

IMPORTAR LAS LIBRERÍAS Y LOS DATASETS

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import datetime
import warnings
warnings.filterwarnings("ignore")
```

IMPORTAR EL DATASET DE VENTAS

```
sales_train_df = pd.read_csv("train.csv")
```

```
sales_train_df.head(5)
# Un millón de observaciones
# 1115 tiendas únicas
# Las ventas de la variable objetivo (la que se intenta predecir)
# Id: ID de transacción (combinación de la tienda y la fecha)
# Store: identificación única de la tienda
# Sales: ventas diarias, esta es la variable objetivo
# Customers: número de clientes de un día dado
# Open: booleano para indicar si la tienda estaba abierta o cerrada (0 = cerrada, 1 = abierta)
# Promo: describe si la tienda tenía algún tipo de promoción ese día o no
# StateHoliday: indica si el día era festivo o no (a = vacaciones públicas, b = vacaciones de Pascua, c = Navidad)
# SchoolHoliday: indica si la tienda, Date) se ve afectado por el cierre de las escuelas públicas
```

	Store	DayOfWeek	Date	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday
0	1	5	2015-07-31	5263	555	1	1	0	1
1	2	5	2015-07-31	6064	625	1	1	0	1
2	3	5	2015-07-31	8314	821	1	1	0	1
3	4	5	2015-07-31	13995	1498	1	1	0	1
4	5	5	2015-07-31	4822	559	1	1	0	1

```
sales_train_df.tail(10)
```

	Store	DayOfWeek	Date	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday
1017199	1106	2	2013-01-01	0	0	0	0	a	1
1017200	1107	2	2013-01-01	0	0	0	0	a	1
1017201	1108	2	2013-01-01	0	0	0	0	a	1
1017202	1109	2	2013-01-01	0	0	0	0	a	1
1017203	1110	2	2013-01-01	0	0	0	0	a	1
1017204	1111	2	2013-01-01	0	0	0	0	a	1
1017205	1112	2	2013-01-01	0	0	0	0	a	1
1017206	1113	2	2013-01-01	0	0	0	0	a	1
1017207	1114	2	2013-01-01	0	0	0	0	a	1
1017208	1115	2	2013-01-01	0	0	0	0	a	1

```
sales_train_df.info()
# 9 columnas en total
# 8 variables objetivas, cada una con 1017209 puntos de datos
# 1 variable objetivo (ventas)

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1017209 entries, 0 to 1017208
Data columns (total 9 columns):
# Column Non-Null Count Dtype
---
0 Store 1017209 non-null int64
1 DayOfWeek 1017209 non-null int64
2 Date 1017209 non-null object
3 Sales 1017209 non-null int64
4 Customers 1017209 non-null int64
5 Open 1017209 non-null int64
6 Promo 1017209 non-null int64
7 StateHoliday 1017209 non-null object
8 SchoolHoliday 1017209 non-null int64
dtypes: float64(0), int64(2), object(5)
memory usage: 69.8+ MB
```

```
sales_train_df.describe()
# Cantidad de ventas promedio por día = 5773 Euros, ventas mínimas por día = 0, ventas máximas por día = 41551
# Número medio de clientes = 633, número mínimo de clientes = 0, número máximo de clientes = 7388
```

	Store	DayOfWeek	Date	Sales	Customers	Open	Promo	StateHoliday
count	1017209+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06	1.017209e+06
mean	5.584297e+02	3.998341e+00	5.773819e+03	6.334136e+02	8.301067e+01	3.815145e-01	1.786467e-01	1.074676e-01
std	3.219087e+02	1.967931e+00	3.849926e+03	4.644177e+02	3.755392e-01	4.857586e-01	3.830564e-01	3.830564e-01
min	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	2.000000e+02	2.000000e+00	3.727000e+03	4.050000e+02	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
50%	5.580000e+02	4.000000e+00	5.744000e+03	6.090000e+02	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
75%	8.380000e+02	6.000000e+00	7.856000e+03	8.370000e+02	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00
max	1.150000e+03	7.000000e+00	4.155100e+04	7.380000e+03	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00

IMPORTAR LA INFORMACIÓN SOBRE LAS TIENDAS

