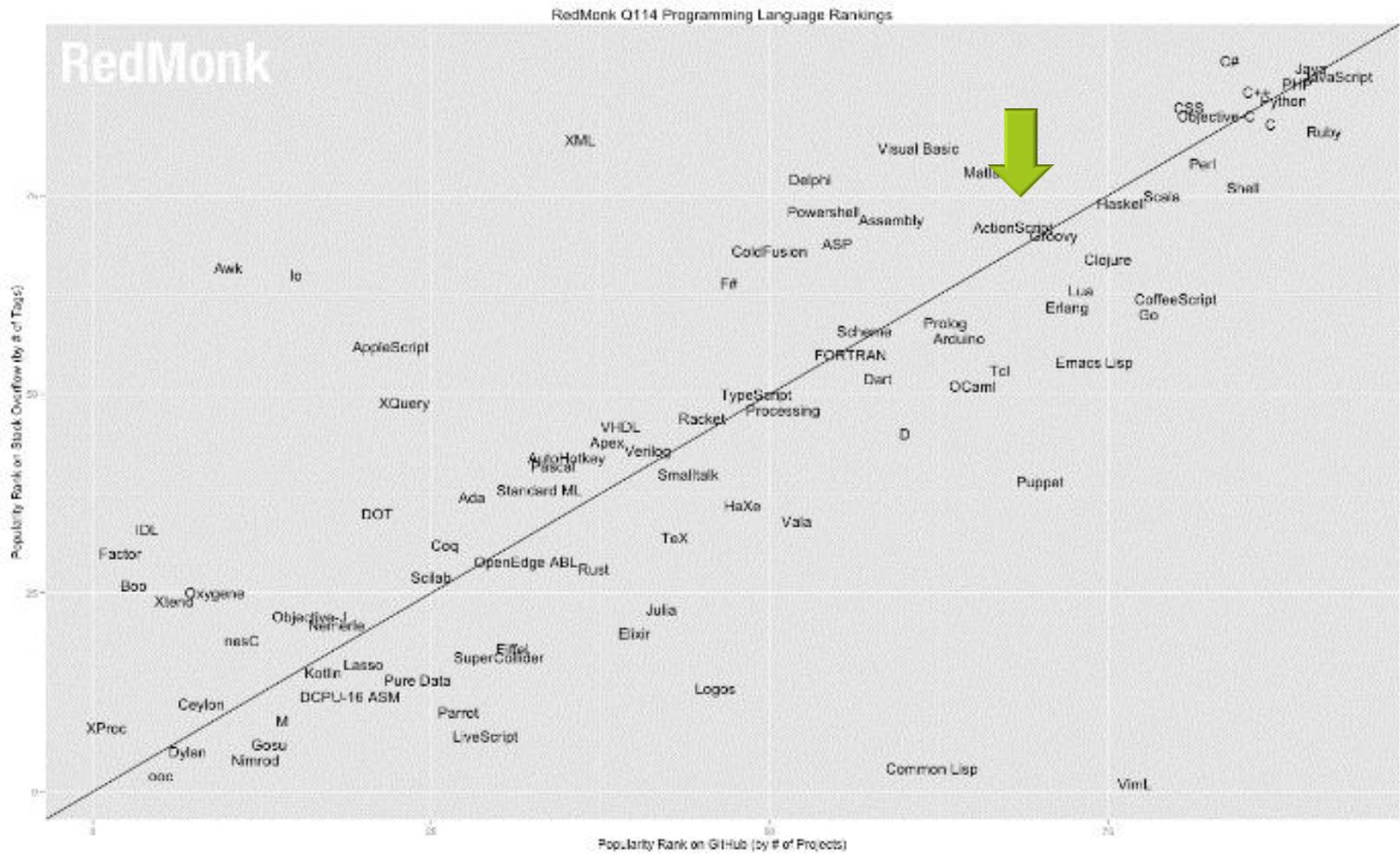
- ¿Por qué aprender R?
- Herramientas de trabajo
- Tipos de datos
- Carga de datos
- Tratamiento de datos ausentes
- Análisis exploratorio

