

Práctica 2: Limpieza y análisis de datos

Maite Gracia

4 de January, 2021

Contents

1	Descripción del dataset	2
2	Integración y selección de los datos de interés a analizar	2
3	Limpieza de los datos	3
3.1	Normalización de los datos	4
3.2	Valores atípicos	7
3.3	Imputación de valores	9
3.4	Selección de datos	9
3.5	Exportación de los datos limpios	11
4	Análisis de los datos	11
4.1	Selección de los grupos de datos a analizar	11
4.2	Normalidad y homocedasticidad	11
4.3	Pruebas estadísticas	13
4.3.1	Contraste de hipótesis de dos muestras	13
4.3.2	Modelo regresión logística simple	14
4.3.2.1	Predicciones	14
4.3.3	Modelo de regresión logística múltiple	15
4.3.3.1	Comparación de las predicciones con las observaciones	17
4.3.3.2	Predicción	17
5	Representación de los resultados	18
6	Conclusiones	20
7	Agradecimientos	20
8	Tabla de contribuciones	20

1 Descripción del dataset

Se ha decidido utilizar un dataset de la web Kaggle para la presente práctica. [enlace](https://www.kaggle.com/kemical/kickstarter-projects) Kickstarter es una plataforma de micro mecenazgo, es decir, gente de todo el mundo ayuda a financiar las ideas y proyectos de pequeñas empresas o particulares.

En la web de Kickstarter [enlace](https://www.kickstarter.com/) se pueden encontrar miles de campañas que buscan financiación para desarrollar productos de todo tipo. Desde películas independientes, a juegos de mesa o ropa, peluches, libros etc. Cada una de estas campañas tendrá un periodo de tiempo en el que cualquiera podrá aportar dinero al proyecto y si se consigue llegar al objetivo de dinero propuesto la campaña será fundada.

Yo personalmente utilicé Kickstarter hace unos años para lanzar una serie de productos lo cual es una de las razones por las que he elegido el presente dataset. El objetivo principal sería poder crear un modelo que predijera que probabilidad tiene cualquier tipo de producto de conseguir recaudar dinero mediante una campaña de Kickstarter antes de ser lanzado.

2 Integración y selección de los datos de interés a analizar

Las variables que componen el dataset son:

- ID: identificador interno de Kickstarter
- name: nombre del proyecto
- category: categoría específica en la que se encuentra el proyecto
- main_category: categoría principal de la campaña
- currency: divisa en la que se creó el proyecto
- deadline: fecha límite
- goal: cantidad de dinero que el creador necesita para completar el proyecto
- launched: fecha lanzamiento
- pledged: cantidad total aportada al proyecto
- state: condición en la que se encuentra el Proyecto (failed, successful, canceled, live, undefined)
- backers: total de mecenas.
- country: país en el que se encuentra el Proyecto.
- usd_pledged: conversión en dólares de la columna pledged hecha por Kickstarter
- usd_pledged_real: conversión en dólares de la columna pledged hecha a través de Fixer.io API
- usd_goal_real: conversión en dólares de la columna goal hecha a través de Fixer.io API

Antes de cargar el archivo en R se hace una inspección de los datos. Al tratarse de un archivo con extensión .csv, hay que cerciorarse del tipo de separador utilizado (en este caso la “,”) y posteriormente se procede a su carga teniendo en cuenta el separador antes mencionado:

```
# Asignamos los datos del fichero cargado a una variable denominada dataSet
dataSet <- read.csv('../data/ks-projects-201801.csv')
nrow(dataSet)
```

```
## [1] 378661
```

```
names(dataSet)
```

```
## [1] "ID"           "name"          "category"      "main_category"
## [5] "currency"     "deadline"      "goal"          "launched"
## [9] "pledged"      "state"         "backers"       "country"
## [13] "usd.pledged"  "usd_pledged_real" "usd_goal_real"
```

Vemos que el dataset original se compone de 378,661 muestras y 15 variables. Ya que se trata de una cantidad de muestras muy elevadas, se ha decidido aplicar una técnica para reducir la cantidad de estas, se empleará la técnica de muestreo aleatorio simple sin sustitución, es decir, se van a extraer 3000 muestras aleatorias

del conjunto de datos, donde la probabilidad de escoger cada una de las muestras será la misma para todas, $1/378,661$.

```
library(sampling)
indices <- sample( 1:nrow( dataSet ), 3000 )
dataSet <- dataSet[ indices, ]
```

A partir de ahora cuando se haga referencia al dataset, estaremos hablando del dataset que contiene las 3000 muestras, no el dataset original.

3 Limpieza de los datos

```
# Muestra de las 5 primeras líneas del dataset completo
head(dataSet, 5)
```

```
##           ID                                     name
## 277557  48211191 Mint Julep Jazz Band's First CD - Durham on Saturday Night
## 140412 1712941959                                     QuikLid
## 10877   1054544730                                     Enamel Pins Set: Twisted Tales
## 254046  361859883      Let Tinned Pineapple Make Your Soapy Dreams Come True!
## 273982  463992846      EDC "NOMAD" Survival Multi-Tool
##           category main_category currency  deadline  goal
## 277557           Jazz           Music    USD 2012-12-31  6000
## 140412 Product Design           Design    USD 2017-07-03 225000
## 10877   Accessories           Fashion    USD 2017-04-12   300
## 254046           Crafts           Crafts    USD 2012-03-31  1000
## 273982 Product Design           Design    USD 2016-03-25   650
##           launched pledged      state backers country  usd.pledged
## 277557 2012-11-01 21:33:49    7460 successful    216     US        7460
## 140412 2017-05-04 19:42:57    1698 canceled     78     US        546
## 10877   2017-03-22 02:04:45     203 canceled      8     US        203
## 254046 2012-02-09 05:53:33     327 failed       18     US        327
## 273982 2016-02-24 06:10:02   11152 successful    234     US       11152
##           usd_pledged_real usd_goal_real
## 277557           7460           6000
## 140412           1698          225000
## 10877           203           300
## 254046           327          1000
## 273982          11152           650
```

```
# Análisis descriptivo del dataset
summary(dataSet)
```

```
##           ID           name           category           main_category
## Min.      :5.141e+05 Length:3000 Length:3000 Length:3000
## 1st Qu.:5.259e+08  Class :character Class :character Class :character
## Median :1.055e+09  Mode  :character Mode  :character Mode  :character
## Mean      :1.060e+09
## 3rd Qu.:1.600e+09
## Max.      :2.147e+09
##
##           currency           deadline           goal           launched
## Length:3000 Length:3000 Min.      :           1 Length:3000
## Class :character Class :character 1st Qu.:    2000 Class :character
## Mode  :character Mode  :character Median :    5000 Mode  :character
```

```
##                               Mean   :   59075
##                               3rd Qu.:   15000
##                               Max.   :100000000
##
## pledged                      state                backers                country
## Min.   :      0.0    Length:3000    Min.   :      0.00    Length:3000
## 1st Qu.:     40.0    Class :character    1st Qu.:      2.00    Class :character
## Median :    625.5    Mode  :character    Median :     12.00    Mode  :character
## Mean   :   9154.1                    Mean   :    102.64
## 3rd Qu.:   4002.8                    3rd Qu.:     54.25
## Max.   :2669009.6                    Max.   :20680.00
##
##  usd.pledged    usd_pledged_real    usd_goal_real
## Min.   :      0    Min.   :      0.0    Min.   :      1
## 1st Qu.:     20    1st Qu.:     40.0    1st Qu.:    2000
## Median :    390    Median :    605.9    Median :    5457
## Mean   :   6181    Mean   :   8879.9    Mean   :   58343
## 3rd Qu.:   3033    3rd Qu.:   4000.0    3rd Qu.:   15000
## Max.   :1068328    Max.   :2669009.6    Max.   :100000000
## NA's   :24
```

```
# Comprobamos si hay NA en el dataset original
sapply(dataSet, function(x) sum(is.na(x)))
```

```
##           ID           name           category    main_category
##           0             0             0             0
## currency    deadline           goal           launched
##           0             0             0             0
## pledged      state           backers           country
##           0             0             0             0
##  usd.pledged usd_pledged_real    usd_goal_real
##           24             0             0
```