```
store_info_df = pd.read_csv("store.csv")
# store_type: categoría que indica el tipo de tienda (a, b, c, d)
# Assortment: a = básico, b = extra, c = extendido
# CompetitionDistance (en metros): distancia a la tienda de la competencia más cercana
# CompetitionOpenSinceMonth: fecha en que abrió la competencia
# Promo2: Promoción de una promoción continuada y consecutiva en algunas tiendas (0 = la tienda no participa, 1 = Participa)
# Promo2SinceWeek (Año/Semana): fecha en la que la tienda empieza a participar en la Promo2
# PromoInterval: describe los intervalos consecutivos donde la Promo2 empieza, indicando los meses en los que se repite
```

```
store_info_df.head(5)
```

	Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear
0	1	c	a	1270.0	9.0	2008.0	0	NaN	NaN
1	2	a	a	570.0	11.0	2007.0	0	NaN	13.0
2	3	a	a	14130.0	12.0	2006.0	1	NaN	14.0
3	4	c	c	620.0	9.0	2009.0	0	NaN	NaN
4	5	a	a	29910.0	4.0	2015.0	0	NaN	NaN

```
# Hay que tener en cuenta que el data frame anterior incluye las transacciones registradas por día (en millones)
# Este data frame solo incluye información sobre las 1115 tiendas exclusivas que forman parte de este estudio
```

```
store_info_df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1115 entries, 0 to 1114
Data columns (total 10 columns):
# Column Non-Null Count Dtype
---
0 Store 1115 non-null int64
1 StoreType 1115 non-null object
2 Assortment 1115 non-null object
3 CompetitionDistance 1112 non-null float64
4 CompetitionOpenSinceMonth 761 non-null int64
5 CompetitionOpenSinceYear 761 non-null int64
6 Promo2 1115 non-null int64
7 Promo2SinceWeek 571 non-null float64
8 Promo2SinceYear 571 non-null float64
9 PromoInterval 1115 non-null object
dtypes: float64(3), int64(2), object(5)
memory usage: 87.2+ KB
```