# ¿Por qué aprender R? ► Lenguaje más usado para análisis de datos



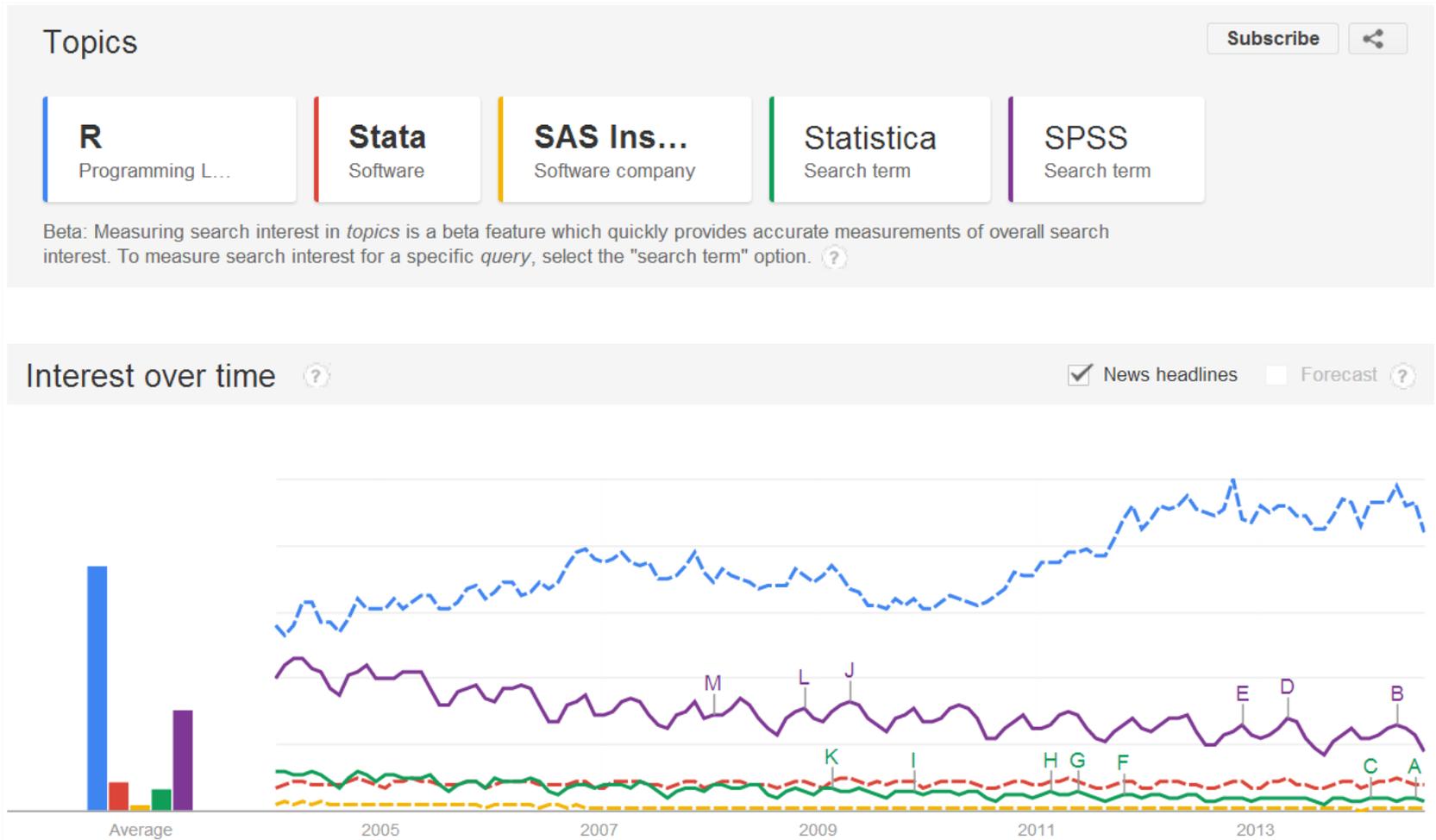
Fuente: KDnuggets poll - Languages for analytics/data mining (Aug 2014)

