



TASCA 12 - IT Academy

Curs Data Science

```
In [137]: import pandas as pd
import numpy as np
from numpy import random
import datetime as pd
from datetime import datetime
import statistics as stat
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib.pyplot import figure
from pyod.models.knn import KNN
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
import sklearn
from sklearn.decomposition import PCA
import sklearn.neighbors
from sklearn.neighbors import neighbors_graph
from sklearn.datasets import make_blobs
from sklearn.model_selection import train_test_split

pd.options.display.float_format = '{:,.0f}'.format
```

Exercici 1.

Parteix el conjunt de dades `DelayedFlights.csv` en `train` i `test`. Estudia els dos conjunts per separat, a nivell descriptiu. [\[X\] Fet](#)

```
In [138]: file=pd.read_csv("DelayedFlights.csv", sep=",", encoding='utf8')

In [139]: file.head()

Out[139]:
```

Unnamed: 0	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	...	TaxiIn	TaxiOut	Canc
0	2008	1	3	4	2003	1955	2211	2225	WN	...	4	8	
1	2008	1	3	4	754	735	1002	1000	WN	...	5	10	
2	2008	1	3	4	628	620	804	750	WN	...	3	17	
3	2008	1	3	4	1829	1755	1959	1925	WN	...	3	10	
4	2008	1	3	4	1940	1915	2121	2110	WN	...	4	10	

5 rows × 30 columns

```
In [140]: file.describe().round(2)

Out[140]:
```

Unnamed: 0	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	FlightNum	...	Distance	Taxi
count	1936758	1936758	1936758	1936758	1936758	1936758	1936758	1929648	1936758	...	1936758	1929648
mean	3341651	2008	6	16	4	1519	1467	1610	1634	2184	...	766
std	2060605	0	3	9	2	450	425	548	465	1945	...	571
min	0	2008	1	1	1	1	0	1	0	1	...	11
25%	1517452	2008	3	8	2	1203	1135	1316	1325	610	...	338
50%	3242558	2008	6	16	4	1545	1510	1715	1705	1543	...	606
75%	4972467	2008	9	23	6	1900	1815	2030	2014	3422	...	998
max	7009727	2008	12	31	7	2400	2359	2400	2400	9742	...	4962

8 rows × 25 columns

```
In [141]: file.shape

Out[141]: (1936758, 30)

In [142]: file.columns

Out[142]: Index(['Unnamed: 0', 'Year', 'Month', 'DayOfMonth', 'DayOfWeek', 'DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime', 'UniqueCarrier', 'FlightNum', 'ActualElapsedTime', 'CRSElapsedTime', 'AirTime', 'ArrDelay', 'DepDelay', 'Origin', 'Dest', 'Distance', 'TaxiIn', 'TaxiOut', 'Cancelled', 'Diverted', 'CarrierDelay', 'NASDelay', 'WeatherDelay', 'LateAircraftDelay'], dtype='object')
```

A continuació, vemós que la columna 'Unnamed' no aporta nada, no es más que un contador de filas, así que la eliminamos.

```
In [143]: file['Unnamed: 0'].value_counts()

Out[143]:
```

Unnamed: 0	count
0	1
6846151	1
2232597	1
5967571	1
2242832	1
...	...
4495774	1
692169	1
3917269	1
3919316	1
2047	1

