In [1]:
```python
import pandas as pd
import numpy as np
from sklearn import metrics
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

In [2]:
```python
import warnings
warnings.filterwarnings("ignore")
```

In [3]:
```python
from sklearn.datasets import load_boston
boston = load_boston()
```

In [4]:
```python
data = pd.DataFrame(boston.data)
```

In [5]:
```python
data.head()
```

Out[5]:

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|----|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 | 5.33 |

In [6]:
```python
#Adding the feature names to the dataframe
data.columns = boston.feature_names
data.head()
```

Out[6]:

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B |
|---|------|----|-------|------|-----|----|-----|-----|-----|-----|---------|---|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 | 396.90 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 | 396.90 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 | 392.83 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 | 394.63 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 | 396.90 |

In [7]:
```python
#Adding target variable to dataframe
data['PRICE'] = boston.target
```

In [8]: 
```python
#Check the shape of dataframe
data.shape
```

Out[8]: (506, 14)

In [9]: 
```python
data.columns
```

Out[9]: 
```
Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', '
RAD', 'TAX',
        'PTRATIO', 'B', 'LSTAT', 'PRICE'],
      dtype='object')
```

In [10]: 
```python
data.dtypes
```

Out[10]: 
```
CRIM        float64
ZN          float64
INDUS       float64
CHAS        float64
NOX         float64
RM          float64
AGE         float64
DIS         float64
RAD         float64
TAX         float64
PTRATIO     float64
B           float64
LSTAT       float64
PRICE       float64
dtype: object
```

In [11]: 
```python
# Identifying the unique number of values in the dataset
data.nunique()
```

Out[11]: 
```
CRIM       504
ZN          26
INDUS       76
CHAS         2
NOX         81
RM         446
AGE        356
DIS        412
RAD          9
TAX         66
PTRATIO     46
B          357
LSTAT      455
PRICE      229
dtype: int64
```

In [12]: 
```python
# Check for missing values
data.isnull().sum()
```

Out[12]: 
```
