

PRAC2 - Neteja i validació de les dades

Rubén Vasallo González

25 de December 2018

- 1 Descripción del dataset
 - 1.1 Quina pregunta/problema pretén respondre?
- 2 Integració i selecció de les dades d'interès a analitzar.
- 3 Neteja de les dades.
 - 3.1 Les dades contenen zeros o elements buits? Com gestionaries aquests casos?
 - 3.2 Identificació i tractament de valors extrems.

1 Descripción del dataset

El conjunt de dades que anem a analitzar s'ha extret de la competició de la web de Kaggle anomenada "Titanic: Machine Learning from Disaster" (<https://www.kaggle.com/c/titanic>) que permet iniciar-se en les competicions que aquesta pagina web te. Aquest conjunt de dades (o dataset) te 891 observacions que representen algunes de les característiques dels passatgers, junt amb un atribut que indica si aquest van sobreviure al famós enfonsament o no. A continuació detallem el significat de cada atribut i, si es te, els possibles valors que poden tenir (en cas de ser atributs qualitatius)

Variable	Definició	Key
PassengerId	Identificador del passatger	
Survived	Sobreviu	0 = No, 1 = Yes
pclass	Classe de Tiquet	1 = 1st, 2 = 2nd, 3 = 3rd
Name	Nom del passatger	
Sex	Sexe	male, female
Age	Edat en anys	
SibSp	# de germans / esposes a bord del Titanic	
Parch	# de pares / fills a bord del Titanic	
Ticket	Numero de Tiquet	
Fare	Tarifa del passatger	
Cabin	Numero de Cabina	
Embarked	Port d'embarcament	C = Cherbourg, Q = Queenstown, S = Southampton

1.1 Quina pregunta/problema pretén respondre?

L'enfonsament del Titanic va ser un dels naufragis mes famosos de l'història. A dia d'avui encara presenta grans misteris i preguntes sense resolució. Una de les raons per les quals el naufragi va comportar aquesta pèrdua de vida va ser que no hi havia prou barques salvavides per als passatgers i la tripulació. Ates a la norma moral de que en cas d'enfonsament, les dones i els nens haurien estat els primers en ser evacuats (https://es.wikipedia.org/wiki/Mujeres_y_ni%C3%B1os_primeros), a l'hora de la veritat, això no va ser així, i tot i que hi va haver algun element de sort en la supervivència de l'enfonsament, la història ens ha demostrat que alguns grups de persones tenien més probabilitats de sobreviure que altres, com ara la classe alta.

L'objectiu d'aquest dataset es el estudi de les característiques dels passatgers que van sobreviure per tractar de crear un model d'aprenentatge automàtic que sigui capaç de predir si nous passatgers amb característiques semblants sobreviurien o no.

2 Integració i selecció de les dades d'interès a analitzar.

Per assolir l'objectiu primer carregarem el dataset i li farem una ullada als 6 primers resultats.

```
# Carregar el fitxer de dades en R
titanicdataset <- read.csv2(file = "train.csv", header = TRUE, sep = ",", quote = "\"", fill
= TRUE, encoding="UTF-8")
# breu resumen dels atributs
head(titanicdataset)
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Braund, Mr. Owen Harris	male	22	1	0	A/5 21171	7.25		S
2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	Heikkinen, Miss. Laina	female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1	0	113803	53.1	C123	S
5	0	3	Allen, Mr. William Henry	male	35	0	0	373450	8.05		S
6	0	3	Moran, Mr. James	male		0	0	330877	8.4583		Q

Veiem que tenim 12 atributs i que no tots tenen valors assignats. Cridem a la funció summary per veure un resum dels valors del dataset sencer.

