

# IMPORTAR LIBRERÍAS Y DATASETS

```
In [ ]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import zipfile
import cv2
import tensorflow as tf
from tensorflow.python.keras import Sequential
from tensorflow.keras.layers import Input, Add, Dense, Activation, ZeroPadding2D, BatchNormalization, Flatten,
from tensorflow.keras.models import Model, load_model
from tensorflow.keras.initializers import GlorotUniform
from tensorflow.keras.utils import plot_model
from tensorflow.keras.callbacks import ReduceLROnPlateau, EarlyStopping, ModelCheckpoint, LearningRateScheduler
from sklearn.preprocessing import StandardScaler, normalize
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.metrics import silhouette_score
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.metrics.pairwise import cosine_similarity
import plotly.express as px
import plotly.graph_objects as go

from google.colab import files #libreria para cargar ficheros directamente en Colab
!cat librerias.txt
```

```
In [ ]: # Procedemos a montar el drive usando los siguientes comandos:
from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive
```

```
In [ ]: sales_df = pd.read_csv('/content/drive/MyDrive/01_Oriol/01_Machine Learning/07_Masterclass-IA-Moderna/03_Market/sales_df.csv')
```

2819

10373

29

100.00

1

3978.51

1/31/2005  
0:00

Shipped

1

1

2005

2820

10386

43

100.00

4

5417.57

3/1/2005  
0:00

Resolved

1

3

2005

2821

10397

34

62.24

1

2116.16

3/28/2005  
0:00

Shipped

1

3

2005

2822

10414

47

65.52

9

3079.44

5/6/2005  
0:00

On Hold

2

5

2005

2823 rows x 25 columns

# Veamos los tipos de datos

sales\_df.dtypes

ORDERNUMBER

int64

QUANTITYORDERED

int64

PRICEEACH

float64

ORDERLINENUMBER

int64

SALES

float64

ORDERDATE

object

STATUS

object

QTR\_ID

int64

MONTH\_ID

int64

YEAR\_ID

int64

PRODUCTLINE

object

MSRP

int64

PRODUCTCODE

object

CUSTOMERNAME

object

PHONE

object

ADDRESSLINE1

object

ADDRESSLINE2

object

CITY

object

STATE

object

POSTALCODE

object

COUNTRY

object

TERRITORY

object

CONTACTLASTNAME

object

CONTACTFIRSTNAME

object

DEASIZE

object

dtype: object

# Convertir la fecha del pedido en formato de fecha y hora

sales\_df['ORDERDATE'] = pd.to\_datetime(sales\_df['ORDERDATE'])

# Comprobar el tipo de datos

sales\_df.dtypes

ORDERNUMBER

int64

QUANTITYORDERED

int64

PRICEEACH

float64

ORDERLINENUMBER

int64

SALES

float64

ORDERDATE

datetime64[ns]

STATUS

object

QTR\_ID

int64

MONTH\_ID

int64

YEAR\_ID

int64

PRODUCTLINE

object

MSRP

int64

PRODUCTCODE

object

CUSTOMERNAME

object

PHONE

object

ADDRESSLINE1

object

ADDRESSLINE2

object

CITY

object

STATE

object

POSTALCODE

object

COUNTRY

object

TERRITORY

object

CONTACTLASTNAME

object

CONTACTFIRSTNAME

object

dtype: object

```
In [ ]: # Veamos los tipos de datos
sales_df.dtypes
```

sales\_df.dtypes

ORDERNUMBER	int64
QUANTITYORDERED	int64
PRICEEACH	float64
ORDERLINENUMBER	int64
SALES	float64
ORDERDATE	object
STATUS	object
QTR_ID	int64
MONTH_ID	int64
YEAR_ID	int64
PRODUCTLINE	object
MSRP	int64
PRODUCTCODE	object
CUSTOMERNAME	object
PHONE	object
ADDRESSLINE1	object
ADDRESSLINE2	object
CITY	object
STATE	object
POSTALCODE	object
COUNTRY	object
TERRITORY	object
CONTACTLASTNAME	object
CONTACTFIRSTNAME	object
DEALSIZE	object
dtype:	object

```
In [ ]: # Convertir la fecha del pedido en formato de fecha y hora
sales_df['ORDERDATE'] = pd.to_datetime(sales_df['ORDERDATE'])
```

```
In [ ]: # Comprobar el tipo de datos
sales_df.dtypes
```