# ¿Por qué aprender R? ► 15º lenguaje más activo en GitHub/Stack Overflow



Fuente: RedMonk Programming Languages Ranking 2014

# ¿Por qué aprender R? ▶ Es el software de analítica de datos más usado



Fuente: [Google Trends](#) (Abrir en navegador)

¿Por qué aprender R? ► Conocimiento de R valorado más que ningún otro

AVERAGE SALARY FOR <b>High Paying Skills and Experience</b>		
SKILL	2013	YR/YR CHANGE
<b>R</b>	\$ 115,531	n/a
<b>NoSQL</b>	\$ 114,796	1.6%
<b>MapReduce</b>	\$ 114,396	n/a
<b>PMBok</b>	\$ 112,382	1.3%
<b>Cassandra</b>	\$ 112,382	n/a
<b>Omnigraffle</b>	\$ 111,039	0.3%
<b>Pig</b>	\$ 109,561	n/a
<b>SOA (Service Oriented Architecture)</b>	\$ 108,997	-0.5%
<b>Hadoop</b>	\$ 108,669	-5.6%
<b>Mongo DB</b>	\$ 107,825	-0.4%

## ¿Por qué aprender R? ► Además ...

- R es *Open Source* (multiplataforma, libre, abierto, etc.)
- Gran número de paquetes disponibles
- Extensa comunidad de usuarios
- Ciclo completo de trabajo:
  - Implementación de algoritmos
  - Preparación de datos
  - Análisis de resultados
  - Generación de documentación

## Herramientas de trabajo

- Consola de R
- Editor de texto (Emacs, Vim, Notepad++)
- RStudio
- RStudio Server

## Herramientas de trabajo ► R

- Binarios disponibles para Linux, OS X y Windows
- Descarga desde <http://www.r-project.org/>
- Disponible en repositorios Linux

```
francisco@Ubuntu14LTS: ~  
francisco@Ubuntu14LTS:~$  
francisco@Ubuntu14LTS:~$ sudo apt-get install r-base r-base-dev
```

