

Spaceship-Titanic

Juan Luis González Rodríguez & Rocío González Martínez

2022-05-31

Índice

1 Contexto	2
1.1 Descripción del Dataset	2
1.2 ¿Por qué es importante y qué pregunta/problema pretende responder?	3
2 Integración y selección de los datos de interes.	4
3 Limpieza de los Datos.	5
3.1 Tratamiento valores nulos.	5
3.2 Tratamiento valores extremos.	6

```
# Package names
packages <- c("tidyr", "dplyr", "ggplot2", "keras", "reshape2", "tidyverse")
# Install packages not yet installed
installed_packages <- packages %in% rownames(installed.packages())
if (any(installed_packages == FALSE)) {
  install.packages(packages[!installed_packages])
}
# Packages loading
invisible(lapply(packages, library, character.only = TRUE))
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

## Warning: package 'keras' was built under R version 4.1.3

##
## Attaching package: 'reshape2'
```

```
## The following object is masked from 'package:tidyr':
##
## smiths

## -- Attaching packages ----- tidyverse 1.3.1 --

## v tibble 3.1.5      v stringr 1.4.0
## v readr 2.0.2      v forcats 0.5.1
## v purrr 0.3.4

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()
```

1 Contexto

1.1 Descripción del Dataset

El dataset *Spaceship Titanic* [1]. Este dataset es parte de la competición homónima y tiene por objetivo crear un algoritmo para predecir qué pasajeros han desaparecido al colisionar una nave espacial denominada Titanic con una anomalía espaciotemporal. Con el conjunto de datos, se pretende predecir si el pasajero ha desaparecido o no, para enviar a un equipo a rescatarlo. Para ello, se facilitan 2 ficheros (separados por entrenamiento y test), Se usará el fichero *train.csv* en uno para limpiar todos los registros y posteriormente se usará este para entrenar al modelo. Con el fichero test podremos probar el modelo (no incluye la variable objetivo).

Descripción de **Train.csv**: Conjunto de datos con información de unos 8 700 pasajeros. Este consta de los campos que se especifican más abajo.

Nombre	Tipo	Descripción
PassengerId	chr	Identificador de cada pasajero. El formato es gggg_pp (gggg hace referencia al grupo de pasajeros y pp al número dentro del grupo). Normalmente los miembros del grupo son familia.
HomePlanet	factor	Platena de origen del pasajero.
CryoSleep	logical	Indica si el pasajero está en animación suspendida durante el viaje o no.
Cabin	chr	Indican la cabina del pasajero. El formato es “plataforma/numero/lado”. Lado será P o S
Destination	factor	Indica el nombre del planeta de destino del pasajero.
Age	integer	Indica la edad biológica del pasajero en años en el momento del viaje.

Nombre	Tipo	Descripción
VIP	logical	Indica si el pasajero ha pagado por un servicio VIP o no
RoomService, FoodCouert, ShopingMall, Spa, VRDeck	numeric	Indica la cantidad de dinero que el pasajero ha gastado en cada uno de los servicios
Name	chr	Indica el nombre y apellido del pasajero
Transported	logical	Variable objetivo, indica si el pasajero ha sido transportado a otra dimensión o no (es decir si ha desaparecido).