```
summary(titanicdataset)
```

```
## PassengerId      Survived      Pclass
## Min.   : 1.0      Min.   :0.0000      Min.   :1.000
## 1st Qu.:223.5      1st Qu.:0.0000      1st Qu.:2.000
## Median :446.0      Median :0.0000      Median :3.000
## Mean   :446.0      Mean   :0.3838      Mean   :2.309
## 3rd Qu.:668.5      3rd Qu.:1.0000      3rd Qu.:3.000
## Max.   :891.0      Max.   :1.0000      Max.   :3.000
##
##                               Name      Sex      Age
## Abbing, Mr. Anthony          : 1   female:314      :177
## Abbott, Mr. Rossmore Edward  : 1   male  :577      24      : 30
## Abbott, Mrs. Stanton (Rosa Hunt) : 1                               22      : 27
## Abelson, Mr. Samuel          : 1                               18      : 26
## Abelson, Mrs. Samuel (Hannah Wizosky): 1                               19      : 25
## Adahl, Mr. Mauritz Nils Martin : 1                               28      : 25
## (Other)                      :885      (Other):581
## SibSp      Parch      Ticket      Fare
## Min.   :0.000      Min.   :0.0000      1601      : 7      8.05      : 43
## 1st Qu.:0.000      1st Qu.:0.0000      347082     : 7      13        : 42
## Median :0.000      Median :0.0000      CA. 2343: 7      7.8958     : 38
## Mean   :0.523      Mean   :0.3816      3101295 : 6      7.75      : 34
## 3rd Qu.:1.000      3rd Qu.:0.0000      347088 : 6      26        : 31
## Max.   :8.000      Max.   :6.0000      CA 2144 : 6      10.5      : 24
##                               (Other) :852      (Other):679
## Cabin      Embarked
##           :687      : 2
## B96 B98     : 4      C:168
## C23 C25 C27: 4      Q: 77
## G6          : 4      S:644
## C22 C26     : 3
## D           : 3
## (Other)     :186
```

De tots els atributs, veiem que els de identificador del passatger, nom del passatger i numero de tiquet son irrelevantes a l'hora de predir si un passatger sobreviu o no, per el que els descartem del dataset.

```
titanicdataset$PassengerId <- NULL
titanicdataset$Name <- NULL
titanicdataset$Ticket <- NULL
```

També podem veure que el programa R ha detectat gaire be tots els atributs com a numèrics (Excepte l'atribut sexe i edat). Realment en aquest dataset tots els valors es poden considerar com a categòrics ja que la informació a la que fa referencia està acotat (no es infinit)

```
sapply(titanicdataset, function(x) class(x))
```

```
## Survived      Pclass      Sex      Age      SibSp      Parch      Fare
## "integer" "integer" "factor" "factor" "integer" "integer" "factor"
## Cabin Embarked
## "factor" "factor"
```

Transformem els atributs a categòrics.

```
titanicdataset[, 'Pclass']<- factor(titanicdataset[, 'Pclass'])
titanicdataset[, 'Survived']<- factor(titanicdataset[, 'Survived'])
titanicdataset[, 'SibSp']<- factor(titanicdataset[, 'SibSp'])
titanicdataset[, 'Parch']<- factor(titanicdataset[, 'Parch'])
```

3 Neteja de les dades.

3.1 Les dades contenen zeros o elements buits? Com gestionaries aquests casos?

Primer comprovem si tenim cap observació amb atributs a null.

```
sapply(titanicdataset, function(x) sum(is.na(x)))
```

```
## Survived    Pclass    Sex    Age    SibSp    Parch    Fare    Cabin
##          0          0      0      0        0        0        0        0
## Embarked
##          0
```

No tenim cap, a continuació comprovem si tenim cap observació amb atributs a 0.

```
sapply(titanicdataset, function(x) sum(x == 0))
```

```
## Survived    Pclass    Sex    Age    SibSp    Parch    Fare    Cabin
##        549          0      0      0        608        678        15        0
## Embarked
##          0
```

En aquest cas tenim valors a 0 en els atributs Survived, SibSp Parch i Fare

Analitzem cadascun d'ells.

- Survived te sentit que tingui valors a 0 ja que son tots aquells passatgers que no van sobreviure.
- SibSp te sentit que tingui valors a 0 ja que son tots aquells passatgers que no tenien cap germà o esposa a bord.
- Parch te sentit que tingui valors a 0 ja que son tots aquells passatgers que no tenien cap pare, mare o fill a bord.
- Fare te sentit que tingui valors a 0 ja que podrien ser tots aquells passatgers que van ser invitats, no van pagar cap tarifa per viatjar a bord o eren tripulants.

No modifiquem pas cap valor.

A continuació comprovem si tenim cap atribut amb valor string buit.

