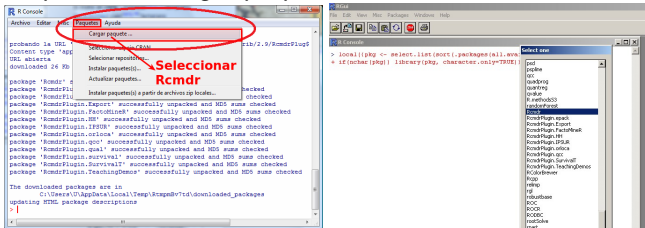


# Práctica 6. Regresión lineal

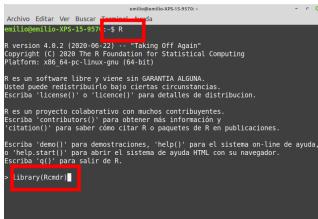
Departamento de Estadística  
Universidad de Oviedo

# Cargar el programa R y el paquete Rcommander

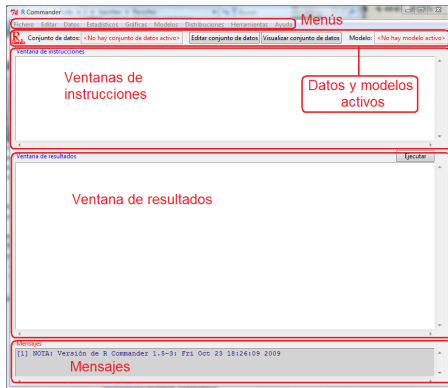
- Iniciamos el programa R.
- Cargamos el paquete RCommander. Dos opciones:
  - 1 Menú *Paquetes* → *Cargar paquete* → Seleccionamos **Rcmdr**.



- 2 Escribimos **library(Rcmdr)** en la consola y pulsamos retorno de carro.



# Cargar Rcommander

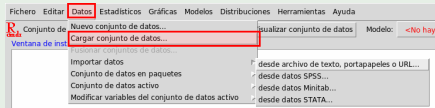


El fichero `acero.rda` se encuentra en el Campus Virtual. Hay que haberlo descargado previamente.

## Cargar la base de datos `acero.rda`

Datos

- ➔ Cargar conjunto de datos
- ➔ Seleccionar **`acero.rda`**



```
> load("/home/emilio/clases/acero.rda")
```

NOTA: El conjunto de datos `acero` tiene 117 filas y 20 columnas.

- Para visualizar la base de datos:

R Commander interface showing the 'acero' dataset loaded. The 'Visualizar conjunto de datos' button is highlighted with a red box. A red arrow points from this button to the data preview window below.

	consumo	pr.tbc	pr.cc	pr.ca	pr.galv1	pr.galv2	pr.pint	linea	averia	hora	r
1	135.311	6840	830	0	579	1401	0	1	Si	1	
2	84.082	443	903	58	611	1635	717	1	No	2	
3	131.615	7270	572	36	982	1969	243	1	No	3	
4	90.460	5031	694	122	896	1568	0	1	No	4	
5	120.043	9365	1054	157	403	1480	0	1	No	5	
6	153.678	9281	1003	172	605	1525	473	1	Si	6	
7	99.089	3223	1118	0	643	1424	732	1	No	7	
8	226.375	10490	1077	179	737	1333	93	1	No	8	
9	140.068	7394	1204	167	580	924	247	1	No	9	

Aparece una ventana con los datos disponibles. Moviendo el cursor hacia la izquierda o hacia abajo podemos recorrer toda la base de datos.

Preview window showing a grid of data. Annotations indicate: "Tiene 20 variables (columnas)" and "Dispone de 117 observaciones (filas)".

# Variables de la base de datos `acero`

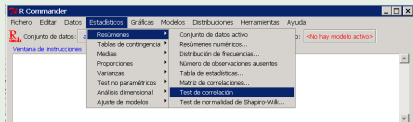
- 1 `consumo` Consumo energético de la empresa (Megavatios/hora).
- 2 `pr.tbc` Producción del tren de bandas calientes (Toneladas de acero).
- 3 `pr.cc` Producción de colada continua (Toneladas de acero).
- 4 `pr.ca` Producción del convertidor de acero (Toneladas de acero).
- 5 `pr.galv1` Producción de galvanizado de tipo I (Tns. de acero).
- 6 `pr.galv2` Producción de galvanizado de tipo II (Tns. de acero).
- 7 `pr.pint` Producción de chapa pintada (Tns. de acero).
- 8 `linea` Línea de producción empleada (A o B).
- 9 `turno` Turno de mañana (M), tarde (T), noche (N).
- 10 `temperatura` Temperatura del sistema: Alta, Media y Baja.
- 11 `pres.aver` Presencia de averías: hubo Averías (A), no hubo averías (NoA).
- 12 `nun.aver` Número de averías detectadas.
- 13 `sistema` Activación de un sistema de detección de sobrecalentamiento: encendido (ON), apagado (OFF).
- 14 ...