La estructura del dataset es la siguiente:

```
df <- read.csv("~/MASTER CIENCIA DE DATOS/Tipologia y ciclo de vida de los datos/Practicas/Práctica2/Ejemplo1.csv",
               colClasses=c("HomePlanet"="factor",
                             "CryoSleep"="logical",
                             "Destination"="factor",
                             "VIP"="logical",
                             "Transported"="logical"))
df$Age <- as.integer(df$Age)
str(df)
```

```
## 'data.frame': 8693 obs. of 14 variables:
## $ PassengerId : chr "0001_01" "0002_01" "0003_01" "0003_02" ...
## $ HomePlanet : Factor w/ 4 levels "", "Earth", "Europa", ...: 3 2 3 3 2 2 2 2 2 3 ...
## $ CryoSleep : logi FALSE FALSE FALSE FALSE FALSE FALSE ...
## $ Cabin : chr "B/0/P" "F/0/S" "A/0/S" "A/0/S" ...
## $ Destination : Factor w/ 4 levels "", "55 Cancri e", ...: 4 4 4 4 4 3 4 4 4 2 ...
## $ Age : int 39 24 58 33 16 44 26 28 35 14 ...
## $ VIP : logi FALSE FALSE TRUE FALSE FALSE FALSE ...
## $ RoomService : num 0 109 43 0 303 0 42 0 0 0 ...
## $ FoodCourt : num 0 9 3576 1283 70 ...
## $ ShoppingMall: num 0 25 0 371 151 0 3 0 17 0 ...
## $ Spa : num 0 549 6715 3329 565 ...
## $ VRDeck : num 0 44 49 193 2 0 0 NA 0 0 ...
## $ Name : chr "Maham Ofracculy" "Juanna Vines" "Altark Susent" "Solam Susent" ...
## $ Transported : logi FALSE TRUE FALSE FALSE TRUE TRUE ...
```

1.2 ¿Por qué es importante y qué pregunta/problema pretende responder?

El objetivo que se persigue con el proyecto es el de, partiendo del conjunto de datos anteriormente comentado, desarrollar un modelo supervisado que permita responder a la pregunta: **¿Ha desaparecido el pasajero que se indica?**

Con ello, la tripulación podrá dirigir los esfuerzos de una manera más eficiente y maximizar las vidas salvadas.

2 Integración y selección de los datos de interes.

Solo hay 1 fichero de origen, por lo que no hay que combinar los datos de diferentes fuentes.

Como ya se tiene a los usuarios identificados a los usuarios en base a los identificadores, no es necesario almacenar sus nombres de cara al análisis. Por otro lado, de los campos *Passenger_id* y *Cabin* se pueden extraer aún más campos como el grupo y número dentro del grupo en el primer caso y la plataforma, número de cabina y lado en el segundo.

Se elimina la variable *Name* y se crean las nuevas variables derivadas.

```
df <- select(df, -Name)

df <- df %>%
  mutate(PassengerGroup=
    as.character(sapply(strsplit(PassengerId,"_"), `[, 1])) %>%
  mutate(PassengerNumInGroup=
    as.factor(sapply(strsplit(PassengerId,"_"), `[, 2])) %>%
  mutate(CabinPlatform =
    as.factor(sapply(strsplit(Cabin,"/"), `[, 1])) %>%
  mutate(CabinNumber =
    as.integer(sapply(strsplit(Cabin,"/"), `[, 2])) %>%
  mutate(CabinSide =
    as.factor(sapply(strsplit(Cabin,"/"), `[, 3]))

df <- select(df, -Cabin)
```

Tras crear las nuevas variables derivadas se elimina *Cabin* porque ya tenemos su información separada. *PassengerId* no se eliminará porque sirve para identificar los registros. Se muestra un resumen de los campos con la función *summary*.

```
summary(df)
```

##	PassengerId	HomePlanet	CryoSleep	Destination
##	Length:8693	: 201	Mode :logical	: 182
##	Class :character	Earth :4602	FALSE:5439	55 Cancr i e :1800
##	Mode :character	Europa:2131	TRUE :3037	PSO J318.5-22: 796
##		Mars :1759	NA's :217	TRAPPIST-1e :5915
##				
##				
##	Age	VIP	RoomService	FoodCourt
##	Min. : 0.00	Mode :logical	Min. : 0.0	Min. : 0.0
##	1st Qu.:19.00	FALSE:8291	1st Qu.: 0.0	1st Qu.: 0.0
##	Median :27.00	TRUE :199	Median : 0.0	Median : 0.0
##	Mean :28.83	NA's :203	Mean : 224.7	Mean : 458.1
##	3rd Qu.:38.00		3rd Qu.: 47.0	3rd Qu.: 76.0
##	Max. :79.00		Max. :14327.0	Max. :29813.0
##	NA's :179		NA's :181	NA's :183
##	ShoppingMall	Spa	VRDeck	Transported
##	Min. : 0.0	Min. : 0.0	Min. : 0.0	Mode :logical
##	1st Qu.: 0.0	1st Qu.: 0.0	1st Qu.: 0.0	FALSE:4315
##	Median : 0.0	Median : 0.0	Median : 0.0	TRUE :4378

```
## Mean : 173.7 Mean : 311.1 Mean : 304.9
## 3rd Qu.: 27.0 3rd Qu.: 59.0 3rd Qu.: 46.0
## Max. :23492.0 Max. :22408.0 Max. :24133.0
## NA's :208 NA's :183 NA's :188
## PassengerGroup PassengerNumInGroup CabinPlatform CabinNumber
## Length:8693 01 :6217 F :2794 Min. : 0.0
## Class :character 02 :1412 G :2559 1st Qu.: 167.2
## Mode :character 03 : 571 E : 876 Median : 427.0
## 04 : 231 B : 779 Mean : 600.4
## 05 : 128 C : 747 3rd Qu.: 999.0
## 06 : 75 (Other): 739 Max. :1894.0
## (Other): 59 NA's : 199 NA's :199
## CabinSide
## P :4206
## S :4288
## NA's: 199
##
##
##
##
```

