

```
In [1]: import numpy as np
import matplotlib.pyplot as plt

import pandas as pd
import seaborn as sns

%matplotlib inline

In [2]: from sklearn.datasets import load_boston
boston_dataset = load_boston()

~/home/ubuntu/.local/lib/python3.10/site-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function load_boston is deprecated; 'load_boston' is deprecated in 1.0 and will be removed in 1.2.

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\t", skiprows=22, header=None)
data = np.loadtxt(raw_df.values[:,2:], dtype=float, delimiter=',')
target = raw_df.values[:,22]

Alternative datasets include the California housing dataset (i.e. :func:`sklearn.datasets.fetch_california_housing`) and the Ames housing dataset. You can load the datasets as follows::

from sklearn.datasets import fetch_california_housing
housing = fetch_california_housing()

for the California housing dataset and::

from sklearn.datasets import fetch_ohio
housing = fetch_ohio(name='house_prices', as_frame=True)

for the Ames housing dataset.
warnings.warn(msg, category=FutureWarning)
```

```
In [3]: print(boston_dataset.keys())

dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename', 'data_module'])

In [4]: boston_dataset.DESCR

Out[4]: "... boston_dataset:\n\nBoston house prices dataset\n\n-----\n\nData Set Characteristics: ** \n\n :Number of Instances: 506 \n\n :Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.\n\n :Attribute Information (in order)\n :CRIM per capita crime rate by town\n :ZN proportion of residential land zoned for lots over 25,000 sq.ft.