## Herramientas de trabajo ► Consola de R

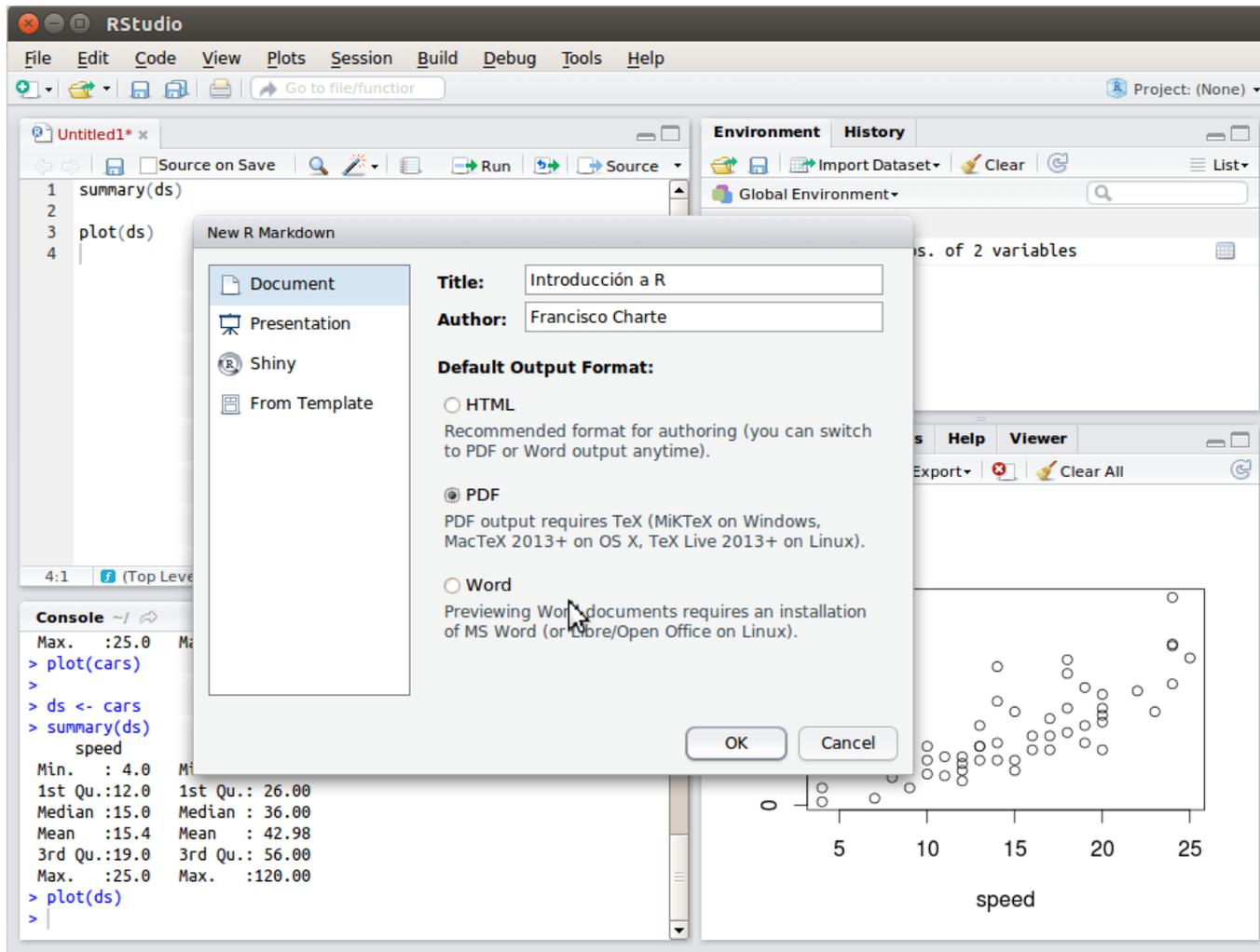
- Trabajo interactivo mediante línea de comandos

```
francisco@Ubuntu14LTS: ~  
francisco@Ubuntu14LTS:~$ R  
  
R version 3.0.2 (2013-09-25) -- "Frisbee Sailing"  
Copyright (C) 2013 The R Foundation for Statistical Computing  
Platform: x86_64-pc-linux-gnu (64-bit)  
  
R es un software libre y viene sin GARANTIA ALGUNA.  
Usted puede redistribuirlo bajo ciertas circunstancias.  
Escriba 'license()' o 'licence()' para detalles de distribución.  
  
R es un proyecto colaborativo con muchos contribuyentes.  
Escriba 'contributors()' para obtener más información y  
'citation()' para saber cómo citar R o paquetes de R en publicaciones.  
  
Escriba 'demo()' para demostraciones, 'help()' para el sistema on-line de ayuda,  
o 'help.start()' para abrir el sistema de ayuda HTML con su navegador.  
Escriba 'q()' para salir de R.  
  
> 2 + 3  
[1] 5  
>  
>
```

## Herramientas de trabajo ► RStudio

- Binarios disponibles para Linux, OS X y Windows
- Descarga desde <http://www.rstudio.com/>
- Licencia *Open Source* y comercial
- IDE estándar para trabajar con R
- **Será la herramienta que usemos en el curso**

# Herramientas de trabajo ► RStudio: IDE completo para trabajar con R



The screenshot displays the RStudio IDE interface. A 'New R Markdown' dialog box is open in the center, with the 'Document' tab selected. The dialog contains the following information:

- Title:** Introducción a R
- Author:** Francisco Charte
- Default Output Format:**
  - HTML
  - PDF
  - Word

The background shows the RStudio workspace with the following elements:

- Source Editor:** Contains R code:
 