CRIM        0
ZN          0
INDUS       0
CHAS        0
NOX         0
RM          0
AGE         0
DIS         0
RAD         0
TAX         0
PTRATIO     0
B           0
LSTAT       0
PRICE       0
dtype: int64
```

In [13]: 
```python
# See rows with missing values
data[data.isnull().any(axis=1)]
```

Out[13]:

| CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | PRICE |
|------|----|-------|------|-----|----|----|-----|-----|-----|---------|---|-------|-------|

In [14]: 
```python
# Viewing the data statistics
data.describe()
```
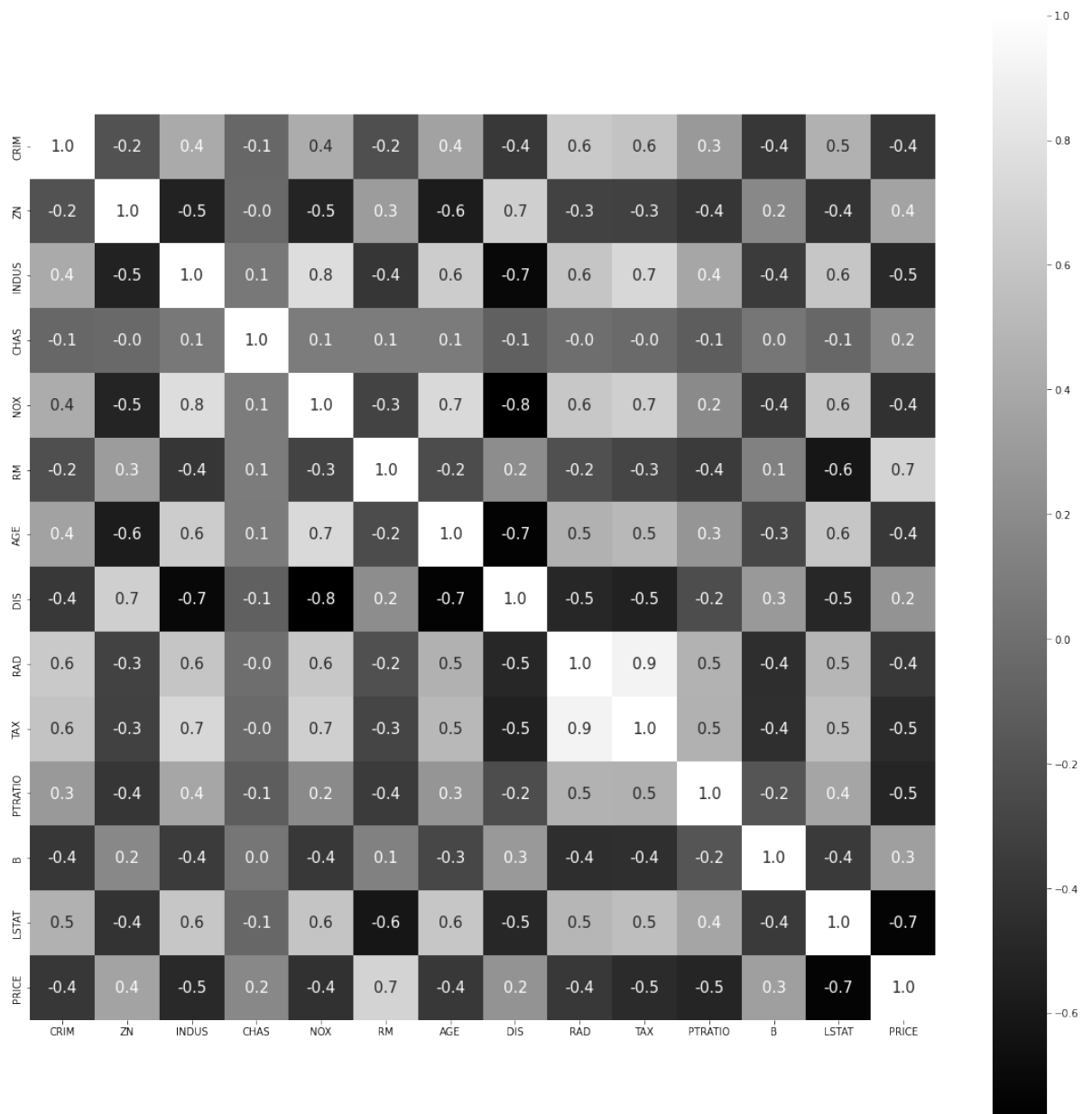
Out[14]:

|       | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE |
|-------|------|----|-------|------|-----|----|----|
| count | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 |
| mean | 3.613524 | 11.363636 | 11.136779 | 0.069170 | 0.554695 | 6.284634 | 68.574901 |
| std | 8.601545 | 23.322453 | 6.860353 | 0.253994 | 0.115878 | 0.702617 | 28.148861 |
| min | 0.006320 | 0.000000 | 0.460000 | 0.000000 | 0.385000 | 3.561000 | 2.900000 |
| 25% | 0.082045 | 0.000000 | 5.190000 | 0.000000 | 0.449000 | 5.885500 | 45.025000 |
| 50% | 0.256510 | 0.000000 | 9.690000 | 0.000000 | 0.538000 | 6.208500 | 77.500000 |
| 75% | 3.677083 | 12.500000 | 18.100000 | 0.000000 | 0.624000 | 6.623500 | 94.075000 |
| max | 88.976200 | 100.000000 | 27.740000 | 1.000000 | 0.871000 | 8.780000 | 100.000000 |

In [15]: 
```python
# Finding out the correlation between the features
corr = data.corr()
corr.shape
```

Out[15]: (14, 14)

In [16]:
```python
# Plotting the heatmap of correlation between features
plt.figure(figsize=(20,20))
sns.heatmap(corr, cbar=True, square= True, fmt='.1f', annot=True, a
```

Out[16]: <AxesSubplot:>

|          | CRIM | ZN   | INDUS | CHAS | NOX  | RM   | AGE  | DIS  | RAD  | TAX  | PTRATIO | B    | LSTAT | PRICE |
|----------|------|------|-------|------|------|------|------|------|------|------|---------|------|-------|-------|
| CRIM     | 1.0  | -0.2 | 0.4   | -0.1 | 0.4  | -0.2 | 0.4  | -0.4 | 0.6  | 0.6  | 0.3     | -0.4 | 0.5   | -0.4  |
| ZN       | -0.2 | 1.0  | -0.5  | -0.0 | -0.5 | 0.3  | -0.6 | 0.7  | -0.3 | -0.3 | -0.4    | 0.2  | -0.4  | 0.4   |
