

Practica 2: Data Cleaning

Jose Luis Rivas Calduch y Mariano Jiménez Barca

11 de diciembre de 2020

Contents

1. Descripción del dataset.	1
2. Integración y selección de los datos de interés a analizar.	3
3. Limpieza de los datos.	3
4. Análisis de los datos	3
4.1. Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar).	3
4.2. Comprobación de la normalidad y homogeneidad de la varianza.	3
4.3. Aplicación de pruebas estadísticas para comparar los grupos de datos.	3
5. Representación de los resultados a partir de tablas y gráficas.	3
6. Resolución del problema.	3
Bibliografía	3
Agradecimientos	3

1. Descripción del dataset.

El data set objeto de estudio esta tomado de la plataforma *Kaggle*. Esta es una comunidad en línea de científicos de datos y profesionales del aprendizaje automático, actualmente es una subsidiaria de *Google LLC*.

El nombre del data set es *Heart Failure Prediction* (ECV) que son la principal causa de muerte a nivel mundial, cobrando un estimado de 17,9 millones de vidas cada año, lo que representa el 31% de todas las muertes en todo el mundo. La insuficiencia cardíaca es un evento común causado por las enfermedades cardiovasculares y este conjunto de datos contiene 12 características que se pueden usar para predecir la mortalidad por insuficiencia cardíaca.

La mayoría de las enfermedades cardiovasculares se pueden prevenir abordando los factores de riesgo conductuales como el consumo de tabaco, la dieta poco saludable y la obesidad, la inactividad física y el consumo nocivo de alcohol utilizando estrategias para toda la población.

Las personas con enfermedad cardiovascular o que se encuentran en alto riesgo cardiovascular (debido a la presencia de uno o más factores de riesgo como hipertensión, diabetes, hiperlipidemia o enfermedad ya establecida) necesitan una detección y manejo precoces donde un modelo de aprendizaje automático puede ser de gran ayuda.

Tipos de variables

```
sapply(rawData, class)
```

```
##           age           anaemia creatinine_phosphokinase
##      "numeric"         "integer"         "integer"
##      diabetes ejection_fraction high_blood_pressure
##      "integer"         "integer"         "integer"
##      platelets serum_creatinine serum_sodium
##      "numeric"         "numeric"         "integer"
##      sex smoking time
##      "integer"         "integer"         "integer"
##      DEATH_EVENT
##      "integer"
```

Descripción de las variables

```
str(rawData)
```

```
## 'data.frame': 299 obs. of 13 variables:
## $ age : num 75 55 65 50 65 90 75 60 65 80 ...
## $ anaemia : int 0 0 0 1 1 1 1 1 0 1 ...
## $ creatinine_phosphokinase: int 582 7861 146 111 160 47 246 315 157 123 ...
## $ diabetes : int 0 0 0 0 1 0 0 1 0 0 ...
## $ ejection_fraction : int 20 38 20 20 20 40 15 60 65 35 ...
## $ high_blood_pressure : int 1 0 0 0 0 1 0 0 0 1 ...
## $ platelets : num 265000 263358 162000 210000 327000 ...
## $ serum_creatinine : num 1.9 1.1 1.3 1.9 2.7 2.1 1.2 1.1 1.5 9.4 ...
## $ serum_sodium : int 130 136 129 137 116 132 137 131 138 133 ...
## $ sex : int 1 1 1 1 0 1 1 1 0 1 ...
## $ smoking : int 0 0 1 0 0 1 0 1 0 1 ...
## $ time : int 4 6 7 7 8 8 10 10 10 10 ...
## $ DEATH_EVENT : int 1 1 1 1 1 1 1 1 1 1 ...
```

Resumen descriptivo de las variables

```
summary(rawData)
```

```
##      age      anaemia creatinine_phosphokinase diabetes
## Min.   :40.00 Min.   :0.0000 Min.   : 23.0 Min.   :0.0000
## 1st Qu.:51.00 1st Qu.:0.0000 1st Qu.: 116.5 1st Qu.:0.0000
## Median :60.00 Median :0.0000 Median : 250.0 Median :0.0000
## Mean   :60.83 Mean   :0.4314 Mean   : 581.8 Mean   :0.4181
## 3rd Qu.:70.00 3rd Qu.:1.0000 3rd Qu.: 582.0 3rd Qu.:1.0000
## Max.   :95.00 Max.   :1.0000 Max.   :7861.0 Max.   :1.0000
## ejection_fraction high_blood_pressure platelets serum_creatinine
## Min.   :14.00 Min.   :0.0000 Min.   : 25100 Min.   :0.500
## 1st Qu.:30.00 1st Qu.:0.0000 1st Qu.:212500 1st Qu.:0.900
## Median :38.00 Median :0.0000 Median :262000 Median :1.100
## Mean   :38.08 Mean   :0.3512 Mean   :263358 Mean   :1.394
## 3rd Qu.:45.00 3rd Qu.:1.0000 3rd Qu.:303500 3rd Qu.:1.400
## Max.   :80.00 Max.   :1.0000 Max.   :850000 Max.   :9.400
## serum_sodium sex smoking time
```

```
## Min.      :113.0   Min.      :0.0000   Min.      :0.0000   Min.      : 4.0
## 1st Qu.   :134.0   1st Qu.   :0.0000   1st Qu.   :0.0000   1st Qu.   : 73.0
## Median    :137.0   Median    :1.0000   Median    :0.0000   Median    :115.0
## Mean      :136.6   Mean      :0.6488   Mean      :0.3211   Mean      :130.3
## 3rd Qu.   :140.0   3rd Qu.   :1.0000   3rd Qu.   :1.0000   3rd Qu.   :203.0
## Max.      :148.0   Max.      :1.0000   Max.      :1.0000   Max.      :285.0
## DEATH_EVENT
## Min.      :0.0000
## 1st Qu.   :0.0000
## Median    :0.0000
## Mean      :0.3211
## 3rd Qu.   :1.0000
## Max.      :1.0000
```

2. Integración y selección de los datos de interés a analizar.

3. Limpieza de los datos.

4. Análisis de los datos

4.1. Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar).

4.2. Comprobación de la normalidad y homogeneidad de la varianza.

4.3. Aplicación de pruebas estadísticas para comparar los grupos de datos.

5. Representación de los resultados a partir de tablas y gráficas.

6. Resolución del problema.

Bibliografía

Subirats Maté, Laila; Pérez Trenard, Diego O.; Calvo González, Mireia (2019) Introducción al ciclo de la vida de los datos. UOC *Subirats Maté, Laila; Calvo González, Mireia (2019)* Web scraping. UOC *Subirats Maté, Laila; Pérez Trenard, Diego O.; Calvo González, Mireia (2019)* Introducción a limpieza y análisis de los datos. UOC *Hernández Orallo, José; Ramírez Quintana, M José; Ferri Ramírez, Cesar (2004)* Introducción a la Minería de Datos. PEARSON. *Gironés Roig, Jordi; Casas Roma, Jordi; Minguillon Alfonso, Julia; Caichuelas Quiles, Ramon (2017)* Minería de datos: Modelos y algoritmos. UOC.

Agradecimientos

Cita Davide Chicco, Giuseppe Jurman: Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone. BMC Medical Informatics and Decision Making 20, 16 (2020). (link)

License CC BY 4.0