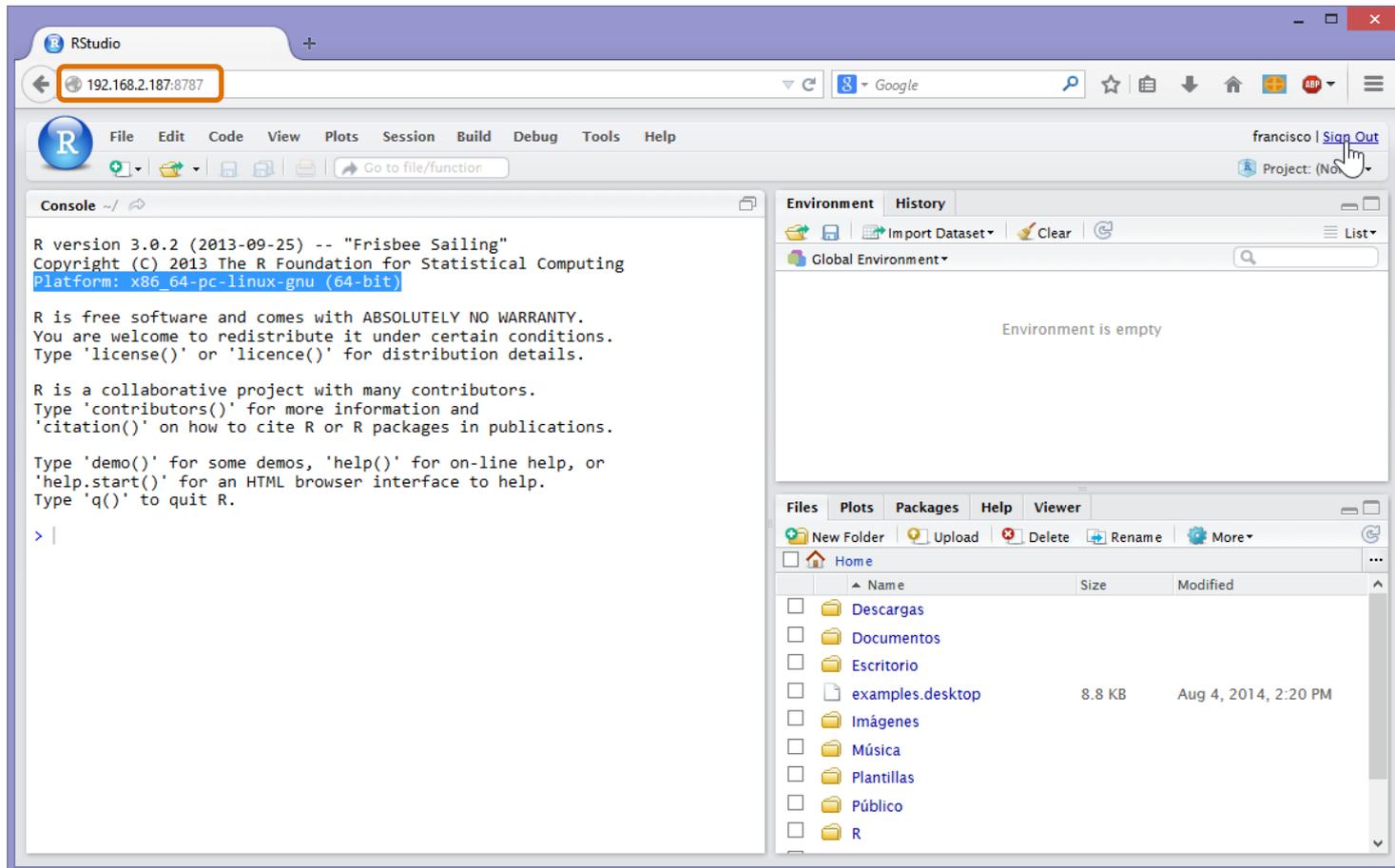
```
1 summary(ds)
2
3 plot(ds)
4
```
- Console:** Shows the execution of the code:
 

```
> plot(cars)
>
> ds <- cars
> summary(ds)
  speed
Min.   : 4.0  1st Qu.: 12.0  Median : 15.0  Mean   : 15.4  3rd Qu.: 19.0  Max.   : 25.0
  miles 1st Qu.: 26.00 Median : 36.00 Mean   : 42.98 3rd Qu.: 56.00 Max.   :120.00
> plot(ds)
>
```
- Environment/History:** Shows the 'Global Environment'.
- Viewer:** Displays a scatter plot of 'speed' vs. 'miles' (partially visible).

