

PRAC2: Neteja i anàlisi de les dades

Marc Ferrer Margarit (mferrermargarit@uoc.edu) i Marc Ramos Bruach (mramosbru@uoc.edu)

5/15/2021

1. Descripció del dataset. Perquè és important i quina pregunta/problema pretén respondre? (0.5) (MF)

El dataset que farem servir per a aquesta pràctica conté els retards i cancel·lacions dels vols 2015 que surten de l'aeroport de San Francisco. Normalment els motius principals dels retards de vol són relacionats amb el temps, però en alguns casos també hi ha retards de vols relacionats amb les companyies aèries o aeroports. Aquest document examina i mostra les causes de retard i cancel·lació en diversos aspectes. Així doncs aquest dataset és important per saber quines són les principals causes dels retards que s'hi han produït durant el 2015.

Les preguntes que volem respondre amb aquest dataset són quins dies de la setmana els quals es produeixen més retards, quines aerolínies són les tenen els retards i veure si la distància dels vols influeix en els retards que es produeixen.

Aquest dataset ha sigut obtingut a partir d'una pràctica anterior realitzada durant el màster de Data Science de la UOC. També es pot obtenir el dataset complet, amb totes les dades dels vols (aprox. 600 MB) al següent enllaç: <https://www.kaggle.com/usdot/flight-delays?select=flights.csv>

2. Integració i selecció de les dades d'interès a analitzar. (0.5) (MF)

Per veure les dades que conté el dataset el carregarem i mostrarem les columnes que conté i la mida del dataset:

```
flights <- read.csv("../data/flights.csv")
str(flights)

## 'data.frame': 145952 obs. of 28 variables:
## $ YEAR           : int  2015 2015 2015 2015 2015 2015 2015 2015 2015 2015 ...
## $ MONTH          : int  12 1 11 5 5 3 4 4 4 6 ...
## $ DAY            : int  31 8 27 30 23 29 11 18 25 2 ...
## $ DAY_OF_WEEK    : int  4 4 5 6 6 7 6 6 6 2 ...
## $ AIRLINE         : chr  "F9" "00" "AA" "00" ...
## $ FLIGHT_NUMBER   : int  668 6287 2207 5647 6508 6289 304 304 317 5459 ...
## $ ORIGIN_AIRPORT  : chr  "SFO" "SFO" "SFO" "SFO" ...
## $ DESTINATION_AIRPORT: chr  "DEN" "SBA" "DFW" "SBA" ...
## $ SCHEDULED_DEPARTURE: int  2000 1350 1200 1840 2230 831 2135 2135 2145 900 ...
## $ DEPARTURE_TIME   : int  1926 1318 1129 1814 2205 809 2113 2113 2123 838 ...
## $ DEPARTURE_DELAY   : int  -34 -32 -31 -26 -25 -22 -22 -22 -22 ...
## $ TAXI_OUT          : int  14 52 43 36 39 25 23 15 20 35 ...
## $ WHEELS_OFF        : int  1940 1410 1212 1850 2244 834 2136 2128 2143 913 ...
## $ SCHEDULED_TIME    : int  150 71 208 75 92 63 84 84 124 66 ...
```

```

## $ ELAPSED_TIME      : int  139 99 247 87 110 71 95 84 136 75 ...
## $ AIR_TIME          : int  109 44 177 47 64 41 67 65 111 36 ...
## $ DISTANCE          : int  967 262 1464 262 421 193 421 421 679 190 ...
## $ WHEELS_ON          : int  2229 1454 1709 1937 2348 915 2243 2233 2334 949 ...
## $ TAXI_IN            : int  16 3 27 4 7 5 5 4 5 4 ...
## $ SCHEDULED_ARRIVAL : int  2330 1501 1728 1955 2 934 2259 2259 2349 1006 ...
## $ ARRIVAL_TIME       : int  2245 1457 1736 1941 2355 920 2248 2237 2339 953 ...
## $ ARRIVAL_DELAY      : int  -45 -4 8 -14 -7 -14 -11 -22 -10 -13 ...
## $ DIVERTED           : int  0 0 0 0 0 0 0 0 0 0 ...
## $ CANCELLED          : int  0 0 0 0 0 0 0 0 0 0 ...
## $ CANCELLATION_REASON: chr  "" "" "" ...
## $ AIR_SYSTEM_DELAY    : int  0 0 0 0 0 0 0 0 0 0 ...
## $ LATE_AIRCRAFT_DELAY: int  NA 0 NA NA NA NA NA NA NA ...
## $ WEATHER_DELAY       : int  NA NA NA NA NA NA NA NA NA ...

```

Com podem veure tenim 28 columnnes i un total de 145.592 dades en el dataset obtingut. També podem veure quines són les dades que conté el nostre dataset. Com que volem veure els retards o problemes que poden haver-hi, les causes i en quines aerolínies només cal que seleccionem aquelles columnnes que ens proporcionin aquesta informació. En aquest cas serien:

```

col_interest = c(
  "DAY_OF_WEEK",
  "AIRLINE",
  "DEPARTURE_DELAY",
  "ARRIVAL_DELAY",
  "DISTANCE"
)
print(col_interest)

## [1] "DAY_OF_WEEK"      "AIRLINE"           "DEPARTURE_DELAY" "ARRIVAL_DELAY"
## [5] "DISTANCE"

```

Amb aquestes dades ja podem fer un anàlisi complet per tal de donar resposta a les preguntes proposades.

```
flights <- subset(flights, select=col_interest)
```

3. Neteja de les dades. (2) (MR)

3.1. Les dades contenen zeros o elements buits? Com gestionaries aquests casos? (1)

Amb R és fàcil veure si tenim valors buits (NA) dins les nostres dades:

```
colSums(is.na(flights))
```

	DAY_OF_WEEK	AIRLINE	DEPARTURE_DELAY	ARRIVAL_DELAY	DISTANCE
##	0	0	0	461	0

Veiem alguns valors buits en la variable ARRIVAL_DELAY. Hi ha varíes estratègies per a resoldre problemes amb els elements buits, una tècnica eficaç és aplicar la funció kNN amb la qual omplirem els buits fent servir informació de k veïns més propers. Aquesta opció escollida es basa en que les variables del nostre dataset guarden certa relació i no són completament independents. Tindrem així uns valors aproximats als esperats que és millor que tenir-ne de buits.

```
# kNN on ARRIVAL_DELAY
suppressWarnings(suppressMessages(library(VIM)))
flights = kNN(flights, variable="ARRIVAL_DELAY", k=5)
```

```
# CHECK NA values in dataset.
colSums(is.na(flights))

##      DAY_OF_WEEK          AIRLINE  DEPARTURE_DELAY  ARRIVAL_DELAY
##            0                  0                  0                  0
##      DISTANCE ARRIVAL_DELAY_imp
##            0                  0
```

Veiem que no tenim valors buits en les variables conflictives.

3.2. Identificació i tractament de valors extrems. (1)

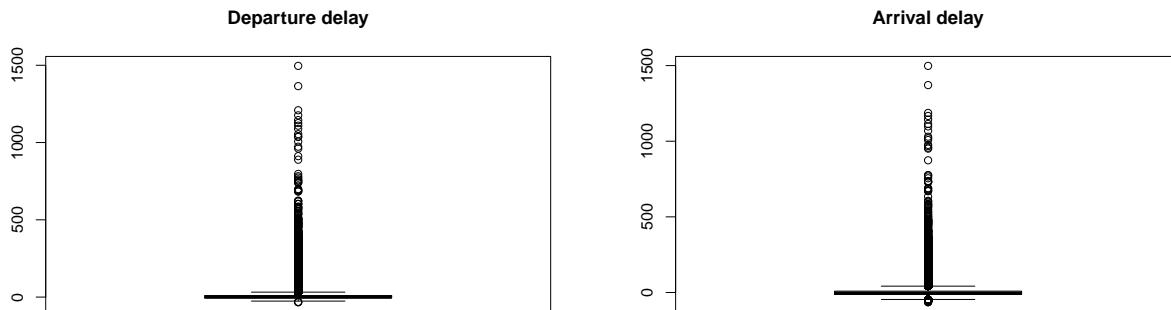
Començarem amb estadística descriptiva bàsica fent ús del summary que ens proporciona R. Aquí podem veure els valors màxims i mínims de cada variable.

```
summary(flights)
```

```
##      DAY_OF_WEEK          AIRLINE  DEPARTURE_DELAY  ARRIVAL_DELAY
##  Min.   :1.000  Length:145952  Min.   :-34.00  Min.   :-65.000
##  1st Qu.:2.000  Class  :character  1st Qu.: -5.00  1st Qu.: -13.000
##  Median :4.000  Mode   :character  Median : -1.00  Median : -4.000
##  Mean   :3.925                               Mean   : 11.19  Mean   : 5.785
##  3rd Qu.:6.000                               3rd Qu.: 10.00  3rd Qu.:  9.000
##  Max.   :7.000                               Max.   :1496.00  Max.   :1498.000
##      DISTANCE  ARRIVAL_DELAY_imp
##  Min.   : 77  Mode  :logical
##  1st Qu.: 414 FALSE:145491
##  Median : 679 TRUE :461
##  Mean   :1201
##  3rd Qu.:2139
##  Max.   :2704
```

