

Manejo de datos con R

Oscar Perpiñán Lamigueiro

25 de Enero de 2013

Contenidos

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de formato

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

Fuentes de datos

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

- ▶ The R Datasets Package
- ▶ Enlaces en Bibsonomy

Contenidos

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de formato

setwd, getwd, dir

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

```
getwd()  
old <- setwd("~/R/intro")  
dir()  
dir(pattern='.R')  
dir('data')
```

► Con un fichero local

```
download.file('http://oscarperpinan.github.com/  
spacetime-vis/data/C02_GNI_BM.csv',  
              destfile='data/C02_GNI_BM.csv')
```

```
C02 <- read.table('data/C02_GNI_BM.csv', header=TRUE,  
                  sep=',')  
head(C02)
```

► Directamente de un enlace URL

```
C02 <- read.table('http://oscarperpinan.github.com/  
spacetime-vis/data/C02_GNI_BM.csv',  
                  header=TRUE, sep=',')  
head(C02)
```

read.csv, read.csv2

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

- `read.csv` y `read.csv2` son como `read.table` con valores por defecto para encabezado y separadores

```
C02 <- read.csv('data/C02_GNI_BM.csv')
```

```
head(C02)
```

```
tail(C02)
```

```
summary(C02)
```

```
names(C02)
```

Contenidos

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de formato

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato


```
chromo <- data.frame(chromosome = gl(3, 10,  
  labels = c('A', 'B', 'C')),  
  probeset = gl(3, 10,  
  labels = c('X', 'Y', 'Z')),  
  ensg = gl(3, 10,  
  labels = c('E1', 'E2', 'E3')),  
  symbol = gl(3, 10,  
  labels = c('S1', 'S2', 'S3')),  
  XXA_00 = rnorm(30),  
  XXA_36 = rnorm(30),  
  XXB_00 = rnorm(30))
```

```
table(chromo$chromosome, df$XXA_00 > 0)  
table(chromo$probeset, df$XXA_00 > -1 & df$XXA_00 <  
1)
```

```
xtabs(XXA_00 > 1 ~ chromosome + probeset,  
      data=chromo)
```

```
tapply(CO2$X2000, CO2$Indicator.Name,
       FUN=mean)
```

CO2 emissions (kg per PPP \$ of GDP)	4.777875e-01
CO2 emissions (metric tons per capita)	7.580861e+00
GNI per capita, PPP (current international \$)	1.981000e+04
GNI, PPP (current international \$)	2.078196e+12

```
tapply(CO2$X2000, CO2[,c("Indicator.Name", "Country.
Name")],
FUN=mean)
```

Indicator.Name	Country.Name	
	Brazil	China
CO2 emissions (kg per PPP \$ of GDP)	2.699746e-01	1.140619e+00
CO2 emissions (metric tons per capita)	1.892645e+00	2.696862e+00
GNI per capita, PPP (current international \$)	6.820000e+03	2.340000e+03
GNI, PPP (current international \$)	1.188790e+12	2.948850e+12

Indicator.Name	Country.Name	
	Finland	France
CO2 emissions (kg per PPP \$ of GDP)	3.923481e-01	2.384221e-01
CO2 emissions (metric tons per capita)	1.007322e+01	6.016236e+00
GNI per capita, PPP (current international \$)	2.548000e+04	2.566000e+04
GNI, PPP (current international \$)	1.318800e+11	1.558990e+12

Indicator.Name	Country.Name	
	Germany	Greece
CO2 emissions (kg per PPP \$ of GDP)	3.929031e-01	4.598579e-01
CO2 emissions (metric tons per capita)	1.012147e+01	8.391709e+00
GNI per capita, PPP (current international \$)	2.549000e+04	1.832000e+04
GNI, PPP (current international \$)	2.095450e+12	2.000130e+11

Indicator.Name	Country.Name	
	India	Norway
CO2 emissions (kg per PPP \$ of GDP)	7.448517e-01	2.391275e-01
CO2 emissions (metric tons per capita)	1.125975e+00	8.641315e+00
GNI per capita, PPP (current international \$)	1.500000e+03	3.565000e+04
GNI, PPP (current international \$)	1.575930e+12	1.601000e+11

Indicator.Name	Country.Name	
	Spain	United States
CO2 emissions (kg per PPP \$ of GDP)	3.428950e-01	5.568755e-01
CO2 emissions (metric tons per capita)	7.312922e+00	1.953626e+01
GNI per capita, PPP (current international \$)	2.115000e+04	3.569000e+04

aggregate

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

```
aggregate(X2000 ~ Indicator.Name,  
          data=C02, FUN=mean)
```

```
aggregate(cbind(X2000, X2001) ~ Indicator.Name,  
          data=C02, FUN=mean)
```

```
aggregate(X2000 ~ Indicator.Name + Country.Name,  
          data=C02, FUN=mean)
```

aggregate

```
aggregate(cbind(X2000, X2001) ~  
  Indicator.Name + Country.Name,  
  data=C02, FUN=mean)
```

```
aggregate(cbind(X2000, X2001) ~  
  Indicator.Name + Country.Name,  
  data=C02, FUN=mean)
```

```
aggregate(cbind(X2000, X2001) ~  
  Indicator.Name + Country.Name,  
  subset=(Country.Name %in% c('United_States',  
    'China')),  
  data=C02, FUN=mean)
```

aggregate

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

```
aggregate(cbind(XXA_00, XXA_36, XXB_00) ~  
          ensg + chromosome + symbol,  
          data = chromo, FUN = mean)  
  
aggregate(cbind(XXA_00, XXA_36, XXB_00) ~ ensg ,  
          data = chromo, FUN = mean)
```

Contenidos

Manejo de datos
con R

Oscar Perpiñán
Lamigueiro

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de formato

Fuentes de datos

Lectura de datos

Datos agregados

Cambio de
formato

- Primero escogemos un subconjunto

```
C02China <- subset(C02,  
  subset=(Country.Name=='China' &  
    Indicator.Name=='C02_emissions_  
      (kg_per_PPP_$_of_GDP)'),  
  select=-c(Country.Name, Country.Code,  
    Indicator.Name, Indicator.  
      Code))
```

- Pasamos de formato wide a long

```
stack(C02China)
```

reshape: wide a long

► Primer intento

```
C02long <- reshape(C02,  
                    varying=list(names(C02)[5:16]),  
                    direction='long')  
head(C02long)
```

► Añadimos argumentos

```
C02long <- reshape(C02,  
                    varying=list(names(C02)[5:16]),  
                    timevar='Year', v.names='Value',  
                    times=2000:2011,  
                    direction='long')  
head(C02long)
```

reshape: long a wide

- ▶ Primero escogemos las columnas de interés

```
C02subset <- C02long[c("Country.Name",  
                      "Indicator.Name",  
                      "Year", "Value")]
```

- ▶ Ahora cambiamos formato

```
C02wide <- reshape(C02subset,  
                  idvar=c('Country.Name','Year'),  
                  timevar='Indicator.Name',  
                  direction='wide')
```

- ▶ Y ponemos nombres al gusto

```
names(C02wide)[3:6] <- c('C02.PPP', 'C02.capita',  
                        'GNI.PPP', 'GNI.capita')
```

```
head(C02wide)
```