
Heart Disease

PORTAFOLIOS DE EQUIPO

Caso de estudio 1

Descripción del Contexto





Principal causa de muerte en todo el mundo

17,5 millones

Muertes anuales



Principales factores de riesgo

Colesterol

Hipertensión

Obesidad

Mala alimentación

Inactividad física

Consumo alcohol

Consumo de tabaco

**Familiares que las
padezcan**



Las 10 más comunes

1 Infarto

2 Cardiopatía isquémica

3 Insuficiencia cardíaca

4 Muerte súbita

5 Miocardiopatías

6 Miocardiopatía dilatada

7 Miocardiopatía hipertrófica

8 Miocardiopatía restrictiva

9 Valvulopatías

10 Valvulopatía mitral

Preprocesamiento de los datasets



Preprocesamiento



1

Analizar los datasets e **identificar**
cada instancia

cleveland.data	
1	1 0 63 1 -9 -9 -9 1 1
2	-9 1 145 1 233 -9 50 20 1
3	1 -9 1 2 2 3 81 0 1
4	0 0 0 0 1 10.5 6 13 1
5	150 60 190 90 145 85 0 0 1
6	2 3 3 -9 172 0 -9 -9 -9 1
7	-9 -9 -9 6 -9 -9 -9 2 1
8	16 81 0 1 1 1 -9 1 1
9	-9 1 -9 1 1 1 1 1 1
10	1 1 -9 -9 name 1
11	2 0 67 1 -9 -9 -9 1 1
12	-9 4 160 1 286 -9 40 40 1
13	0 -9 1 2 3 5 81 0 1
14	1 0 0 0 1 9.5 6 13 1
15	108 64 160 90 160 90 1 0 1
16	1.5 2 -9 185 3 -9 -9 -9 1
17	-9 -9 -9 3 -9 -9 -9 2 1
18	5 81 2 1 2 2 -9 2 1
19	-9 1 -9 1 1 1 1 1 1
20	1 1 -9 -9 name 1
21	3 0 67 1 -9 -9 -9 1 1
22	-9 4 120 1 229 -9 20 35 1
23	0 -9 1 2 2 19 81 0 1
24	1 0 0 0 1 8.5 6 10 1
25	129 78 140 80 120 80 1 0 1
26	2 6 2 -9 150 2 -9 -9 -9 1
27	-9 -9 -9 7 -9 -9 -9 2 1
28	20 81 1 1 1 1 -9 1 1
29	-9 1 -9 2 2 1 1 1 1
30	7 3 -9 -9 name 1



2

Agregar **encabezados** de atributos y ordenar **instancias** por filas

```
cleveland.data
1 ekgday ekggyr dig prop nitr pro diuretic proto thaldur thaltime met thalach thalrest tpeakbps tpeakbpd dummy
2 -9 -9 -9 -9 -9 -9 6 -9 -9 -9 216 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
3 -9 -9 -9 -9 -9 -9 3 -9 -9 -9 25 81 2 1 2 2 -9 2 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
4 -9 -9 -9 -9 -9 -9 7 -9 -9 -9 220 81 1 1 1 1 -9 1 -9 1 -9 2 2 1 1 17 3 -9 -9 name113
5 9 -9 -9 -9 -9 -9 3 -9 -9 -9 24 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
6 -9 -9 -9 -9 3 -9 -9 -9 218 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
7 0 -9 -9 -9 -9 -9 -9 3 -9 -9 -9 310 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
8 -9 -9 -9 -9 -9 3 -9 -9 -9 22 81 3 1 2 1 -9 1 -9 2 -9 2 1 1 1 17 1 -9 -9 name113
9 -9 -9 -9 -9 -9 3 -9 -9 -9 721 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
10 9 -9 -9 -9 -9 7 -9 -9 -9 73 81 2 1 2 1 -9 1 -9 1 -9 2 1 1 1 67 2 -9 -9 name113
11 -9 -9 -9 -9 -9 -9 7 -9 -9 -9 76 81 1 1 1 1 -9 1 -9 1 -9 2 1 1 1 11 1 -9 -9 name113
12 0 -9 -9 -9 -9 -9 -9 6 -9 -9 -9 71 81 0 1 1 1 -9 1 -9 1 -9 1 1 1 1 11 1 -9 -9 name113
```




3

Importar los datasets a  **rapidminer**

A



cleveland

282 instancias
75 atributos

C



long-beach-va

200 instancias
75 atributos

B



hungarian

294 instancias
75 atributos

D



switzerland

123 instancias
75 atributos

899 instancias, **75** atributos



4

Eliminar atributos por cantidad de **missing values**

✗ **pncaden** (100%)

✗ **dm** (89,4%)

✗ **famhist** (46,9%)

✗ **restwm** (96,7%)

✗ **exerwm** (99,4%)

✗ **ca** (66,7%)

✗ **exeref** (99,8%)

✗ **restef** (96,9%)

✗ **diag** (62,1%)

✗ **om2** (63,6%)

✗ **ramus** (63,1%)

✗ **thal** (53,1%)

12 atributos **eliminados**



5

Eliminar atributos **irrelevantes** (1/2)



id



ekgmo



dig



pro



name



ekgday



prop



diuretic



ccf



ekgyr



nitr



proto



thalach



thalrest



met



thaldur



tpeakbpd



thalttime



tpeakbps



oldpeak



5

Eliminar atributos **irrelevantes** (2/2)

✗ slope	✗ rldv5e	✗ exerckm
✗ rldv5	✗ restckm	✗ cmo
✗ painloc	✗ cday	✗ cyr

29 atributos **eliminados**



6

Eliminar atributos por **falta de información**

✗ thalsev	✗ earlobe	✗ lvx2	✗ lvx4
✗ thalpul	✗ lvx1	✗ lvx3	✗ lvf
✗ cathef	✗ junk	✗ dummy	✗ xhypo

12 atributos **eliminados**



7

Generación del atributo **smoke_gen**

SI $\text{cigs} > 0 \parallel \text{years} > 0 \parallel \text{smoke} == 1$

1

SINO SI $!\text{missing}(\text{cigs}) \parallel !\text{missing}(\text{years}) \parallel !\text{missing}(\text{smoke})$

0

SINO

Vecino más cercano

× cigs **×** years **×** smoke

+ smoke_gen

Atributos Utilizados

Total : 20



✓	age	✓	trestbpd
✓	sex	✓	exang
✓	painexer	✓	lmt
✓	relres	✓	ladprox
✓	cp	✓	laddist
✓	trestbps	✓	cxmain
✓	htn	✓	om1
✓	chol	✓	rcadist
✓	fbs	✓	smoke_gen
✓	restecg	✓	rcaprox



8

Imputación de missing values

K - NN

con $K = 1$



9

Normalización de los datos

Transformación Z

Resultado de la regresión





Confianza

87,98%

+/- 2,55%



Preguntas?

—

Muchas gracias!