## Herramientas de trabajo ► >> tareasHabituales.R <<

- Acceso a la documentación
  - `help('source')`, `vignette('grid')`, `demo('image')`
- Ruta de trabajo
  - `getwd()`, `setwd()`
- Espacio de trabajo
  - `save()`, `save.image()`, `load()`
- Instalación de paquetes
  - `install.packages()`, `library()`

# Herramientas de trabajo ► RStudio Server, accesible desde un navegador

A screenshot of a web browser displaying the RStudio Server interface. The browser's address bar shows the URL '192.168.2.187:8787'. The RStudio interface includes a menu bar (File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help), a toolbar with icons for file operations, and a 'Go to file/function' search bar. The main area is divided into three panes: a Console pane on the left showing the R startup message, an Environment pane on the top right showing 'Global Environment' and 'Environment is empty', and a Files pane on the bottom right showing a file explorer view of the 'Home' directory. The Files pane lists folders like 'Descargas', 'Documentos', 'Escritorio', 'Imágenes', 'Música', 'Plantillas', 'Público', and 'R', along with a file 'examples.desktop' (8.8 KB, modified Aug 4, 2014, 2:20 PM). The browser's user interface includes a search bar, navigation icons, and a user profile 'francisco' with a 'Sign Out' link.

# Tipos de datos

▣ Simples

▣ Vectores

▣ Matrices

▣ Factors

▣ Data Frames

▣ Listas

} >> tiposDatos.R <<

} >> tiposDatosII.R <<

## Tipos de datos ► Simples

- numeric
  - Enteros :: 1024, -3
  - Punto flotante :: 3.1415927
  - Notación exponencial :: 3.85e6
  - Otros :: Inf, NaN
  
- integer :: `as.integer(numeric)`
  
- complex :: 1+2i
  
- character :: 'R', "Hola"
  
- logical :: TRUE, FALSE, NA

## Tipos de datos ► Variables

### □ Asignación

□ `a = 1024` | `a <- 1024` | `1024 -> a`

### □ Obtención de clase y tipo

□ `class(a)` # numeric | `typeof(a)` # double

### □ Comprobación de tipo

□ `is.numeric(a)`, `is.character(a)`, `is.integer(a)`,  
`is.infinite(a)`, `is.na(a)`

### □ Objetos en el espacio de trabajo

□ `ls()`, `rm(var)`, `str(var)`,  
`save(var, file = arch)`, `save.image()`, `load()`

## Tipos de datos ► Vectores

### □ Definición

- `diasMes <- c(31,29,31,30,31,30,31,31,30,31,30,31)`
- `dias <- c('Lun','Mar','Mié','Jue','Vie','Sáb','Dom')`
- `quincena <- 16:30`
- `semanas <- seq(1, 365, 7)`
- `rep(T, 5)`

### □ Obtención número de elementos

- `length(dias)`

### □ Acceso a elementos

- `dias[2]` # Segundo elemento del vector
- `dias[-2]` # Todos los elementos menos el segundo
- `dias[c(3, 7)]` # Elementos 3 y 7

## Tipos de datos ► Vectores

- Generación de valores aleatorios
  - Establecimiento de la semilla: `set.seed(4242)`
  - Distribución normal
    - `rnorm(100, mean = 10, sd = 3)`
  - Distribución uniforme
    - `runif(6, min = 1, max = 49)`
  - Otras distribuciones
    - `rbinom()`, `rlogis()`, `rpois()`, etc.
  