Calcule la matriz de correlaciones de la emisión de óxido nitroso ( $\text{N}_2\text{O}$ ), óxidos de monóxido de carbono ( $\text{CO}$ ), dióxido de carbono ( $\text{CO}_2$ ), mezcla de óxidos de nitrógeno ( $\text{NO}_x$ ) y dióxido de azufre ( $\text{SO}_2$ )

Estadísticos

↳ Resúmenes

↳ Matriz de correlaciones...

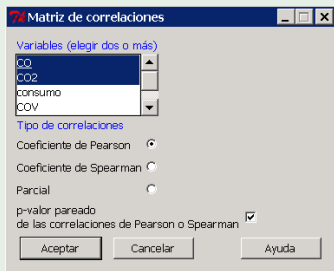


Seleccionamos  $\text{CO}$ ,  $\text{CO}_2$ ,  $\text{N}_2\text{O}$ ,  $\text{NO}_x$  y  $\text{SO}_2$ .

↳ Coeficiente de Pearson

↳ p-valor pareado de las correlaciones de Pearson o Spearman

↳ Aceptar



(Sigue →)

## Matriz de correlaciones:

	CO	CO2	N2O	NOx	SO2
CO	1.00	0.94	0.82	0.52	0.04
CO2	0.94	1.00	0.85	0.57	-0.03
N2O	0.82	0.85	1.00	0.53	0.01
NOx	0.52	0.57	0.53	1.00	-0.13
SO2	0.04	-0.03	0.01	-0.13	1.00

¿Qué variable es la que tiene más relación con N2O?

	CO	CO2	N2O	NOx	SO2
CO	1.00	0.94	0.82	0.52	0.04
CO2	0.94	1.00	0.85	0.57	-0.03
<b>N2O</b>	<b>0.82</b>	<b>0.85</b>	<b>1.00</b>	<b>0.53</b>	<b>0.01</b>
NOx	0.52	0.57	0.53	1.00	-0.13
SO2	0.04	-0.03	0.01	-0.13	1.00



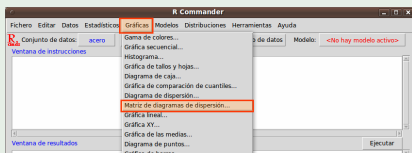
Dibuje la matriz de diagrama de dispersión de la emisión de óxido nitroso ( $\text{N}_2\text{O}$ ), óxidos de monóxido de carbono ( $\text{CO}$ ), dióxido de carbono ( $\text{CO}_2$ ), mezcla de óxidos de nitrógeno ( $\text{NO}_x$ ) y dióxido de azufre ( $\text{SO}_2$ )

Gráficas

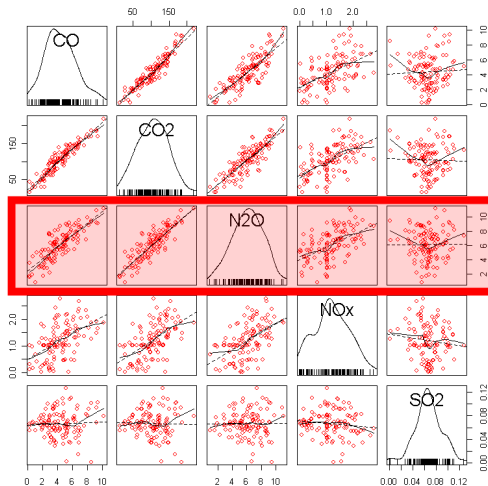
➡ Matriz de diagrama de dispersión...

Seleccionamos  $\text{CO}$ ,  $\text{CO}_2$ ,  $\text{N}_2\text{O}$ ,  $\text{NO}_x$  y  $\text{SO}_2$ .