sales\_df.dtypes

ORDERNUMBER	int64
QUANTITYORDERED	int64
PRICEEACH	float64
ORDERLINENUMBER	int64
SALES	float64
ORDERDATE	datetime64[ns]
STATUS	object
QTR_ID	int64
MONTH_ID	int64
YEAR_ID	int64
PRODUCTLINE	object
MSRP	int64
PRODUCTCODE	object
CUSTOMERNAME	object
PHONE	object
ADDRESSLINE1	object
ADDRESSLINE2	object
CITY	object
STATE	object
POSTALCODE	object
COUNTRY	object
TERRITORY	object
CONTACTLASTNAME	object
CONTACTFIRSTNAME	object
DEALSIZE	object
dtype:	object

```
In [ ]: # Comprobar el número de elementos no nulos del data frame
sales_df.info()
```

```
In [ ]: # Comprobar el número de elementos nulos del data frame
sales_df.isnull().sum()
```

sales\_df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   ORDERNUMBER        2823 non-null    int64
 1   QUANTITYORDERED    2823 non-null    int64
 2   PRICEEACH           2823 non-null    float64
 3   ORDERLINENUMBER    2823 non-null    int64
 4   SALES               2823 non-null    float64
 5   ORDERDATE          2823 non-null    datetime64[ns]
 6   STATUS             2823 non-null    object
 7   QTR_ID             2823 non-null    int64
 8   MONTH_ID          2823 non-null    int64
 9   YEAR_ID            2823 non-null    int64
10  PRODUCTLINE        2823 non-null    object
11  MSRP               2823 non-null    int64
12  PRODUCTCODE        2823 non-null    object
13  CUSTOMERNAME       2823 non-null    object
14  PHONE              2823 non-null    object
15  ADDRESSLINE1       2823 non-null    object
16  ADDRESSLINE2       2823 non-null    object
17  CITY               2823 non-null    object
18  STATE              2823 non-null    object
19  POSTALCODE         2823 non-null    object
20  COUNTRY            2823 non-null    object
21  TERRITORY          2823 non-null    object
22  CONTACTLASTNAME    2823 non-null    object
23  CONTACTFIRSTNAME   2823 non-null    object
24  DEALSIZE           2823 non-null    object
dtypes: datetime64[ns](1), float64(2), int64(7), object(15)
memory usage: 551.5+ KB
```

```
In [ ]: # Dado que hay muchos valores nulos en 'addressline2', 'state', 'postal code' y 'territorio', podemos eliminar la información geográfica del pedido.
# También podemos eliminar la ciudad, la dirección, el número de teléfono, el nombre del contacto, el apellido y el email del cliente.
sales_df = sales_df.drop(['addressline2', 'state', 'postal code', 'territory', 'city', 'state', 'contactlast', 'contactfirst', 'email', 'dealsize'], axis=1)
```

```
In [ ]: # Obtener las observaciones únicas de cada columna
sales_df.nunique()
```

sales\_df.nunique()

QUANTITYORDERED	58
PRICEEACH	1016
ORDERLINENUMBER	18
SALES	2763
ORDERDATE	252
STATUS	6
QTR_ID	4
MONTH_ID	12
YEAR_ID	3
PRODUCTLINE	7
MSRP	80
PRODUCTCODE	109
COUNTRY	19
DEALSIZE	3
dtype:	int64

## ANALISIS EXPLORATORIO DE LOS DATOS Y LIMPIEZA

```
In [ ]: sales_df['COUNTRY'].value_counts().index
```

```
In [ ]: Index(['USA', 'Spain', 'France', 'Australia', 'UK', 'Italy', 'Finland', 'Morway', 'Singapore', 'Canada', 'Denmark', 'Germany', 'Sweden', 'Austria', 'Japan', 'Belgium', 'Switzerland', 'Philippines', 'Ireland'], dtype=object)
```

```
In [ ]: sales_df['COUNTRY'].value_counts()
```

sales\_df['COUNTRY'].value\_counts()

USA	1004
Spain	342
France	314
Australia	185
UK	144
Italy	113
Finland	92
Norway	85
Singapore	79
Canada	70
Denmark	63
Germany	62
Sweden	57
Austria	55
Japan	52
Belgium	33
Switzerland	31
Philippines	26
Ireland	16
Name: COUNTRY, dtype: int64	