3.1 Normalización de los datos

Basándonos en la estadística descriptiva de la muestra y en la descripción de cada variable podemos ver que todas las variables menos ID son de tipo carácter. Para poder analizar de forma eficaz los datos haremos las siguientes conversiones:

- Variables category, main_category, currency y country van a convertirse a tipo factor para poder agrupar proyectos.

```
dataSet$category <- as.factor(dataSet$category)
dataSet$main_category <- as.factor(dataSet$main_category)
dataSet$currency <- as.factor(dataSet$currency)
dataSet$country <- as.factor(dataSet$country)
dataSet$state <- as.factor(dataSet$state)
```

```
# Valores que toman las variables currency y country
unique(dataSet$currency)
```

```
## [1] USD EUR AUD GBP CAD NZD NOK SGD SEK MXN CHF HKD DKK
## Levels: AUD CAD CHF DKK EUR GBP HKD MXN NOK NZD SEK SGD USD
```

```
unique(dataSet$country)
```

```
## [1] US FR NL BE N,0" AU GB CA IT ES NZ NO SG DE SE
```

```
## [16] MX CH HK IE DK AT
## Levels: AT AU BE CA CH DE DK ES FR GB HK IE IT MX N,0" NL NO NZ SE SG US
```

Vemos que country tiene un carácter especial en algunos de los casos, vamos a sustituirlos por NA y más adelante imputaremos estos valores basándonos en la variable currency.

```
dataSet$country[dataSet$country == 'N,0"'] <- NA
```

- Las variables deadline y launched se convertirán a tipo Date.

```
dataSet$deadline <- as.Date(dataSet$deadline, '%Y-%m-%d')
dataSet$launched <- as.Date(dataSet$launched, '%Y-%m-%d')
```

- goal, pledged y usd.pledged van a pasar a ser tipo numérico.

```
dataSet$goal <- as.numeric(dataSet$goal)
dataSet$pledged <- as.numeric(dataSet$pledged)
dataSet$usd_pledged <- as.numeric(dataSet$usd_pledged)
dataSet$usd_pledged_real <- as.numeric(dataSet$usd_pledged_real)
dataSet$usd_goal_real <- as.numeric(dataSet$usd_goal_real)
```

- La variable state, como ya se ha explicado, detalla el estado en el que acabó o estaba en ese momento la campaña. Vemos que hay 5 estados failed, successful, canceled, suspended y undefined. Se ha decidido añadir una nueva columna status, que contendrá dos valores, 0 si el proyecto no ha sido fundado y 1 si el proyecto ha recaudado los fondos suficientes. Esta nueva variable se creará comprobando si el dinero recaudado para el proyecto es superior al goal propuesto y si el estado de dicha muestra es distinto de live o undefined, es decir, es un proyecto que ya ha terminado. Si no se hace esta comprobación el modelo final no será tan preciso ya que se puede dar el caso en el que un proyecto que esté live no haya conseguido llegar al objetivo todavía, pero puede que termine siendo exitoso.

```
dataSet['status'] <- as.factor(
  ifelse(
    (dataSet$pledged > dataSet$goal) &
    (dataSet$state != 'live' | dataSet$state != 'undefined'),
    1, 0))
```

- Se va a añadir una columna nueva euros_pledged que contendrá la conversión de usd_pledged_real a euros. Se utilizará la conversión 1€ = 1.23\$ a 4 de enero de 2021.

```
dataSet['euros_pledged'] <- as.numeric(
  format(as.numeric(dataSet$usd_pledged_real)/1.23), nsmall = 1)
```

- Se va a añadir una columna nueva euros_goal que contendrá la conversión de goal a euros. Se utilizará la conversión 1€ = 1.23\$ a 4 de enero de 2021.

```
dataSet['euros_goal'] <- round(as.numeric(
  format(as.numeric(dataSet$goal)/1.23), nsmall = 1), 2)
```

- Se va a añadir una columna nueva proyect_length de tipo numérico, que contendrá el total de días que el proyecto ha estado abierto a financiación. Esta nueva columna será resultado de la diferencia entre la columna deadline y launched.

```
dataSet['proyect_length'] <- as.numeric(dataSet$deadline - dataSet$launched)
```

```
# Muestra set de datos
head(dataSet, 5)
```

```
##              ID                                     name
## 277557 48211191 Mint Julep Jazz Band's First CD - Durham on Saturday Night
## 140412 1712941959                                     QuikLid
```

```
## 10877 1054544730 Enamel Pins Set: Twisted Tales
## 254046 361859883 Let Tinned Pineapple Make Your Soapy Dreams Come True!
## 273982 463992846 EDC "NOMAD" Survival Multi-Tool
## category main_category currency deadline goal launched
## 277557 Jazz Music USD 2012-12-31 6000 2012-11-01
## 140412 Product Design Design USD 2017-07-03 225000 2017-05-04
## 10877 Accessories Fashion USD 2017-04-12 300 2017-03-22
## 254046 Crafts Crafts USD 2012-03-31 1000 2012-02-09
## 273982 Product Design Design USD 2016-03-25 650 2016-02-24
## pledged state backers country usd.pledged usd_pledged_real
## 277557 7460 successful 216 US 7460 7460
## 140412 1698 canceled 78 US 546 1698
## 10877 203 canceled 8 US 203 203
## 254046 327 failed 18 US 327 327
## 273982 11152 successful 234 US 11152 11152
## usd_goal_real usd_pledged status euros_pledged euros_goal project_length
## 277557 6000 7460 1 6065.0410 4878.05 60
## 140412 225000 1698 0 1380.4880 182926.80 60
## 10877 300 203 0 165.0407 243.90 21
## 254046 1000 327 0 265.8537 813.01 51
## 273982 650 11152 1 9066.6670 528.46 30
```

```
summary(dataSet)
```

```
## ID name category
## Min. :5.141e+05 Length:3000 Product Design: 173
## 1st Qu.:5.259e+08 Class :character Documentary : 141
## Median :1.055e+09 Mode :character Music : 122
## Mean :1.060e+09 Tabletop Games: 107
## 3rd Qu.:1.600e+09 Food : 102
## Max. :2.147e+09 Shorts : 97
## (Other) :2258
## main_category currency deadline goal
## Film & Video: 518 USD :2350 Min. :2009-07-16 Min. : 1
## Music : 411 GBP : 267 1st Qu.:2013-06-28 1st Qu.: 2000
## Publishing : 330 EUR : 139 Median :2015-02-01 Median : 5000
## Games : 259 CAD : 112 Mean :2014-11-13 Mean : 59075
## Design : 235 AUD : 69 3rd Qu.:2016-05-20 3rd Qu.: 15000
## Technology : 235 NZD : 17 Max. :2018-02-12 Max. :100000000
## (Other) :1012 (Other): 46
## launched pledged state backers
## Min. :2009-06-03 Min. : 0.0 canceled : 314 Min. : 0.00
## 1st Qu.:2013-05-22 1st Qu.: 40.0 failed :1574 1st Qu.: 2.00
## Median :2014-12-28 Median : 625.5 live : 20 Median : 12.00
## Mean :2014-10-10 Mean : 9154.1 successful:1059 Mean : 102.64
## 3rd Qu.:2016-04-16 3rd Qu.: 4002.8 suspended : 11 3rd Qu.: 54.25
## Max. :2018-01-02 Max. :2669009.6 undefined : 22 Max. :20680.00
##
## country usd.pledged usd_pledged_real usd_goal_real
## US :2334 Min. : 0 Min. : 0.0 Min. : 1
## GB : 266 1st Qu.: 20 1st Qu.: 40.0 1st Qu.: 2000
## CA : 112 Median : 390 Median : 605.9 Median : 5457
## AU : 68 Mean : 6181 Mean : 8879.9 Mean : 58343
## DE : 37 3rd Qu.: 3033 3rd Qu.: 4000.0 3rd Qu.: 15000
## (Other): 159 Max. :1068328 Max. :2669009.6 Max. :100000000
```

```
## NA's : 24 NA's :24
##   usd_pledged      status   euros_pledged      euros_goal
## Min. :    0.0  0:1943   Min. :    0.0   Min. :    1
## 1st Qu.:   40.0  1:1057   1st Qu.:   32.5   1st Qu.:   1626
## Median :   605.9      Median :   492.6   Median :   4065
## Mean :   8879.9      Mean :   7219.5   Mean :   48028
## 3rd Qu.:  4000.0      3rd Qu.:  3252.0   3rd Qu.:  12195
## Max. :2669009.6      Max. :2169927.0   Max. :81300810
##
## proyect_length
## Min. : 1.00
## 1st Qu.:30.00
## Median :30.00
## Mean :34.03
## 3rd Qu.:36.00
## Max. :91.00
##
```

