

```
import numpy as np
import pandas as pd
from sklearn import read_csv
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
from matplotlib.pyplot import figure
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import r2_score
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing importMinMaxScaler
from sklearn.ensemble import RandomForestRegressor
import seaborn as sns
import math
from sklearn.cluster import KMeans
import statistics as stat
from sklearn.neural_network import MLPRegressor
from sklearn.decomposition import PCA
```

Siguemos el hilo argumental del **ejercicio 9.1** y **ejercicio 9.2** para trabajarcon el ejercicio 10. Primero cargamos el Data Set.

```
In [3]: df = pd.read_csv("DelayedFlights.csv") # este es el conjunto de datos proporcionado en el ejercicio
df.head(10)
```

	Unnamed: 0	Year	Month	DayofMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	...	TaxiIn	TaxiOut	Canc
0	0	2008	1	3	4	2003.0	1955	2211.0	2225	WN	...	4.0	8.0	
1	1	2008	1	3	4	754.0	735	1002.0	1000	WN	...	5.0	10.0	
2	2	2008	1	3	4	628.0	620	804.0	750	WN	...	3.0	17.0	
3	4	2008	1	3	4	1829.0	1755	1959.0	1925	WN	...	3.0	10.0	
4	5	2008	1	3	4	1940.0	1915	2117.0	2110	WN	...	4.0	10.0	
5	6	2008	1	3	4	1937.0	1830	2037.0	1940	WN	...	3.0	7.0	
6	10	2008	1	3	4	706.0	700	916.0	915	WN	...	5.0	19.0	
7	11	2008	1	3	4	1640.0	1510	1845.0	1725	WN	...	6.0	8.0	
8	15	2008	1	3	4	1029.0	1020	1021.0	1010	WN	...	6.0	9.0	
9	16	2008	1	3	4	1452.0	1425	1640.0	1625	WN	...	7.0	8.0	

10 rows × 30 columns

1. Tratamiento de variables

En este apartado vamos a hacer tratamiento de variables para luego aplicar a programas de Clustering

Hemos dejado los enlaces del **ejercicio 9.1** y **ejercicio 9.2** en que se explica los razonamientos para seleccionar una variable u otra,

aunque volveremos a explicar los motivos por los que se seleccionan variables, unas sí u otras no, no entraremos tan en detalle

a) Variable Unnamed 0 y Year, básicamente son un índice y el año de vuelos del 2008. Year es una constante. Así que las eliminamos.

b) las variables "UniqueCarrier", "FlightNum", "Month", "DayofMonth", "DayOfWeek", "Cancelled", "CancellationCode", "Diverted", "Origin", "Dest" son eliminadas por siguiente motivo. *KMeans funciona calculando distancias en un espacio dimensional, de volumen R^n . Si implementamos variables categóricas, por ejemplo binarias, o con más clases como los meses, los puntos van a estar concentrados en hiperplanos de dimensión $n-1$. Si las variables no categóricas, están normalizadas o estandarizadas, la distancia entre hiperplanos será por lo general mayor, que la distancia entre puntos dentro de un hiperplano de dimensión $n-1$.*

c) También eliminamos "ActualElapsedTime", "CRSElapsedTime", "AirTime" ya que tienen un 0.95 de correlación o más con Distance y aportan la misma información. ActualElapsed es el tiempo esperado del vuelo(desamboro, salida, vuelo, más aterrizaje). CRSElapsedTime es el mismo tiempo previsto, mientras que AirTime es el total que el avión está en el aire y Distance la distancia recorrida en millas.

```
In [4]: df[["ActualElapsedTime", "CRSElapsedTime", "AirTime", "Distance" ]].corr()
```