D'aquí veiem alguns casos interessants, volem veure els que tenen mínims i màxims que s'allunyen clarament dels quartils (1r i 3r). Amb un boxplot podrem veure quants valors són extrems dins d'aquestes variables.

```
boxplot(flights$DEPARTURE_DELAY, main="Departure delay")
boxplot(flights$ARRIVAL_DELAY, main="Arrival delay")
```



Amb un gràfic de caixes podem veure clarament si tenim *outliers* o valors extrems a les dades. R representa els valors extrems com a cercles més enllà del rang interquartil. Aquests valors són normals ja que pot ser que els vols hagin tingut gran retards. La raó per la qual quasi no podem veure la caixa (a prop de zero) és perquè la gran majoria de vols no tenen retard i els outliers coincideixen amb els vols que en tenen.

4. Anàlisi de les dades. (2.5) (MR)

4.1. Selecció dels grups de dades que es volen analitzar/comparar (planificació dels anàlisis a aplicar). (0.5)

- Com volem analitzar el retard, considerarem que els vols han tingut retard si la suma del retard de sortida més el d'arribada supera els 10 minuts. En cas que el retard hagi sigut de l'avió mirarem que sigui per sobre de 15 minuts.

```
flights$TOTAL_DELAY = flights$DEPARTURE_DELAY + flights$ARRIVAL_DELAY
flights <- within(flights, {
  DELAYED <- NA
  DELAYED[TOTAL_DELAY > 10] <- 1
  DELAYED[TOTAL_DELAY <= 10] <- 0
})
})
```

- Analitzarem els vols també per aerolínia i per dia de la setmana. Categoritzem les variables.

```
flights$AIRLINE = factor(flights$AIRLINE)
flights$DAY_OF_WEEK_FAC = factor(flights$DAY_OF_WEEK)
levels = c("Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun")
levels(flights$DAY_OF_WEEK_FAC) <- levels

flights$DELAYED = factor(flights$DELAYED)
levels = c("No", "Yes")
levels(flights$DELAYED) <- levels
```

- Separarem en grups de vols llargs (més de 500 mi) de vols curts (menys de 500 mi).

```
long_index = flights$DISTANCE > 500
flights.long = flights[long_index,]
flights.short = flights[-long_index,]
```

4.2. Comprovació de la normalitat i homogeneïtat de la variància. (1)

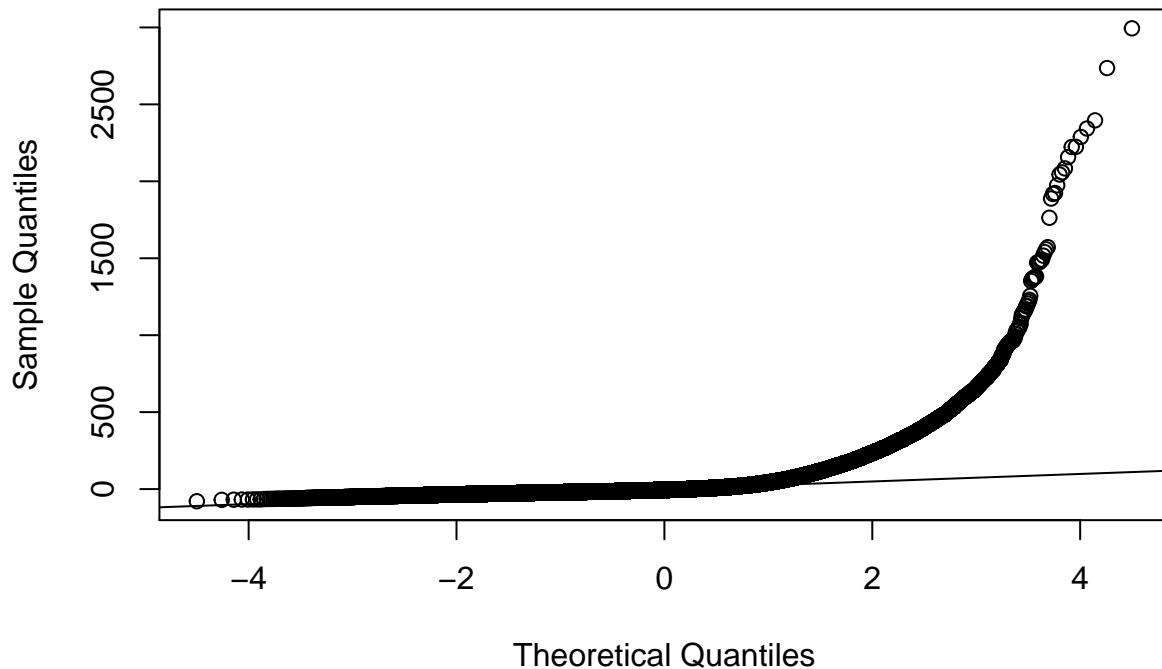
Per fer la comprovació de la normalitat de les dades un dels mètodes més habituals és fer servir la funció de Shapiro-Wilk. Mirem en les variables numèriques:

```
shapiro.test(flights[c(1:5000), "TOTAL_DELAY"])

##
##  Shapiro-Wilk normality test
##
## data: flights[c(1:5000), "TOTAL_DELAY"]
## W = 0.96812, p-value < 2.2e-16

qqnorm(flights$TOTAL_DELAY)
qqline(flights$TOTAL_DELAY)
```

Normal Q-Q Plot



Aquest test és més restrictiu que el test de Kolmogorov-Smirnov. Només podem presentar 5000 mostres en el test. Com el p valor que resulta és inferior a 0.05 es rebutja l'hipòtesi nul · la i considera que la distribució no és normal. El gràfic QQ-plot ens indica que les dades tenen una distribució normal quan els punts segueixen la línia dibuixada. En aquest cas, veiem que divergeixen pels valors a partir de 2 i per tant no podem assumir normalitat.

Veiem el test de l'homoscedasticitat per asegurar igualtat de variàncies. Aplicarem el test de Fligner-Killeen, que es tracta de l'alternativa no paramètrica, utilitzada quan les dades no compleixen amb la condició de normalitat.

```
fligner.test(TOTAL_DELAY~DISTANCE, data = flights)

##
##  Fligner-Killeen test of homogeneity of variances
##
## data:  TOTAL_DELAY by DISTANCE
## Fligner-Killeen:med chi-squared = 1650.5, df = 80, p-value < 2.2e-16

fligner.test(TOTAL_DELAY~DAY_OF_WEEK, data = flights)

##
##  Fligner-Killeen test of homogeneity of variances
##
## data:  TOTAL_DELAY by DAY_OF_WEEK
## Fligner-Killeen:med chi-squared = 1212.3, df = 6, p-value < 2.2e-16
```

Atès que les proves resulten en un p-valor inferior al nivell de significació (< 0.05), es rebutja la hipòtesi nul · la d'homoscedasticitat i es conclou que les variables presenta variàncies estadísticament diferents per als diferents grups de spray.