Name: Unnamed: 0, Length: 1936758, dtype: int64

```
In [144]: file.drop(['Unnamed: 0'], axis=1, inplace=True)

In [145]:
```

```
# Miramos los valores NaN que tienen las columnas
total_nan_values = file.isna().sum()
print ("Total Number of NaN values: "+'\n'+str(total_nan_values))

Total Number of NaN values:
Year                                0
Month                              0
DayOfMonth                         0
DayOfWeek                         0
DepTime                           0
CRSDepTime                        0
ArrTime                           0
CRSArrTime                        0
UniqueCarrier                     0
FlightNum                         0
TailNum                           5
ActualElapsedTime                 8387
CRSElapsedTime                   198
AirTime                          8387
ArrDelay                         8387
DepDelay                         0
Origin                           0
Dest                             0
Distance                         0
TaxiIn                           7110
TaxiOut                          455
Cancelled                         0
CancellationCode                 0
CarrierDelay                     689270
WeatherDelay                     689270
NASDelay                         689270
SecurityDelay                    689270
LateAircraftDelay                689270
dtype: int64
```

```
In [146]: # Vamos a tratar los valores NaN de dos maneras: las columnas que representa una hora concreta ('ArrTime') o el número de minutos, las vamos a sustituir los NaN's por la media aritmética.

file = file.dropna(subset=['ArrTime', 'TailNum'])

In [147]: # Sustituimos los datos NaN por la media aritmética

file=file.fillna(file.mean())

In [148]: file.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1929645 entries, 0 to 1936757
Data columns (total 23 columns):
 #   Column        Dtype
---  -
 0   Year          int64
 1   Month         int64
 2   DayOfMonth    int64
 3   DayOfWeek     int64
 4   DepTime       float64
 5   CRSDepTime    int64