➡ Aceptar



(Sigue →)



En la tercera fila, el N<sub>2</sub>O aparece en representado en el eje de ordenadas, y el resto de variables en el de abscisas.  
 ¿Qué diagrama de dispersión de la tercera fila muestra un patrón más claro de relación? Si bien usualmente no se puede responder de forma concluyente a esta pregunta a través de estos gráficos, se ve claramente en este caso cómo no parece haber relación con SO<sub>2</sub>, no hay una relación lineal clara con NO<sub>x</sub> y la relación lineal es fuerte con CO y CO<sub>2</sub>.

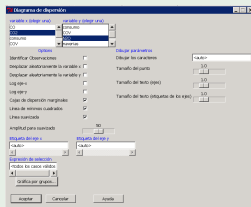
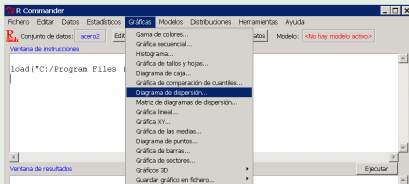
Dibuje el diagrama de dispersión con la emisión de óxido nitroso ( $N_2O$ ) en el eje de ordenadas y la emisión de dióxido de carbono ( $CO_2$ ) en el eje de abcisas.

Gráficas

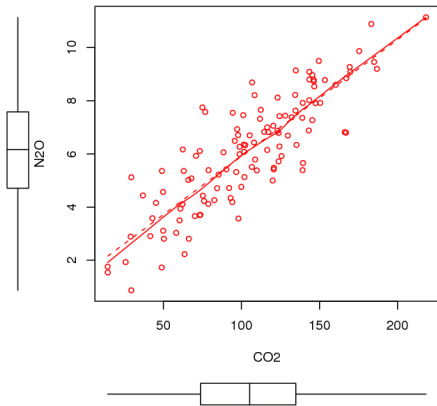
➡ Diagrama de dispersión...

Seleccionamos:  **$CO_2$**  y  **$N_2O$**   
En Opciones, marcar **Cajas de dispersión marginales**,  
**Línea de mínimos cuadrados**  
y  
**Líneas suavizadas**

➡ Aceptar



(Sigue →)



El eje de abscisas representa la emisión de  $\text{CO}_2$  (con un gráfico de cajas) y el de ordenadas la de  $\text{N}_2\text{O}$  (con un gráfico de cajas). Se observa una relación creciente entre ambas magnitudes.

En el gráfico aparecen dos líneas. Una es la recta de regresión (el modelo más simple) y la otra la línea de regresión no paramétrica (el mejor ajuste posible a los datos, respecto de mínimos cuadrados). Si ambas líneas son muy similares, el ajuste lineal resulta adecuado.

En este caso la línea recta sigue muy bien el comportamiento de la línea no paramétrica, por lo que el modelo lineal parece ajustar bien los datos.

# Estime la emisión de $N_2O$ a partir de la emisión de $CO_2$ . Llame al modelo como **RegModel.1**

Determinar los coeficientes  $\beta_0$ ,  $\beta_1$  y  $\sigma_\epsilon$  tales que

$$N_2O = \beta_0 + \beta_1 \cdot CO_2 + \epsilon, \quad \epsilon \sim \text{Normal}(0, \sigma_\epsilon)$$

donde  $\epsilon$  denota el error aleatorio.

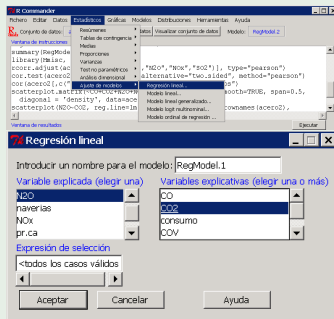
Estadísticos

- ➔ Ajuste de modelos
- ➔ Regresión lineal

Variable explicada: **N2O**

- ➔ Variables explicativas: **CO2**
- ➔ Aceptar

(El modelo por defecto se llama **RegModel.1**)



(Sigue →)

```

> RegModel.1 <- lm(N2O~CO2, data=acero)
> summary(RegModel.1)
Call:
lm(formula = N2O ~ CO2, data = acero)

Residuals:
    Min       1Q   Median       3Q      Max
-2.2585 -0.7287  0.0404  0.6511  2.9353

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.526865   0.280149   5.45 2.91e-07 ***
CO2          0.043850   0.002491  17.60 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.111 on 115 degrees of freedom
Multiple R-squared:  0.7293,    Adjusted R-squared:  0.7269
F-statistic: 309.8 on 1 and 115 DF,  p-value: < 2.2e-16