4.3. Aplicació de proves estadístiques per comparar els grups de dades.

En funció de les dades i de l'objectiu de l'estudi, aplicar proves de contrast d'hipòtesis, correlacions, regressions, etc. Aplicar almenys tres mètodes d'anàlisi diferents. (1)

Comprovarem si podem entendre el retard de sortida es pot entendre a partir del retard d'arribada i amb el de l'avió a través d'una regressió lineal. Podem afegir la distància i el dia de la setmana per millorar el model?

```
# Agafem el valor de referencia de dia de la setmana dilluns
day_monday = relevel(factor(flights$DAY_OF_WEEK), ref = 1)

# Agafem com a aeroport de referencia el AA.
airline_aa=relevel(factor(flights$AIRLINE), ref = 'AA')

model1 <- lm(DEPARTURE_DELAY ~ ARRIVAL_DELAY, data=flights)
model2 <- lm(DEPARTURE_DELAY ~ ARRIVAL_DELAY + DISTANCE, data=flights)
model3 <- lm(DEPARTURE_DELAY ~ ARRIVAL_DELAY + day_monday, data=flights)
model4 <- lm(DEPARTURE_DELAY ~ ARRIVAL_DELAY + airline_aa, data=flights)
model5 <- lm(DEPARTURE_DELAY ~ ARRIVAL_DELAY + airline_aa + DISTANCE + day_monday, data=flights)

tabla.coeficientes <-matrix(c(1,summary(model1)$r.squared,2,summary(model2)$r.squared,3,summary(model3)$r.squared), nrow=1, byrow=TRUE)

colnames(tabla.coeficientes) <-c("Model", "R^2")

# Veiem tots els R^2 per cada model:
print(tabla.coeficientes)

##      Model      R^2
## [1,]    1 0.9083035
## [2,]    2 0.9128394
## [3,]    3 0.9083851
## [4,]    4 0.9137738
## [5,]    5 0.9161822

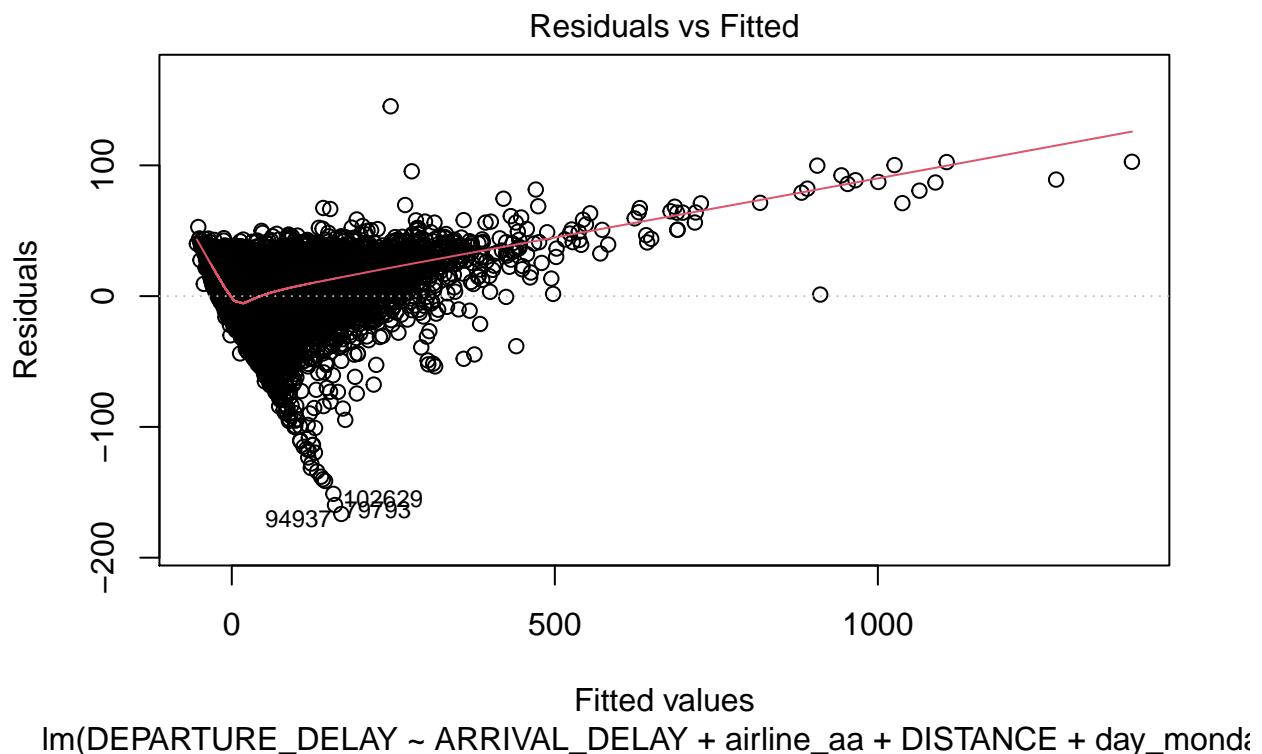
summary(model5)

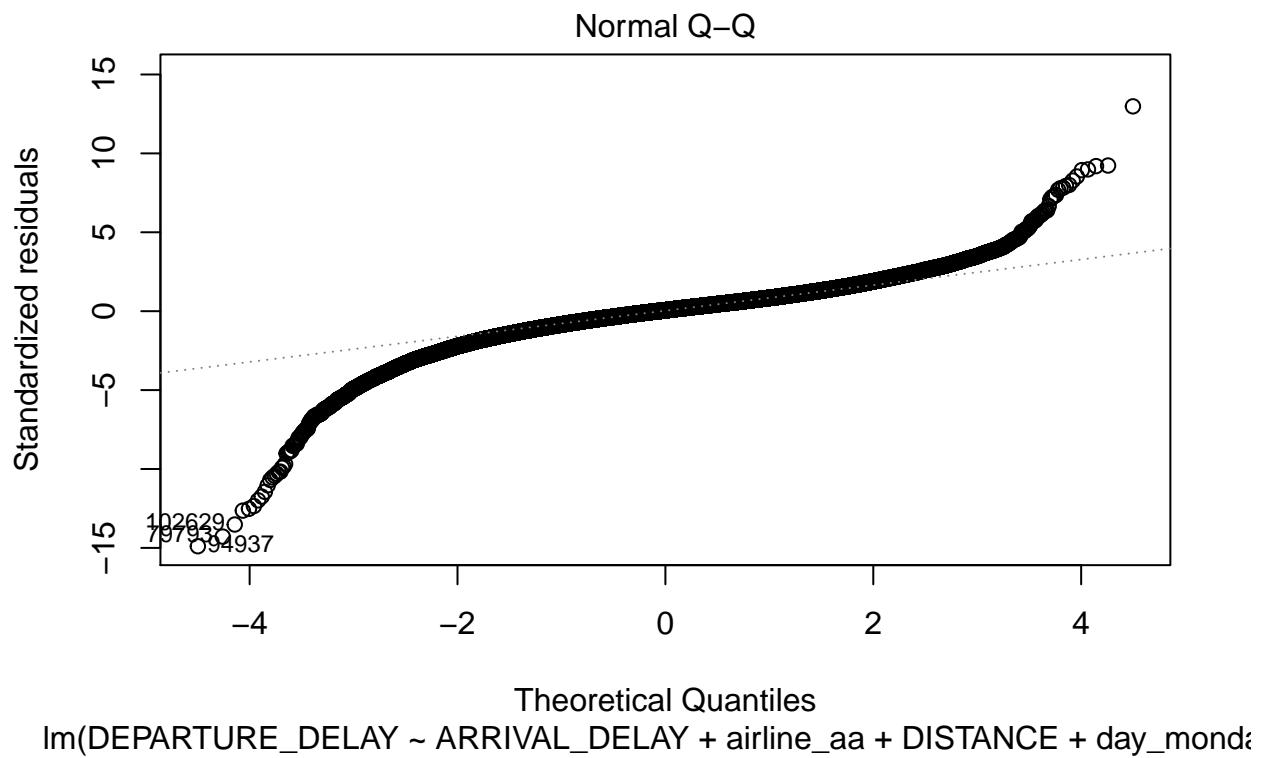
##
## Call:
## lm(formula = DEPARTURE_DELAY ~ ARRIVAL_DELAY + airline_aa + DISTANCE +
##     day_monday, data = flights)
##
## Residuals:
##       Min     1Q   Median     3Q    Max 
## -166.577 -5.776   0.683   6.467 145.158 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 1.625e+00 1.403e-01 11.584 < 2e-16 ***
## ARRIVAL_DELAY 9.265e-01 7.375e-04 1256.386 < 2e-16 ***
## airline_aaAS -9.813e-01 1.909e-01 -5.140 2.75e-07 ***
## airline_aaB6  1.327e+00 1.898e-01  6.993 2.69e-12 ***
## airline_aaDL  9.224e-01 1.528e-01  6.035 1.59e-09 ***
## airline_aaF9 -5.432e+00 2.745e-01 -19.789 < 2e-16 ***
## airline_aaHA -1.258e+01 4.465e-01 -28.178 < 2e-16 ***
## airline_aaOO -3.142e-01 1.273e-01 -2.468 0.013572 *
```