	ActualElapsedTime	CRSElapsedTime	AirTime	Distance
ActualElapsedTime	1.000000	0.971122	0.976660	0.952980
CRSElapsedTime	0.971122	1.000000	0.986086	0.981759
AirTime	0.976660	0.986086	1.000000	0.980294
Distance	0.952980	0.981759	0.980294	1.000000

```
In [5]: df=df.drop(["Unnamed: 0", "Year", "UniqueCarrier", "FlightNum", "TailNum", "Month", "DayofMonth", "DayOfWeek", "Cancelled", "CancellationCode", "Diverted", "Origin", "Dest", "ActualElapsedTime", "CRSElapsedTime", "AirTime"], axis=1)
```

	DepTime	CRSDepTime	ArrTime	CRSArrTime	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
0	2003.0	1955	2211.0	2225	-14.0	8.0	810	4.0	8.0	NaN	NaN	NaN	NaN	NaN
1	754.0	735	1002.0	1000	2.0	19.0	810	5.0	10.0	NaN	NaN	NaN	NaN	NaN
2	628.0	620	804.0	750	14.0	8.0	515	3.0	17.0	NaN	NaN	NaN	NaN	NaN
3	1829.0	1755	1959.0	1925	34.0	34.0	515	3.0	10.0	2.0	0.0	0.0	0.0	0.0
4	1940.0	1915	2117.0	2110	11.0	25.0	688	4.0	10.0	NaN	NaN	NaN	NaN	NaN
5	1937.0	1830	2037.0	1940	57.0	67.0	1591	3.0	7.0	10.0	0.0	0.0	0.0	0.0
6	706.0	700	916.0	915	1.0	6.0	828	5.0	19.0	NaN	NaN	NaN	NaN	NaN
7	1640.0	1510	1845.0	1725	80.0	94.0	828	6.0	8.0	8.0	0.0	0.0	0.0	0.0
8	1029.0	1020	1021.0	1010	11.0	9.0	162	6.0	9.0	NaN	NaN	NaN	NaN	NaN
9	1452.0	1425	1640.0	1625	15.0	27.0	1489	7.0	8.0	3.0	0.0	0.0	0.0	0.0

Como vamos a empezar a transformar variables, vamos a hacer previamente un muestreo.

```
In [6]: df2=df.sample(390000,random_state=55)
```

Como vimos en los ejercicios 9.1 y 9.2, los valores NaN de ArrDelay, así como de otras variables coinciden con aquellos que el vuelo había desviado o cancelado, valores imposibles de valor (los NaN) por el mismo concepto de cancelación o desvío. por lo que eliminamos los valores NaN de ArrDelay

```
In [7]: df3=df2.dropna( subset=["ArrDelay"]).reset_index(drop=True)
df3.isna().sum()
```

	DepTime	CRSDepTime	ArrTime	CRSArrTime	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
DepTime	0													
CRSDepTime	0													
ArrTime	0													
CRSArrTime	0													
ArrDelay	0													
DepDelay	0													
Distance	0													
TaxiIn	0													
TaxiOut	0													
CarrierDelay	137592													
WeatherDelay	137592													
NASDelay	137592													
SecurityDelay	137592													
LateAircraftDelay	137592													
dtype:	int64													

```
In [8]: df3.shape
```

```
Out[8]: (388357, 14)
```

0.1. Transformación de variables

En esta parte vamos a convertir las variables DepTime, CRSDepTime, ArrTime, CRSArrTime en la función ciclica.

Estas cuatro variables vienen en formato horario hhmm. Lo que haremos será contar todos los minutos transcurridos durante el día, siendo 0 minutos a las 00:00 y 1440 los minutos transcurridos durante el día a las 23:59.

Más adelante, en el apartado 0.4, transformaremos estas variables en ciclicas.

```
In [9]: # Primero de todo convertiro las variables horarias en formato hora y para eso tienen que haber 4 dígitos, que
# ignorará con ceros
# Primer tengo que convertir en entero las variables DepTime y ArrTime en enteros para evitar los decimales
df3['DepTime'] = df3['DepTime'].astype(int)
```

```
In [10]: df3['ArrTime'] = df3['ArrTime'].astype(int)
```

```
In [11]: #relleno por la izquierda con ceros
df3['DepTime'] = df3['DepTime'].astype(str).str.zfill(4)
df3['CRSDepTime'] = df3['CRSDepTime'].astype(str).str.zfill(4)
df3['ArrTime'] = df3['ArrTime'].astype(str).str.zfill(4)
df3['CRSArrTime'] = df3['CRSArrTime'].astype(str).str.zfill(4)
df3.head()
```