 6   ArrTime       float64
 7   CRSArrTime    int64
 8   UniqueCarrier object
 9   FlightNum     int64
10  TailNum       object
11  ActualElapsedTime float64
12  CRSElapsedTime float64
13  AirTime       float64
14  ArrDelay      float64
15  DepDelay      float64
16  Origin        object
17  Dest          object
18  Distance      int64
19  TaxiIn        float64
20  TaxiOut       float64
21  Cancelled     int64
22  CancellationCode object
23  Diverted      int64
24  CarrierDelay  float64
25  WeatherDelay  float64
26  NASDelay      float64
27  SecurityDelay float64
28  LateAircraftDelay float64
dtypes: float64(14), int64(10), object(5)
memory usage: 441.7+ MB
```

```
In [149]: # Pasamos a GRAFICAR: a través de la matriz de correlación podemos ver cómo algunas de las 29 variables del conjunto de datos están relacionadas, es decir, pueden predecirse linealmente a partir de las demás.

In [150]: # Matriz de correlación #

corrmat = file.corr()
f, ax = plt.subplots(figsize=(12, 9))
sns.heatmap(corrmat, vmax=.8, square=True);
plt.show()
```



```
In [151]: # Scatterplot #

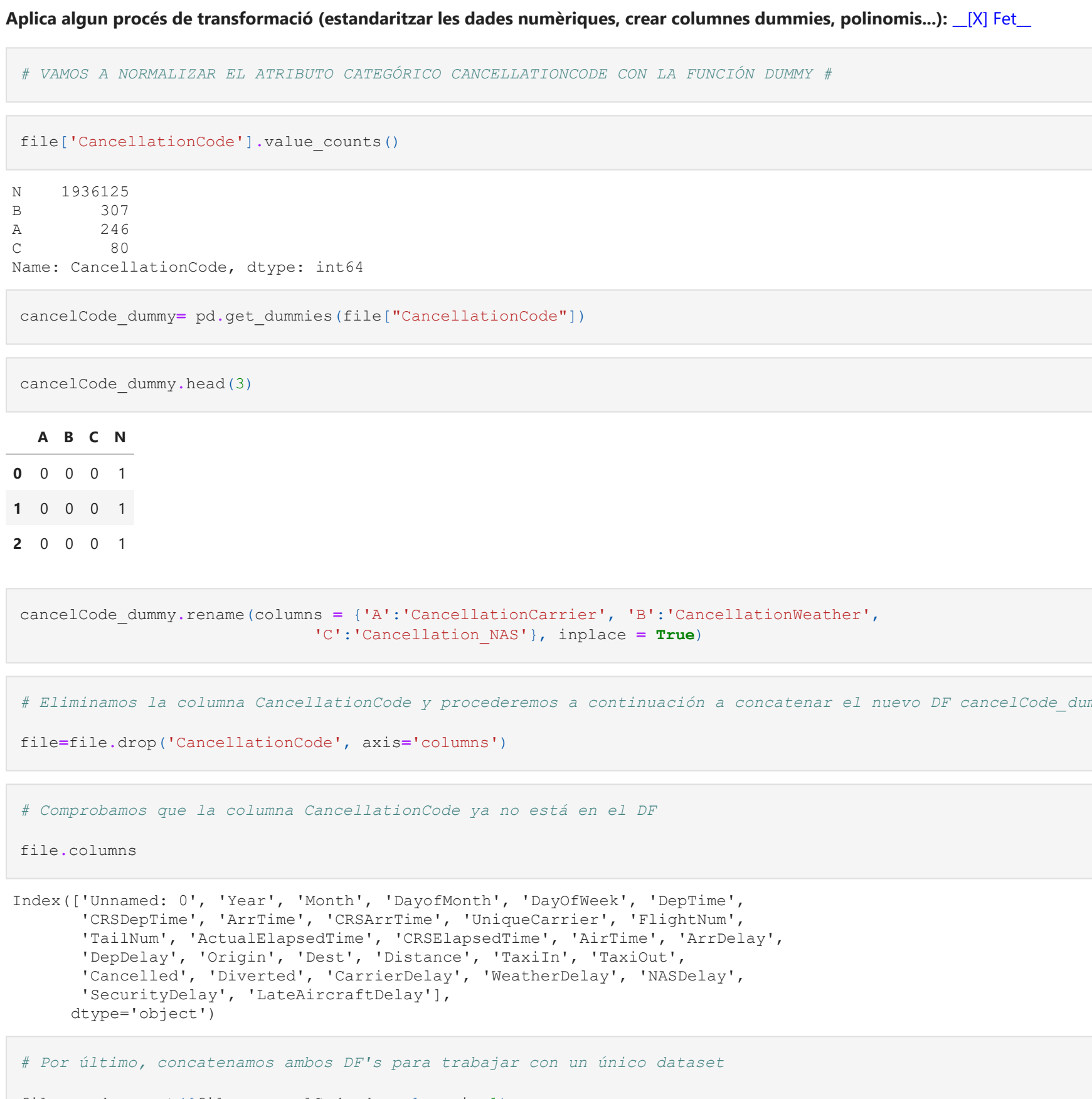
# La matriz de correlación anterior identifica las tres variables principales que generan los retrasos: Late / Carrier Delay and NAS Delay. El scatterplot corrobora este supuesto, mostrando cómo esas tres variables crean los retrasos durante el año.

# Antes de graficar, creamos la columna Status que representa si el vuelo llegó en hora (0), con un poco retrasado (1) o con mucho retrasado (2), desviado (3) o cancelado (4)

for dataset in file:
    file.loc[file['ArrDelay'] <= 15, 'Status'] = 0
    file.loc[file['ArrDelay'] >= 15, 'Status'] = 1
    file.loc[file['ArrDelay'] >= 60, 'Status'] = 2
    file.loc[file['Diverted'] == 1, 'Status'] = 3
    file.loc[file['Cancelled'] == 1, 'Status'] = 4

sns.set()
DelayedFlights = file[(file.Status >= 1) & (file.Status < 3)]
cols = ['ArrDelay', 'CarrierDelay', 'LateAircraftDelay', 'NASDelay', 'WeatherDelay']
sns.pairplot(DelayedFlights[cols], size = 2.5)
plt.show()
```