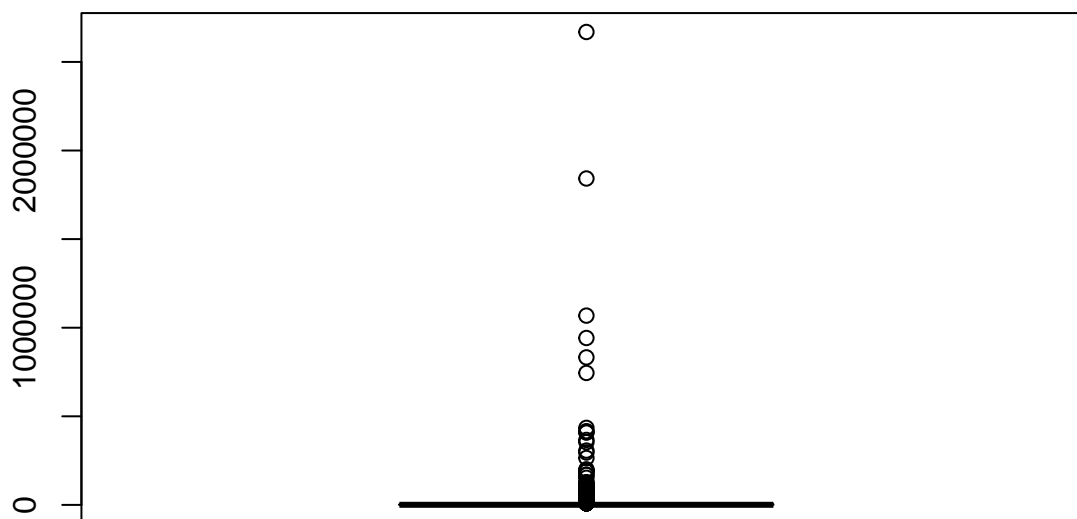
3.2 Valores atípicos

Volviendo a la estadística descriptiva vemos que la diferencia entre la media y el máximo y mínimo valor de muestras de la variable pledged y proyect_length es bastante significativa, lo que puede indicar la presencia de outliers. Vamos a comprobar si tenemos outliers mediante diagrama de cajas.

```
# Importamos la librería ggplot2
library(ggplot2)

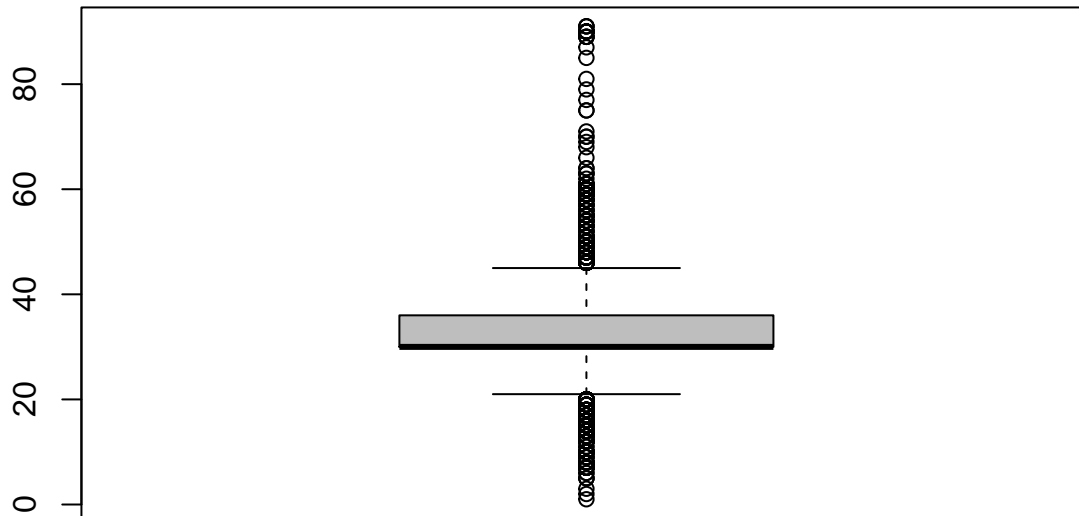
# Diagrama de cajas para la variable pledged y proyect_length
boxplot(dataSet$pledged, main="Box plot", col="gray")
```

Box plot



```
boxplot(dataSet$project_length, main="Box plot", col="gray")
```

Box plot



Vemos que ambas variables tienen valores extremos por lo que vamos a analizar para determinar cómo proceder con ellos.

```
tail(sort(dataSet$pledged), 10)
```

```
## [1] 408160.1 411141.0 415922.0 434805.0 744628.0 832523.0 941966.9
## [8] 1068328.0 1842141.7 2669009.6
```

```
tail(sort(dataSet$project_length), 10)
```

```
## [1] 90 90 90 90 90 90 90 90 91 91 91
```

Vemos que hay bastante diferencia entre la media de la variable pledged y los valores más altos, pero haciendo un poco de investigación online se ha encontrado que algún proyecto ha llegado a recaudar más de 20,000,000\$, [enlace](<https://www.marketwatch.com/story/10-kickstarter-products-that-raised-the-most-money-2017-06-22-10883052>). Por este motivo se ha decidido aceptar dichos outliers y tratarlos como datos válidos.

Haciendo un poco de investigación sobre la normativa de Kickstarter, se ha encontrado [enlace](<https://help.kickstarter.com/hc/en-us/articles/115005128434-What-is-the-maximum-project-duration-#:~:text=Projects%20on%20Kickstarter%20can%20last,at%2030%20days%20or%20less.>) que, hoy en día, la duración máxima por proyecto es de 60 días. También en este otro artículo se explica que hasta el año 2011 la duración máxima era de 90 días [enlace](<https://www.kickstarter.com/blog/shortening-the-maximum-project-length>).