Cabe destacar que en *HomePlanet* y en *Destination* hay campos con valores vacíos que no se han considerado como NA's. Por otro lado, Hay algunos campos que presenta NA's que podrán tratarse o desestimarse. También se observan valores extremos en algunos campos.

3 Limpieza de los Datos.

En este apartado se tratará de mejorar la calidad de los datos presentes en base a la falta de calidad. Por límite de extensión del proyecto, nos centraremos en el tratamiento de outliers y de valores nulos.

3.1 Tratamiento valores nulos.

Se remapean los campos en blanco de los campos *HomePlanet* y *Destination* por el valor *Unknown*. Con esto, no perdemos información y evitamos confundir a las personas que interpreten los resultados.

```
levels(df$HomePlanet) <- c("Unknown", "Earth", "Europa", "Mars")
levels(df$Destination) <- c("Unknown", "55 Cancri e", "PSO J318.5-22", "TRAPPIST-1e")
```

Se muestran la cantidad de valores nulos que tiene cada campo.

```
sapply(df, function(x) sum(length(which(is.na(x)))))
```

```
## PassengerId HomePlanet CryoSleeper Destination
## 0 0 217 0
## Age VIP RoomService FoodCourt
## 179 203 181 183
## ShoppingMall Spa VRDeck Transported
## 208 183 188 0
## PassengerGroup PassengerNumInGroup CabinPlatform CabinNumber
## 0 0 199 199
## CabinSide
## 199
```

Son relativamente pocos registros en comparación con el total que constan en el dataset. Por lo que se decide con contar con estos registros para entrenar al modelo predictivo.

```
nrow(df)
```

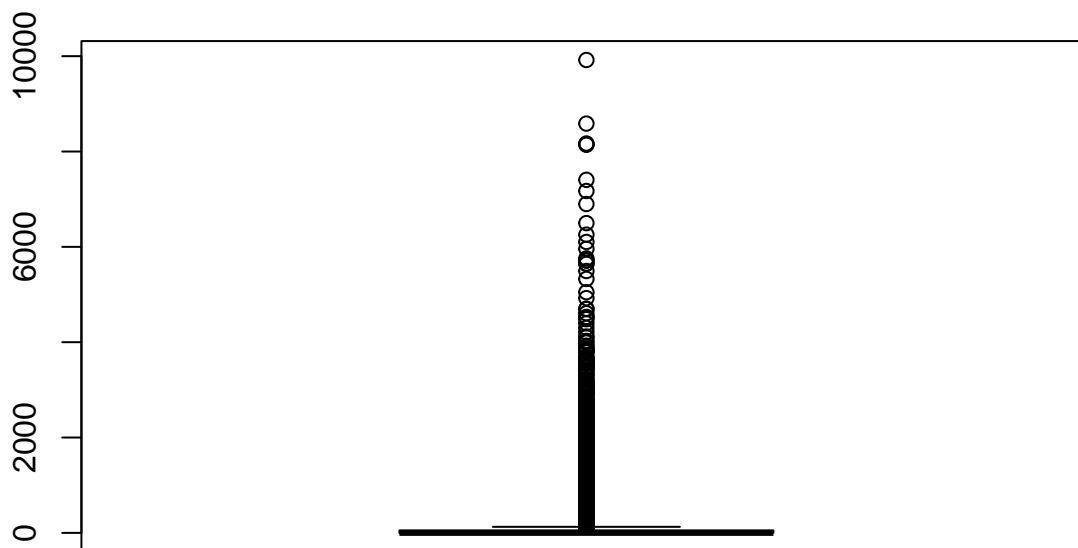
```
## [1] 8693
```

```
df <- na.omit(df)  
nrow(df)
```