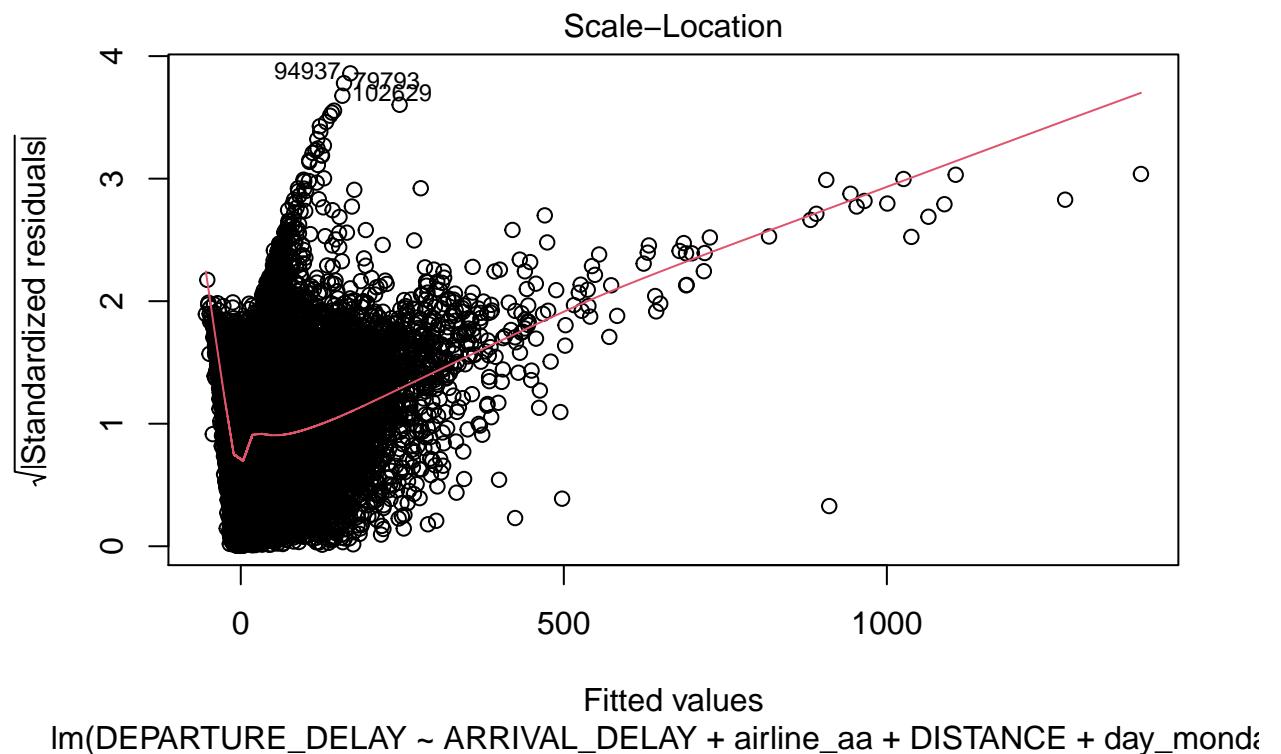
```

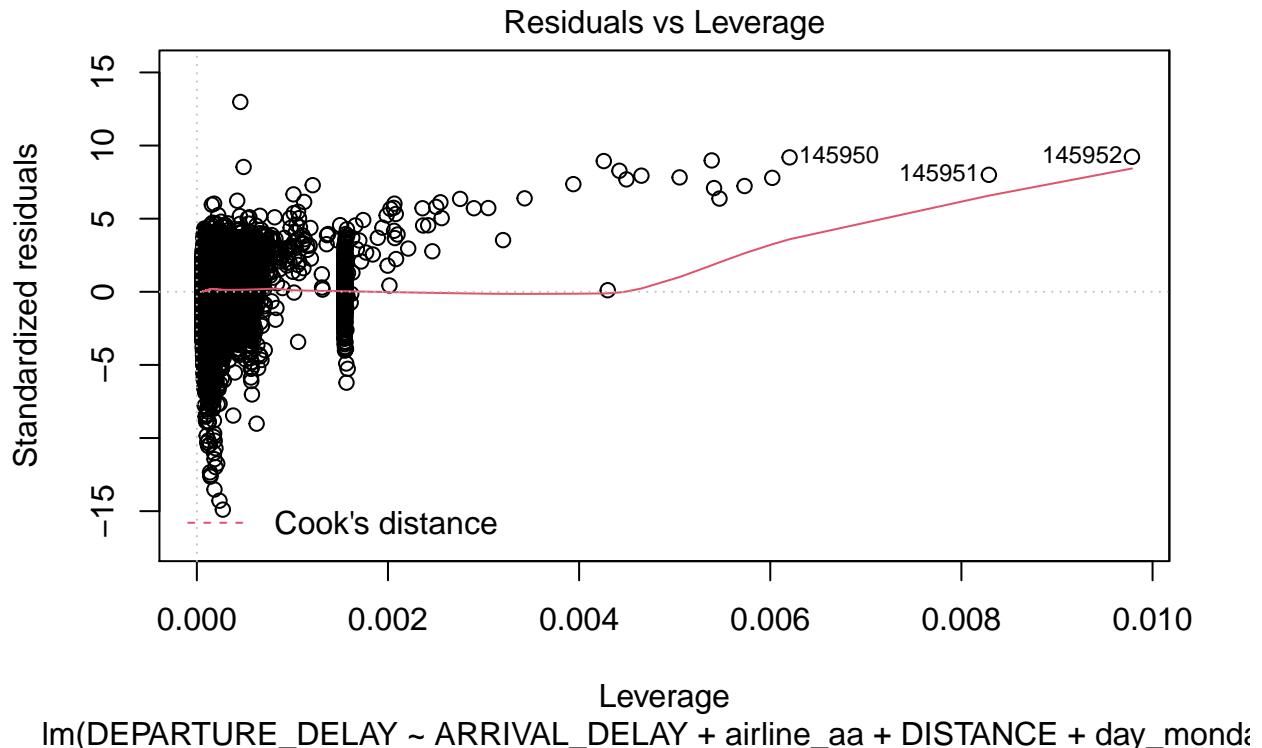
## airline_aaUA 3.320e+00 1.147e-01 28.948 < 2e-16 ***
## airline_aaUS -3.436e+00 2.422e-01 -14.187 < 2e-16 ***
## airline_aaVX -1.400e+00 1.362e-01 -10.279 < 2e-16 ***
## airline_aaWN 2.819e+00 1.446e-01 19.491 < 2e-16 ***
## DISTANCE 2.529e-03 4.007e-05 63.132 < 2e-16 ***
## day_monday2 4.695e-02 1.079e-01 0.435 0.663547
## day_monday3 3.003e-01 1.075e-01 2.793 0.005217 **
## day_monday4 -3.079e-01 1.073e-01 -2.868 0.004129 **
## day_monday5 8.740e-02 1.077e-01 0.812 0.417043
## day_monday6 1.180e+00 1.138e-01 10.372 < 2e-16 ***
## day_monday7 4.184e-01 1.087e-01 3.849 0.000119 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.18 on 145933 degrees of freedom
## Multiple R-squared: 0.9162, Adjusted R-squared: 0.9162
## F-statistic: 8.862e+04 on 18 and 145933 DF, p-value: < 2.2e-16
plot(model5)