Por ello se ha decidido que cualquier duración significativamente mayor de 90 días se va a tratar como outlier y se reemplazará por NA para posteriormente imputarlo con un valor de 90.

```
dataSet$project_length <- ifelse(dataSet$project_length > 90, NA, dataSet$project_length)
index <- which(is.na(dataSet$project_length))
index
```

```
## [1] 1223 1466 2935
```

Muestra de todas las variables y sus valores NA's.

```
# Comprobamos si quedan NA's
sapply(dataSet, function(x) sum(is.na(x)))
```

```
## ID name category main_category
```



```
##           0           0           0           0
##      currency      deadline      goal      launched
##           0           0           0           0
##      pledged      state      backers      country
##           0           0           0           24
##      usd.pledged usd_pledged_real usd_goal_real usd_pledged
##           24           0           0           0
##      status      euros_pledged      euros_goal      project_length
##           0           0           0           3
```

3.3 Imputación de valores

- Cómo se ha mencionado anteriormente los valores NA de la variable `project_length` se van a reemplazar por 90 ya que es el máximo número de días que un proyecto puede estar recaudando dinero.

```
dataSet$project_length[index] <- 90
```

- Imputación de valores perdidos para la variable `country` en base a la variable `currency`.

```
idx <- which(is.na(dataSet$country))
# encontrar las combinaciones únicas de country y currency pero no cuando country
# es NA
uniques <- unique(dataSet[c('country', 'currency')])
uniques <- uniques[!is.na(uniques$country),]

# reemplazar los NA's de country con los valores únicos asociados con currency
na.country <- which(is.na(dataSet$country))
na.currency <- dataSet$currency[na.country]
dataSet$country[idx] <- uniques$country[match(na.currency, uniques$currency)]
```

Por último comprobamos si quedan NA's en los datos.

```
# Comprobamos si quedan NAs
sapply(dataSet, function(x) sum(is.na(x)))
```

```
##           ID           name      category      main_category
##           0           0           0           0
##      currency      deadline      goal      launched
##           0           0           0           0
##      pledged      state      backers      country
##           0           0           0           0
##      usd.pledged usd_pledged_real usd_goal_real usd_pledged
##           24           0           0           0
##      status      euros_pledged      euros_goal      project_length
##           0           0           0           0
```

3.4 Selección de datos

A continuación, vamos a detallar que atributos hemos descartado y cuales hemos decidido sean imprescindibles para el análisis:

- Se ha decidido borrar del dataset la columna de `usd_pledged`, esta representa la conversión a dólares por parte de Kickstarter del atributo `pledged`, pero se han descubierto bastantes inconsistencias. El creador del dataset, por este mismo motivo, decidió incluir un nuevo atributo con una conversión más precisa de `pledged`, que es la que vamos a usar.
- También se ha decidido descartar la variable `usd_goal_pledged` porque no resulta significativa para el estudio.

- Por otro lado, vamos a prescindir de la variable state. Como se ha mencionado anteriormente, un Kickstarter es satisfactorio si el proyecto consigue recaudar el dinero marcado como objetivo en el tiempo estimado, por lo que no es necesario para nuestro estudio si dicho proyecto se ha cancelado, o se ha suspendido o sigue activo. Se puede dar el caso por ejemplo que un proyecto llegue al objetivo económico marcado dentro de tiempo, pero el organizador, por cualquier motivo decida suspenderlo. En ese caso el proyecto aparecerá como cancelado, pero desde el punto de vista del objetivo del proyecto, la recaudación ha sido satisfactoria.

```
# Quitar columnas usd_pledged y usd_goal_pledged del dataset
drops <- c('usd.pledged', 'usd_goal_pledged', 'state')
dataSet <- dataSet[ , !(names(dataSet) %in% drops)]
# Análisis descriptivo del dataset limpio
summary(dataSet)
```

```
##          ID          name          category
## Min.      :5.141e+05 Length:3000      Product Design: 173
## 1st Qu.:5.259e+08   Class :character Documentary   : 141
## Median :1.055e+09   Mode  :character Music          : 122
## Mean      :1.060e+09      Tabletop Games: 107
## 3rd Qu.:1.600e+09      Food          : 102
## Max.      :2.147e+09      Shorts        : 97
##                                     (Other)       :2258
##          main_category    currency    deadline    goal
## Film & Video: 518   USD      :2350   Min.      :2009-07-16   Min.      :      1
## Music          : 411   GBP      : 267   1st Qu.:2013-06-28   1st Qu.:    2000
## Publishing     : 330   EUR      : 139   Median :2015-02-01   Median :    5000
## Games          : 259   CAD      : 112   Mean    :2014-11-13   Mean     :   59075
## Design         : 235   AUD      : 69   3rd Qu.:2016-05-20   3rd Qu.:   15000
## Technology     : 235   NZD      : 17   Max.    :2018-02-12   Max.     :100000000
## (Other)        :1012 (Other): 46
##          launched          pledged          backers          country
## Min.      :2009-06-03   Min.      :      0.0   Min.      :      0.00   US      :2350
## 1st Qu.:2013-05-22   1st Qu.:    40.0   1st Qu.:      2.00   GB      : 267
## Median :2014-12-28   Median :   625.5   Median :    12.00   CA      : 112
## Mean      :2014-10-10   Mean      : 9154.1   Mean      : 102.64   AU      : 69
## 3rd Qu.:2016-04-16   3rd Qu.: 4002.8   3rd Qu.:    54.25   DE      : 37
## Max.      :2018-01-02   Max.      :2669009.6   Max.      :20680.00   NL      : 25
##                                     (Other): 140
##          usd_pledged_real    usd_goal_real    usd_pledged    status
## Min.      :      0.0   Min.      :      1   Min.      :      0.0   0:1943
## 1st Qu.:    40.0   1st Qu.:    2000   1st Qu.:    40.0   1:1057
## Median :   605.9   Median :   5457   Median :   605.9
## Mean      : 8879.9   Mean      : 58343   Mean      : 8879.9
## 3rd Qu.: 4000.0   3rd Qu.: 15000   3rd Qu.: 4000.0
## Max.      :2669009.6   Max.      :100000000   Max.      :2669009.6
##
##          euros_pledged    euros_goal    proyect_length
## Min.      :      0.0   Min.      :      1   Min.      : 1.00
## 1st Qu.:    32.5   1st Qu.:   1626   1st Qu.:30.00
## Median :   492.6   Median :   4065   Median :30.00
## Mean      : 7219.5   Mean      : 48028   Mean      :34.03
## 3rd Qu.: 3252.0   3rd Qu.: 12195   3rd Qu.:36.00
## Max.      :2169927.0   Max.      :81300810   Max.      :90.00
##
```

3.5 Exportación de los datos limpios

Una vez el procesamiento de los datos ha finalizado, se genera un archivo csv con nombre “ks-projects-201801_clean.csv”, que contendrá el dataset con 3000 muestras limpias.

```
# Exportación de los datos limpios en .csv
# write.csv(dataSet, '../data/ks-projects-201801_sample_clean.csv')
dataSet <- read.csv('../data/ks-projects-201801_sample_clean.csv')
attach(dataSet)
```

4 Análisis de los datos

4.1 Selección de los grupos de datos a analizar

De todo el conjunto de datos, se han seleccionado los siguientes atributos para poder ser analizados creyendo que son estos los que aportarán más valor al análisis posterior:

- main_category: recoge las 15 principales categorías presentes.

```
unique(main_category)
```

```
## [1] "Fashion"      "Publishing"    "Film & Video" "Music"         "Comics"
## [6] "Food"         "Art"           "Design"        "Crafts"        "Theater"
## [11] "Games"        "Dance"         "Journalism"    "Technology"    "Photography"
```

- project_length: tiempo de duración de cada proyecto, expresado en días.
- euros_goal: conversión a € de la variable goal.
- country: país en el que se publicó el proyecto Kickstarter. La variable currency representa la moneda de dicho país por lo que nos resulta redundante.
- backers: cantidad total de mecenas del proyecto.
- status: estado del proyecto, 0 no ha conseguido el objetivo, 1 si lo ha conseguido.