	DepTime	CRSDepTime	ArrTime	CRSArrTime	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
0	0918	0905	1040	1035	9.0	13.0	414	3.0	12.0	NaN	NaN	NaN	NaN	NaN
1	1421	1400	1500	1450	10.0	21.0	192	1.0	7.0	NaN	NaN	NaN	NaN	NaN
2	2250	2153	0013	2330	43.0	57.0	432	5.0	17.0	0.0	43.0	0.0	0.0	0.0
3	0049	0054	0521	0530	-9.0	9.0	1235	4.0	17.0	NaN	NaN	NaN	NaN	NaN
4	1854	1845	2017	2010	7.0	9.0	488	4.0	11.0	NaN	NaN	NaN	NaN	NaN

```
In [12]: # las convierto en formato horario( Nota: en un principio lo pasé a formato horario por si lo necesitaba para
# pasar luego segundos a minutos y con la reconversión ya comentada aplica la función minutos()
def minutos(x):
    xx=split(sep=":")
    seg= sum(int(x[i])*60+int(x[i+1]))
    return seg
```

	DepTime	CRSDepTime	ArrTime	CRSArrTime	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
0	091800	090500	104000	103500	9.0	13.0	414	3.0	12.0	NaN	NaN	NaN	NaN	NaN
1	142100	140000	150000	145000	10.0	21.0	192	1.0	7.0	NaN	NaN	NaN	NaN	NaN
2	225000	215300	001300	233000	43.0	57.0	432	5.0	17.0	0.0	43.0	0.0	0.0	0.0
3	004900	004000	052100	053000	-9.0	9.0	1235	4.0	17.0	NaN	NaN	NaN	NaN	NaN
4	185400	184500	201700	201000	7.0	9.0	488	4.0	11.0	NaN	NaN	NaN	NaN	NaN
...
388357	085900	084500	110100	102000	41.0	14.0	1333	5.0	15.0	14.0	0.0	0.0	0.0	27.0
388353	145900	145000	165200	165000	2.0	9.0	594	2.0	14.0	NaN	NaN	NaN	NaN	NaN
388354	213200	203000	224300	213000	73.0	62.0	255	5.0	14.0	11.0	0.0	0.0	0.0	11.0
388355	194200	175000	211800	193500	103.0	112.0	583	3.0	14.0	83.0	0.0	0.0	0.0	0.0
388356	081500	074500	093200	090500	27.0	30.0	427	5.0	13.0	0.0	0.0	0.0	0.0	27.0

388357 rows × 14 columns

```
In [13]: # creamos la función minutos, que divide la hora hhmm con un Split, en una lista ("hh","mm"), reconvierte hh y
# mm luego segundos a minutos y con la reconversión ya comentada aplica la función minutos()
def minutos(x):
    xx=split(sep=":")
    seg= sum(int(x[i])*60+int(x[i+1]))
    return seg
```

```
In [14]: dfhora=df3[['DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime']]
dfhora_DT=dfhora['DepTime'].apply(minutos)
dfhora_CRSD=dfhora['CRSDepTime'].apply(minutos)
dfhora_ARR=dfhora['ArrTime'].apply(minutos)
dfhora_CRSA=dfhora['CRSArrTime'].apply(minutos)
```

```
In [15]: df4=df3.drop(['DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime'], axis=1)
```

```
In [16]: # ahora añadimos las cuatro columnas nuevas
df5=pd.concat([df4, dfhora_DT, dfhora_CRSD, dfhora_ARR], axis=1)
df5.columns
```

```
Out[16]: Index(['ArrDelay', 'DepDelay', 'Distance', 'TaxiIn', 'TaxiOut', 'CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay', 'LateAircraftDelay', 'X10', 'X11'], dtype='object')
```

```
In [17]: df5[['DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime']].describe()
```