```
In [ ]: # Visualizar el recuento de elementos en una columna determinada
df = pd.DataFrame(sales_df['COUNTRY'].value_counts().index, columns=['COUNTRY'])
```

```
In [ ]: # Creamos a esta función para cualquier columna determinada, como 'COUNTRY'
def barplot_visualization(column):
    fig = plt.figure(figsize=(12, 6))
    fig = px.bar(df[column], title=f'{column} value counts', color=sales_df[column].value_counts().index)
    fig.show()
```

```
In [ ]: barplot_visualization('STATUS')
```

barplot\_visualization('STATUS')

Shipped	2055
On Hold	1
Resolved	77

```
In [ ]: barplot_visualization('PRODUCTLINE')
```

barplot\_visualization('PRODUCTLINE')

Motorcycles	95
Ships	54
Ships	54
Ships	54

```
In [ ]: barplot_visualization('DEALSIZE')
```

barplot\_visualization('DEALSIZE')

0	1004
1	342
2	314
3	185
4	144
5	113
6	92
7	85
8	79
9	70
10	63
11	62
12	57
13	55
14	52
15	33
16	31
17	26
18	16

```
In [ ]: # Agregar variables ficticias para reemplazar variables categóricas
def dummies(x):
    dummy = pd.get_dummies(sales_df[x])
    sales_df.drop(columns=x, inplace=True)
    return pd.concat([sales_df, dummy], axis=1)
```

```
In [ ]: # Obtenemos variables ficticias para la columna 'Sales'
sales_df = dummies('SALES')
```

sales\_df

	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	QTR_ID	MONTH_ID	YEAR_ID	PRODUCTLINE	MSRP	PRO
0	30	95.70		2	2871.00	2003-02-24	1	2	2003	Motorcycles	95
1	34	81.35		5	2765.90	2003-05-07	2	5	2003	Motorcycles	95
2	41	94.74		2	3884.34	2003-07-01	3	7	2003	Motorcycles	95
3	45	83.26		6	3746.70	2003-08-25	3	8	2003	Motorcycles	95
4	49	100.00		14	5205.27	2003-10-10	4	10	2003	Motorcycles	95
...	...	...	...	...	...	...	...	...	...	...	...
2818	20	100.00		15	2244.40	2004-12-02	4	12	2004	Ships	54
2819	29	100.00		1	3978.51	2005-01-31	1	1	2005	Ships	54
2820	43	100.00		4	5417.57	2005-03-01	1	3	2005	Ships	54
2821	34	62.24		1	2116.16	2005-03-28	1	3	2005	Ships	54
2822	47	65.52		9	3079.44	2005-05-06	2	5	2005	Ships	54

2823 rows x 13 columns

```
In [ ]: sales_df = dummies('PRODUCTLINE')
```

```
In [ ]: sales_df
```

sales\_df

	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	QTR_ID	MONTH_ID	YEAR_ID	MSRP	PRODUCTCODE	DE
0	30	95.70		2	2871.00	2003-02-24	1	2	2003	95	510,1678
1	34	81.35		5	2765.90	2003-05-07	2	5	2003	95	510,1678
2	41	94.74		2	3884.34	2003-07-01	3	7	2003	95	510,1678
3	45	83.26		6	3746.70	2003-08-25	3	8	2003	95	510,1678
4	49	100.00		14	5205.27	2003-10-10	4	10	2003	95	510,1678
...	...	...	...	...	...	...	...	...	...	...	...
2818	20	100.00		15	2244.40	2004-12-02	4	12	2004	54	572,3212
2819	29	100.00		1	3978.51	2005-01-31	1	1	2005	54	572,3212
2820	43	100.00		4	5417.57	2005-03-01	1	3	2005	54	572,3212
2821	34	62.24		1	2116.16	2005-03-28	1	3	2005	54	572,3212
2822	47	65.52		9	3079.44	2005-05-06	2	5	2005	54	572,3212

2823 rows x 13 columns