```

La columna `Estimate` proporciona los valores de las estimaciones de los coeficientes, con lo que el modelo de regresión lineal simple que mejor se ajusta a estos datos es:

$$N2O = 1.526865_{\sigma=0.280149} + 0.043850_{\sigma=0.002491} \cdot CO2 + \epsilon, \quad \epsilon \sim Normal(0, 1.111)$$

Así pues, por cada unidad que se incremente la emisión de `CO2`, la emisión promedio de `N2O` se espera que se incremente en 0.043850 unidades [Compruebe los intervalos de confianza asociados: el intervalo de confianza al 95% para este efecto es  $0.04385 \pm 1.96 \cdot 0.002491$ ].

Hemos obtenido que el coeficiente de determinación es:  $R^2 = 72.69\%$ , que estima el porcentaje de la variación de la variable dependiente que es explicado por la regresión. En este caso, el 72.69% de la variación de la emisión de `N2O` se debe a la emisión de `CO2`.

Recordemos que los valores que proporciona la recta de regresión para un valor dado de la variable explicativa pueden interpretarse como:

- predicciones del valor de la variable explicada; o
- estimaciones de su media.

Tanto para estas predicciones como para estas estimaciones, podemos proporcionar intervalos de confianza al nivel que se considere apropiado, normalmente al 95%.

Ejemplo:

[interval="prediction"] ¿Cuál es la previsión de gasto medio en alcohol de Juan, que tiene 20 años?

[interval="confidence"] ¿Cuál es la previsión de gasto medio en alcohol de quienes tienen 20 años?

Utilice el modelo `RegModel.1` para estudiar los valores de emisión de  $\text{N}_2\text{O}$  en las horas en las que se emiten 110t/h de  $\text{CO}_2$  (estimación puntual).

Escriba en la Ventana de instrucciones:

```
predict (RegModel.1, data.frame (CO2=c (110) ) )
```

Las salidas de este procedimiento nos indican que la estimación puntual, obtenida a partir de este modelo de regresión lineal simple, de la emisión de  $\text{N}_2\text{O}$  que se producirá en una hora en la que se hayan emitido 110t de  $\text{CO}_2$  es de 6'350341.



Utilice el modelo `RegModel.1` para estudiar los valores de emisión de  $N_2O$  en las horas en las que se emiten 110t/h (toneladas/hora) de  $CO_2$  (estimación por **intervalo**).

Escriba en la Ventana de instrucciones:

```
predict (RegModel.1, data.frame (CO2=c (110) ) , interval="prediction")
```

```
> predict (RegModel.1, data.frame (CO2=c (110) ) , interval="prediction")
```

```
      fit      lwr      upr  
1 6.350341 4.140992 8.55969
```

El valor 6.350341 es la estimación puntual de la emisión de  $N_2O$  que se obtiene a partir de este modelo. Es el valor que se obtiene en la recta de regresión para  $CO_2=110$ .

Se tiene una confianza del 95% de que la emisión de  $N_2O$  estará entre 4.140992t/h y 8.55969t/h para una hora en la que haya una emisión de 110t/h de  $CO_2$ .

Utilice el modelo `RegModel.1` para estudiar los valores de emisión de  $\text{N}_2\text{O}$  en las horas en las que se emiten 110t/h de  $\text{CO}_2$  (estimación por **intervalo** con una confianza del 99%).

Escriba en la Ventana de instrucciones:

```
predict(RegModel.1,data.frame(CO2=c(110)),interval="prediction",level=0.99)
```

	fit	lwr	upr
1	6.350341	3.428878	9.271804

Utilice el modelo `RegModel.1` para estudiar los valores de emisión *media* de  $\text{N}_2\text{O}$  en las horas en las que se emiten 110 y 100 t/h de  $\text{CO}_2$  (estimación por intervalo para el **promedio**).

Escriba en la Ventana de instrucciones:

```
predict(RegModel.1,data.frame(CO2=c(110,100)),interval="confidence")
```

	fit	lwr	upr
1	6.350341	6.145248	6.555434
2	5.911843	5.707193	6.116494

Con una confianza del 95% la emisión *media* de  $\text{N}_2\text{O}$  en aquellas horas en la que se emiten 110t/h de  $\text{CO}_2$  está entre 6.145248 y 6.555434.

Y, respectivamente, la emisión media cuando la  $\text{CO}_2$  es 100 está entre 5.707193 y 6.116494.