```







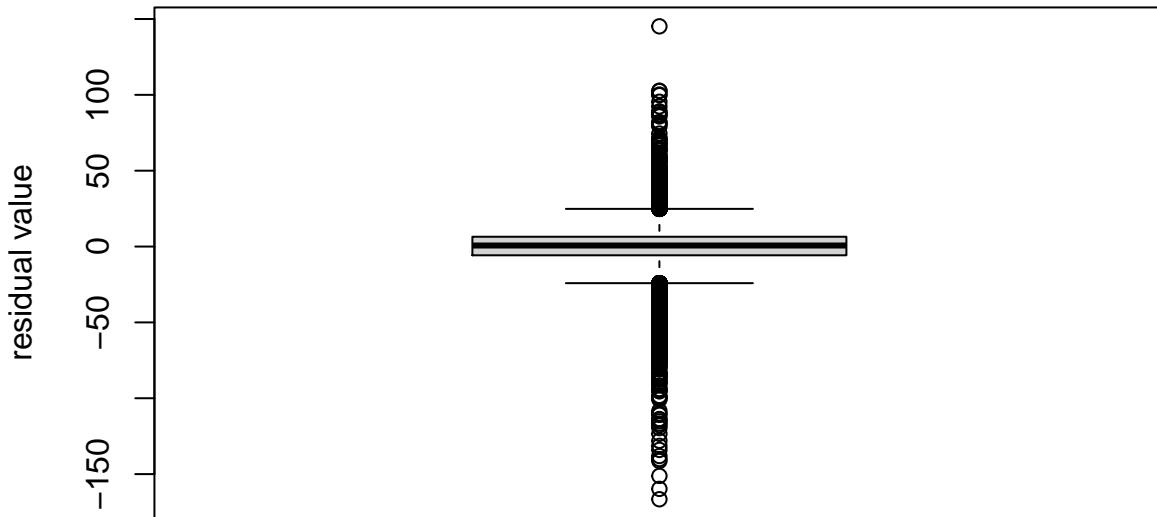


Obtenim una R-squared superior a 0.9 per a tots els models. El model 5 sembla que té una millor capacitat descriptiva, però incloure tots els paràmetres tampoc es tradueix una millora significativa. Amb la variable ARRIVAL_DELAY ja tenim un model prou robust.

Amb un valor p inferior al valor de significació podem dir que es pot explicar la variable DEPARTURE_DELAY amb ARRIVAL_DELAY. Aquesta correlació ens indica que sabent el retard del vol d'arribada podriem predir el retard del següent vol de sortida. En la gràfica de QQ-plot dels residuals veiem que segueixen una distribució aproximadament normal. EN un boxplot hauríem de veure la mediana al voltant de zero.

```
boxplot(model5[['residuals']], main='Boxplot: Residuals', ylab='residual value')
```

Boxplot: Residuals



Amb un model de regressió logística veurem si és possible obtenir el retard a partir del dia de la setmana i de la distància. Preveurem només el retard de sortida:

```

flights <- within(flights, {
  delay_SFO <- NA
  delay_SFO[DEPARTURE_DELAY >= 15] <- 1
  delay_SFO[DEPARTURE_DELAY < 15] <- 0
  })

flights$delay_SFO <- factor(flights$delay_SFO, levels = c(0, 1))
str(flights$delay_SFO)

##  Factor w/ 2 levels "0","1": 1 1 1 1 1 1 1 1 1 1 ...
model5 <- glm(flights$delay_SFO ~ day_monday + flights$DISTANCE, family="binomial")
summary(model5)

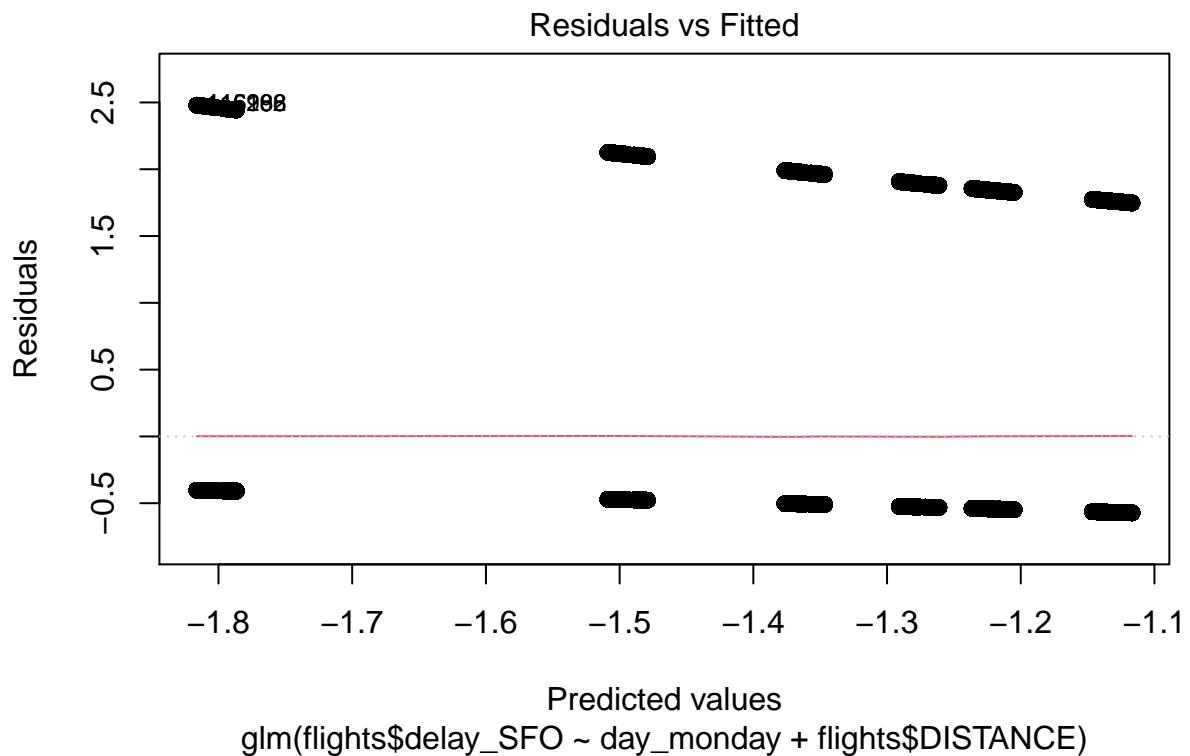
##
## Call:
## glm(formula = flights$delay_SFO ~ day_monday + flights$DISTANCE,
##      family = "binomial")
##
## Deviance Residuals:
##       Min        1Q     Median        3Q       Max
## -0.7525  -0.7185  -0.6782  -0.5526   1.9834
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)

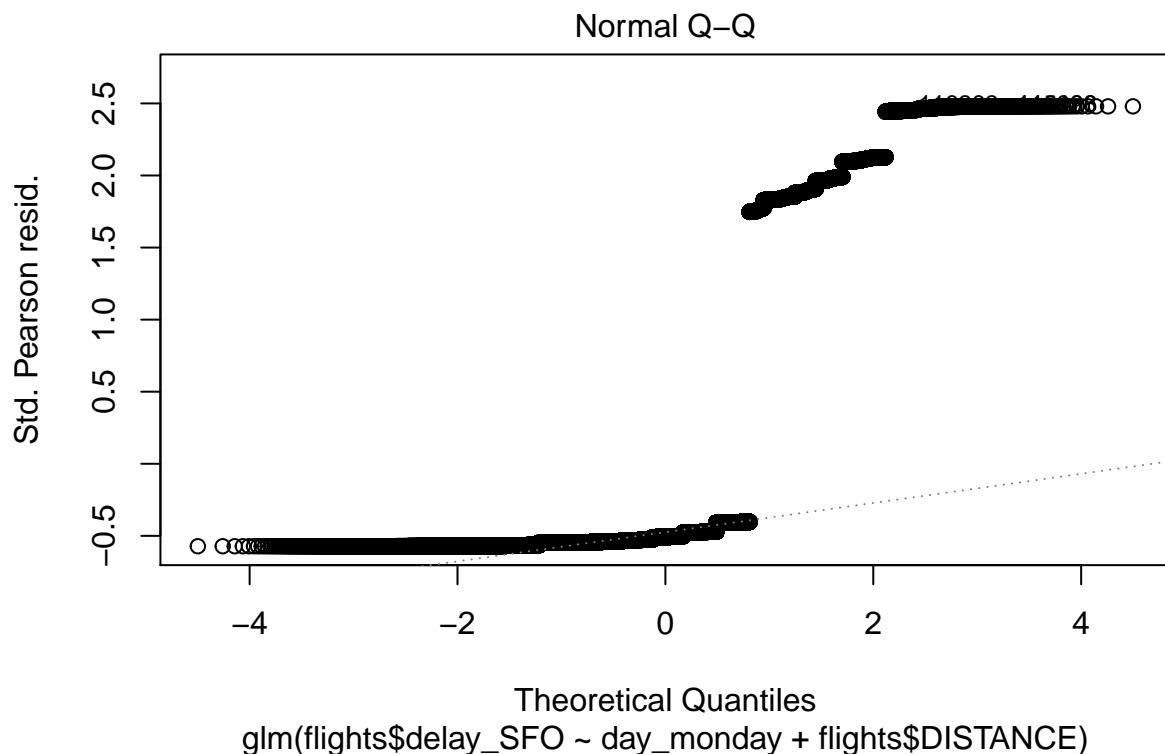
```