c:\Users\marta\appdata\local\programa\python\python37\lib\site-packages\seaborn\axisgrid.py:1969: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)



```
In [152]: # Una vez realizada una pequeña exploración inicial, pasamos a dividir el dataset con el método train_test_split()

train, test = train_test_split(file, test_size = 0.30, random_state=4)

print("Ejemplos usados para Train: ", len(train))
print("Ejemplos usados para Test: ", len(test))

Ejemplos usados para Train: 1350751
Ejemplos usados para Test: 578894

In [153]: train.describe().round(2)

Out[153]:
```

	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	FlightNum	ActualElapsedTime	...	TaxiIn
count	1350751	1350751	1350751	1350751	1350751	1350751	1350751	1350751	1350751	1349869	...	1350751
mean	2008	6	16	4	1519	1468	1610	1634	2184	732	...	7
std	0	3	9	2	450	425	548	465	1944	173	...	5
min	2008	1	1	1	1	0	1	0	1	15	...	0
25%	2008	3	8	2	1203	1135	1316	1325	611	80	...	4
50%	2008	6	16	4	1545	1510	1715	1705	1543	116	...	6
75%	2008	9	23	6	1900	1815	2030	2014	3422	165	...	8
max	2008	12	31	7	2400	2359	2400	2359	9740	1114	...	200

8 rows × 25 columns

```
In [154]: test.describe().round(2)

Out[154]:
```

	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	FlightNum	ActualElapsedTime	...	TaxiIn
count	578894	578894	578894	578894	578894	578894	578894	578894	578894	578499	...	578894
mean	2008	6	16	4	1519	1468	1610	1634	2185	133	...	7
std	0	3	9	2	451	425	548	465	1946	72	...	5
min	2008	1	1	1	1	0	1	0	1	14	...	0
25%	2008	3	8	2	1203	1135	1315	1325	610	80	...	4
50%	2008	6	16	4	1545	1510	1715	1706	1543	116	...	6
75%	2008	9	23	6	1901	1815	2031	2015	3423	165	...	8
max	2008	12	31	7	2400	2359	2400	2359	9741	790	...	200

8 rows × 25 columns

Exercici 2.

Aplica algun procés de transformació (estandaritzar les dades numèriques, crear columnes dummies, polinomis...): [\[X\] Fet](#)

```
In [155]: # VAMOS A NORMALIZAR EL ATRIBUTO CATEGÓRICO CANCELLATIONCODE CON LA FUNCIÓN DUMMY #

In [159]: file['CancellationCode'].value_counts()

Out[159]:
```

CancellationCode	count
0	1936125
1	307
2	24
3	80
4	0

Name: CancellationCode, dtype: int64

```
In [160]: cancelCode_dummy = pd.get_dummies(file['CancellationCode'])

In [161]: cancelCode_dummy.head(3)

Out[161]:
```

A	B	C	N
0	0	0	1
1	0	0	0
2	0	0	1

```
In [162]: cancelCode_dummy.rename(columns = {'A':'CancellationCarrier', 'B':'CancellationWeather', 'C':'Cancellation_NAS'}, inplace = True)

In [163]: # Eliminamos la columna CancellationCode y procederemos a continuación a concatenar el nuevo DF cancelCode_dummy con el dataset original

file=file.drop('CancellationCode', axis='columns')

In [168]: # Comprobamos que la columna CancellationCode ya no está en el DF

file.columns

Out[168]: Index(['Unnamed: 0', 'Year', 'Month', 'DayOfMonth', 'DayOfWeek', 'DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime', 'FlightNum', 'ActualElapsedTime', 'CRSElapsedTime', 'AirTime', 'ArrDelay', 'DepDelay', 'Origin', 'Dest', 'Distance', 'TaxiIn', 'TaxiOut', 'Cancelled', 'Diverted', 'CarrierDelay', 'NASDelay', 'WeatherDelay', 'LateAircraftDelay', 'CancellationCarrier', 'CancellationWeather', 'Cancellation_NAS', 'N'], dtype='object')
```

```
In [169]: # Por último, concatenamos ambos DF's para trabajar con un único dataset

file = pd.concat([file, cancelCode_dummy], axis=1)
file.head(3)

Out[169]:
```

Unnamed: 0	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	...	Diverted	CarrierDelay
0	2008	1	3	4	2003	1955	2211	2225	WN	...	0	NaN
1	2008	1	3	4	754	735	1002	1000	WN	...	0	NaN
2	2008	1	3	4	628	620	804	750	WN	...	0	NaN