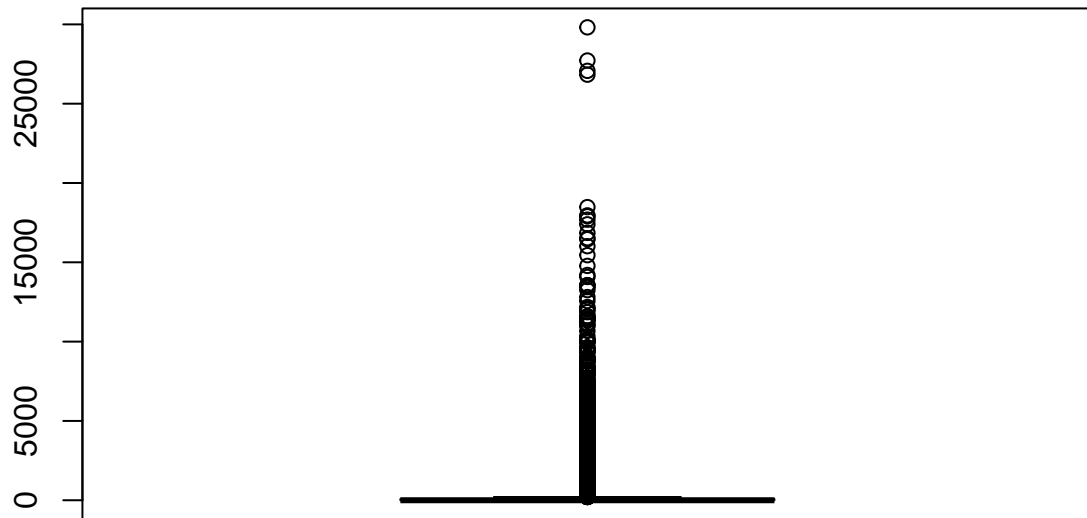
```
## [1] 7084
```

3.2 Tratamiento valores extremos.

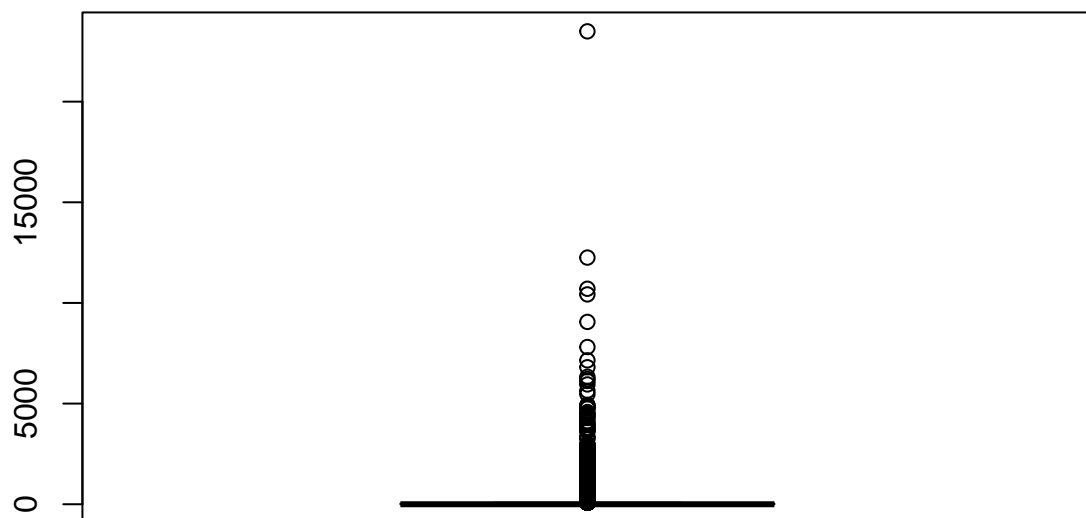
```
attach(df)  
boxplot(RoomService)
```