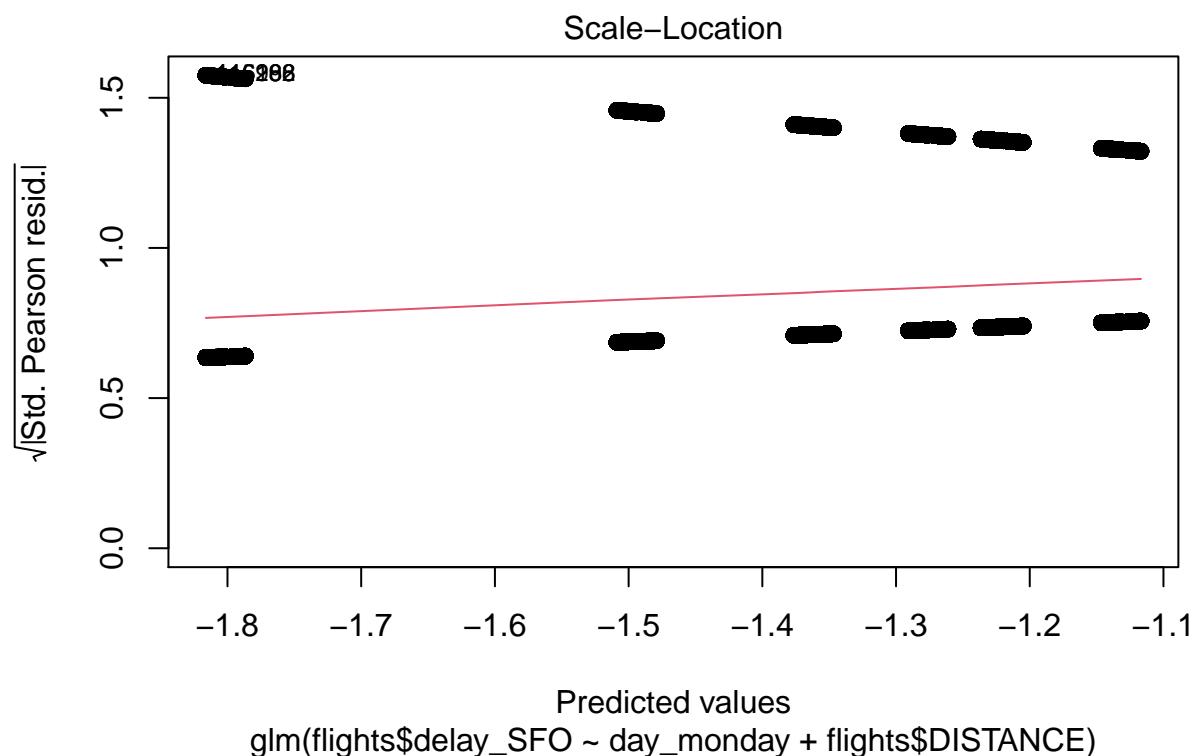
```

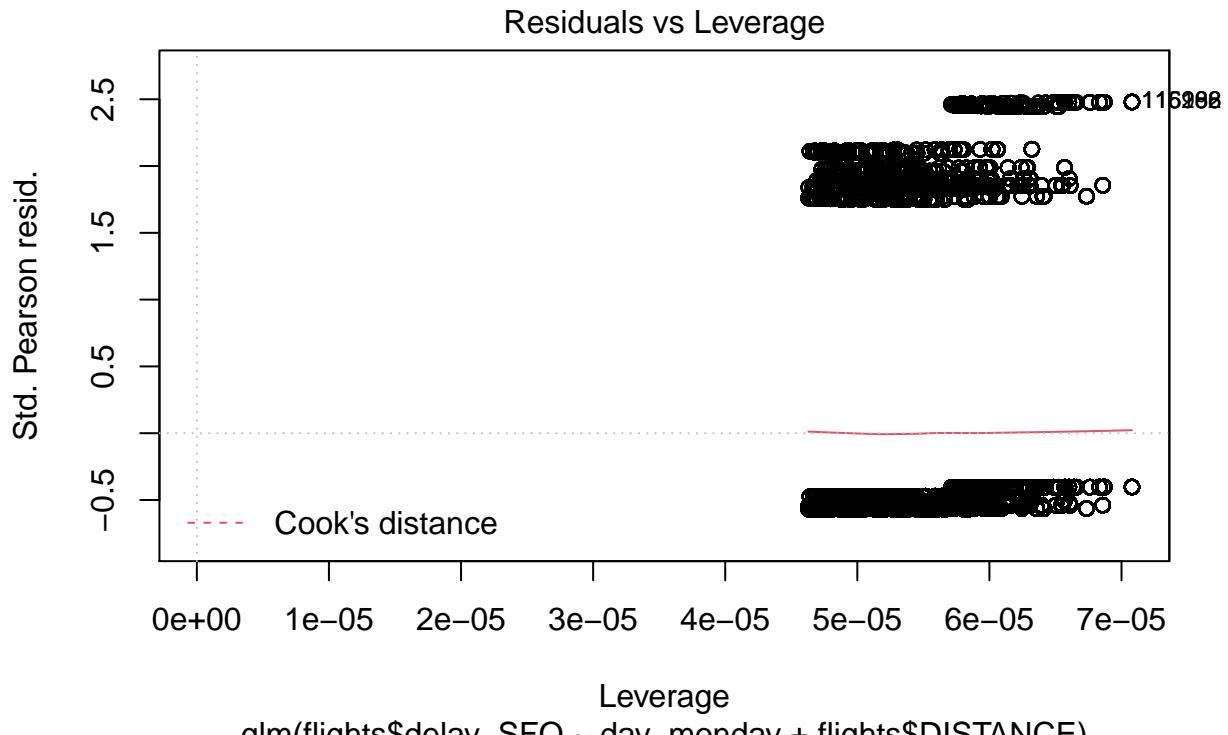
## (Intercept) -1.116e+00 1.798e-02 -62.083 < 2e-16 ***
## day_monday2 -2.300e-01 2.321e-02 -9.911 < 2e-16 ***
## day_monday3 -3.623e-01 2.361e-02 -15.348 < 2e-16 ***
## day_monday4 -8.811e-02 2.260e-02 -3.899 9.66e-05 ***
## day_monday5 -1.444e-01 2.286e-02 -6.317 2.67e-10 ***
## day_monday6 -6.698e-01 2.677e-02 -25.021 < 2e-16 ***
## day_monday7 -8.999e-02 2.290e-02 -3.929 8.52e-05 ***
## flights$DISTANCE -1.124e-05 7.213e-06 -1.559 0.119
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 149532 on 145951 degrees of freedom
## Residual deviance: 148655 on 145944 degrees of freedom
## AIC: 148671
##
## Number of Fisher Scoring iterations: 4
plot(model5)

```









Podem fer un segon test i respondre a la pregunta, “els vols llargs tenen més retards que els curts?”, o dit d’altra manera, “la mitjana dels retards dels vols llargs és major a la dels vols curts?”. Així plantegem el següent test de contrast d’hipòtesis de dos mostres sobre la diferència de mitjanes, amb valor de significació $\alpha = 0.05$. La hipòtesi nul · la és que el retard dels vols curts és el mateix que el dels vols llargs. La hipòtesi alternativa, és que existeix una diferència (bilateral).

```
flights.short.late <- flights.short$TOTAL_DELAY
flights.long.late <- flights.long$TOTAL_DELAY
t.test(flights.short.late, flights.long.late, alternative = "less")
```

```
##
##  Welch Two Sample t-test
##
## data:  flights.short.late and flights.long.late
## t = 1.7347, df = 186827, p-value = 0.9586
## alternative hypothesis: true difference in means is less than 0
## 95 percent confidence interval:
##      -Inf 1.142494
## sample estimates:
## mean of x mean of y
## 16.97594 16.38950
```

Amb el valor de p superior a 0.05 no podem rebutjar la hipòtesis nul · la i per tant no podem concloure que siguin diferents.