	DepTime	CRSDepTime	ArrTime	CRSArrTime
count	388357.000000	388357.000000	388357.000000	388357.000000
mean	922.620877	891.358868	978.678528	992.171942
std	269.858220	255.01740	327.865454	278.288941
min	1.000000	0.000000	1.000000	0.000000
25%	723.000000	695.000000	797.000000	805.000000
50%	944.000000	910.000000	1035.000000	1025.000000
75%	1140.000000	1095.000000	1230.000000	1213.000000
max	1440.000000	1439.000000	1440.000000	1439.000000

```
In [18]: df5.isna().sum()
```

	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
ArrDelay	0									
DepDelay	0									
Distance	0									
TaxiIn	0									
TaxiOut	0									
CarrierDelay	137592									
WeatherDelay	137592									
NASDelay	137592									
SecurityDelay	137592									
LateAircraftDelay	137592									
dtype:	int64									

0.2 Valores NaN

En esta sección vamos a completar las variables del motivo del retraso que tiene varios NaN. Como ya pudimos observar en el ejercicio 9.2, los vuelos con retrasos de menos de 14 minutos, tienen valores NaN

Carrier Delay es el retraso de la compañía

WeatherDelay es el retraso por las condiciones climatológicas SecurityDelay es el retraso por cuestiones de seguridad

LateAircraftDelay es el retraso de la misma aerovne.

Nas delay son los retrasos causado por el Sistema Nacional del Espacio Aéreo (NAS)

por lo que vamos a asignar 0 a los valores NaN

```
In [20]: df_delay=df5[['ArrDelay','CRSDepTime','CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay','LateAircraftDelay']]
df5_0m=df3.drop(['CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay', 'LateAircraftDelay'], axis=1)
delay_not_NaN=df_delay[['CarrierDelay', 'WeatherDelay', 'NASDelay', 'SecurityDelay', 'LateAircraftDelay']].fill(0)
```

```
In [20]: df5=pd.concat([df5_0, delay_not_NaN], axis=1)
df5.head(10)
```

```

rmdelmodels del data get --del las dos nuevas variables
df5b=df5.drop(['DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime'], axis=1)
df6=pd.concat([df5b, x10,x11], axis =1)
df6.head(10)

```

	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay	X10	X11
0	9.0	13.0	414	3.0	12.0	0.0	0.0	0.0	0.0	0.0	0.0	9
1	10.0	21.0	192	1.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	10
2	43.0	57.0	432	5.0	17.0	0.0	43.0	0.0	0.0	0.0	0.0	43
3	-9.0	9.0	1235	4.0	17.0	0.0	0.0	0.0	0.0	0.0	0.0	-9
4	7.0	9.0	488	4.0	11.0	0.0	0.0	0.0	0.0	0.0	0.0	7
5	134.0	105.0	316	15.0	39.0	0.0	0.0	87.0	0.0	47.0	134	105
6	-15.0	8.0	1235	5.0	16.0	0.0	0.0	0.0	0.0	0.0	0.0	-15
7	32.0	40.0	304	5.0	13.0	32.0	0.0	0.0	0.0	0.0	0.0	32
8	24.0	17.0	1589	5.0	29.0	17.0	0.0	7.0	0.0	0.0	0.0	24
9	49.0	36.0	787	7.0	16.0	6.0	0.0	13.0	0.0	30.0	49	36

```

df6.shape
(388357, 12)

```

En el siguiente paso vamos a analizar la multicolinealidad. Para ver si alguna variable tiene una multicolinealidad muy elevada $VIF = 1/R^2$ donde R^2 es el coeficiente de determinación de la regresión lineal