| INDUS    | 0.4  | -0.5 | 1.0   | 0.1  | 0.8  | -0.4 | 0.6  | -0.7 | 0.6  | 0.7  | 0.4     | -0.4 | 0.6   | -0.5  |
| CHAS     | -0.1 | -0.0 | 0.1   | 1.0  | 0.1  | 0.1  | 0.1  | -0.1 | -0.0 | -0.0 | -0.1    | 0.0  | -0.1  | 0.2   |
| NOX      | 0.4  | -0.5 | 0.8   | 0.1  | 1.0  | -0.3 | 0.7  | -0.8 | 0.6  | 0.7  | 0.2     | -0.4 | 0.6   | -0.4  |
| RM       | -0.2 | 0.3  | -0.4  | 0.1  | -0.3 | 1.0  | -0.2 | 0.2  | -0.2 | -0.3 | -0.4    | 0.1  | -0.6  | 0.7   |
| AGE      | 0.4  | -0.6 | 0.6   | 0.1  | 0.7  | -0.2 | 1.0  | -0.7 | 0.5  | 0.5  | 0.3     | -0.3 | 0.6   | -0.4  |
| DIS      | -0.4 | 0.7  | -0.7  | -0.1 | -0.8 | 0.2  | -0.7 | 1.0  | -0.5 | -0.5 | -0.2    | 0.3  | -0.5  | 0.2   |
| RAD      | 0.6  | -0.3 | 0.6   | -0.0 | 0.6  | -0.2 | 0.5  | -0.5 | 1.0  | 0.9  | 0.5     | -0.4 | 0.5   | -0.4  |
| TAX      | 0.6  | -0.3 | 0.7   | -0.0 | 0.7  | -0.3 | 0.5  | -0.5 | 0.9  | 1.0  | 0.5     | -0.4 | 0.5   | -0.5  |
| PTRATIO  | 0.3  | -0.4 | 0.4   | -0.1 | 0.2  | -0.4 | 0.3  | -0.2 | 0.5  | 0.5  | 1.0     | -0.2 | 0.4   | -0.5  |
| B        | -0.4 | 0.2  | -0.4  | 0.0  | -0.4 | 0.1  | -0.3 | 0.3  | -0.4 | -0.4 | -0.2    | 1.0  | -0.4  | 0.3   |
| LSTAT    | 0.5  | -0.4 | 0.6   | -0.1 | 0.6  | -0.6 | 0.6  | -0.5 | 0.5  | 0.5  | 0.4     | -0.4 | 1.0   | -0.7  |
| PRICE    | -0.4 | 0.4  | -0.5  | 0.2  | -0.4 | 0.7  | -0.4 | 0.2  | -0.4 | -0.5 | -0.5    | 0.3  | -0.7  | 1.0   |

In [17]:
```python
# Spliting target variable and independent variables
X = data.drop(['PRICE'], axis = 1)
y = data['PRICE']
```

In [18]:
```python
# Splitting to training and testing data

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size
```

In [19]: 
```python
# Import library for Linear Regression
from sklearn.linear_model import LinearRegression
```

In [20]: 
```python
# Create a Linear regressor
lm = LinearRegression()

# Train the model using the training sets
lm.fit(X_train, y_train)
```

Out[20]: LinearRegression()

In [21]: 
```python
# Value of y intercept
lm.intercept_
```

Out[21]: 36.357041376594964

In [22]: 
```python
#Converting the coefficient values to a dataframe
coeffcients = pd.DataFrame([X_train.columns,lm.coef_]).T
coeffcients = coeffcients.rename(columns={0: 'Attribute', 1: 'Coeff
coeffcients
```

Out[22]:

|    | Attribute | Coefficients |