```
In [ ]: sales_df = dummies('DEALSIZE')
```

sales\_df

	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	QTR_ID	MONTH_ID	YEAR_ID	MSRP	PRODUCTCODE	Aus
0	30	95.70		2	2871.00	2003-02-24	1	2	2003	95	510,1678
1	34	81.35		5	2765.90	2003-05-07	2	5	2003	95	510,1678
2	41	94.74		2	3884.34	2003-07-01	3	7	2003	95	510,1678
3	45	83.26		6	3746.70	2003-08-25	3	8	2003	95	510,1678
4	49	100.00		14	5205.27	2003-10-10	4	10	2003	95	510,1678
...	...	...	...	...	...	...	...	...	...	...	...
2818	20	100.00		15	2244.40	2004-12-02	4	12	2004	54	572,3212
2819	29	100.00		1	3978.51	2005-01-31	1	1	2005	54	572,3212
2820	43	100.00		4	5417.57	2005-03-01	1	3	2005	54	572,3212
2821	34	62.24		1	2116.16	2005-03-28	1	3	2005	54	572,3212
2822	47	65.52		9	3079.44	2005-05-06	2	5	2005	54	572,3212

2823 rows x 13 columns

```
In [ ]: # Agrupamos los datos según la fecha del pedido
sales_df_group = sales_df.groupby(by=['ORDERDATE']).sum()
```

sales\_df\_group

ORDERDATE	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	QTR_ID	MONTH_ID	YEAR_ID	MSRP	PRODUCTCODE	Australia
2003-01-06	151	288.78	10	12133.25	4	4	8012	363	174.0	0
2003-01-09	142	284.96	10	11432.34	4	4	8012	372	181.0	0
2003-01-10	80	150.14	3	6864.05	2	2	4006	155	37.0	0
2003-01-29	541	1417.54	136	54702.00	16	16	32048	1695	723.0	0
2003-01-31	443	1061.89	91	44621.96	13	13	26039	1365	720.0	0
...	...	...	...	...	...	...	...	...	...	...
2005-05-13	259	561.18	21	31821.90	12	30	12830	728	101.0	0
2005-05-17	509	1269.43	105	59475.10	28	70	28070	1669	462.0	0
2005-05-29	607	1148.40	94	51233.18	30	75	30075	1328	797.0	13
2005-05-30	187	542.16	18	14578.75	14	35	14035	618	265.0	0
2005-05-31	696	1561.40	112	78918.03	38	95	38095	2065	899.0	0

252 rows x 11 columns

```
In [ ]: fig = px.line(x=sales_df_group.index, y=sales_df_group.SALES, title='Sales')
fig.show()
```