4.2 Normalidad y homocedasticidad

A la hora de identificar los métodos de análisis más adecuados se debe conocer antes las características de los datos, por ejemplo, si estos siguen una distribución normal o si presentan homocedasticidad. Por ello vamos a comprobar que las variables numéricas elegidas siguen una distribución normal o presentan homogeneidad de la varianza.

- Test de normalidad

Se va a utilizar el test Shapiro-Wilk, asumiendo un intervalo de confianza del 95%. Esto quiere decir que si el p-valor es menor o igual que el nivel de significancia con un valor de 0.05, entonces podemos rechazar la presunción de normalidad, es decir, la variable no sigue una distribución normal.

```
shapiro.test(project_length)
```

```
##
## Shapiro-Wilk normality test
##
## data:  project_length
## W = 0.8273, p-value < 2.2e-16
```

```
shapiro.test(euros_goal)
```

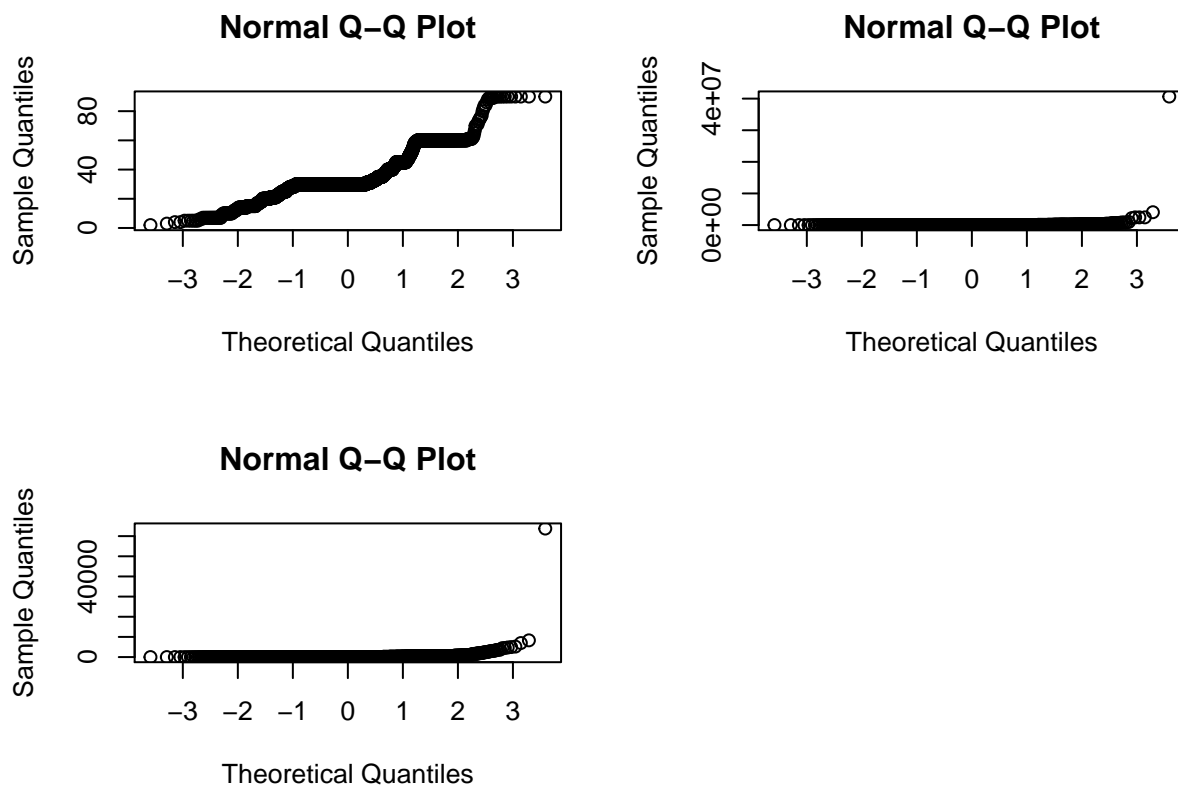
```
##
## Shapiro-Wilk normality test
```

```
##
## data: euros_goal
## W = 0.015699, p-value < 2.2e-16
```

```
shapiro.test(backers)
```

```
##
## Shapiro-Wilk normality test
##
## data: backers
## W = 0.041695, p-value < 2.2e-16
```

```
# Representación de la distribución
par(mfrow=c(2,2))
qqnorm(proyect_length)
qqnorm(euros_goal)
qqnorm(backers)
```



Se puede apreciar que los datos no siguen una distribución normal ya que en el total de las comprobaciones el p-valor del Test de Shapiro-Wilk el $p\text{-value} < 2.2e-16$ además que visualmente se se puede corroborar esto, por lo tanto podemos decir que las variables no siguen una distribución normal.

- Test de homocedasticidad

Ya que hemos comprobado que nuestros datos no siguen una distribución normal ($p\text{-value} < 2.2e-16$ en todos los casos), para el test de homocedasticidad utilizaremos el test de Fligner-Killeen. La hipótesis nula asume la igualdad de varianzas, por lo que p-values inferiores al nivel de significancia (0.05), indicarán heterocedasticidad.

Para ello comprobaremos distintos grupos de datos entre sí:

```

fligner.test(proyect_length ~ euros_goal, data = dataSet)

##
##  Fligner-Killeen test of homogeneity of variances
##
## data:  proyect_length by euros_goal
## Fligner-Killeen:med chi-squared = 331.81, df = 376, p-value = 0.951
fligner.test(euros_goal ~ backers, data = dataSet)

##
##  Fligner-Killeen test of homogeneity of variances
##
## data:  euros_goal by backers
## Fligner-Killeen:med chi-squared = 630.07, df = 389, p-value = 1.105e-13
fligner.test(proyect_length ~ backers, data = dataSet)

##
##  Fligner-Killeen test of homogeneity of variances
##
## data:  proyect_length by backers
## Fligner-Killeen:med chi-squared = 300.29, df = 389, p-value = 0.9997

```

De este análisis podemos observar dos casos, para las variables `proyect_length-euros_goal` y `proyect_length-backers` el test de Fligner-Killeen da un p-value mayor que 0.05 (0.951 y 0.9997 respectivamente), por lo que se asume homocedasticidad.

Por otro lado, la prueba para `euros_goal` y `backers` se resuelve con un p-value 1.105e-13, por lo que en este caso si se puede rechazar la hipótesis nula de homocedasticidad y se concluye que la variable `euros_goal` presenta varianzas estadísticamente diferentes para los diferentes grupos de `backers`.

4.3 Pruebas estadísticas

4.3.1 Contraste de hipótesis de dos muestras

Para comprobar si existe relación entre las variables `status` y `main_category`, es decir, proyecto fundado exitosamente y tipo categoría vamos a aplicar el test no paramétrico khi cuadrado mediante la función `chisq.test()`.