Amb un anàlisis de correlació podem investigar quina de les variables influeix més en el retard dels vols.

```

library(dplyr)

## Warning: package 'dplyr' was built under R version 4.0.3
##
## 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

columns_to_compare = c("DISTANCE", "DAY_OF_WEEK")
corr_matrix <-matrix(nc = 2, nr = 0)
colnames(corr_matrix) <-c("estimate", "p-value")

# Calculem el coeficient de correlació per a cada camp respecte la variable quantitativa TOTAL_DELAY
for (i in columns_to_compare) {
  spearman_test =cor.test(flights[,i], flights$TOTAL_DELAY, method = "spearman", exact = FALSE )
  corr_coef = spearman_test$estimate
  p_val = spearman_test$p.value
  # Add row to matrix
  pair =matrix(ncol = 2, nrow = 1)
  pair[1][1] = corr_coef
  pair[2][1] = p_val
  corr_matrix <-rbind(corr_matrix, pair)
  rownames(corr_matrix)[nrow(corr_matrix)] <-i
}
print(corr_matrix)

##           estimate      p-value
## DISTANCE    -0.03683230 5.338833e-45
## DAY_OF_WEEK -0.02937156 3.132281e-29

```

La correlació és molt baixa en qualsevol cas i no veiem cap relació directa entre les variables explicatives i el retard.

5. Representació dels resultats a partir de taules i gràfiques. (2) (MF)

Un cop ja tenim les dades netes farem un anàlisis descriptiu per tal de veure com són les dades i els seus valors. En primer lloc mostrem un resum dels valors que conté la nostra taula:

```
summary(flights, digits = 1)
```

	DAY_OF_WEEK	AIRLINE	DEPARTURE_DELAY	ARRIVAL_DELAY	DISTANCE
## Min.	:1	UA :45135	Min. : -34	Min. : -65	Min. : 77
## 1st Qu.	:2	OO :34223	1st Qu.: -5	1st Qu.: -13	1st Qu.: 414
## Median	:4	VX :15848	Median : -1	Median : -4	Median : 679
## Mean	:4	WN :13790	Mean : 11	Mean : 6	Mean : 1201
## 3rd Qu.	:6	AA :12062	3rd Qu.: 10	3rd Qu.: 9	3rd Qu.: 2139
## Max.	:7	DL : 9659	Max. :1496	Max. :1498	Max. :2704

```

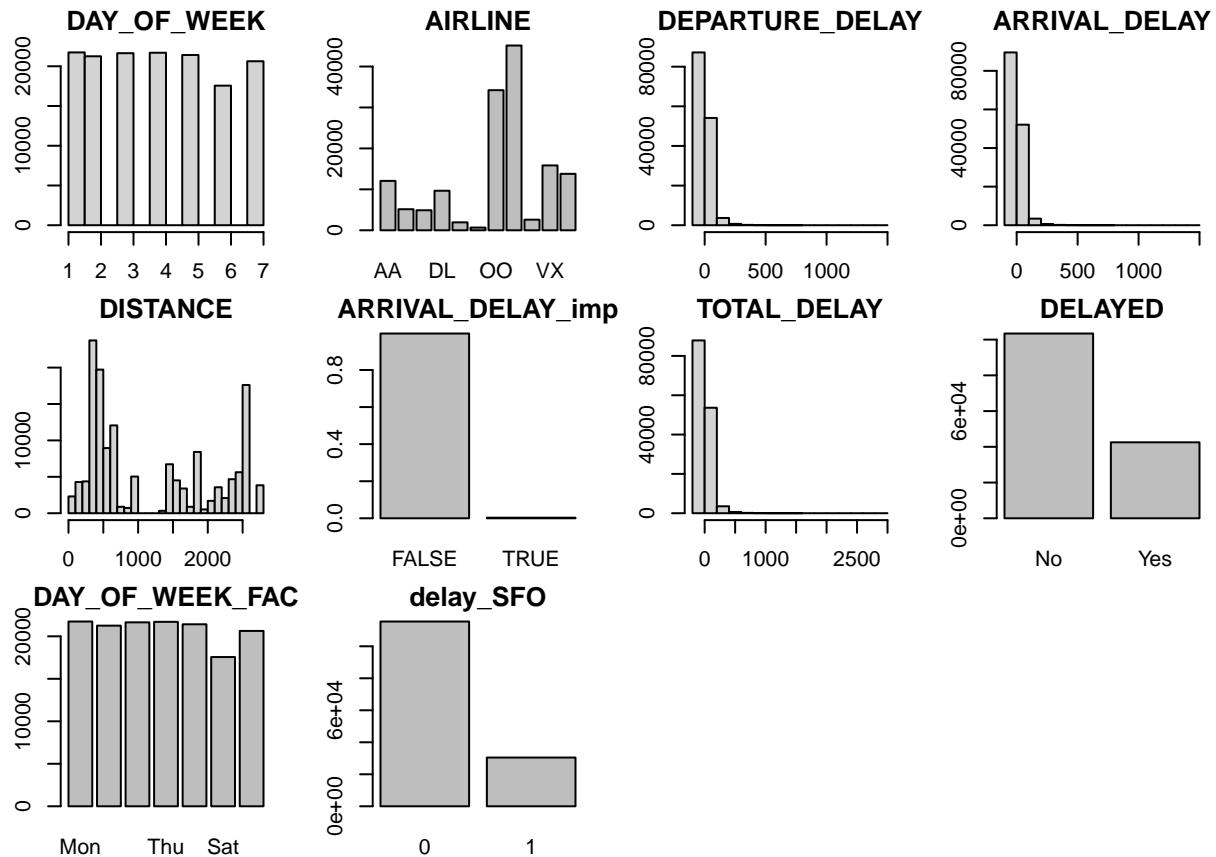
##                               (Other) :15235
## ARRIVAL_DELAY_imp  TOTAL_DELAY    DELAYED      DAY_OF_WEEK_FAC delay_SFO
## Mode :logical      Min.   : -79  No :103425  Mon:21742       0:115488
## FALSE:145491       1st Qu.: -16 Yes: 42527  Tue:21254       1: 30464
## TRUE :461          Median : -5
##                      Mean   : 17
##                      3rd Qu.: 17
##                      Max.   :2994
##
```

Com podem observar en la taula generada, tenim les diferents variables de la taula i com es troben distribuïdes i els seus valors més importants. A continuació representarem gràficament totes les variables en histogrames per tal de visualitzar molt millor aquests valors.

```

par(mfrow=c(3,4), mar=c(2,2,2,2))
hist(flights$DAY_OF_WEEK, main = "DAY_OF_WEEK")
barplot(summary(flights$AIRLINE), main = "AIRLINE")
hist(flights$DEPARTURE_DELAY, main = "DEPARTURE_DELAY")
hist(flights$ARRIVAL_DELAY, main = "ARRIVAL_DELAY")
hist(flights$DISTANCE, main = "DISTANCE")
barplot(prop.table(table(flights$ARRIVAL_DELAY_imp)), main = "ARRIVAL_DELAY_imp")
hist(flights$TOTAL_DELAY, main = "TOTAL_DELAY")
barplot(summary(flights$DELAYED), main = "DELAYED")
barplot(summary(flights$DAY_OF_WEEK_FAC), main = "DAY_OF_WEEK_FAC")
barplot(summary(flights$delay_SFO), main = "delay_SFO")

```



Un cop ja tenim un resum de les variables, ens centrarem en les variables més importants i les quals ens

ajudaran a respondre les preguntes més endavant. Les variables són les aerolínies i els dies de la setmana. D'aquesta forma, a continuació mostrarem les taules d'aquestes variables segons els retards.