- Operaciones sobre vectores
  - No necesaria la iteración por los elementos
  - Posibilidades de paralelización

## Tipos de datos ► Matrices

### □ Definición

- `mes <- matrix(1:35, nrow = 5)`
- `mes <- matrix(1:35, ncol = 7, byrow = T)`

### □ Obtención número de elementos

- `length(mes) | nrow(mes) | ncol(mes)`

### □ Acceso a elementos

- `mes[2, ]` # Segunda fila completa
- `mes[, 2]` # Segunda columna completa
- `mes[2, 5]` # Quinta columna de la segunda fila
- `fix(mes)` # Edición de elementos en la matriz

## Tipos de datos ► Factors

### □ Definición

- `herramientas <- factor('Consola', 'RStudio')`
- `fdias <- factor(días)`
- `tam <- ordered(c('Ligero', 'Medio', 'Pesado'))`

### □ Obtención niveles

- `nlevels(fdias)`
- `levels(días)`

### □ Relación de orden (factors ordenados)

- `tam[2] < tam[1] # FALSE`

## Tipos de datos ► Data Frames

### □ Definición

- `df <- data.frame(vect1, ..., vectN)`

- `df <- data.frame(matrix)`

- `df <- data.frame(col1 = tipo(N), ..., colN = tipo(N))`

### □ Ejemplo

- `df <- data.frame(Dia = fdias,  
                  Estimado = rep(c(T, F), 7),  
                  Lectura = rnorm(14, 5))`

### □ Obtención número de elementos

- `nrow(mes)`

- `ncol(mes)`

## Tipos de datos ► Data Frames

### ■ Selección y proyección de datos

- `df[5, ]` # 5ª fila
- `df[ ,3]` # 3ª columna
- `df[c(-3,-6), ]` # Menos 3ª y 6ª fila
- `df$Lectura` # 3ª columna
- `df$Lectura[5]` # 5ª fila de 3ª col.
- `df[, c('Dia','Lectura')]` # Columnas 1 y 3
- `df[df$Estimado == F, ]` # Filas condición
- # Selección de filas y columnas  
`df[df$Estimado == F & df$Lectura > 3,  
c('Dia', 'Lectura')]`

## Tipos de datos ► Data Frames

- Agregar nuevas filas
  - `df[15, 1] <- 'vie'`
  - `df$Dia[15] <- 'vie'`
  - `rbind(df, data.frame(Dia=fdias[1], Est=F, Lect=5))`
  - `rbind(df[1:9, ], nuevaFila, df[10:14, ])`
  
- Agregar nuevas columnas
  - `df$Ajustado <- df$Lectura + rnorm(15, 2)`
  - `cbind(df, Ajustado = df$Lectura + rnorm(15, 2))`
  - `cbind(df[ ,c(1,3)], nuevaCol, df$Estimado)`
  
- Nombres de filas y columnas
  - `names(df)` # vector con nombres de columnas
  - `rownames(df)` # vector con nombres de filas

## Tipos de datos ► Listas

### □ Definición

- `l1st <- list(3.1415927, 'Hola', TRUE, fdias[4])`
- `l1st <- list(fdias, mes, df)`

### □ Información sobre la lista

- `length(l1st)`
- `names(l1st)`

### □ Acceso a los elementos

- `l1st[[1]]`
- `l1st[['PI']]`
- `l1st$PI`

## Carga de datos ► >> cargaDatos.R <<

- ▣ Lectura de archivos CSV
- ▣ Importación de hojas de cálculo Excel
- ▣ Carga de datasets en formato ARFF
- ▣ Obtención de datos de otras fuentes

## Carga de datos ► csv

- `datosCSV <- read.table(  
 file = "miArchivo.csv",  
 header = T,  
 sep = ",",  
 dec = ".",  
 quote = "\"")`
- `read.csv("miArchivo.csv") # sep="," , dec="."`
- `read.csv2("miArchivo.csv") # sep=";" , dec=","`

## Carga de datos ► Excel

- Múltiples posibilidades
  - Exportar desde Excel a CSV
  - Copiar datos al portapapeles
  - Leer archivo Excel desde R
  
- Paquetes R para trabajar con archivos Excel
  - XLConnect
    - `datos <- readWorksheetFromFile('archivo.xls', sheet=n)`
  - xlsx
    - `datos <- read.xlsx('archivo.xlsx', sheetName = n, rango)`
  
- `vignette(paquete)` # Abrir el manual asociado

## Carga de datos ► ARFF

- Paquete `foreign`
  - Funciones para leer múltiples formatos de archivo
  - `read.arff('dataset.arff')`
  