Las hipótesis nula y alternativa quedarían de la siguiente manera:

- Hipótesis nula, H_0 : el éxito del proyecto y la categoría en la que se encuentre son variables independientes.
- Hipótesis alternativa, H_1 : existe relación entre la categoría en la que se encuentra un proyecto y el éxito de este.

```

chisqTable <- table( status, main_category )
chisq.test(chisqTable, correct=FALSE)

##
##  Pearson's Chi-squared test
##
## data:  chisqTable
## X-squared = 129.27, df = 14, p-value < 2.2e-16

```

Vemos que el p-value resultante del test es $< 2.2e-16$ por lo que podemos rechazar la hipótesis nula y afirmar, con un 95% de confianza que existe relación significativa entre la categoría en la que se encuentra un proyecto y su éxito o no.

4.3.2 Modelo regresión logística simple

En este primero modelo de regresión logística simple se quiere analizar la probabilidad de que un proyecto consiga recaudar los fondos propuestos en base a la longitud establecida.

Como se ha comentado anteriormente, Kickstarter decidió disminuir la duración máxima de los proyectos a 60 días y de hecho recomiendan configurar la campaña para que dure 30 días o menos, ya que la probabilidad de éxito disminuye conforme se alarga el proyecto. Es por ello por lo que queremos comprobar si esto es cierto mediante este modelo.

```
# Ajuste de un modelo logístico
modelLogisticSimple <- glm(status ~ project_length, family = "binomial")
summary(modelLogisticSimple)
```

```
##
## Call:
## glm(formula = status ~ project_length, family = "binomial")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2891  -0.9566  -0.8299   1.4157   2.1746
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.316066   0.119579   2.643  0.00821 **
## project_length -0.028687   0.003456  -8.302 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 3861.9  on 2999  degrees of freedom
## Residual deviance: 3785.5  on 2998  degrees of freedom
## AIC: 3789.5
##
## Number of Fisher Scoring iterations: 4
confint(object = modelLogisticSimple, level = 0.95 )

##              2.5 %      97.5 %
## (Intercept)    0.08325544  0.55217237
## project_length -0.03555386 -0.02200222
```

4.3.2.1 Predicciones Vamos a comparar probabilidad de que un proyecto sea exitoso cuando la duración de este es de 25 días y cuando es de 45 días y 60 días.

```
# Predicción para 25 días
round(predict(modelLogisticSimple, data.frame(project_length = 25), type="response"), 2)

##      1
## 0.4

# Predicción para 45 días
round(predict(modelLogisticSimple, data.frame(project_length = 45), type="response"), 2)

##      1
## 0.27
```

```
# Predicción para 60 días
round(predict(modelLogisticSimple, data.frame(proyect_length = 60), type="response"), 2)
```

```
## 1
## 0.2
```

Podemos dar por cierto la afirmación de que a más largo el proyecto menos probabilidad de éxito hay ya que hemos obtenido que, para los proyectos de 25 días hay un 40% de probabilidad, para los de 45 días un 27% y para los de 60 días un 20%, se aprecia la tendencia a la baja.

4.3.3 Modelo de regresión logística múltiple

Vamos a crear un primer modelo predictivo de regresión logística para predecir la expectativa de que un proyecto sea exitoso antes de lanzarlo. Para ello tendremos como variable respuesta status, y como variables explicativas usaremos: main_category, proyect_length, backers y euros_goal.

Vamos a especificar el nivel base de referencia para la variable cualitativa:

- Para la variable main_category, la categoría 'Dance'.

```
status <- as.factor(status)
main_category <- as.factor(main_category)

# Nivel de referencia
main_category <- relevel(main_category, ref = 'Dance')

modelLogistic = glm(formula = status ~ main_category + proyect_length +
                    backers + euros_goal, family = binomial(link = logit))
summary(modelLogistic)
```

```
##
## Call:
## glm(formula = status ~ main_category + proyect_length + backers +
##     euros_goal, family = binomial(link = logit))
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -8.4904  -0.6031  -0.1044   0.2609   4.5555
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.477e-01  4.834e-01   0.719  0.471938
## main_categoryArt    -3.655e-01  4.922e-01  -0.743  0.457720
## main_categoryComics -7.867e-01  5.451e-01  -1.443  0.148975
## main_categoryCrafts -1.148e+00  5.869e-01  -1.957  0.050380 .
## main_categoryDesign -1.055e+00  5.247e-01  -2.011  0.044309 *
## main_categoryFashion -8.089e-01  5.148e-01  -1.571  0.116159
## main_categoryFilm & Video -3.991e-01  4.798e-01  -0.832  0.405569
## main_categoryFood    -1.096e+00  5.229e-01  -2.096  0.036111 *
## main_categoryGames    -1.782e+00  5.191e-01  -3.432  0.000598 ***
## main_categoryJournalism -9.212e-01  6.578e-01  -1.401  0.161357
## main_categoryMusic    -1.391e-01  4.821e-01  -0.288  0.773001
## main_categoryPhotography -8.319e-01  5.632e-01  -1.477  0.139665
## main_categoryPublishing -9.017e-01  4.889e-01  -1.845  0.065096 .
## main_categoryTechnology -8.363e-01  5.280e-01  -1.584  0.113219
## main_categoryTheater    -5.243e-02  5.366e-01  -0.098  0.922166
```

```
## project_length          -1.605e-02  4.669e-03  -3.438 0.000587 ***
## backers                 3.895e-02  1.862e-03  20.919 < 2e-16 ***
## euros_goal             -2.472e-04  1.537e-05 -16.088 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 3861.9 on 2999 degrees of freedom
## Residual deviance: 1965.7 on 2982 degrees of freedom
## AIC: 2001.7
##
## Number of Fisher Scoring iterations: 12
```

```
# Odds Ratio
```

```
exp(coefficients(modelLogistic))
```

```
## (Intercept)          main_categoryArt      main_categoryComics
## 1.4158301           0.6938535           0.4553503
## main_categoryCrafts  main_categoryDesign  main_categoryFashion
## 0.3171124           0.3480801           0.4453674
## main_categoryFilm & Video  main_categoryFood      main_categoryGames
## 0.6709450           0.3342483           0.1683337
## main_categoryJournalism  main_categoryMusic  main_categoryPhotography
## 0.3980378           0.8701811           0.4352336
## main_categoryPublishing  main_categoryTechnology  main_categoryTheater
## 0.4058637           0.4333262           0.9489238
## project_length          backers          euros_goal
## 0.9840785           1.0397144           0.9997528
```

```
exp(confint.default(modelLogistic, level = 0.95))
```

```
## 2.5 % 97.5 %
## (Intercept) 0.54897116 3.6515121
## main_categoryArt 0.26443859 1.8205840
## main_categoryComics 0.15643816 1.3254047
## main_categoryCrafts 0.10036935 1.0019022
## main_categoryDesign 0.12445863 0.9734943
## main_categoryFashion 0.16236441 1.2216476
## main_categoryFilm & Video 0.26198089 1.7183207
## main_categoryFood 0.11993823 0.9314956
## main_categoryGames 0.06085578 0.4656290
## main_categoryJournalism 0.10965741 1.4448096
## main_categoryMusic 0.33827593 2.2384542
## main_categoryPhotography 0.14431909 1.3125659
## main_categoryPublishing 0.15569239 1.0580177
## main_categoryTechnology 0.15395575 1.2196465
## main_categoryTheater 0.33150745 2.7162477
## project_length 0.97511453 0.9931250
## backers 1.03592755 1.0435152
## euros_goal 0.99972269 0.9997829
```

Observando el sumario del modelo podemos decir que las variables backers, project_length y euros_goal son significativas al tener p-valores menores que el nivel de significancia 0.05 ($< 2e-16$, 0.000587 y $< 2e-16$ respectivamente). La variable project_length, como ya sabíamos afecta de forma negativa al éxito del proyecto con un coeficiente asociado igual a -1.605e-02. Por otra parte, la variable backers afecta a al éxito

del proyecto de manera positiva $3.895e-02$. La variable `euros_goal` (objetivo económico marcado al que se tiene que llegar antes de que se termine la campaña), afecta de forma negativa a razón de $2.472e-04$.