3 rows × 33 columns

```
In [170]: file.columns

Out[170]: Index(['Unnamed: 0', 'Year', 'Month', 'DayOfMonth', 'DayOfWeek', 'DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime', 'UniqueCarrier', 'FlightNum', 'ActualElapsedTime', 'CRSElapsedTime', 'AirTime', 'ArrDelay', 'DepDelay', 'Origin', 'Dest', 'Distance', 'TaxiIn', 'TaxiOut', 'Cancelled', 'Diverted', 'CarrierDelay', 'NASDelay', 'WeatherDelay', 'LateAircraftDelay', 'CancellationCarrier', 'CancellationWeather', 'Cancellation_NAS', 'N'], dtype='object')
```

```
In [171]: # VAMOS A NORMALIZAR LOS ATRIBUTOS NUMÉRICOS ARRDELAY Y DEPDELAY CON STANDARDSCALER #

In [171]: file_ss=file[['ArrDelay', 'DepDelay']].copy()

In [172]: ss = StandardScaler()
file_transformed = ss.fit_transform(file_ss)
file_transformed

Out[172]: array([[ -0.98970118, -0.65886773],
       [-0.70793514, -0.45288482],
       [-0.49661061, -0.65886773],
       [ 1.00027146,  0.68938404],
       [-0.5846625 , -0.60269057],
       [-0.83120778, -0.67759345])

In [173]: file_transformed = pd.DataFrame(file_transformed)

In [174]: file_transformed.head(3)

Out[174]:
```

	0	1
0	-1	-1
1	-1	-0
2	-0	-1

```
In [175]: file_transformed.rename(columns = {'0':'ArrDelay', '1':'DepDelay'}, inplace = True)

In [176]: file_transformed.head(3)

Out[176]:
```

	ArrDelay	DepDelay
0	-1	-1
1	-1	-0
2	-0	-1

```
In [177]: file_ok=file._drop(['ArrDelay', 'DepDelay'], axis='columns')

In [178]: file_scaled=pd.concat([file_ok, file_transformed], axis=1)
file_scaled.head()

Out[178]:
```

Unnamed: 0	Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	...	WeatherDelay	NASDelay
0	2008	1	3	4	2003	1955	2211	2225	WN	...	NaN	NaN
1	2008	1	3	4	754	735	1002	1000	WN	...	NaN	NaN
2	2008	1	3	4	628	620	804	750	WN	...	NaN	NaN
3	2008	1	3	4	1829	1755	1959	1925	WN	...	0	0
4	2008	1	3	4	1940	1915	2121	2110	WN	...	NaN	NaN

5 rows × 33 columns

```
In [179]: file_scaled.columns

Out[179]: Index(['Unnamed: 0', 'Year', 'Month', 'DayOfMonth', 'DayOfWeek', 'DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime', 'UniqueCarrier', 'FlightNum', 'ActualElapsedTime', 'CRSElapsedTime', 'AirTime', 'ArrDelay', 'DepDelay', 'Origin', 'Dest', 'Distance', 'TaxiIn', 'TaxiOut', 'Cancelled', 'Diverted', 'CarrierDelay', 'NASDelay', 'WeatherDelay', 'LateAircraftDelay', 'CancellationCarrier', 'CancellationWeather', 'Cancellation_NAS', 'N', 'ArrDelay', 'DepDelay'], dtype='object')
```

```
In [180]: file_scaled.shape

Out[180]: (1936758, 33)
```

Exercici 3.

Resumeix les noves columnes generades de manera estadística i gràfica: [\[X\] Fet](#)

```
In [181]: # COLUMNAS DUMMIES #

In [181]: # Las columnas generadas con dummies son columnas con datos binarios, es decir, con valores 0 o 1. Por tanto, y a continuación, su resumen estadístico se basa en establecer su valor mín, que obviamente es cero, el valor máx y el número de valores. Vemos que no existe media, ni desviación estándar, ni cuantiles en las columnas generadas.

file_scaled['CancellationWeather'].describe()

Out[181]:
```

	count	mean	std	min	25%	50%	75%	max
0	1936758	0	0	0	0	0	0	1
1	307	0	0	0	0	0	0	1
2	24	0	0	0	0	0	0	1
3	80	0	0	0	0	0	0	1
4	0	0	0	0	0	0	0	1