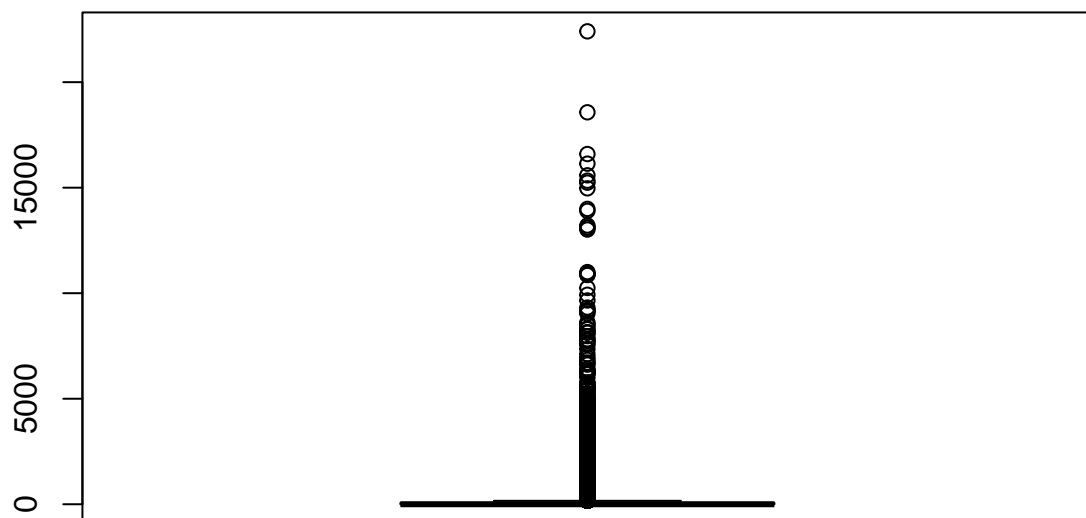
```
boxplot(FoodCourt)
```



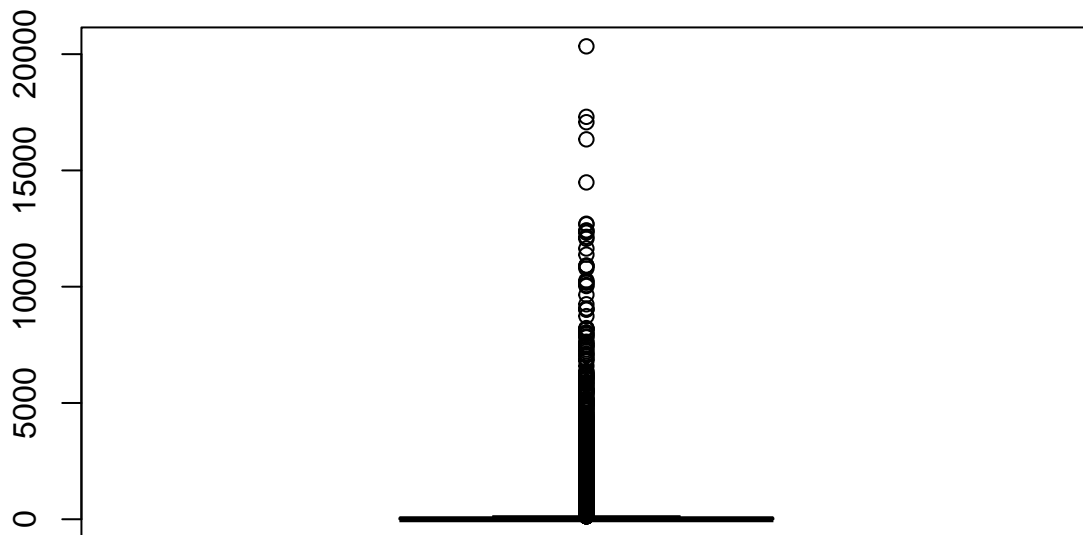
```
boxplot(ShoppingMall)
```



```
boxplot(Spa)
```

```
boxplot(VRDeck)
```



```
detach(df)
```

Aunque encontramos valores muy alejados de los valores centrales. No se consideraran como valores extremos. Se considerarán valores atípicos pero que son representativos de la variedad de nuestra muestra y por tanto formarán parte de los datos para entrenar al modelo. No se eliminará ningún valor extremo.