\n :INDUS proportion of non-retail business acres per town\n :CHAS Charles River dummy variable (= 1 if tract b\n :NOX nitric oxides concentration (parts per 10 million)\n :AGE average number of rooms per dwelling\n :DIS weighted distances to five Boston emplo\n :RAD index of accessibility to radial highways\n :TAX full-value property-tax rate per $10,000\n :PTRATIO pupil-teacher ratio by town\n :B 1000(BK - 0.63)^2 where BK is the proportion of black people by town\n :LST lower status of the population\n :MEDV Median value of owner-occupied homes in $1000's\n\n :Missing Attribute Values: None\n\n :Creator: Harrison, D. and Rubinfeld, D.I.\n\nThis is a copy of UCI ML housing dataset. \nhttps://archive.ics.uci.edu/ml/machine-learning-databases/housing/\n\nThis dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.\n\nThe Boston house-price data of Harrison, D. and Rubinfeld, D.I. "Hedonic prices and the demand for clean air", J. Environ. Economics & Management, vol. 5, 81-102, 1978. Used in Belsley, Kuh & Welsch, "Regression diagnostics: Identifying Influential Data and Sources of Collinearity", Wiley, 1980. N.B. Various transformations are used in the table on pages 244-261 of the latter.\n\nThe Boston house-price data has been used in many machine learning papers that address regression problems.\n\n :Source: Refer to Belsley, Kuh & Welsch, "Regression diagnostics: Identifying Influential Data and Sources of Collinearity", Wiley, 1980. 244-261. \n - Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Coference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.\n"
```

```
In [5]: boston = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
boston.head()

Out[5]:
```

```
In [6]: boston["MEDV"] = boston_dataset.target

In [7]: boston.isnull().sum()
```

```
Out[7]:
```

```
In [8]: sns.set(rc={'figure.figsize':(11.7,8.27)})
sns.distplot(boston["MEDV"], bins=30)
plt.show()

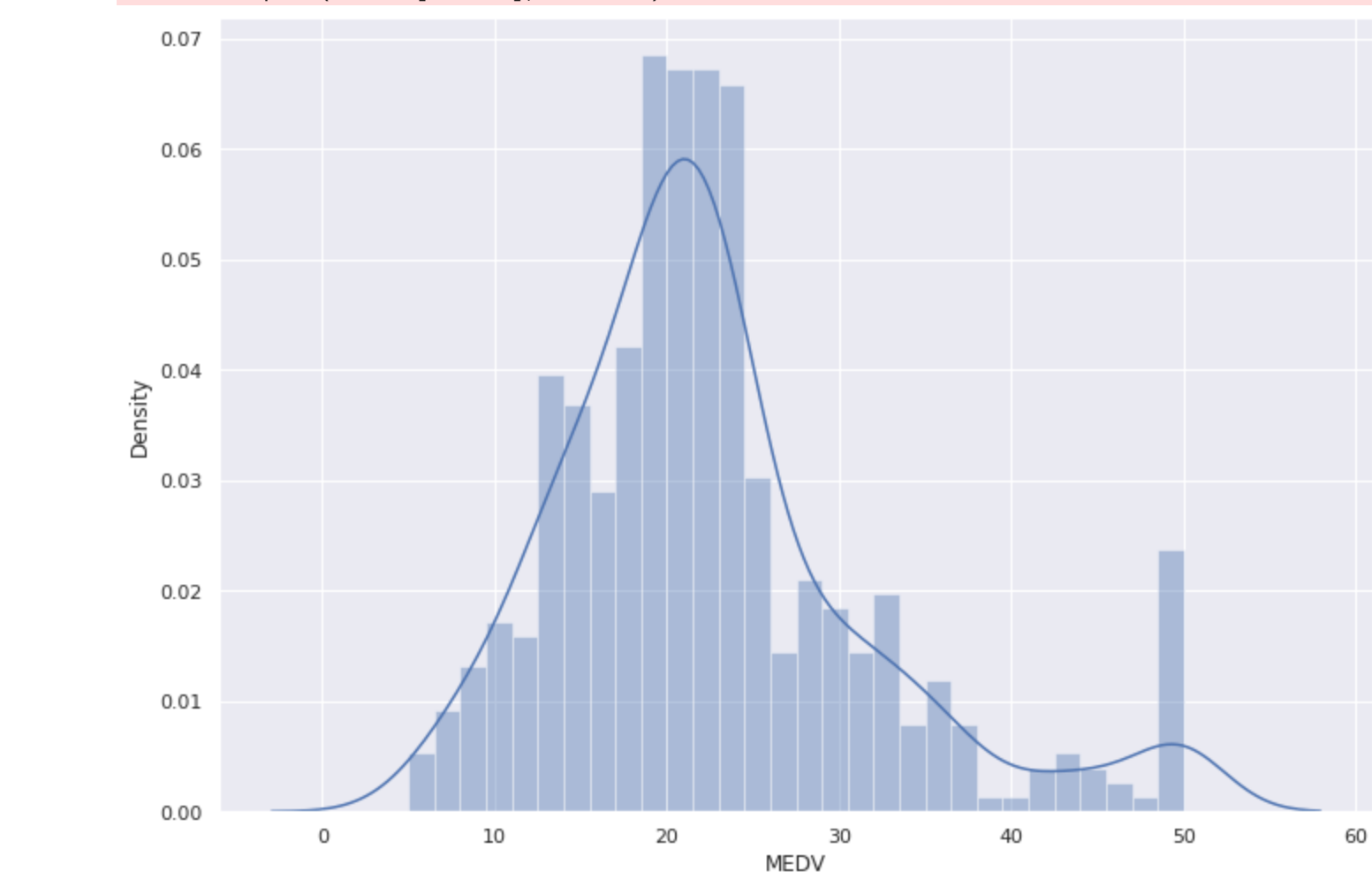
~/tmp/ipykernel_6673/3962179664.py:2: UserWarning:

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