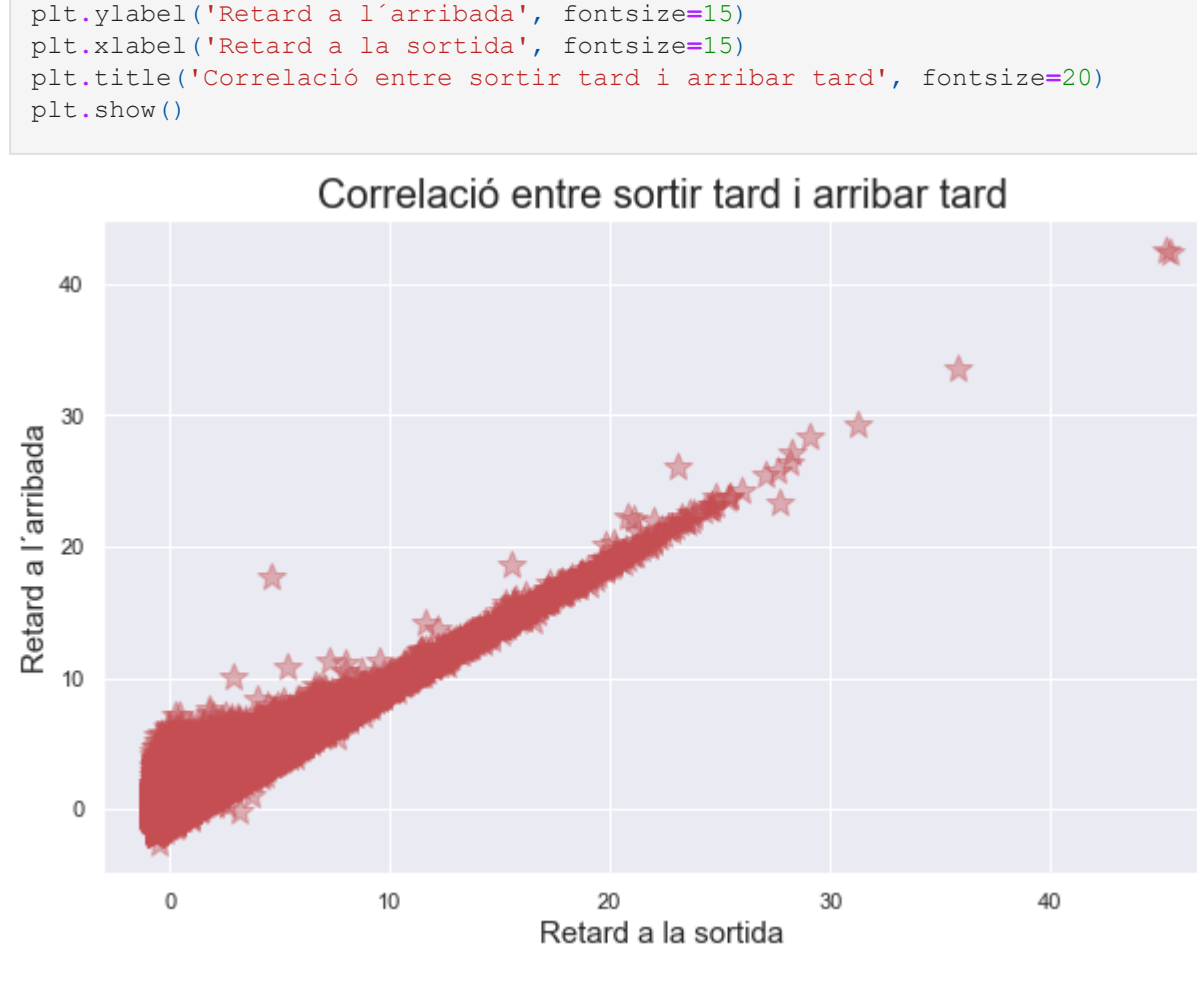
Name: CancellationWeather, dtype: float64

```
In [182]: # Gráficamente, al representar una columna dummy, vemos que simplemente se representan valores absolutos (0 y 1). He querido comparar las cancelaciones causadas por problemas meteorológicos (X) versus el total de cancelaciones (Y).

x = file_scaled['CancellationWeather']
y = file_scaled['Cancelled']

plt.subplots(figsize=(10, 7))
plt.xlabel("Cancelaciones por causas climatológicas", size=16)
plt.ylabel("Total de cancelaciones", size=16)
plt.legend('figsize=(10, 7)')
myexplode = [0, 0.2, 0, ]
fig=plt.gcf()
plt.autotitle("%i.16%", labels=mylabels, explode=myexplode)
plt.legend('figsize=(10, 7)')
plt.show()
```

```
Out[182]: [matplotlib.lines.Line2D at 0x24757f9b48b]
```



```
In [171]: # Para representar un gráfico que nos proporcione información con las columnas dummies que hemos creado, por el momento, vamos a utilizar un gráfico de barras. Vemos que no existe media, ni desviación estándar, ni cuantiles en las columnas generadas.

In [185]: # Contamos el número de vuelos totales cancelados con código de "Motivo de cancelación" y los pasamos a DF

NAS=file_scaled['Cancellation_NAS'].value_counts().to_list()
Weather=file_scaled['CancellationWeather'].value_counts().to_list()
Carrier=file_scaled['CancellationCarrier'].value_counts().to_list()

In [186]: # A continuación pasamos los datos de cada motivo de cancelación a lista, y hacemos una lista total que incluye los vuelos de cancelados por cada causa, y por último, graficamos.

NAS=NAS[1:]
Weather=Weather[1:]

In [187]: Weather

Out[187]: [0]
```

```
In [188]: Carrier

Out[188]: [246]
```

```
In [189]: DF=NAS+Weather+Carrier

In [194]: # PLOT #

mylabels = ["Problemas del National Air System", "Causas Climatológicas", "Problemas de la Aerolínea"]
plt.title('Motivos de cancelación de los vuelos', fontsize=20)
plt.legend('figsize=(10, 7)')
myexplode = [0, 0.2, 0, ]
fig=plt.gcf()
plt.autotitle("%i.16%", labels=mylabels, explode=myexplode)
plt.legend('figsize=(10, 7)')
plt.show()
```



```
In [191]: # COLUMNAS ESTANDARIZADAS CON STANDARDSCALER #

In [191]: # Cuando realizamos estandarización de valores numéricos, a diferencia que ocurre con la normalización de valores, con dummies, vemos que el resumen estadístico no muestra que se traten de variables con valores binarios, sino que sí media, desviación estándar, los cuantiles, valores mínimos y máximos, etc. y podremos operar con ellos como si se tratara de valores numéricos.

In [215]: pd.options.display.float_format = '{:,.2f}'.format

file_scaled['ArrDelay'].describe()

Out[215]:
```