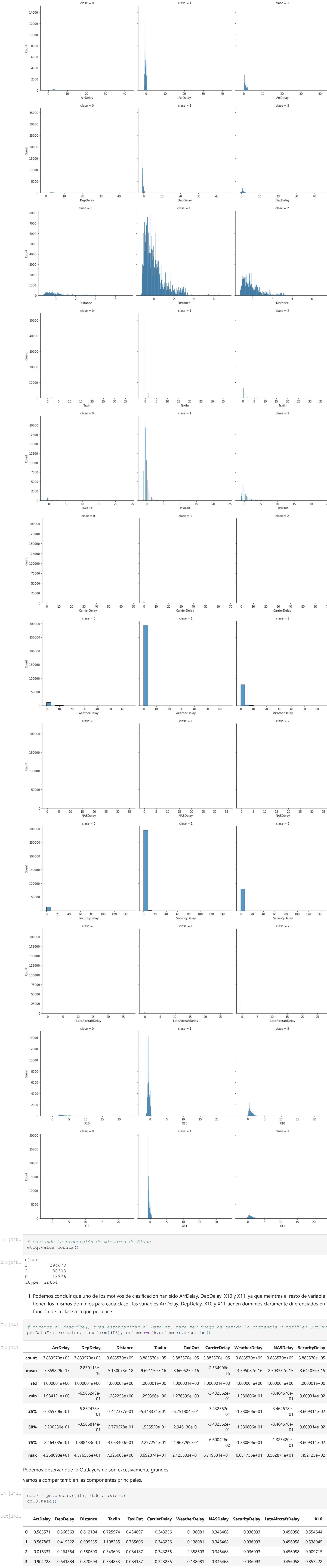
```

from statsmodels.stats.outliers_influence import variance_inflation_factor

def vif(X):
    vifdf = pd.DataFrame()
    vifdf['variables'] = X.columns
    vifdf['vif'] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
    return vifdf

round(vif(df6), 2)

```

In [140]:
contando la proporción de miembros de Clase
etiq.value_counts()

Out[140]:
class
1 294678
2 80303
0 13376
dtype: int64

1. Podemos concluir que uno de los motivos de clasificación han sido ArrDelay, DepDelay, X10 y X11, ya que mientras el resto de variable tienen los mismos dominios para cada clase, las variables ArrDelay, DepDelay, X10 y X11 tienen dominios claramente diferenciados en función de la clase a la que pertence

In [141]:
miramos el describe() tras estandarizar el DataSet, para ver como ha tenido la distancia y posibles Outlays
pd.DataFrame(Scaler.transform(df6), columns=df6.columns).describe()