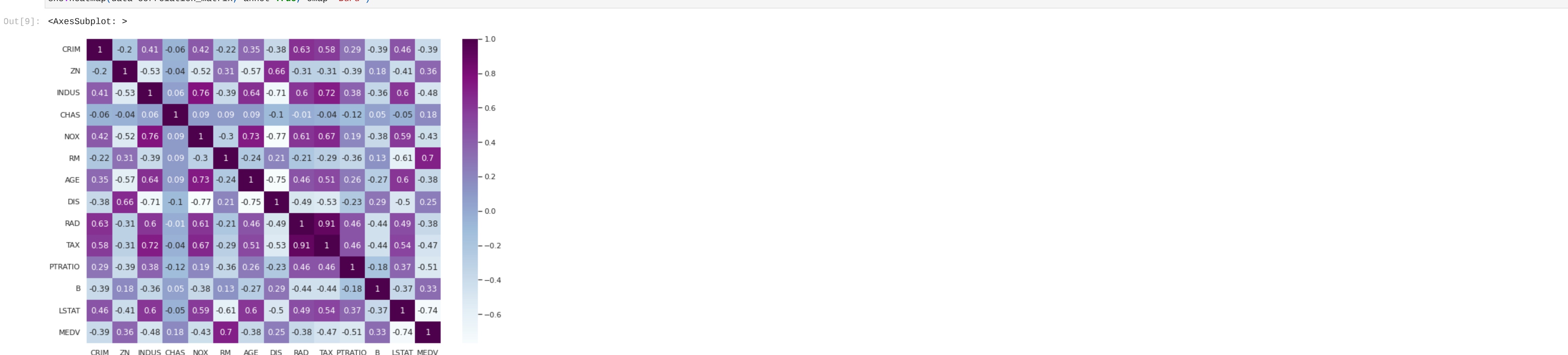
Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/9e44147ed2974457ad8372758bbe5751

sns.distplot(boston["MEDV"], bins=30)
```

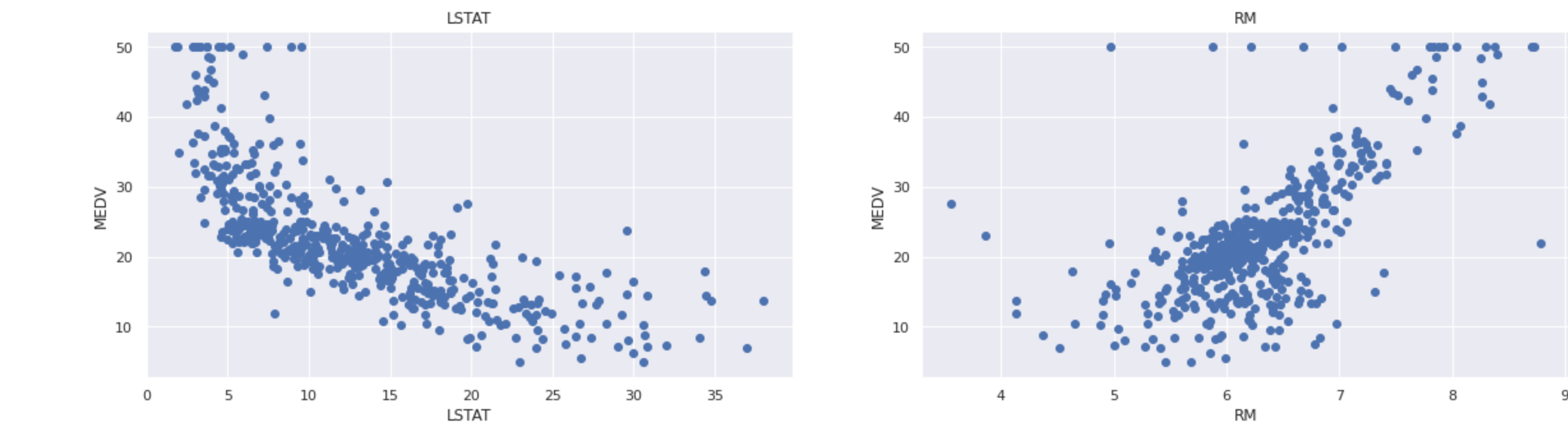


```
In [9]: correlation_matrix = boston.corr().round(2)
sns.heatmap(data=correlation_matrix, annot=True, cmap='BuPu')
```



```
In [10]: plt.figure(figsize=(20, 5))
features = ['LSTAT', 'RM']
target = boston['MEDV']

for i, col in enumerate(features):
    plt.subplot(1, len(features) - 1, i+1)
    x = boston[col]
    y = target
    plt.scatter(x, y, markers='o')
    plt.title(col)
    plt.xlabel(col)
    plt.ylabel('MEDV')
```



```
In [11]: X= pd.DataFrame(np.c_[boston['LSTAT'], boston['RM']], columns = ['LSTAT', 'RM'])
Y = boston['MEDV']
```

```
In [12]: from sklearn.model_selection import train_test_split

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state=5)
print(X_train.shape)
print(X_test.shape)
print(Y_train.shape)
print(Y_test.shape)

(494, 2)
(102, 2)
(494,)
(102,)
```

```
In [13]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

lin_model = LinearRegression()
lin_model.fit(X_train, Y_train)
```

```
Out[13]: LinearRegression()

LinearRegression()
```

```
In [14]: y_train_predict = lin_model.predict(X_train)
rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
r2 = r2_score(Y_train, y_train_predict)

print("The model performance for training set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

y_test_predict = lin_model.predict(X_test)
rmse = (np.sqrt(mean_squared_error(Y_test, y_test_predict)))
r2 = r2_score(Y_test, y_test_predict)

print("The model performance for testing set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

NameError                                Traceback (most recent call last)
~/tmp/ipykernel_6673/3269519552.py in <module>
      1 y_train_predict = lin_model.predict(X_train)
      2 rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
----> 3 r2 = r2_score(Y_train, y_train_predict)
      4
      5 print("The model performance for training set")

NameError: name 'r2_score' is not defined
```

```
In [15]: y_train_predict = lin_model.predict(X_train)
rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
r2 = r2_score(Y_train, y_train_predict)

print("The model performance for training set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

y_test_predict = lin_model.predict(X_test)
rmse = (np.sqrt(mean_squared_error(Y_test, y_test_predict)))
r2 = r2_score(Y_test, y_test_predict)

print("The model performance for testing set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

NameError                                Traceback (most recent call last)
~/tmp/ipykernel_6673/2569517694.py in <module>
      1 y_train_predict = lin_model.predict(X_train)
      2 rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
----> 3 r2 = r2_score(Y_train, y_train_predict)
      4
      5 print("The model performance for training set")

NameError: name 'r2_score' is not defined
```

```
In [16]: plt.scatter(Y_train, y_pred)
plt.xlabel('Prices')
plt.ylabel('Predicted prices')
plt.title("Prices vs Predicted prices")
plt.show()

NameError                                Traceback (most recent call last)
~/tmp/ipykernel_6673/246851517.py in <module>
----> 1 plt.scatter(Y_train, y_pred)
      2 plt.xlabel('Prices')
      3 plt.ylabel('Predicted prices')
      4 plt.title("Prices vs Predicted prices")
      5 plt.show()

NameError: name 'Y_train' is not defined
```

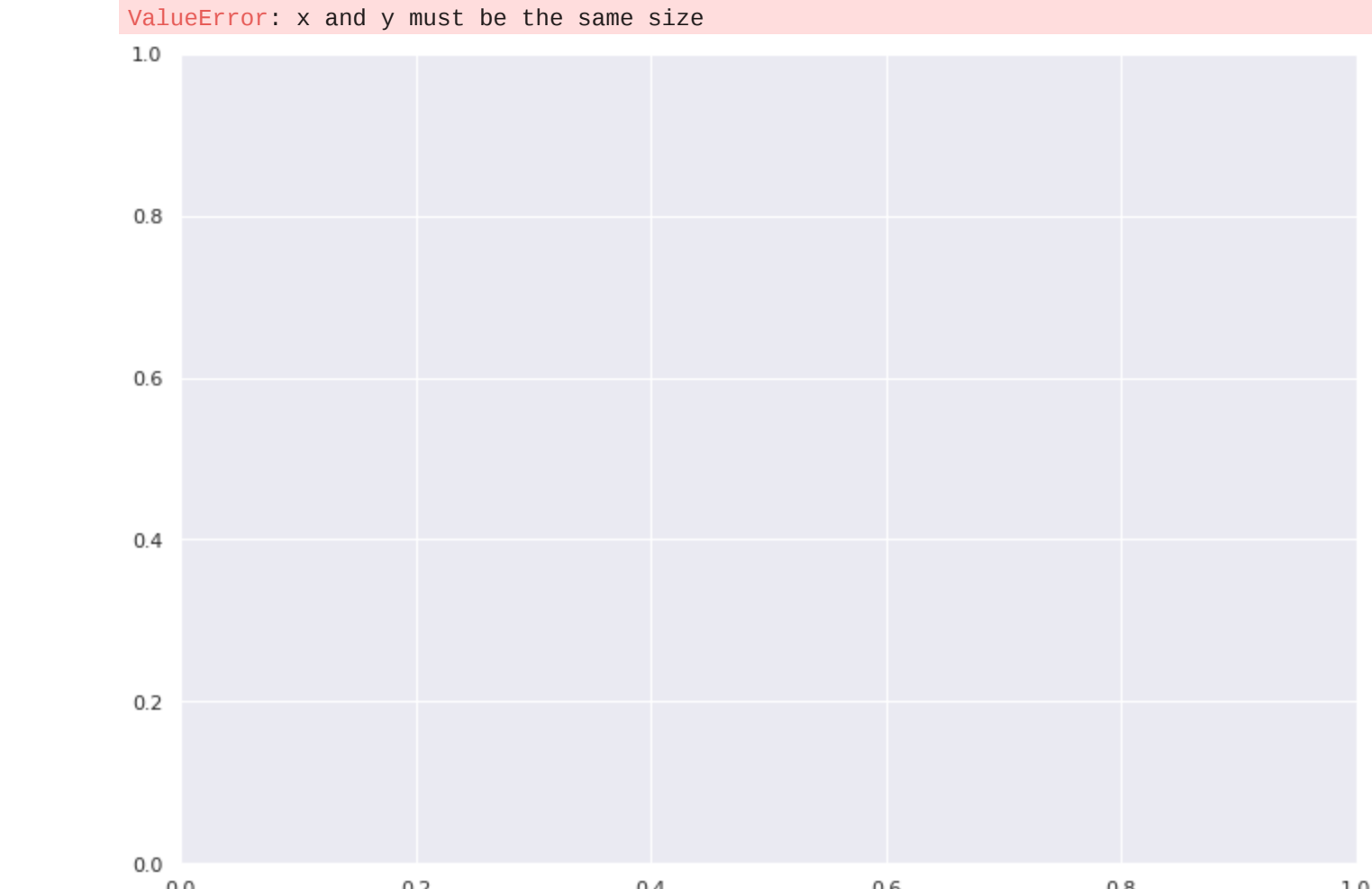
```
In [17]: plt.scatter(Y_train, Y_test)
plt.xlabel('Prices')
plt.ylabel('Predicted prices')
plt.title("Prices vs Predicted prices")
plt.show()
```

```
ValueError                                Traceback (most recent call last)
~/tmp/ipykernel_6673/3194198346.py in <module>
----> 1 plt.scatter(Y_train, Y_test)
      2 plt.xlabel('Prices')
      3 plt.ylabel('Predicted prices')
      4 plt.title("Prices vs Predicted prices")
      5 plt.show()

~/local/lib/python3.10/site-packages/matplotlib/pyplot.py in scatter(x, y, s, c, marker, cmap, norm, vmin, vmax, alpha, linewidths, edgecolors, plotnonfinite, data, **kwargs)
   2788     vmin=vmin, vmax=vmax, alpha=alpha, linewidths=linewidths,
   2789     edgecolors=edgecolors, plotnonfinite=False, data=None, **kwargs):
-> 2790     _ret = gca().scatter(
   2791         x, y, s=s, c=c, marker=marker, cmap=cmap, norm=norm,
   2792         vmin=vmin, vmax=vmax, alpha=alpha, linewidths=linewidths,

~/local/lib/python3.10/site-packages/matplotlib/_api/_init_.py in inner(ax, data, *args, **kwargs)
   1421     def inner(ax, *args, data=None, **kwargs):
   1422         if data is None:
-> 1423             return func(ax, *map(sanitize_sequence, args), **kwargs)
   1424
   1425     bound = new_sig.bind(ax, *args, **kwargs)

~/local/lib/python3.10/site-packages/matplotlib/axes/_axes.py in scatter(self, x, y, s, c, marker, cmap, norm, vmin, vmax, alpha, linewidths, edgecolors, plotnonfinite, **kwargs)
   4518     y = np.ma.ravel(y)
   4519     if x.size != y.size:
-> 4520         raise ValueError("x and y must be the same size")
   4521
   4522     if s is None:
ValueError: x and y must be the same size
```



```
In [18]: y_pred = lin_model.predict(X_train)
```

```
In [19]: plt.scatter(Y_train, y_pred)
plt.xlabel('Prices')
plt.ylabel('Predicted prices')
plt.title("Prices vs Predicted prices")
plt.show()
```