En cuanto a la variable `main_category`, hemos definido antes que el nivel de referencia es 'Dance', si nos fijamos en los OR (odds-ratio) vemos que para el resto las categorías el OR está por debajo de la unidad, lo que indica que es menos probable que un proyecto sea éxito si pertenece a cualquier otra categoría que si pertenece a Dance.

Cuanto más se aleja el valor del odds-ratio de la unidad, más fuerte es la relación entre la variable dependiente e independiente, por ello vemos que los más alejados de la unidad corresponden a categoría Games (0.168). Si calculamos la inversa tenemos que para Technology $1/0.168 = 5.952$, esto se podría interpretar como, si el proyecto pertenece a la categoría Dance, los odds de éxito son 5.952 veces mayor que si pertenece a la categoría Games.

Basándonos en el intervalo de confianza de las tres variables independientes podemos decir que la que más impacto tiene sobre el éxito del proyecto es `main_category`, es decir, la categoría en la que encuentra el proyecto. Y dentro de `main_category`, podemos decir que la categoría que influye de manera más positiva al éxito de un proyecto es Dance.

A modo de nota, se ha probado a introducir la variable `country`, pero empeoraba el modelo considerablemente, el AIC pasaba a superar los 3700. Además de que todos los p-values para las distintas variables dummy creadas para representar los distintos países estaban por encima de 0.05.

4.3.3.1 Comparación de las predicciones con las observaciones Para conocer un poco más del modelo crearemos la matriz de confusión y veremos qué porcentaje de observaciones de entrenamiento es capaz de clasificar correctamente el modelo.

```
predicciones <- ifelse(test = modelLogistic$fitted.values > 0.5, yes = 1, no = 0)
# predicciones
matriz_confusion <- table(modelLogistic$model$status, predicciones,
                          dnn = c("observaciones", "predicciones"))
matriz_confusion
```

```
##           predicciones
## observaciones    0     1
##              0 1896    72
##              1  262   770
```

El modelo es capaz de clasificar correctamente $((1896 + 770) / (1896 + 770 + 262 + 72)) = 0.8886$, es decir, el 88.86% de las observaciones de entrenamiento, por lo que se puede considerar bueno.

4.3.3.2 Predicción A continuación, vamos a aplicar nuestro modelo para predecir que probabilidad tiene de ser exitoso un proyecto que pertenece a la categoría Music, cuya duración va a ser de 30 y cuyo objetivo monetario serán 10,000€. Ya que el número de mecenas no se sabe hasta que el proyecto no ha finalizado, vamos a calcularlo para 5 valores distintos (20, 40, 60, 80, 100).

```
predict(modelLogistic, data.frame(main_category = 'Dance', project_length = 30,
                                   backers = seq(from = 20, to = 100, by = 20),
                                   euros_goal = 10000), type = "response")
```

```
##           1           2           3           4           5
## 0.1385835 0.2595742 0.4330886 0.6247271 0.7839073
```

La probabilidad, expresada en porcentaje para los valores mencionados anteriormente iría desde 13.85% si hay 20 backers a un 78.39% si los backers son 100.

5 Representación de los resultados

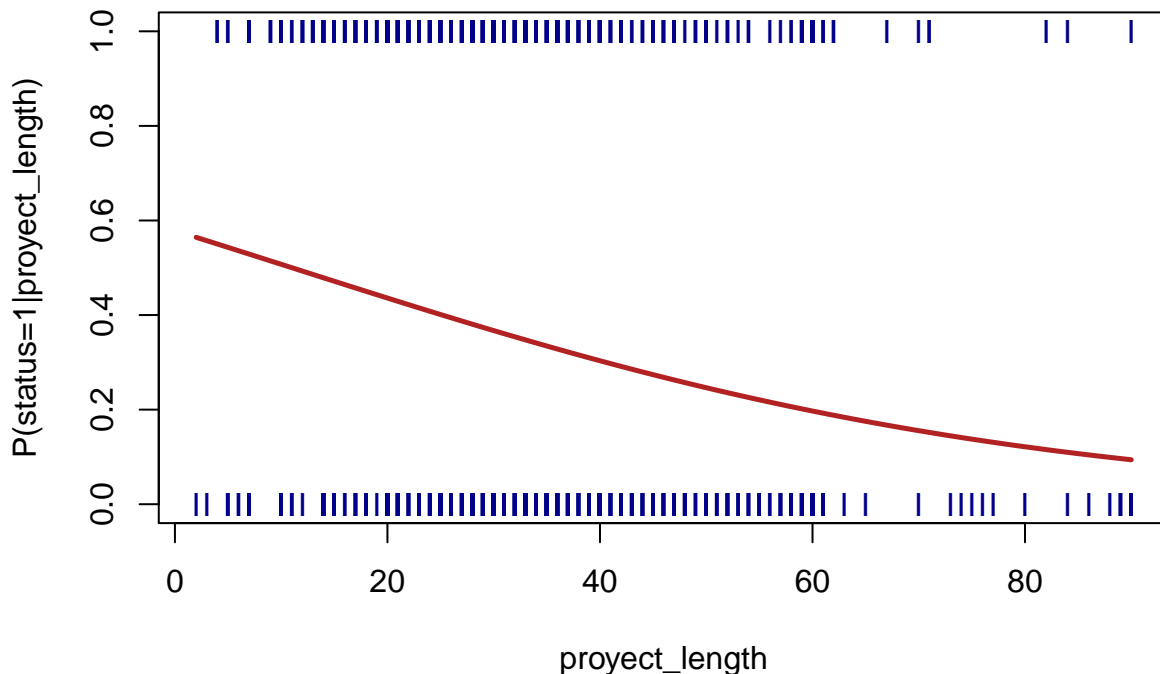
- Visualización predictiva de la variable del éxito de un proyecto en base a la variable `proyect_length`

```
dataSet$status <- as.character(dataSet$status)
dataSet$status <- as.numeric(dataSet$status)

plot(status ~ proyect_length, dataSet, col = "darkblue",
     main = "Modelo regresión logística",
     ylab = "P(status=1|proyect_length)",
     xlab = "proyect_length", pch = "I")

curve(predict(modelLogisticSimple, data.frame(proyect_length = x), type = "response"),
      col = "firebrick", lwd = 2.5, add = TRUE)
```