```
In [ ]: # Podemos eliminar 'ORDERDATE' y quedarnos con el resto de datos relacionados con la fecha como 'MONTH'
sales_df.drop('ORDERDATE', axis=1, inplace=True)
sales_df.shape
```

sales\_df

	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	QTR_ID	MONTH_ID	YEAR_ID	MSRP	PRODUCTCODE	Aus
0	30	95.70		2	2871.00	2003-02-24	1	2	2003	95	510,1678
1	34	81.35		5	2765.90	2003-05-07	2	5	2003	95	510,1678
2	41	94.74		2	3884.34	2003-07-01	3	7	2003	95	510,1678
3	45	83.26		6	3746.70	2003-08-25	3	8	2003	95	510,1678
4	49	100.00		14	5205.27	2003-10-10	4	10	2003	95	510,1678
...	...	...	...	...	...	...	...	...	...	...	...
2818	20	100.00		15	2244.40	2004-12-02	4	12	2004	54	572,3212
2819	29	100.00		1	3978.51	2005-01-31	1	1	2		

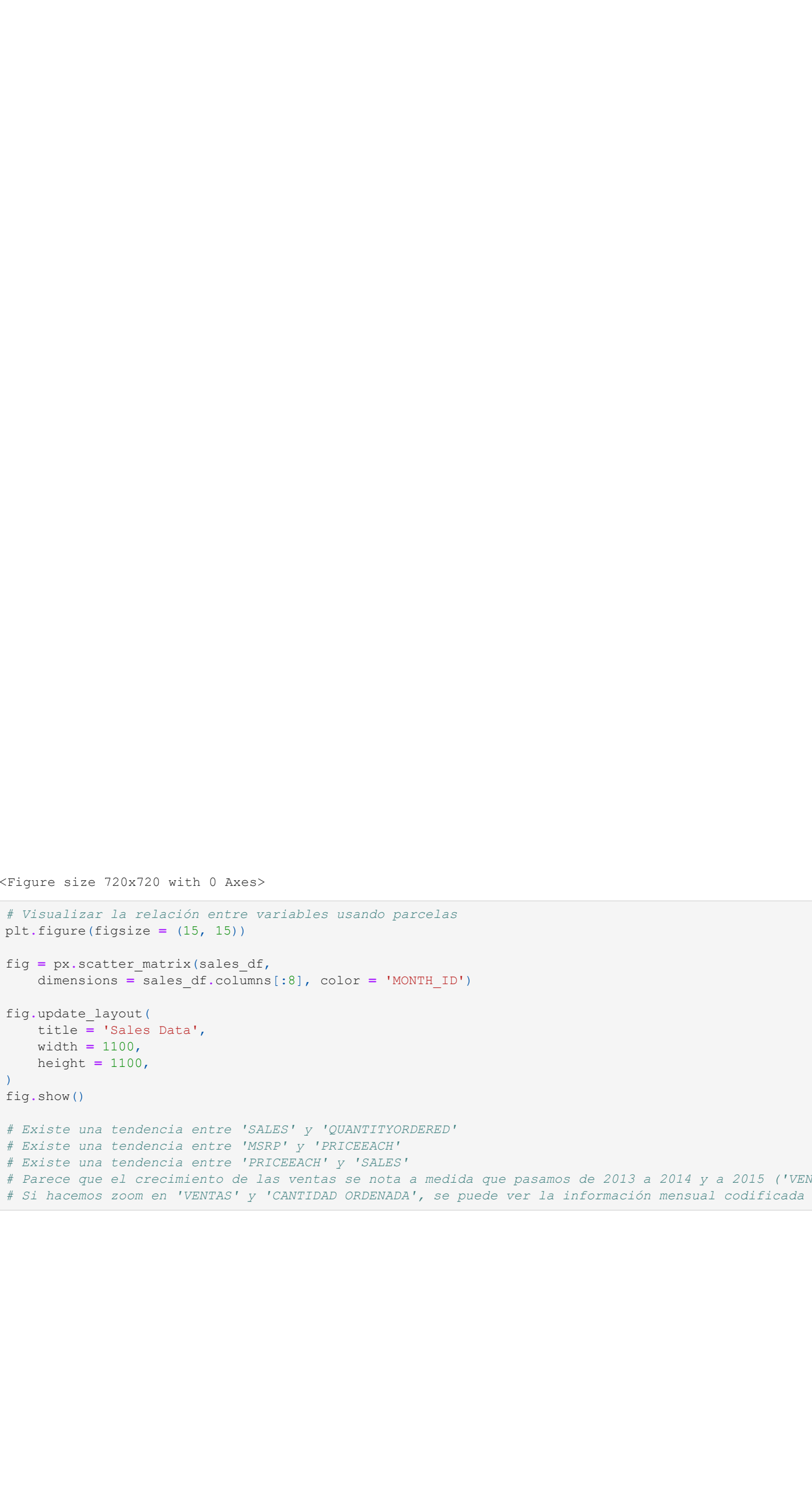


```
# Visualizar la relación entre variables usando parcelas
plt.figure(figsize = (15, 15))

fig = px.scatter_matrix(sales_df,
                        dimensions = sales_df.columns[:8], color = 'MONTH_ID')

fig.update_layout(
    title = 'Sales Data',
    width = 1100,
    height = 1100,
)
fig.show()
```

# Existe una tendencia entre 'SALES' y 'QUANTITYORDERED'  
# Existe una tendencia entre 'MSRP' y 'PRICEEACH'  
# Existe una tendencia entre 'PRICEEACH' y 'SALES'  
# Parece que el crecimiento de las ventas se nota a medida que pasamos de 2013 a 2014 y a 2015 ('VENTAS' frente a 'VENTAS')  
# Si hacemos zoom en 'VENTAS' y 'CANTIDAD ORDENADA', se puede ver la información mensual codificada por colores



<Figure size 720x720 with 0 Axes>

```
In [ ]: # Visualizar la relación entre variables usando parcelas
plt.figure(figsize = (15, 15))

fig = px.scatter_matrix(sales_df,
                        dimensions = sales_df.columns[:8], color = 'MONTH_ID')

fig.update_layout(
    title = 'Sales Data',
    width = 1100,
    height = 1100,
)
fig.show()
```

```
# Existe una tendencia entre 'SALES' y 'QUANTITYORDERED'
# Existe una tendencia entre 'MSRP' y 'PRICEEACH'
# Existe una tendencia entre 'PRICEEACH' y 'SALES'
# Parece que el crecimiento de las ventas se nota a medida que pasamos de 2013 a 2014 y a 2015 ('VENTAS' frente a 'VENTAS')
# Si hacemos zoom en 'VENTAS' y 'CANTIDAD ORDENADA', se puede ver la información mensual codificada por colores
```

```
<Figure size 1080x1080 with 0 Axes>
```

## ENCONTRAR EL NÚMERO ÓPTIMO DE CLUSTERS UTILIZANDO EL MÉTODO DE CODO

```
In [ ]: # Escalamos los datos
scaler = StandardScaler()
sales_df_scaled = scaler.fit_transform(sales_df)
```

```
In [ ]: sales_df_scaled.shape
```

```
Out [ ]: (2823, 37)
```

```
In [ ]: scores = []

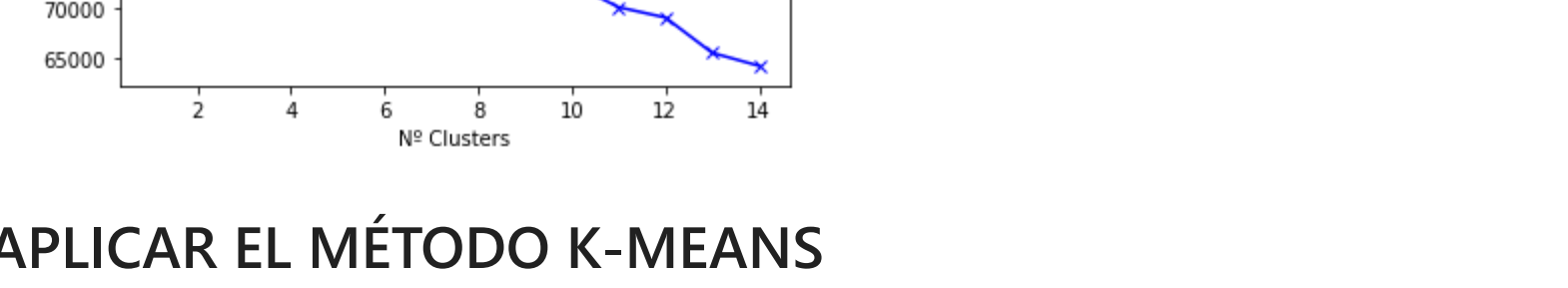
range_values = range(1, 15)

for i in range_values:
    kmeans = KMeans(n_clusters = i)
    kmeans.fit(sales_df_scaled)
    scores.append(kmeans.inertia_) # la inercia es la suma de los cuadrados de las distancias de las observaciones a los centros de los clusters

plt.plot(range_values, scores, 'bx-')
plt.title('Encontrar el número correcto de clusters')
plt.xlabel('Nº Clusters')
plt.ylabel('WCSS')
```

```
plt.show()
```

# A partir de esto podemos observar que, el 5º grupo parece estar cambiando el codo de la curva.  
# Tengamos en cuenta que la curva cambia cada vez que ejecutemos la celda



<Figure size 1080x1080 with 0 Axes>

## APLICAR EL MÉTODO K-MEANS

```
In [ ]: # Agrupar los datos usando k-means
kmeans = KMeans(5)
kmeans.fit(sales_df_scaled)
labels = kmeans.labels_
```

```
In [ ]: labels
```

```
Out [ ]: array([2, 2, 1, ..., 4, 4, 4], dtype=int32)
```

```
In [ ]: kmeans.cluster_centers_.shape
```

```
Out [ ]: (5, 37)
```

```
In [ ]: # Echemos un vistazo a los centros de los clusters
cluster_centers = pd.DataFrame(data = kmeans.cluster_centers_, columns = [sales_df.columns])
```