```
store_info_df.describe()
# De media, la distancia de la competencia es de 5404 metros (5,4 kms)
```

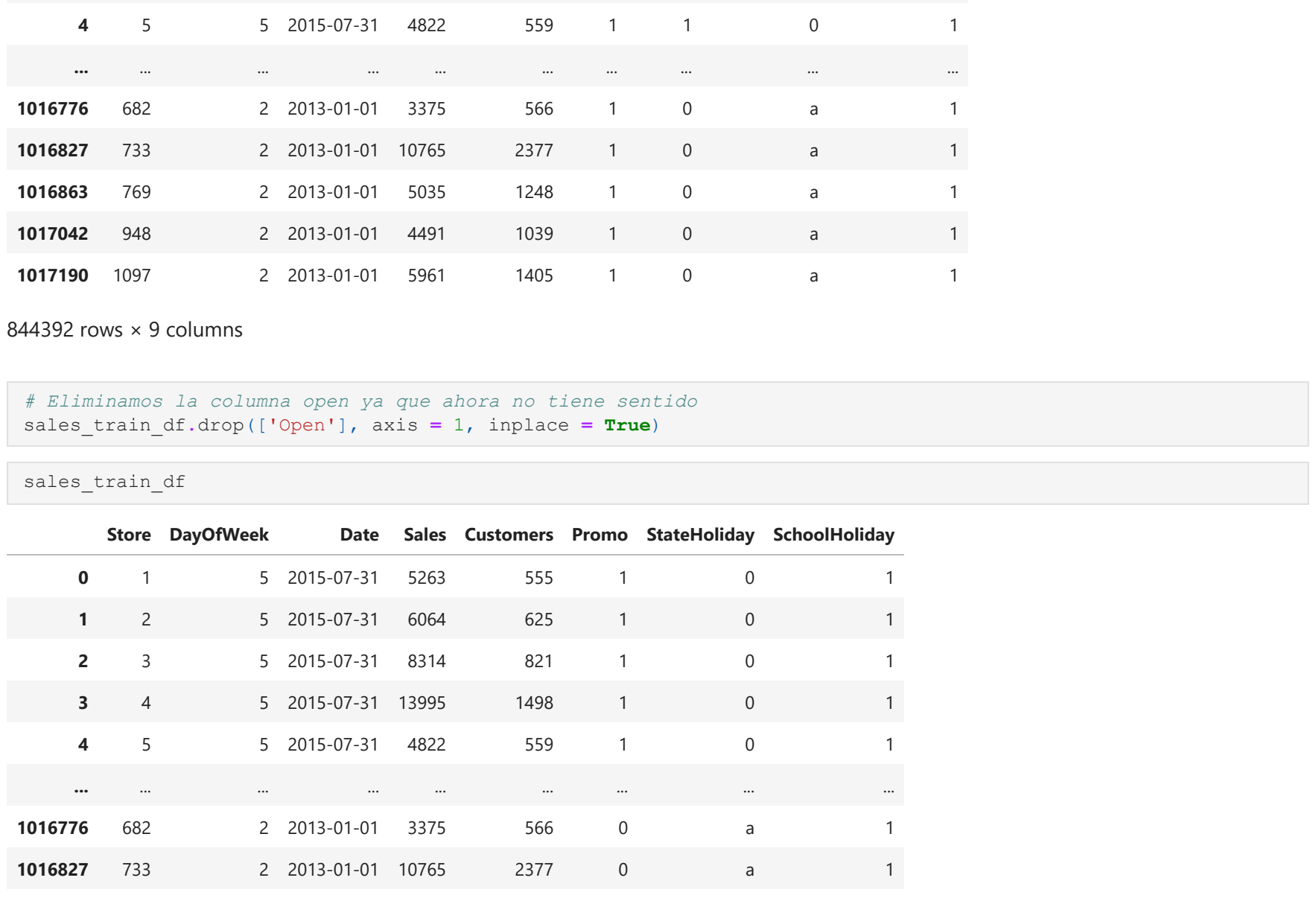
	Store	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek	Promo2SinceYear
count	1115	0.000000	1112.000000	761.000000	761.000000	0.512108	571.000000
mean	558.00000	5040.901079	7.221704	6.199583	0.500078	14.141984	2011.76
std	332.01708	7063.174720	3.123248	3.755392	0.500078	14.141984	1.67
min	1.00000	20.00000	1.000000	1900.000000	0.000000	1.000000	2009.00
25%	279.50000	717.50000	4.000000	2006.000000	0.000000	13.000000	2012.00
50%	558.00000	2325.00000	8.000000	2010.000000	1.000000	22.000000	2010.00
75%	856.50000	6882.50000	10.000000	2013.000000	1.000000	37.000000	2013.00
max	1115.00000	75860.00000	12.000000	2015.000000	1.000000	50.000000	2015.00

EXPLORAR EL DATASET

EXPLORAR EL DATASET DE VENTAS

```
# Verificamos si no faltan datos
sns.heatmap(sales_train_df.isnull(), yticklabels=False, cbar=False, cmap = "Blues")
# Comprobamos que no tenemos datos nulos
```