#4 Análisis de los datos ## 4.1 ## 4.2 ## 4.3 Analisis crudos

```
model <- glm(Transported ~ HomePlanet, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ HomePlanet, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.4619  -1.0522   0.9174   1.1479   1.3079
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -0.0585    0.1530  -0.382   0.702
## HomePlanetEarth -0.2432    0.1565  -1.554   0.120
## HomePlanetEuropa  0.7063    0.1611   4.384 1.17e-05 ***
```

```
## HomePlanetMars      0.1284      0.1619    0.793    0.428
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9560.9  on 7080  degrees of freedom
## AIC: 9568.9
##
## Number of Fisher Scoring iterations: 4
```

```
model <- glm(Transported ~ CryoSleep, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ CryoSleep, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8505  -0.8974   0.6309   0.6309   1.4861
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.7016     0.0313  -22.42  <2e-16 ***
## CryoSleepTRUE  2.2149     0.0609   36.37  <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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 8191.5  on 7082  degrees of freedom
## AIC: 8195.5
##
## Number of Fisher Scoring iterations: 4
```

```
model <- glm(Transported ~ Destination, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ Destination, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3886  -1.1203   0.9799   1.2356   1.2356
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
```

```
## (Intercept)          -0.05264      0.16228  -0.324  0.74563
## Destination55 Cancr e    0.53664      0.17095   3.139  0.00169 **
## DestinationPS0 J318.5-22  0.08935      0.18015   0.496  0.61992
## DestinationTRAPPIST-1e   -0.08306      0.16483  -0.504  0.61433
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9714.7  on 7080  degrees of freedom
## AIC: 9722.7
##
## Number of Fisher Scoring iterations: 4
```

```
model <- glm(Transported ~ Age, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ Age, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.316  -1.185   1.045   1.160   1.405
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.320287   0.053447   5.993 2.06e-09 ***
## Age         -0.010782   0.001655  -6.514 7.32e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9777.6  on 7082  degrees of freedom
## AIC: 9781.6
##
## Number of Fisher Scoring iterations: 3
```

```
model <- glm(Transported ~ VIP, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ VIP, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.186  -1.186   1.168   1.168   1.402
```

```
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.02112    0.02406   0.878 0.379958
## VIPTRUE     -0.53507    0.15982  -3.348 0.000814 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9808.8  on 7082  degrees of freedom
## AIC: 9812.8
##
## Number of Fisher Scoring iterations: 3
```

```
model <- glm(Transported ~ RoomService, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ RoomService, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.299  -1.298   1.061   1.061   3.670
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.805e-01  2.649e-02  10.59  <2e-16 ***
## RoomService -1.757e-03  9.268e-05 -18.95  <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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9139.9  on 7082  degrees of freedom
## AIC: 9143.9
##
## Number of Fisher Scoring iterations: 5
```

```
model <- glm(Transported ~ FoodCourt, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ FoodCourt, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -1.6792 -1.1676 0.8183 1.1872 1.1872
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.307e-02 2.481e-02 -0.930 0.352
## FoodCourt    6.840e-05 1.586e-05 4.313 1.61e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9800.0  on 7082  degrees of freedom
## AIC: 9804
##
## Number of Fisher Scoring iterations: 4
```

```
model <- glm(Transported ~ ShoppingMall, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ ShoppingMall, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.376  -1.178   1.092   1.177   1.177
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  9.190e-04 2.474e-02 0.037 0.970
## ShoppingMall 4.245e-05 3.891e-05 1.091 0.275
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9819.2  on 7082  degrees of freedom
## AIC: 9823.2
##
## Number of Fisher Scoring iterations: 3
```

```