Modelo regresión logística



- Visualización predictiva de la variable del éxito de un proyecto en base a las variables `main_category`, `proyect_length`, `backers` y `euros_goal`. La variable `main_category` será `Dance`, la longitud del proyecto será 30 días. Luego la variable `backers` tomará valores de 10 a 200 de 10 en 10 y por último para la variable `euros_goal` se han propuesto los valores 1000, 5000, 10000, 20000 y 50000.

```
require(gridExtra)
newdata <- data.frame(proyect_length = c(rep(30, 100)),
                     backers = rep(seq(from = 10, to = 200, by = 10), 5),
                     main_category=c(rep('Dance', 100)),
                     euros_goal=c(rep(1000, 20), rep(5000, 20), rep(10000, 20),
                                   rep(20000, 20), rep(50000, 20)))

successPrediction <- predict(modelLogistic, newdata, type="response")
goalPrediction <- as.factor(c(rep(1000, 20), rep(5000, 20), rep(10000, 20),
                             rep(20000, 20), rep(50000, 20)))
backersTotal <- rep(seq(from = 10, to = 200, by = 10), 5)
```

```

predictionDataFrame <- data.frame(goalPrediction, successPrediction, backersTotal)

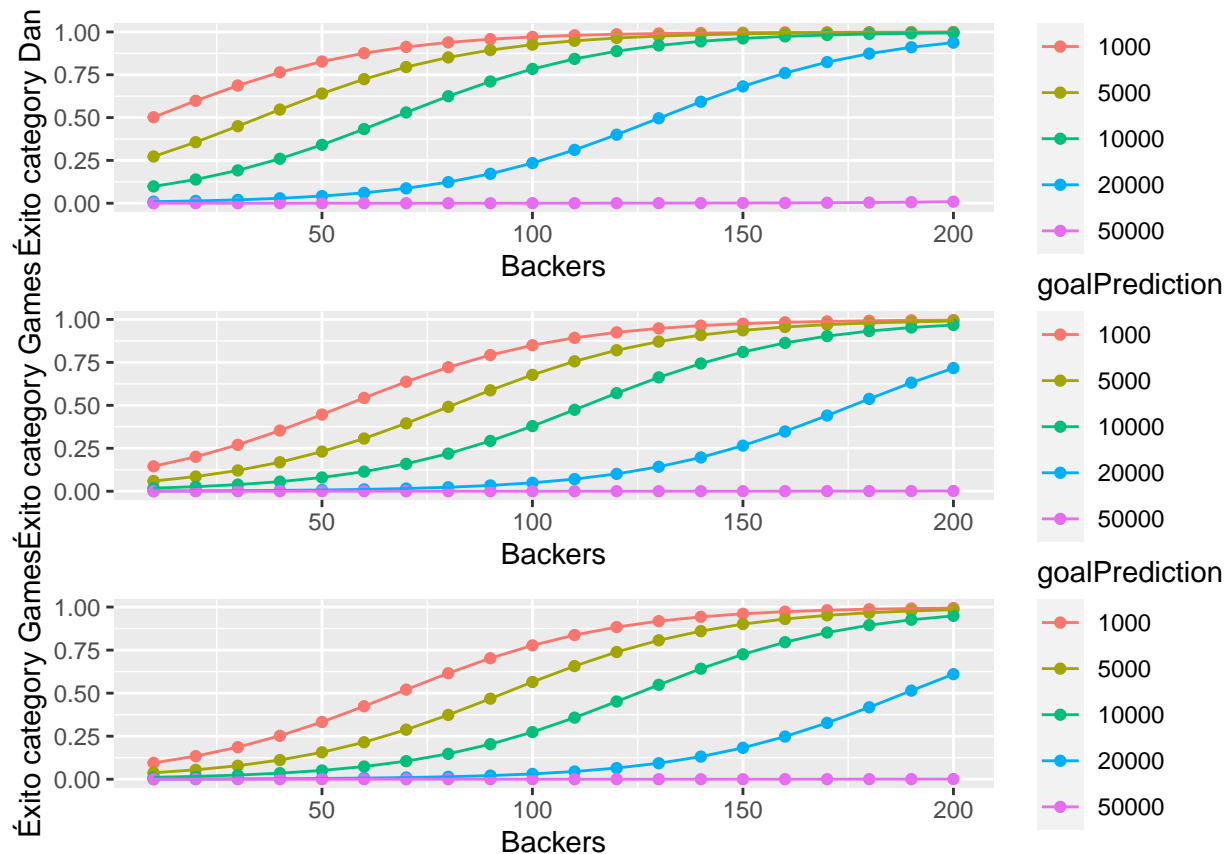
# Gráfica para la categoría Dance
plot1 <- ggplot(predictionDataFrame, aes(x = backersTotal, y = successPrediction,
col = goalPrediction))+geom_point()+geom_line()+
ylab('Éxito category Dance') + xlab('Backers')

# Gráfica para la categoría Games
newdata$main_category <- c(rep('Games', 100))
successPrediction <- predict(modelLogistic, newdata, type="response")
predictionDataFrame <- data.frame(goalPrediction, successPrediction, backersTotal)
plot2 <- ggplot(predictionDataFrame, aes(x = backersTotal, y = successPrediction,
col = goalPrediction))+geom_point()+geom_line()+
ylab('Éxito category Games') + xlab('Backers')

# Gráfica para longitud de proyecto 60 días
newdata$project_length <- c(rep(60, 100))
successPrediction <- predict(modelLogistic, newdata, type="response")
predictionDataFrame <- data.frame(goalPrediction, successPrediction, backersTotal)
plot3 <- ggplot(predictionDataFrame, aes(x = backersTotal, y = successPrediction,
col = goalPrediction))+geom_point()+geom_line()+
ylab('Éxito category Games') + xlab('Backers')

grid.arrange(plot1, plot2, plot3, ncol = 1)

```



Podemos observar como a medida que aumenta la cantidad de mecenas la probabilidad de que el proyecto consiga recaudar el goal propuesto aumenta. Dependiendo de dicho goal, se necesitarán más backers o menos para que la probabilidad aumente más deprisa. Por ejemplo, si el goal propuesto son 5000€, para la categoría

Dance, con apenas 50 backers la probabilidad llegaría al 70%, sin embargo, si el goal son 10000, con los 50 backers de antes la probabilidad estaría en torno al 38%. Vemos también que cuando el objetivo está muy por encima de la media de la variable euros_goal la curva de subida no es tan pronunciada.

Al cambiar la categoría observamos que la probabilidad de éxito disminuye, es decir, se necesitan más mecenas por proyecto para que el proyecto sea exitoso. Y lo mismo ocurre si la duración del proyecto se alarga hasta los 60 días.

6 Conclusiones

Se ha conseguido crear un modelo que predice que probabilidad tiene un proyecto de conseguir el objetivo monetario propuesto en base a la categoría en la que se incluye el proyecto, la longitud de este, el goal fijado y los mecenas. Ciertamente es que los mecenas del proyecto es imposible conocerlos de antemano, pero se pueden estimar ciertos valores que nos ayudarán a saber qué cantidad de backers necesitaremos para que nuestro proyecto sea exitoso.

Claramente la categoría a la que pertenece el mismo es un factor que influye en gran medida al igual que el goal propuesto. Se puede ver en las gráficas que a partir de cierto goal la probabilidad disminuye muchísimo.

De entrada, se puede pensar que a más días dure el proyecto mejor, más probabilidad de conseguir mecenas y llegar al objetivo, pero claramente, mediante los modelos y las gráficas se puede observar que esto no es así, sino todo lo contrario.

7 Agradecimientos

En primer lugar, agradecer y reconocer el trabajo de Mickaël Mouillé [enalce](<https://www.kaggle.com/kemical>), creador del dataset por trabajo para recolectar datos durante tantos años y publicarlos para el uso público.

También agradecer a todas aquellas personas que han publicado sus dudas sobre el dataset en beneficio de todos.

8 Tabla de contribuciones

Contribuciones	Firma
Investigación previa	M.G.
Redacción de las respuestas	M.G.
Desarrollo código	M.G.