|----|-----------|--------------|
| 0  | CRIM      | -0.12257     |
| 1  | ZN        | 0.055678     |
| 2  | INDUS     | -0.008834    |
| 3  | CHAS      | 4.693448     |
| 4  | NOX       | -14.435783   |
| 5  | RM        | 3.28008      |
| 6  | AGE       | -0.003448    |
| 7  | DIS       | -1.552144    |
| 8  | RAD       | 0.32625      |
| 9  | TAX       | -0.014067    |
| 10 | PTRATIO   | -0.803275    |
| 11 | B         | 0.009354     |
| 12 | LSTAT     | -0.523478    |

In [23]: 
```python
# Model prediction on train data
y_pred = lm.predict(X_train)
```
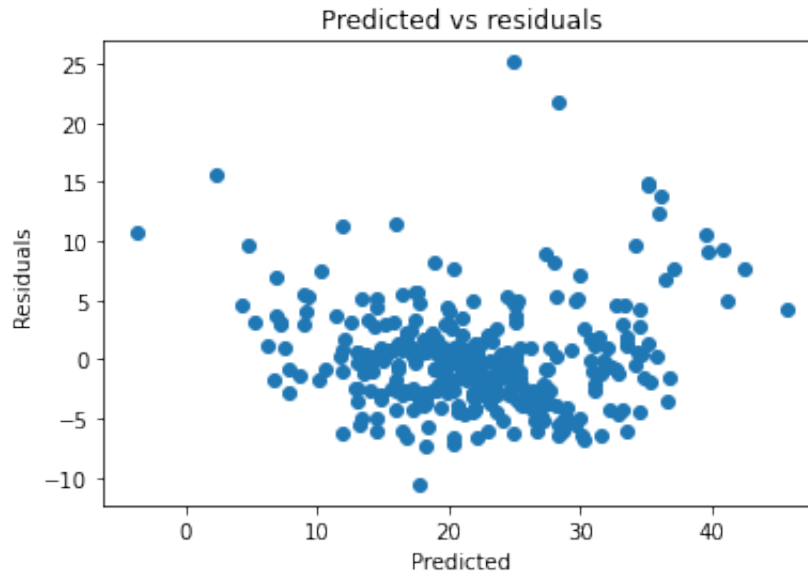
In [24]:
```
Evaluation
R^2:',metrics.r2_score(y_train, y_pred))
Adjusted R^2:',1 - (1-metrics.r2_score(y_train, y_pred))*(len(y_trai
MAE:',metrics.mean_absolute_error(y_train, y_pred))
MSE:',metrics.mean_squared_error(y_train, y_pred))
RMSE:',np.sqrt(metrics.mean_squared_error(y_train, y_pred)))
```

```
R^2: 0.7465991966746854
Adjusted R^2: 0.736910342429894
MAE: 3.0898610949711323
MSE: 19.07368870346903
RMSE: 4.367343437774162
```
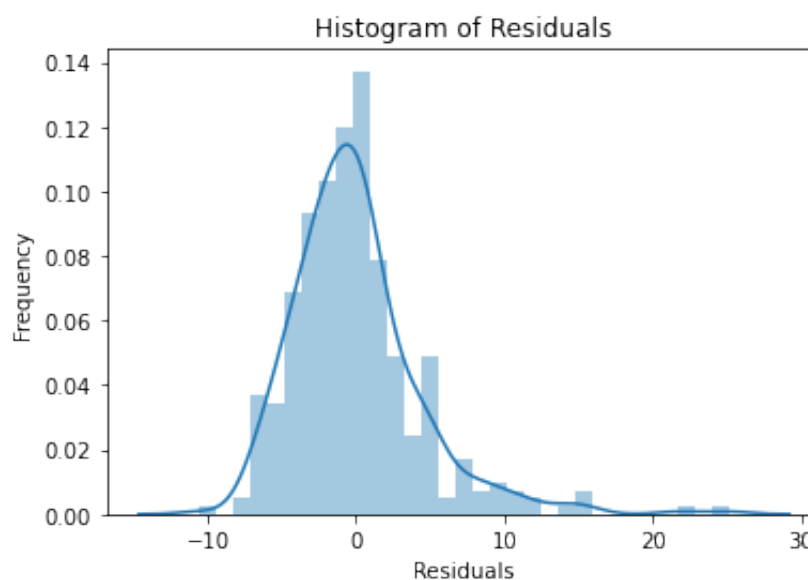
In [25]:
```
# Visualizing the differences between actual prices and predicted v
plt.scatter(y_train, y_pred)
plt.xlabel("Prices")
plt.ylabel("Predicted prices")
plt.title("Prices vs Predicted prices")
plt.show()
```

In [26]:
```python
# Checking residuals
plt.scatter(y_pred,y_train-y_pred)
plt.title("Predicted vs residuals")
plt.xlabel("Predicted")
plt.ylabel("Residuals")
plt.show()
```



In [27]:
```python
# Checking Normality of errors
sns.distplot(y_train-y_pred)
plt.title("Histogram of Residuals")
plt.xlabel("Residuals")
plt.ylabel("Frequency")
plt.show()
```

In [28]:
```python
# Predicting Test data with the model
y_test_pred = lm.predict(X_test)
```

In [29]:
```python
l Evaluation
nreg = metrics.r2_score(y_test, y_test_pred)
'R^2:', acc_linreg)
'Adjusted R^2:',1 - (1-metrics.r2_score(y_test, y_test_pred))*(len(y
'MAE:',metrics.mean_absolute_error(y_test, y_test_pred))
'MSE:',metrics.mean_squared_error(y_test, y_test_pred))
'RMSE:',np.sqrt(metrics.mean_squared_error(y_test, y_test_pred)))
```

```
R^2: 0.7121818377409183
Adjusted R^2: 0.6850685326005699
MAE: 3.859005592370746
MSE: 30.053993307124266
RMSE: 5.482152251362987
```

In [ ]: