

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings("ignore")
```

1. Importe la base de datos a una base en Jupyter Notebook con pandas.

```
In [174]: df=pd.read_csv('base1heart.csv')
Out[174]:
```

	Store	Date	Weekly_Sales	Holiday_Flag	Temperature	Fuel_Price	CP1	Unemployment
0	1	05-02-2010	1643690.90	0	42.31	2.572	211.096358	8.106
1	1	12-02-2010	1641957.44	1	38.51	2.548	211.242170	8.106
2	1	19-02-2010	1611968.17	0	39.93	2.514	211.289143	8.106
3	1	26-02-2010	1409727.59	0	46.63	2.561	211.319643	8.106
4	1	05-03-2010	1554096.68	0	46.50	2.625	211.350143	8.106
...
6430	45	28-09-2012	713173.95	0	64.88	3.997	192.015558	8.664
6431	45	05-10-2012	729455.07	0	64.89	3.985	192.170412	8.667
6432	45	12-10-2012	734464.36	0	54.47	4.000	192.327265	8.667
6433	45	19-10-2012	718125.53	0	56.47	3.969	192.330854	8.667
6434	45	26-10-2012	760281.43	0	58.85	3.882	192.308899	8.667
...
6435	rows > 8 columns							

```
In [175]: df.rename({'Store':'Tienda', 'Date':'Fecha','Weekly_Sales':'VentasSemanales','Holiday_Flag':'SemanaEspecial','Temperature':'Temperatura','Fuel_Price':'PrecioCombustible','CP1':'IPC'})
Out[175]:
```

	Tienda	Fecha	VentasSemanales	SemanaEspecial	Temperatura	PrecioCombustible	IPC	Desempleo
0	1	05-02-2010	1643690.90	0	42.31	2.572	211.096358	8.106
1	1	12-02-2010	1641957.44	1	38.51	2.548	211.242170	8.106
2	1	19-02-2010	1611968.17	0	39.93	2.514	211.289143	8.106
3	1	26-02-2010	1409727.59	0	46.63	2.561	211.319643	8.106
4	1	05-03-2010	1554096.68	0	46.50	2.625	211.350143	8.106
...
6430	45	28-09-2012	713173.95	0	64.88	3.997	192.015558	8.664
6431	45	05-10-2012	729455.07	0	64.89	3.985	192.170412	8.667
6432	45	12-10-2012	734464.36	0	54.47	4.000	192.327265	8.667
6433	45	19-10-2012	718125.53	0	56.47	3.969	192.330854	8.667
6434	45	26-10-2012	760281.43	0	58.85	3.882	192.308899	8.667
...
6435	rows > 8 columns							

VARIABLES DUMMIES FECHA

```
In [176]: df_dummies = pd.get_dummies(df, columns=['Fecha'], dummy_na=True)
```

2. Obtenga los descriptivos resumen de la base de datos e identifique a las variables numéricas y categóricas. ¿Hay algo que le llame la atención?

La variable `SemanaEspecial` parece indicar si una semana es especial o no. La mayoría de las semanas tienen un valor de 0, lo que sugiere que las semanas especiales son menos comunes en los datos.

Las variables `Temperatura`, `PrecioCombustible`, `IPC` y `Desempleo` muestran una variación considerable en sus valores. Esto podría indicar que estas variables tienen un impacto en las ventas semanales y podrían ser importantes para el análisis.

Los valores extremos en algunas de las variables, como la temperatura mínima (2.06) y la tasa de desempleo máxima (14.313), podrían requerir un análisis adicional para comprender su impacto en las ventas.

```
In [177]: df.describe()
Out[177]:
```

	Tienda	VentasSemanales	SemanaEspecial	Temperatura	PrecioCombustible	IPC	Desempleo
count	6435.000000	6.435000e+03	6435.000000	6435.000000	6435.000000	6435.000000	6435.000000
mean	23.000000	1.608959e+06	0.099950	60.662762	3.269607	171.576394	7.999151
std	12.908505	1.640056e+06	0.236569	18.444923	0.408020	39.366712	1.879585
min	1.000000	2.099629e+05	0.000000	27.060000	2.472000	126.064000	3.879000
25%	12.000000	5.537001e+05	0.000000	47.460000	2.933000	131.739000	6.810000
50%	23.000000	9.607400e+05	0.000000	62.670000	3.445000	182.616621	7.874000
75%	34.000000	1.420159e+06	0.000000	74.940000	3.735000	212.742393	8.620000
max	45.000000	3.818606e+06	1.000000	100.140000	4.468000	227.232807	14.313000

4. Evalúe si alguna de las variables contiene datos atípicos (outliers) De ser el caso, detalle cuáles y qué método estadístico aplicarán para corregir

Entendimos generales, estos descriptivos proporcionan una visión general de LOS datos. Algunos puntos de atención podrían incluir la presencia de valores negativos en la temperatura (que podrían ser errores de medición) y la amplia variabilidad en las ventas semanales y otras variables.

Método estadístico la imputación predictiva. Se utilizó modelos de regresión para predecir valores válidos para los casos con valores negativos basados en otras características del conjunto de datos.

3. Evalúe si la base contiene datos perdidos.

No hay datos perdidos en la base. La columna "Tienda" tiene 0 valores faltantes, al igual que las otras columnas, incluyendo "Fecha", "VentasSemanales", "SemanaEspecial", "Temperatura", "PrecioCombustible", "IPC" y "Desempleo". Todos los valores tienen recuentos de 0 faltantes (NaN).

```
In [178]: df.isna().sum()
Out[178]:
Tienda      0
Fecha       0
VentasSemanales  0
SemanaEspecial  0
Temperatura  0
PrecioCombustible  0
IPC         0
Desempleo   0
dtype: int64
In [179]: for i in df.select_dtypes("object").columns:
if df[i].dtype=="object":
df[i].value_counts().plot(kind='bar', title=i)
plt.show()
```



5. Grafique las distribuciones de las variables y a priori comente sobre ellas.

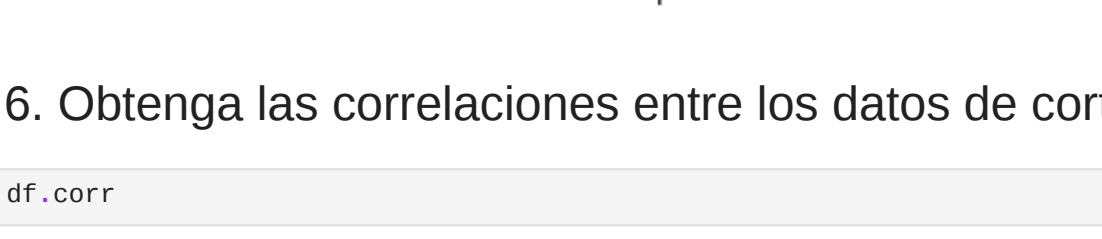
Grupo 1 Las variables en este grupo están muy agrupadas y no muestran valores atípicos. Esto sugiere que los datos en este grupo tienen una distribución bastante uniforme y que la mayoría de las observaciones caen dentro de un rango similar. Esto podría indicar que las variables en este grupo tienen una baja variabilidad y que las observaciones son consistentes en su comportamiento.

Grupo 2 Al igual que en el Grupo 1, las variables en este grupo también están muy agrupadas y no muestran valores atípicos. Esto sugiere una distribución uniforme de los datos en este grupo. Al igual que en el Grupo 1, esto podría indicar una baja variabilidad y consistencia en el comportamiento de las observaciones.

Grupo 3 Una vez más, los datos en este grupo están muy agrupados y no muestran valores atípicos. La distribución uniforme es evidente en este grupo. Esto sugiere una baja variabilidad y consistencia en las observaciones dentro de este grupo.

No se observan valores atípicos y los datos están muy agrupados en los gráficos de distribución, esto indica que las variables en cada grupo tienen comportamientos predecibles y estables.

```
In [179]: x = df["IPC"]
y = df["VentasSemanales"]
In [179]: plt.scatter(x,y)
plt.title('Scatterplot teórico')
plt.xlabel('x - variable independiente')
plt.ylabel('y - variable dependiente')
plt.show()
```



6. Obtenga las correlaciones entre los datos de corte numérico.

```
In [179]: df.corr
Out[179]:
```

		Tienda	Fecha	VentasSemanales	SemanaEspecial	Temperatura
0	1	05-02-2010	1643690.90	0	42.31	
1	12-02-2010	1641957.44	1	38.51		
2	19-02-2010	1611968.17	0	39.93		
3	26-02-2010	1409727.59	0	46.63		
4	05-03-2010	1554096.68	0	46.50		
...	
6430	45	28-09-2012	713173.95	0	64.88	
6431	45	05-10-2012	729455.07	0	64.89	
6432	45	12-10-2012	734464.36	0	54.47	
6433	45	19-10-2012	718125.53	0	56.47	
6434	45	26-10-2012	760281.43	0	58.85	
...	
0	PrecioCombustible	2.572	211.096358	8.106		
1	IPC	2.548	211.242178	8.106		
2	VentasSemanales	2.514	211.289143	8.106		
3	SemanaEspecial	2.561	211.319643	8.106		
4	Temperatura	2.625	211.350143	8.106		
...		
6430	45	192.015558	8.664			
6431	45	192.170412	8.667			
6432	45	192.327265	8.667			
6433	45	192.330854	8.667			
6434	45	192.308899	8.667			
...			
[6435 rows x 8 columns]						

7. Comente que variable escogerán como variable dependiente y que variables introducirán a su modelo.

"VentasSemanales" como variable dependiente (variable objetivo). Al utilizar un modelo de regresión lineal de Scikit-Learn, se introduce todas las demás variables disponibles como variables independientes (características) en el modelo. Esta sera una aproximación de regresión lineal múltiple, donde se intenta modelar la relación entre las ventas semanales y múltiples características.

8. Indique que tipo de modelación realizarán y porqué.

La elección de utilizar un modelo de regresión lineal de Scikit-Learn con los datos mencionados se respalda por las razones que se indican a continuación:

Interpretación de los coeficientes: La regresión lineal proporciona coeficientes directamente interpretados para cada variable independiente. Esto significa que se puede entender cómo cada variable contribuye o se relaciona con la variable dependiente (ventas semanales) en términos de la magnitud y la dirección de la influencia.

Simplicidad y claridad: La regresión lineal es un modelo simple y transparente que se basa en la suposición de una relación lineal entre las variables. Esto hace que sea fácil de entender y comunicar los resultados a partes interesadas no técnicas.

10. Obtenga el modelo definitivo, prediga los valores y comente el grado de ajuste del modelo. Justifique con métricas su respuesta

Error Cuadrático Medio (MSE): El MSE mide la magnitud de los errores cuadrados entre las predicciones y los valores reales. En este caso, tanto el MSE de entrenamiento como el MSE de prueba son relativamente altos, en el orden de cientos de millones. Esto sugiere que el modelo no es muy preciso en la predicción de las ventas semanales, ya que los errores son significativos.

Raíz del Error Cuadrático Medio (RMSE): El RMSE es la raíz cuadrada del MSE y proporciona una medida de la precisión del modelo en la misma unidad que la variable dependiente. Los valores de RMSE también son relativamente altos, en el orden de cientos de miles, lo que indica una precisión limitada.

Error Absoluto Medio (MAE): El MAE mide la magnitud promedio de los errores absolutos entre las predicciones y los valores reales. Los valores de MAE son más bajos que los de MSE y RMSE, pero aún son significativos, en el orden de cientos de miles.

Coefficiente de Determinación (R²): El R² mide la proporción de la varianza en la variable dependiente que es explicada por el modelo. Un R² de 0.14 en el conjunto de entrenamiento y 0.13 en el conjunto de prueba indica que el modelo explica una pequeña fracción de la variación en las ventas semanales y tiene un ajuste muy limitado a los datos.

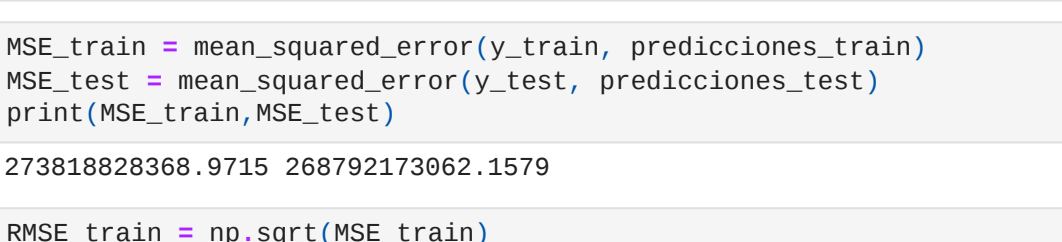
```
In [174]: from sklearn.linear_model import LinearRegression
In [175]: var_cuantitativas = df.select_dtypes("number").columns
var_cualitativas = df.select_dtypes("object").columns
In [176]: from sklearn.preprocessing import LabelEncoder
In [177]: labelencoder = LabelEncoder()
In [178]: df[var_cualitativas] = df[var_cualitativas].apply(labelencoder.fit_transform)
In [179]: X=df[df.columns.difference(['VentasSemanales'])]
y=df.VentasSemanales
In [180]: from sklearn.model_selection import train_test_split
In [181]: X_train, X_test, y_train, y_test=train_test_split(X, y, test_size = 0.10, random_state =123)
In [182]: X_train.shape
Out[182]: (5791, 7)
In [183]: X_test.shape
Out[183]: (644, 7)
In [184]: print(X_train.shape,"",type(X_train))
print(X_train.shape,"",type(X_train))
print(X_test.shape,"",type(X_test))
print(X_test.shape,"",type(X_test))
(5791, 7) <class 'pandas.core.frame.DataFrame'>
(5791, 7) <class 'pandas.core.frame.DataFrame'>
(644, 7) <class 'pandas.core.frame.DataFrame'>
(644, 7) <class 'pandas.core.frame.DataFrame'>
In [185]: modelo_regresion=LinearRegression()
In [186]: modelo_regresion.fit(X_train,y_train)
Out[187]: LinearRegression()
In [188]: predicciones_train = modelo_regresion.predict(X_train)
predicciones_test = modelo_regresion.predict(X_test)
In [189]: len(predicciones_test)
Out[189]: 644
In [190]: len(predicciones_train)
Out[190]: 5791
In [191]: from sklearn.metrics import mean_squared_error, mean_absolute_error
In [192]: MSE_train = mean_squared_error(y_train, predicciones_train)
MSE_test = mean_squared_error(y_test, predicciones_test)
print(MSE_train,MSE_test)
273818828368.9715 268792173982.1579
In [193]: RMSE_train = np.sqrt(MSE_train)
RMSE_test = np.sqrt(MSE_test)
print(RMSE_train,RMSE_test)
523277.0092111553 518451.70755062415
In [194]: MAE_train = mean_absolute_error(y_train, predicciones_train)
MAE_test = mean_absolute_error(y_test, predicciones_test)
print(MAE_train,MAE_test)
438276.9499653293 427167.4888537875
In [195]: from sklearn.metrics import r2_score
In [196]: r_square_train = r2_score(y_train, predicciones_train)
r_square_test = r2_score(y_test, predicciones_test)
print('El R² del subconjunto de entrenamiento es:', r_square_train)
print('El R² del subconjunto de prueba es:', r_square_test)
El R² del subconjunto de entrenamiento es: 0.1427554474366871
El R² del subconjunto de prueba es: 0.1322892098304545
In [197]: modelo_regresion.coef_
Out[197]: array([-21367.25422651, -2337.65182967, 18178.93496641, 62938.74489942, -949.58024181, -15538.64285727, -258.08374993])
In [198]: # Print the Intercept:
print('Intercept:', modelo_regresion.intercept_)
Intercept: 2813147.1906412344
pendiente: [-21367.25422651, -2337.65182967, 18178.93496641, 62938.74489942, -949.58024181, -15538.64285727, -258.08374993]
```

11. Grafique a los valores predicho de modelo vs los valores reales.¿Cómo se ven una vez graficados frente a los valores reales? Argumente su respuesta.

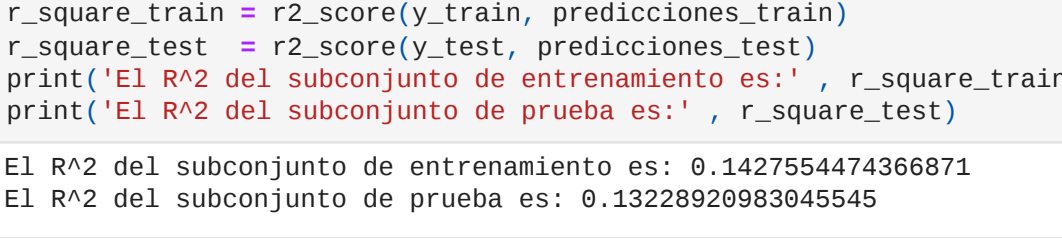
La concentración de los valores predichos en el centro y su alejamiento de los picos de los valores observados en un gráfico de dispersión sugiere una falta de ajuste del modelo a los datos. Este patrón de dispersión indica que el modelo de regresión lineal no está capturando adecuadamente las relaciones subyacentes en los datos.

Cuando los valores predichos se concentran en el centro y están alejados de los picos de los valores observados, significa que el modelo tiende a subestimar las ventas en situaciones de alto rendimiento y sobrestimarlas en situaciones de bajo rendimiento. En otras palabras, el modelo no logra capturar las fluctuaciones en las ventas que ocurren en los extremos de la distribución de datos.

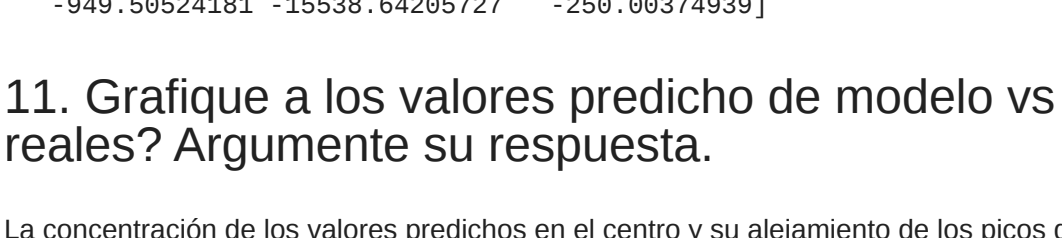
```
In [199]: fig, ax = plt.subplots()
ax.plot(y_train.values, )
ax.plot(predicciones_train)
plt.title('Valores observados vs. predichos en train set')
In [200]: (1098, 0, 1598, 0)
Out[200]:
```



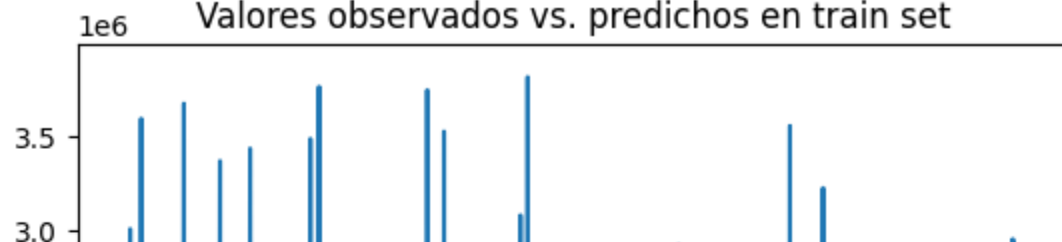
```
In [201]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test)
plt.title('Valores observados vs. predichos en train set')
In [202]: (1098, 0, 1598, 0)
Out[202]:
```



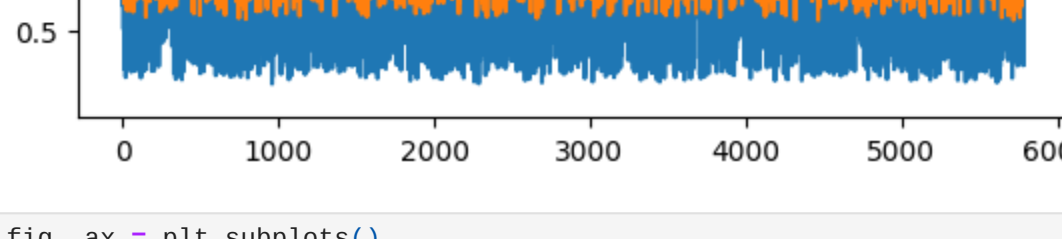
```
In [203]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [204]: (1098, 0, 1598, 0)
Out[204]:
```



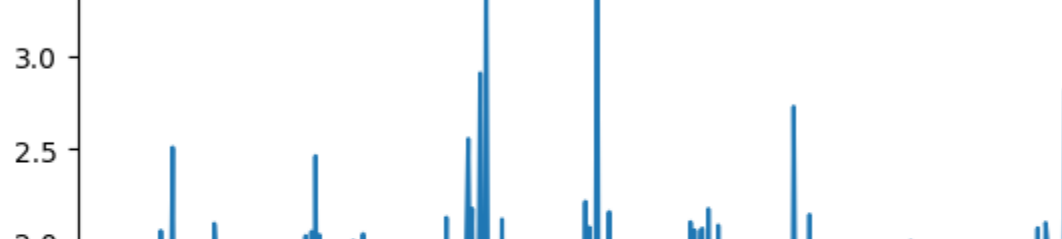
```
In [205]: from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
In [206]: X_train_std = sc.fit_transform(X_train)
X_test_std = sc.fit_transform(X_test)
In [207]: modelo_regresion_std = LinearRegression()
modelo_regresion_std.fit(X_train_std, y_train)
Out[207]: LinearRegression()
In [208]: predicciones_train_std = modelo_regresion_std.predict(X_train_std)
predicciones_test_std = modelo_regresion_std.predict(X_test_std)
In [209]: r2_train=r2_score(y_train, predicciones_train)
r2_test=r2_score(y_test, predicciones_test)
print('El R² del subconjunto de entrenamiento es:', r2_train, 'El R² del subconjunto de prueba es:', r2_test)
El R² del subconjunto de entrenamiento es: 0.1427554474366871 El R² del subconjunto de prueba es: 0.1322892098304545
In [210]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [211]: (1098, 0, 1598, 0)
Out[211]:
```



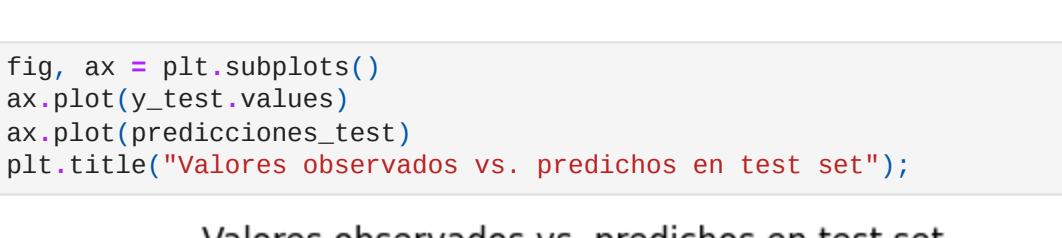
```
In [212]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test_std)
plt.title('Valores observados vs. predichos en test set')
In [213]: (1098, 0, 1598, 0)
Out[213]:
```



```
In [214]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [215]: (1098, 0, 1598, 0)
Out[215]:
```



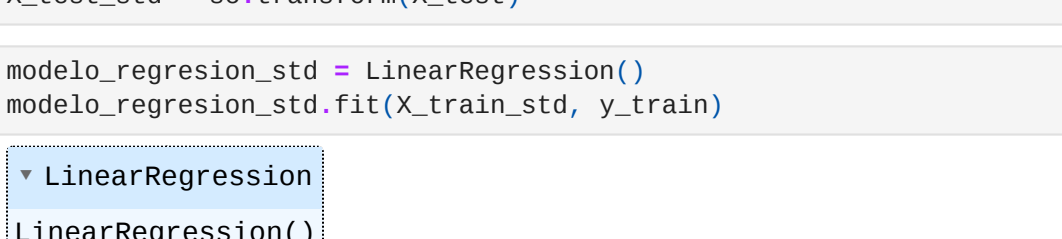
```
In [216]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test_std)
plt.title('Valores observados vs. predichos en test set')
In [217]: (1098, 0, 1598, 0)
Out[217]:
```



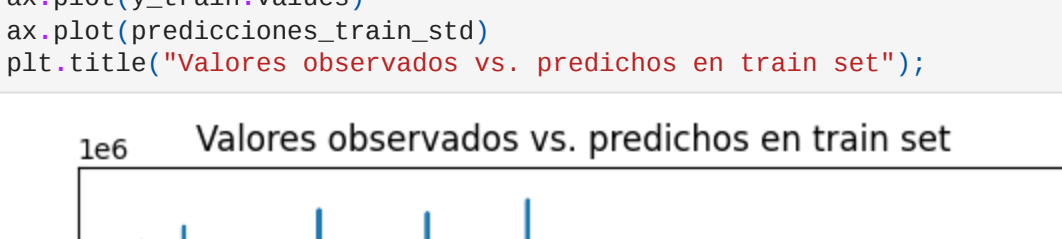
```
In [218]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [219]: (1098, 0, 1598, 0)
Out[219]:
```