	(count)	mean	std	min	25%	50%	75%	max
0	1928371	0.00	1.00	-0.70	-0.32	-0.32	-0.24	42.60
1	307	0.00	1.00	-0.70	-0.32	-0.32	-0.24	42.60
2	24	0.00	1.00	-0.70	-0.32	-0.32	-0.24	42.60
3	80	0.00	1.00	-0.70	-0.32	-0.32	-0.24	42.60
4	0	0.00	1.00	-0.70	-0.32	-0.32	-0.24	42.60

Name: ArrDelay, dtype: float64

```
In [226]: # Fent el gràfic observem que hi ha una correlació lineal entre sortir tard i arribar tard a destí.

file_scaled['ArrDelay']
x=file_scaled['DepDelay']

plt.plot(x,y, 'r', alpha=0.4, marker='x')
plt.xlabel('Retard a l'arribada', fontsize=15)
plt.ylabel('Retard a la sortida', fontsize=15)
plt.title('Correlació entre sortir tard i arribar tard', fontsize=20)
plt.show()
```


Exercici 4.

Exporta el Notebook com a pdf i com a html: [\[X\] Fet](#)

```
In [190]: python -m pip install --user notebook==3.0.4
python -m pip install --user notebook==3.0.4

SyntaxError: invalid syntax

In [191]: jupyter-nbconvert --to pdfviahtml Tasca12.ipynb
```