```
sapply(titanicdataset, function(x) sum(x == ""))
```

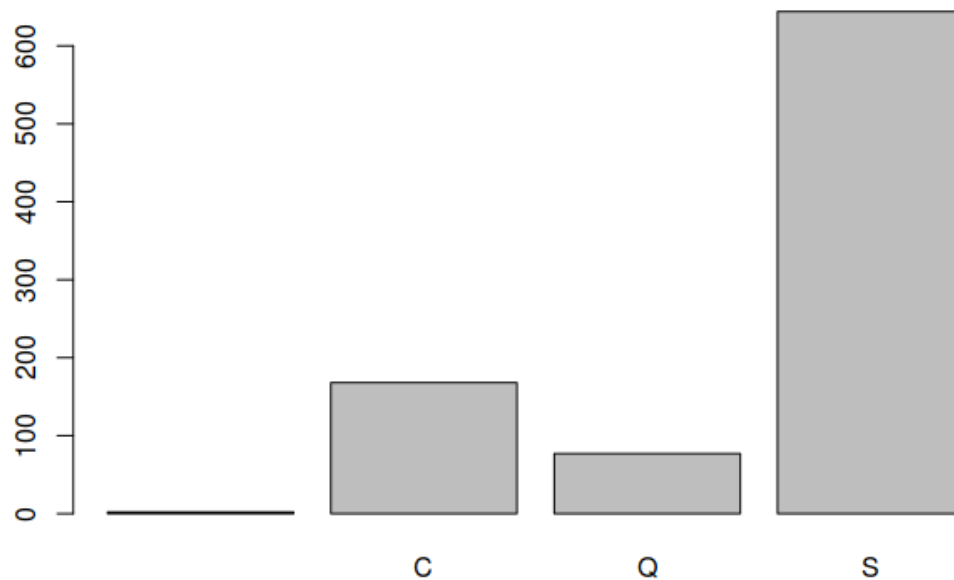
```
## Survived    Pclass    Sex    Age    SibSp    Parch    Fare    Cabin
##          0          0      0    177        0        0        0    687
## Embarked
##          2
```

En aquest cas tenim valors a string buits en els atributs Age, Cabin i Embarked

Analitzem cadascun d'ells.

Observem que tenim en el atribut Embarked valors per informar (Tenim 2 passatgers que no tenen port d'embarcament).

```
plot(titanicdataset$Embarked)
```



Ja que no tenim cap manera de recuperar la informació que ens falta intentarem aproximar el origen dels embarcaments de les observacions que no tenen aquesta informació utilitzant el mètode dels veïns mes propers (utilitzant la funció KNN del paquet VIM), que imputa el valor dels resultats (k veïns) mes propers. (Sempre serà millor aproximar el resultat que descartar-lo ja que tindrem un model amb menor marge d'error)

```
titanicdataset$Embarked <- sapply(titanicdataset$Embarked, function(x) if(x==""){NA}else{x})
)

suppressWarnings(suppressMessages(library(VIM)))

titanicdataset$Embarked <- kNN(titanicdataset)$Embarked
titanicdataset[, 'Embarked'] <- factor(titanicdataset[, 'Embarked'])
```

Observem que tenim el mateix problema amb el atribut Age. Igual que ens passa en el cas aterior, no tenim manera de recuperar aquesta informació per el que intentarem Aproximar-lo mitjançant el mètode dels veïns mes propers.

```
titanicdataset$Age <- sapply(titanicdataset$Age, function(x) if(x==""){NA}else{x} )

suppressWarnings(suppressMessages(library(VIM)))

titanicdataset$Age <- kNN(titanicdataset)$Age
titanicdataset[, 'Age']<- factor(titanicdataset[, 'Age'])
```

Finalment observem que el atribut Cabin li falta 687 valors. Aquest valor es mes complicat d'aproximar ja que no podem calcular quin era la cabina on viatjava els passatgers que no tenen una cabina assignada o si viatjaven sense cabina. En un principi vaig pensar que els passatgers que no tenien cabina assignada es que viatjaven al passadís però analitzant mes les dades, es pot observar que hi ha 40 passatgers de primera classe que no tenen una cabina assignada, i costa de creure que un passatger de primera viatges sense cabina.

```
table(titanicdataset[titanicdataset$Cabin == "",]$Pclass)
```

```
##
## 1 2 3
## 40 168 479
```

Descartem l'atribut.