```
In [220]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test_std)
plt.title('Valores observados vs. predichos en test set')
In [221]: (1098, 0, 1598, 0)
Out[221]:
```



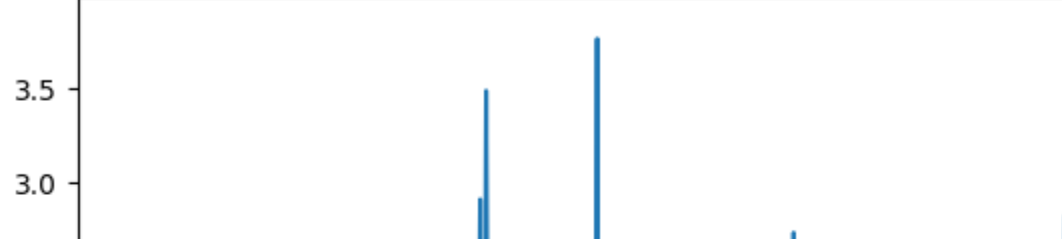
```
In [222]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [223]: (1098, 0, 1598, 0)
Out[223]:
```



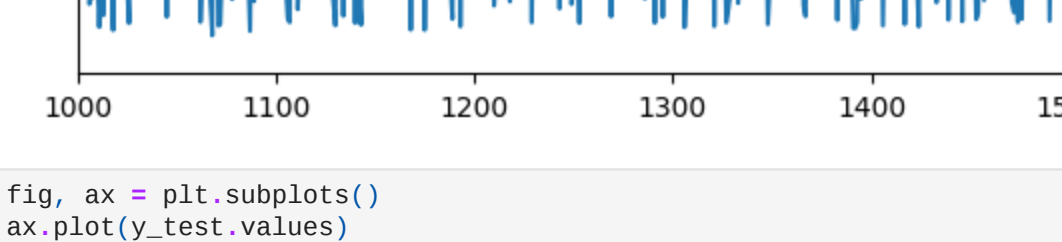
```
In [224]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test_std)
plt.title('Valores observados vs. predichos en test set')
In [225]: (1098, 0, 1598, 0)
Out[225]:
```



```
In [226]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [227]: (1098, 0, 1598, 0)
Out[227]:
```



```
In [228]: fig, ax = plt.subplots()
ax.plot(y_test.values)
ax.plot(predicciones_test_std)
plt.title('Valores observados vs. predichos en test set')
In [229]: (1098, 0, 1598, 0)
Out[229]:
```



```
In [230]: fig, ax = plt.subplots()
ax.plot(y_train.values)
ax.plot(predicciones_train_std)
plt.title('Valores observados vs. predichos en train set')
In [231]: (1098, 0, 1598, 0)
Out[231]:
```