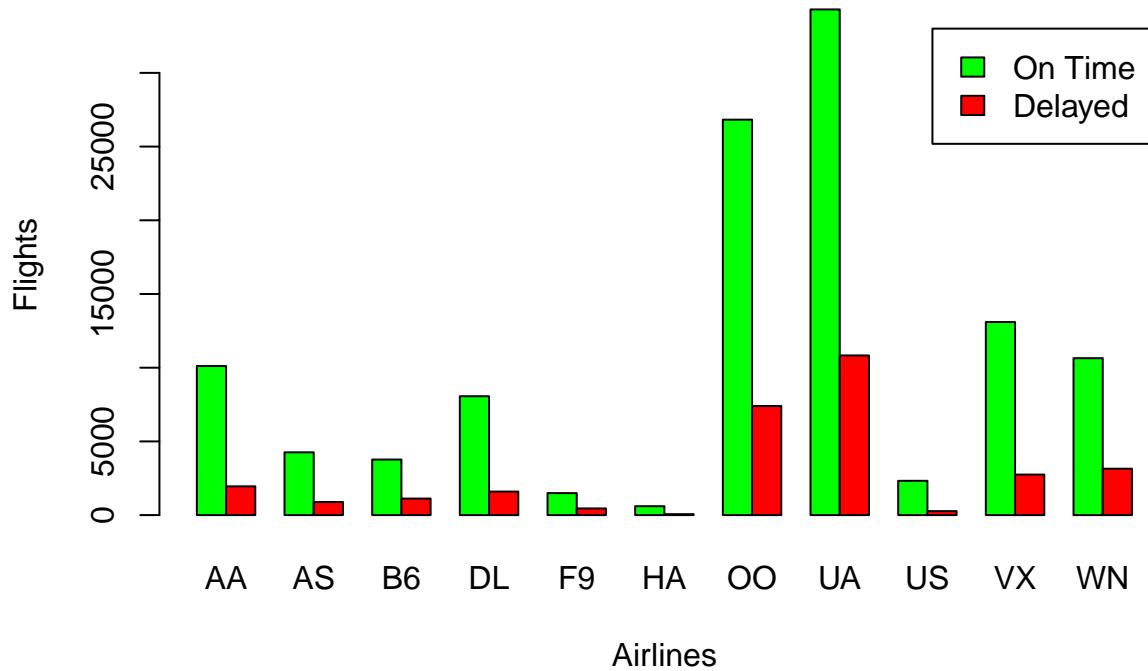
```
# In absolute values


```

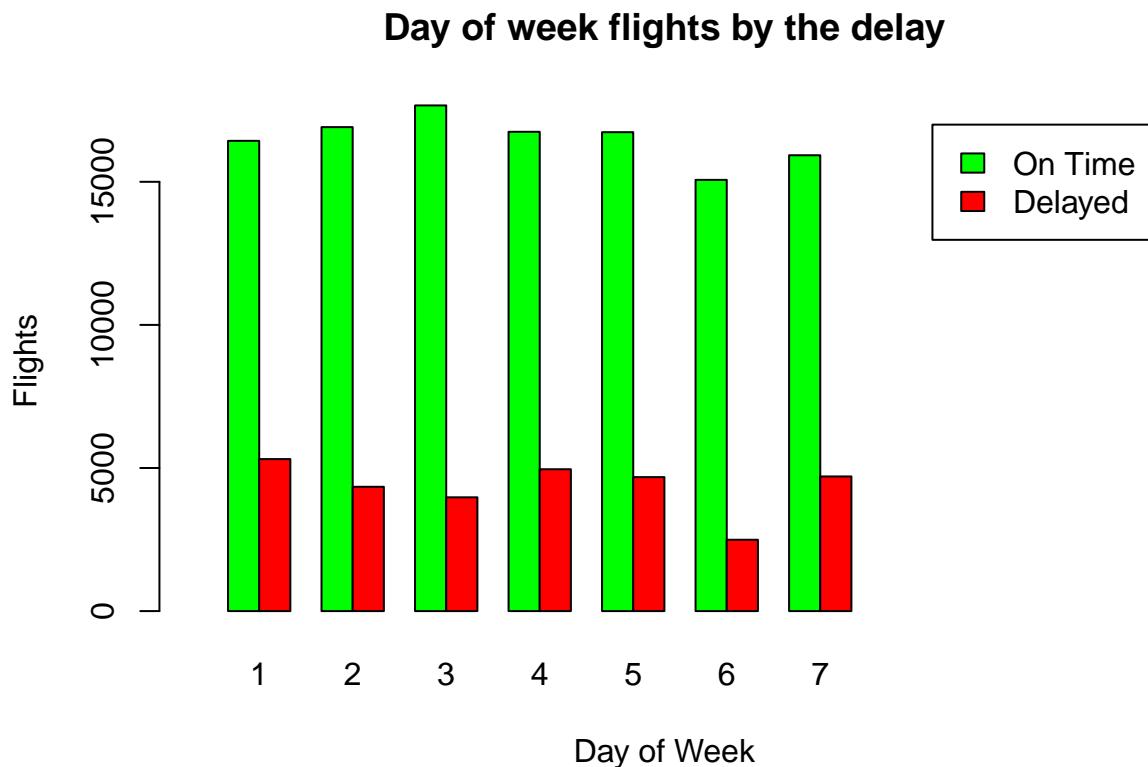
Com podem veure les dues taules ens mostren per a cada una de les variables el nombre de vols que han tingut retard i el nombre de vols que no n'han tingut. Partint d'aquestes taules que hem generat, representarem els valors utilitzant uns gràfics de barres de forma que podem veure més clarament la distribució de tots els valors.

```
table_airlines <- table(flights$delay_SFO, flights$AIRLINE)
table_dayweek <- table(flights$delay_SFO, flights$DAY_OF_WEEK)
barplot(table_airlines, main="Airlines flights by the delay", xlab = "Airlines", ylab = "Flights", col =
```

Airlines flights by the delay



```
barplot(table_dayweek, main="Day of week flights by the delay", xlab = "Day of Week", ylab = "Flights",
```



Com podem veure els gràfics realitzats ens mostren per cada una de les variables que volem el nombre de vols que no han tingut retards, de color verd, i el nombre de vols que han tingut retard, de color vermell, d'aquesta forma podem veure d'una manera molt ràpida i clara com estan distribuïts els valors en les variables escollides.

6. Resolució del problema.

A partir dels resultats obtinguts, quines són les conclusions? Els resultats permeten respondre al problema? (0.5) (MR)

En aquest treball s'han realitzat tres proves estadístiques sobre un conjunt de dades recollides en l'aeroport de San Francisco a l'any 2015. S'ha treballat amb les variables corresponents als retards dels vols i amb els factors que potencialment hi estan relacionats per a solucionar el problema i contestar les preguntes que es plantegen al principi de la pràctica. Al final hem representat els resultats en taules i en gràfics il·lustratius per veure quina informació en podem extreure d'ells.

Després d'analitzar els resultats no podem conculoure que els retards dels vols tinguin a veure amb el dia de la setmana o amb la aerolínia. Tampoc existeix una correlació amb la distància del vol ni influeix que sigui un vol de llarga distància o curta. Els models lineals de predicción de retard en el vol de sortida sabent el vol d'arribada han donat un R-squared per sobre de 0.9 que ens permeten construir un model lineal de predicción sobre els retards de sortida sabent el retard de l'arribada de l'anterior vol.

7. Contribucions al treball

Contribucions	Accuracy
Recerca prèvia	M.R./ M.F.
Redacció de les respostes	M.R./ M.F.
Desenvolupament codi	M.R./ M.F.

8. Codi: Cal adjuntar el codi, preferiblement en R, amb el que s'ha realitzat la neteja, anàlisi i representació de les dades. Si ho preferiu, també podeu treballar en Python. (2) (MF)