```
titanicdataset$Cabin <- NULL
```

Comprovem una vegada més que ja no tenim cap observació amb strings vuits.

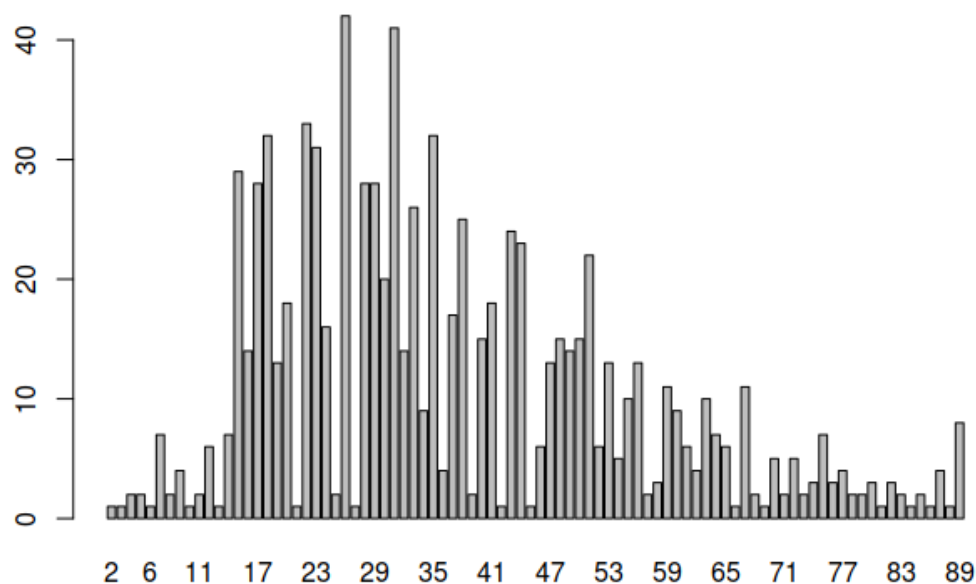
```
sapply(titanicdataset, function(x) sum(x == ""))
```

```
## Survived  Pclass  Sex  Age  SibSp  Parch  Fare Embarked
##          0        0    0    0      0      0      0        0
```

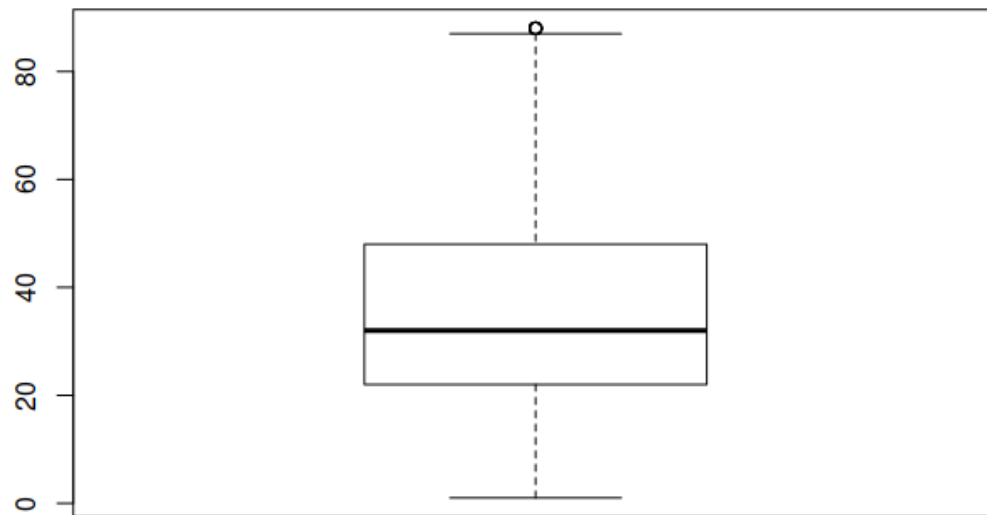
3.2 Identificació i tractament de valors extrems.

Com hem comentat anteriorment, la majoria dels atributs que te aquest dataset son categòrics (tenen valors finits) per el que es difícil detectar valors outliers. En aquest cas podríem intentar detectar outliers visualitzant en el atribut Edat, mostrant la distribució de les dades en gràfics.

```
plot(titanicdataset$Age)
```



```
boxplot(as.numeric(titanicdataset$Age))
```



Podem veure que tenim outliers tant per un extrem com per altre. Això es normal ja que en el Titànic viatjaven tant nens com gent d'avançada edat.

No toquem cap valor.

Amb aquest passos donem per finalitzada la fase de neteja. A continuació fem l'exportació del dataset amb les dades ja processades per utilitzar en les següents fases del projecte.

```
write.csv(titanicdataset, "titanic_data_clean.csv", row.names = FALSE)
```