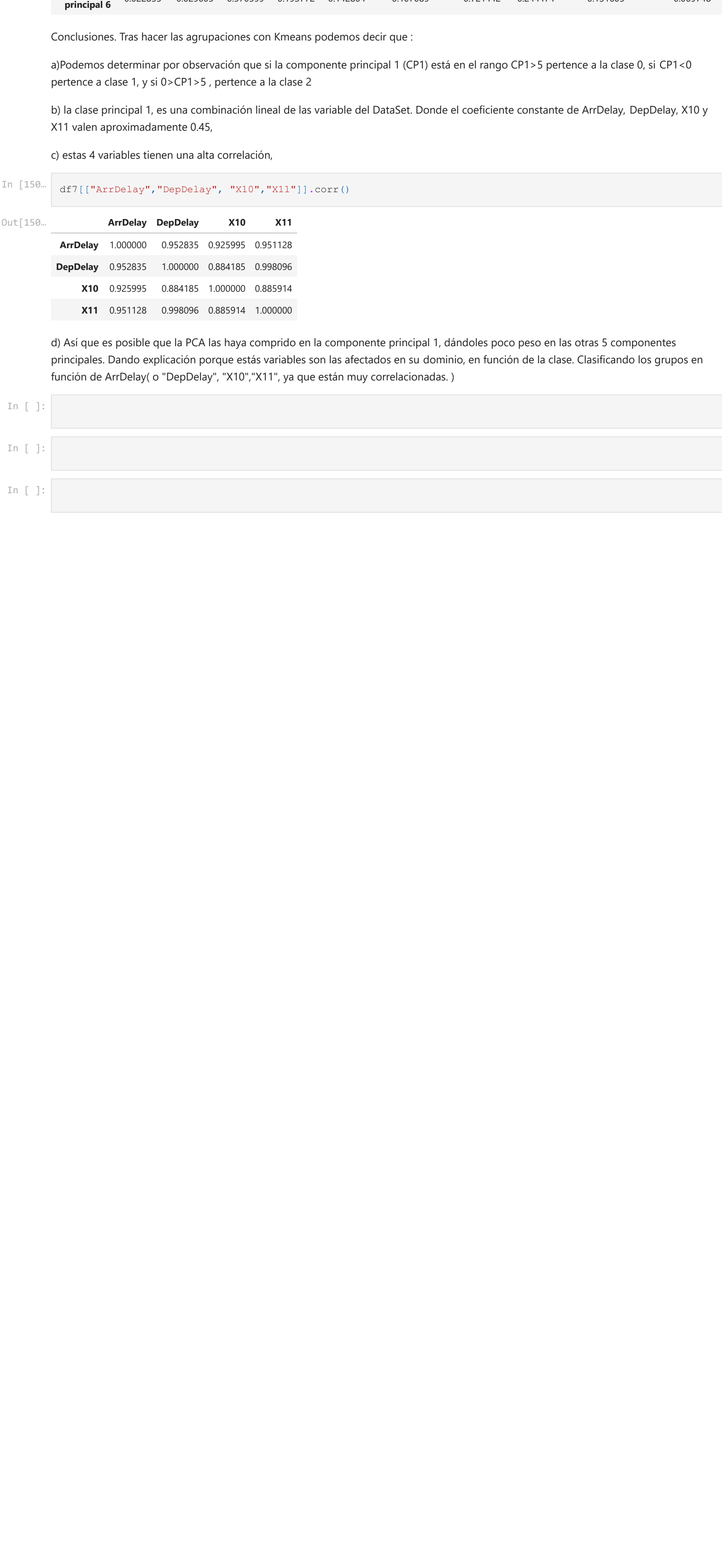
	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay	X10	X11
count	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05	3.883570e+05
mean	-7.859829e-17	-2.830113e-16	-5.150073e-18	-9.691159e-16	-3.660525e-16	-2.534906e-15	4.795082e-16	2.503332e-15	-3.644056e-15	-3.644056e-15	-3.644056e-15	-3.644056e-15
std	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00
min	-1.984121e+00	-6.985243e-01	-1.282255e+00	-1.299396e+00	-1.276599e+00	-3.432562e-01	-1.380806e-01	-3.464678e-01	-3.609314e-02	-3.609314e-02	-3.609314e-02	-3.609314e-02
25%	-5.855706e-01	-5.852433e-01	-7.447377e-01	-5.348334e-01	-5.751804e-01	-3.432562e-01	-1.380806e-01	-3.464678e-01	-3.609314e-02	-3.609314e-02	-3.609314e-02	-3.609314e-02
50%	-3.200230e-01	-3.586814e-01	-2.770278e-01	-1.525520e-01	-2.946130e-01	-3.432562e-01	-1.380806e-01	-3.464678e-01	-3.609314e-02	-3.609314e-02	-3.609314e-02	-3.609314e-02
75%	2.464785e-01	1.888433e-01	4.053400e-01	2.297294e-01	1.963799e-01	-6.600426e-02	-1.380806e-01	-1.325420e-01	-3.609314e-02	-3.609314e-02	-3.609314e-02	-3.609314e-02
max	4.268098e+01	4.576555e+01	7.325003e+00	3.692874e+01	2.425503e+01	6.719531e+01	6.651756e+01	3.562871e+01	1.492125e+02	1.492125e+02	1.492125e+02	1.492125e+02

Podemos observar que lo Outlayers no son excesivamente grandes
vamos a comparar también las componentes principales.