```
Out [ ]:
```

```
0 0.013139 -0.396807 0.160159 -0.334319 0.010234 -0.032711 -0.690094 0.633095 -0.212338 -0.140961 0
1 0.267137 0.599614 -0.027845 0.480034 0.009476 -0.016772 0.440432 -0.275865 0.009749 0.026760 -0
2 -0.471508 -0.766219 0.039814 -0.825614 0.009366 0.003008 -0.577834 0.154705 0.064312 -0.023103 -0
3 1.246891 0.800157 -0.255984 2.571482 -0.081396 0.126876 1.446248 -0.885353 -0.083492 0.044555 -0
4 -0.037171 0.009763 0.078809 -0.271913 -0.043991 0.001642 -0.363043 1.126176 -0.230280 -0.017283 0
```

```
In [ ]: # Para entender lo que significan estos números, realicemos una transformación inversa
cluster_centers = scaler.inverse_transform(cluster_centers)
cluster_centers = pd.DataFrame(data = cluster_centers, columns = [sales_df.columns])
```

```
Out [ ]:
```

```
0 35.220779 75.654673 7.142857 2938.226883 7.129870 2003.792208 72.987013 45.061551 0.012987 -1.3877
1 37.690442 95.753173 6.348521 4437.891103 7.127098 2003.803357 118.412470 53.766234 0.007946 2.3181
2 30.540452 68.203376 6.634389 2033.480287 7.126697 2003.817195 77.497738 58.658824 0.081443 1.6289
3 47.237179 99.798269 5.384615 8289.373141 6.794872 2003.903846 158.826923 25.814103 0.044872 2.5641
4 34.730769 83.855470 6.799145 3053.150128 6.931624 2003.816239 86.128205 89.337607 0.008547 1.7094
```

```
In [ ]: labels.shape # Etiquetas del clúster asociado a cada observación
```

```
Out [ ]: (2823,)
```

```
In [ ]: labels.max()
```

```
Out [ ]: 4
```

```
In [ ]: labels.min()
```

```
Out [ ]: 0
```

```
In [ ]: y_kmeans = kmeans.fit_predict(sales_df_scaled)
```

```
Out [ ]: array([1, 1, 2, ..., 2, 1, 2], dtype=int32)
```

```
In [ ]: y_kmeans.shape
```

```
Out [ ]: (2823,)
```

```
In [ ]: # Agrupar una etiqueta (que clúster) correspondiente a cada punto de datos
sales_df_cluster = pd.concat([sales_df, pd.DataFrame({'cluster': labels})], axis = 1)
sales_df_cluster
```

```
Out [ ]:
```

```
0 30 95.70 2 2871.00 2 2003 95 0 0 0
1 34 81.35 5 2765.90 5 2003 95 0 0 0
2 41 94.74 2 3884.34 7 2003 95 0 0 0
3 45 83.26 8 2003 95 0 0 0
4 49 100.00 14 5205.27 10 2003 95 0 0 0
...
```

```
2818 20 100.00 15 2248.40 12 2004 54 108 0 0
2819 29 100.00 1 3978.51 1 2005 54 108 0 0
2820 43 100.00 4 5417.57 3 2005 54 108 0 0
2821 34 62.24 1 2116.16 3 2005 54 108 0 0
2822 47 65.52 9 3079.44 5 2005 54 108 0 0
```

2823 rows x 38 columns

```
In [ ]: sales_df['ORDERLINENUMBER'] = sales_df['ORDERLINENUMBER'].apply(lambda x: float(x))
```

```
In [ ]: # Representar un histograma para cada característica según el clúster al que pertenece
for i in range_values:
    plt.figure(figsize = (30, 6))
    for j in range(5):
        plt.subplot(1, 5, j+1)
        cluster = sales_df_cluster[sales_df_cluster['cluster'] == j]
        cluster[i].hist()
        plt.title('{}\nCluster - {}'.format(i, j))
    plt.show()
```



```
In [ ]: # Reducir los datos originales a 3 dimensiones usando PCA para visualizar los clusters
pca = PCA(n_components = 3)
principal_comp = pca.fit_transform(sales_df_scaled)
principal_comp
```