```
Out[70]: <matplotlib.axes._subplots.AxesSubplot at 0x1c6aff7b720>
```



```
sales_train_df.hist(bins = 30, figsize=(20,20), color = 'b')
# Promedio de 600 clientes por día, el máximo es 4500 (tenga en cuenta que no podemos ver el valor atípico en
# Los datos se distribuyen por igual en varios días de la semana (~ 150000 observaciones x 7 días = ~ 1.1 mil millones de datos)
# Los datos se distribuyen por igual entre todas las tiendas (sin sesgo)
# La promoción # 1 se ejecutó aproximadamente el 40% del tiempo
# Ventas promedio alrededor de 5000-6000 Euros
# Las vacaciones escolares duran alrededor del 18% del tiempo
```

```
array([[<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF01134A0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF01169F0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF011E280D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>,
<matplotlib.axes._subplots.AxesSubplot object at 0x000001CF001013B0D>],
      dtype=object)
```

```
sales_train_df["Customers"].max()
7388

# Comprobamos cuántas tiendas están abiertas y cerradas
closed_train_df = sales_train_df[sales_train_df["Open"] == 0]
open_train_df = sales_train_df[sales_train_df["Open"] == 1]

# Se reutilizan el número de tiendas que están abiertas y cerradas
print("Total = {}".format(len(sales_train_df)))
print("Número de tiendas abiertas = {}".format(len(open_train_df)))
print("Número de tiendas cerradas = {}".format(len(closed_train_df)))
print("Porcentaje de tiendas cerradas = {}".format((len(closed_train_df)/len(sales_train_df))*100))

Total = 1017209
Número de tiendas abiertas = 844392
Número de tiendas cerradas = 172817
Porcentaje de tiendas cerradas = 16.989330609340426

# Nos quedamos solo con las tiendas abiertas y eliminamos las tiendas cerradas
sales_train_df = sales_train_df[sales_train_df["Open"] == 1]
```

```
sales_train_df
```

```
Out[76]:
```

	Store	DayOfWeek	Date	Sales	Customers	Open	Promo	StateHoliday	SchoolHoliday
0	1	5	2015-07-31	5263	555	1	1	0	1
1	2	5	2015-07-31	6064	625	1	1	0	1
2	3	5	2015-07-31	8314	821	1	1	0	1
3	4	5	2015-07-31	13995	1498	1	1	0	1
4	5	5	2015-07-31	4822	559	1	1	0	1
...
1016776	682	2	2013-01-01	3375	566	1	0	a	1
1016827	733	2	2013-01-01	10765	2377	1	0	a	1
1016863	769	2	2013-01-01	5035	1248	1	0	a	1
1017042	948	2	2013-01-01	4491	1039	1	0	a	1
1017190	1097	2	2013-01-01	5961	1405	1	0	a	1

844392 rows x 9 columns

```
# Eliminamos la columna open ya que ahora no tiene sentido
sales_train_df.drop(["Open"], axis = 1, inplace = True)
```

```
sales_train_df
```

```
Out[78]:
```

	Store	DayOfWeek	Date	Sales	Customers	Promo	StateHoliday	SchoolHoliday
0	1	5	2015-07-31	5263	555	1	0	1
1	2	5	2015-07-31	6064	625	1	0	1
2	3	5	2015-07-31	8314	821	1	0	1
3	4	5	2015-07-31	13995	1498	1	0	1
4	5	5	2015-07-31	4822	559	1	0	1
...
1016776	682	2	2013-01-01	3375	566	0	a	1
1016827	733	2	2013-01-01	10765	2377	0	a	1
1016863	769	2	2013-01-01	5035	1248	0	a	1
1017042	948	2	2013-01-01	4491	1039	0	a	1
1017190	1097	2	2013-01-01	5961	1405	0	a	1

844392 rows x 8 columns

```
# Ventas promedio = 6955 Euros, número promedio de clientes = 762 (ha subido)
```