In [143]:
df10 = pd.concat([df9, df8], axis=1)
df10.head()

	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay	X10	X11
0	-0.585571	-0.566363	-0.612104	-0.725974	-0.434897	-0.343256	-0.138081	-0.346468	-0.036093	-0.456058	-0.554644	-0.554644
1	-0.567867	-0.415322	-0.999535	-1.108255	-0.785606	-0.343256	-0.138081	-0.346468	-0.036093	-0.456058	-0.538045	-0.538045
2	0.016337	0.264364	0.582090	-0.343693	-0.084187	-0.343256	2.358603	-0.346468	-0.036093	-0.456058	0.009715	0.009715
3	-0.904228	-0.641884	0.826094	-0.534833	-0.084187	-0.343256	-0.138081	-0.346468	-0.036093	-0.456058	-0.853422	-0.853422
4	-0.620977	-0.641884	-0.482960	-0.534833	-0.505039	-0.343256	-0.138081	-0.346468	-0.036093	-0.456058	-0.587841	-0.587841

In [145]:
indice2= ["componente principal 1", "componente principal 2", "componente principal 3", "componente principal 4", "componente principal 5", "componente principal 6"]
for i in indice2:
sns.displot(df10, x=i, col="class")



Podemos observar que la variable que determina la clasificación es la **Componente principal 1**, ya que su dominio varía en función de la clase. Vamos a ver como son las combinaciones lineales de las variables para sacar los componentes principales

In [147]:
pd.DataFrame(data = pca2.components_, columns = df7.columns,
index = ["componente principal 1", "componente principal 2", "componente principal 3", "componente principal 4", "componente principal 5", "componente principal 6"])

	ArrDelay	DepDelay	Distance	TaxiIn	TaxiOut	CarrierDelay	WeatherDelay	NASDelay	SecurityDelay	LateAircraftDelay
componente principal 1	0.455679	0.446139	-0.006972	0.072683	0.117151	0.251775	0.127080	0.193445	0.003756	0.261358
componente principal 2	-0.028414	0.144431	-0.137661	-0.346136	-0.585387	0.250336	-0.074603	-0.603122	-0.000460	0.208334
componente principal 3	0.000169	0.007134	0.351833	0.035813	0.035670	0.697176	-0.236844	-0.013206	0.005811	-0.575532
componente principal 4	0.019147	-0.015667	-0.563940	-0.379466	0.157540	0.144160	0.553338	0.032690	-0.112705	-0.413890
componente principal 5	0.003667	0.006027	0.038624	-0.069967	0.018660	-0.024566	0.149734	-0.034522	0.982508	-0.060169
componente principal 6	-0.022855	0.029005	0.570599	0.193772	-0.142804	-0.107089	0.721442	-0.244474	-0.131605	-0.069748

Conclusiones. Tras hacer las agrupaciones con Kmeans podemos decir que :
a)Podemos determinar por observación que si la componente principal 1 (CP1) está en el rango CP1>5 pertenece a la clase 0, si CP1<0 pertenece a clase 1, y si 0<CP1<5, pertenece a la clase 2
b) la clase principal 1, es una combinación lineal de las variables del DataSet. Donde el coeficiente constante de ArrDelay, DepDelay, X10 y X11 valen aproximadamente 0.45,
c) estas 4 variables tienen una alta correlación,

In [150]:
df7[["ArrDelay", "DepDelay", "X10", "X11"]].corr()

	ArrDelay	DepDelay	X10	X11
ArrDelay	1.000000	0.952835	0.925995	0.951128
DepDelay	0.952835	1.000000	0.884185	0.998096
X10	0.925995	0.884185	1.000000	0.885914
X11	0.951128	0.998096	0.885914	1.000000

d) Así que es posible que la PCA las haya comprimido en la componente principal 1, dándoles poco peso en las otras 5 componentes principales. Dando explicación porque estas variables son las afectadas en su dominio, en función de la clase. Clasificando los grupos en función de ArrDelay(o "DepDelay", "X10", "X11", ya que están muy correlacionadas.)

In []:

#

In []:

#

In []:

#