```
Out [ ]: array([[0.48638335, -1.18924267, 1.02227436],
               [-0.78471591, -0.147987, 1.40066622],
               [1.43140417, 0.141865, 0.26694045],
               ...,
               [0.3897192, 4.05662552, -0.11335055],
               [-2.7889184, 2.35317929, 1.52293978],
               [-0.7055867, 3.27631066, -0.52178646]])
```

```
In [ ]: pca_df = pd.DataFrame(data = principal_comp, columns = ['pca1', 'pca2', 'pca3'])
pca_df.head()
```

```
Out [ ]:
```

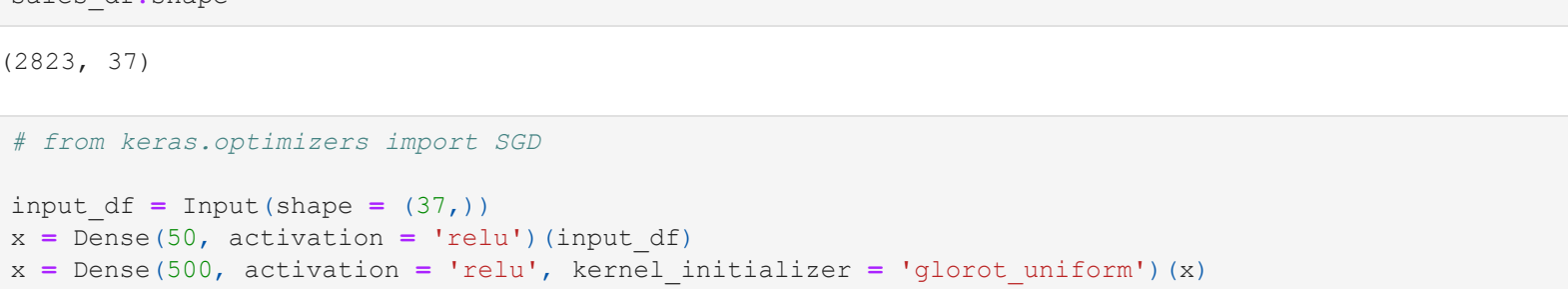
```
0 -0.486383 -1.189243 1.022274
1 -0.784716 -0.14797 1.400666
2 1.431404 0.141865 0.266940
3 1.288338 -0.211897 -0.574293
4 1.975904 -0.123142 -1.062156
```

```
...
```

```
2818 -2.450258 1.083166 -1.601567
2819 -0.158456 3.649761 0.061848
2820 0.389972 4.056626 -0.113351
2821 -2.788918 2.353179 1.522940
2822 -0.705587 3.276311 -0.521786
```

2823 rows x 4 columns

```
In [ ]: # Visualizar los clusters con 3D-Scatterplot
fig = px.scatter_3d(pca_df, x = 'pca1', y = 'pca2', z = 'pca3',
                    color = 'cluster', symbol = 'cluster', size_max = 18, opacity = 0.7)
fig.update_layout(margin=dict(l=0, b=0, t=0, r=0))
```



## APLICAR AUTOENCODERS PARA REALIZAR LA REDUCCIÓN DE DIMENSIONALIDAD

```
In [ ]: sales_df.shape
```

```
Out [ ]: (2823, 37)
```

```
In [ ]: # from keras.optimizers import SGD

input_df = Input(shape = (37,))
x = Dense(50, activation = 'relu', kernel_initializer = 'glorot_uniform')(input_df)
x = Dense(500, activation = 'relu', kernel_initializer = 'glorot_uniform')(x)
x = Dense(2000, activation = 'relu', kernel_initializer = 'glorot_uniform')(x)
encoded = Dense(8, activation = 'relu', kernel_initializer = 'glorot_uniform')(x)
x = Dense(2000, activation = 'relu', kernel_initializer = 'glorot_uniform')(encoded)
x = Dense(500, activation = 'relu', kernel_initializer = 'glorot_uniform')(x)
decoded = Dense(37, kernel_initializer = 'glorot_uniform')(x)
```

```
# autoencoder
autoencoder = Model(input_df, decoded)
```

```
# encoder - utilizado para reducir la dimensión
encoder = Model(input_df, encoded)
```

```
autoencoder.compile(optimizer = 'adam', loss='mean_squared_error')
```

```
In [ ]: autoencoder.fit(sales_df_scaled, sales_df_scaled, batch_size = 128, epochs = 500, verbose = 3)
```