model <- glm(Transported ~ Spa, data=df, family = binomial(link = 'logit'))
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ Spa, family = binomial(link = "logit"),
##      data = df)
```

```
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.301  -1.297   1.059   1.059   3.642
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.2858366  0.0265178   10.78  <2e-16 ***
## Spa         -0.0016860  0.0000961  -17.54  <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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9072.7  on 7082  degrees of freedom
## AIC: 9076.7
##
## Number of Fisher Scoring iterations: 6
```

```
model <- glm(Transported ~ VRDeck, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ VRDeck, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.284  -1.283   1.075   1.075   3.677
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.460e-01  2.622e-02   9.382  <2e-16 ***
## VRDeck       -1.384e-03  8.463e-05 -16.350  <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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9212.7  on 7082  degrees of freedom
## AIC: 9216.7
##
## Number of Fisher Scoring iterations: 6
```

```
#model <- glm(Transported ~ PassengerGroup, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ VRDeck, family = binomial(link = "logit"),
```

```
##      data = df)
##
## Deviance Residuals:
##      Min        1Q      Median        3Q        Max
## -1.284   -1.283    1.075    1.075    3.677
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.460e-01  2.622e-02   9.382  <2e-16 ***
## VRDeck      -1.384e-03  8.463e-05 -16.350  <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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9212.7  on 7082  degrees of freedom
## AIC: 9216.7
##
## Number of Fisher Scoring iterations: 6

model <- glm(Transported ~ PassengerNumInGroup, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ PassengerNumInGroup, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min        1Q      Median        3Q        Max
## -1.4270   -1.1340    0.9468    1.2215    1.2435
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.10303    0.02812  -3.664 0.000249 ***
## PassengerNumInGroup02  0.31833    0.06572   4.844 1.27e-06 ***
## PassengerNumInGroup03  0.50849    0.09924   5.124 2.99e-07 ***
## PassengerNumInGroup04  0.67294    0.15323   4.392 1.13e-05 ***
## PassengerNumInGroup05  0.33042    0.19753   1.673 0.094370 .
## PassengerNumInGroup06  0.31925    0.25110   1.271 0.203582
## PassengerNumInGroup07  0.42845    0.36506   1.174 0.240538
## PassengerNumInGroup08 -0.05112    0.55706  -0.092 0.926881
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9757.4  on 7076  degrees of freedom
## AIC: 9773.4
##
## Number of Fisher Scoring iterations: 4
```



```
model <- glm(Transported ~ CabinPlatform, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ CabinPlatform, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6117  -1.0712   0.7983   1.1453   1.6651
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.07276    0.13493  -0.539 0.589721
## CabinPlatformB  1.05290    0.16051   6.560 5.39e-11 ***
## CabinPlatformC  0.83219    0.16018   5.195 2.04e-07 ***
## CabinPlatformD -0.23106    0.16879  -1.369 0.171016
## CabinPlatformE -0.51336    0.15510  -3.310 0.000934 ***
## CabinPlatformF -0.18238    0.14126  -1.291 0.196674
## CabinPlatformG  0.14870    0.14180   1.049 0.294331
## CabinPlatformT -1.02585    1.16256  -0.882 0.377554
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9482.4  on 7076  degrees of freedom
## AIC: 9498.4
##
## Number of Fisher Scoring iterations: 4
```

```
model <- glm(Transported ~ CabinSide, data=df, family = binomial(link = 'logit'))
summary(model)
```

```
##
## Call:
## glm(formula = Transported ~ CabinSide, family = binomial(link = "logit"),
##      data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.272  -1.090   1.086   1.086   1.267
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.20808    0.03403  -6.115 9.66e-10 ***
## CabinSideS   0.42721    0.04780   8.937 < 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: 9820.4  on 7083  degrees of freedom
## Residual deviance: 9740.0  on 7082  degrees of freedom
## AIC: 9744
##
## Number of Fisher Scoring iterations: 3
```

- [1] Kaggle, “Spaceship titanic.” kaggle.com; Kaggle, 2022.Available: <https://www.kaggle.com/competitions/spaceship-titanic/overview>