- Paquete `RWeka`
  - Interfaz completa entre R y Weka
    - Leer y escribir archivos ARFF
    - Acceso a algoritmos de clasificación, agrupamiento, etc.
  - `read.arff('dataset.arff')`

## Carga de datos ► Otras fuentes

### □ Portapapeles

- `read.delim('clipboard')`
- `write.table(datos, 'clipboard')`

### □ Desde URL

- `conn <- getURL('http://url/datos')` # Conexión abierta
- `datos <- read(conn)`
- `conn <- getURL('https://url/datos')` # Conexión cifrada
- `datos <- read(textConnection(conn))`

### □ Datasets integrados

- `data()` # Lista de todos los datasets integrados
- `summary(iris)`

## Tratamiento de datos ausentes ► >> `tratamientoNulos.R` <<

- Problemática
  - Datos ausentes (*missing values*) dificultan múltiples operaciones
- Detectar existencia de valores ausentes
  - `is.na(variable)`
  - `na.fail(variable)`
- Eliminar valores ausentes
  - `na.omit(variable)`
  - `complete.cases(variable)`
  - `variable[is.na(variable)] <- valor`
- Operar con presencia de valores ausentes
  - `mean(variable, na.rm = T)`
  - `lm(x ~ y, variable, na.action = na.omit)`

## Análisis exploratorio ► Información general >> análisisExploratorio.R <<

- Estructura interna de la variable
  - `str(variable)`
  
- Resumen del contenido
  - `summary(variable)`
  
- Exploración del contenido
  - `head(variable)` | `tail(variable)`
  - `variable[filas, columnas]`
  - `variable$columna`
  - `variable$columna[which(condición)]`
    - `iris$Sepal.Length[which(iris$Species == 'versicolor')]`

## Análisis exploratorio ► Estadística descriptiva

### ■ Funciones básicas (operan sobre vectores)

- `mean` # media
- `median` # mediana
- `var` # varianza
- `sd` # desviación estándar
- `max` # máximo valor
- `min` # mínimo valor
- `range` # rango de valores
- `quantile` # cuantiles

### ■ Para estructuras complejas

- `lapply(iris[,1:4], mean)` # Aplicar a cada columna
- `describe(variable)` # Paquete Hmisc

## Análisis exploratorio ► Agrupamiento de datos

- Tabla de contingencia con número de combinaciones
  - Longitud de sépalo según especie  
`table(iris$Sepal.Length, iris$Species)`
  - Valoración de vendedores según moneda  
`table(ebay$sellerRating, ebay$currency)`
- Agrupamiento y selección
  - Separar los casos por especie de flor  
`split(iris, iris$Species)`
  - Obtener elevación, pendiente y clase de filas que cumplan condición  
`subset(covertypes, slope > 45 & soil_type == '1',  
select = c(elevation, slope, class))`

## Análisis exploratorio ► Ordenación de datos

- Ordenar un vector obteniendo otro
  - `sort(valores)`
  
- Obtener la posición para cada valor
  - `order(valores)`
  
- Generar un ranking a partir de los valores
  - `rank(valores)`
  - `rank(valores, ties.method = 'average')`

## Análisis exploratorio ► Particionamiento de datos

- Tomando el orden en que aparecen en el dataset
  - División entre entrenamiento (75%) y prueba (25%)

```
nTraining <- as.integer(nrow(iris) * .75)
training  <- iris[1:nTraining, ]
test      <- iris[(nTraining+1):nrow(iris), ]
```

- Tomando un subconjunto aleatorio
  - Misma proporción anterior

```
set.seed(4242)          # Asegurar reproducibilidad
indices <- sample(1:nrow(iris), nTraining)
training <- iris[indices, ]
test     <- iris[-indices, ]
```



# CURSOS DE VERANO 2014

**APROXIMACIÓN PRÁCTICA A LA CIENCIA DE DATOS Y BIG DATA: HERRAMIENTAS KNIME, R, HADOOP Y MAHOUT**

**Introducción a R**

**Francisco Charte Ojeda**