```
sales_train_df.describe()
```

	Store	DayOfWeek	Date	Sales	Customers	Promo	StateHoliday
count	844392.000000	844392.000000	844392.000000	844392.000000	844392.000000	844392.000000	844392.000000
mean	558.4297	3.998341	5773.819	633.4136	762.728395	0.446352	0.193580
std	322.0171914	1.723689	31904.214680	401.227674	0.497114	0.395103	0.395103
min	1.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	280.000000	2.000000	4859.000000	519.000000	0.000000	0.000000	0.000000
50%	558.000000	3.000000	6369.000000	676.000000	0.000000	0.000000	0.000000
75%	837.000000	5.000000	8360.000000	893.000000	1.000000	0.000000	0.000000
max	1115.000000	7.000000	41551.000000	7388.000000	1.000000	1.000000	1.000000

EXPLORAR LOS DATOS DE LA INFORMACIÓN DE LAS TIENDAS

```
# Comprobamos si falta algún dato en el data frame de información de la tienda
sns.heatmap(store_info_df.isnull(), yticklabels=False, cbar=False, cmap = "Blues")
```

```
Out[80]: <matplotlib.axes._subplots.AxesSubplot at 0x1c6f80b1c670>
```

```
# Verificamos los valores faltantes en la 'CompetitionDistance'
# Solo faltan 3 filas
store_info_df[store_info_df["CompetitionDistance"].isnull()]

Store StoreType Assortment CompetitionDistance CompetitionOpenSinceMonth CompetitionOpenSinceYear Promo2 Promo2SinceWeek
290 291 d a NaN NaN NaN NaN 0 NaN
621 622 a c NaN NaN NaN NaN 0 NaN
878 879 d a NaN NaN NaN NaN 1 50
```

```
store_info_df[store_info_df["Promo2"] == 0]
```

	Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek
11	12	a	c	1070.0	NaN	NaN	NaN	1
12	13	d	a	310.0	NaN	NaN	NaN	1
15	16	a	c	3270.0	NaN	NaN	NaN	0
18	19	a	c	3240.0	NaN	NaN	NaN	1
21	22	a	a	1040.0	NaN	NaN	NaN	1
...
1095	1096	a	c	1130.0	NaN	NaN	NaN	1
1099	1100	a	a	540.0	NaN	NaN	NaN	1
1112	1113	a	c	9260.0	NaN	NaN	NaN	0
1113	1114	a	c	870.0	NaN	NaN	NaN	0
1114	1115	d	c	5350.0	NaN	NaN	NaN	1

354 rows x 10 columns

```
store_info_df[store_info_df["Promo2"] == 0]
```

	Store	StoreType	Assortment	CompetitionDistance	CompetitionOpenSinceMonth	CompetitionOpenSinceYear	Promo2	Promo2SinceWeek
0	1	c	a	1270.0	9.0	2008.0	0	NaN
3	4	c	c	620.0	9.0	2009.0	0	NaN
4	5	a	a	29910.0	4.0	2015.0	0	NaN
5	6	a	a	310.0	12.0	2013.0	0	NaN
...
1107	1108	a	a	540.0	9.0	2004.0	0	NaN
1109	1110	c	c	900.0	9.0	2010.0	0	NaN
1111	1112	c	c	980.0	9.0	2006.0	0	NaN
1112	1113	a	c	1260.0	NaN	NaN	0	NaN
1113	1114	a	c	870.0	NaN	NaN	0	NaN

544 rows x 10 columns

```
# Parece que si 'promo2' es cero, 'promo2SinceWeek' y 'promo2SinceYear' y la información de 'PromoInterval' se ve afectada
# Hay 354 filas con valores de 'CompetitionDistance' que faltan, llenándolas con el valor promedio de la columna
# Establecemos estos valores en ceros
store_info_df["CompetitionDistance"] = store_info_df["CompetitionDistance"].mean().fillna(0, inplace = True)

for str in str_cols:
    store_info_df[str].fillna(0, inplace = True)
```

```
sns.heatmap(store_info_df.isnull(), yticklabels=False, cbar=False, cmap = "Blues")
```

```
Out[85]: <matplotlib.axes._
```