

Import all the important modules

```
In [1]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

Import the data frame with all the important columns which effect the rate of house price in US

```
In [2]: df = pd.read_csv(r"C:\Users\sagar\Downloads\Home_price_index.csv")
```

Check the shape of the data frame

```
In [3]: df.shape
```

```
Out[3]: (92, 9)
```

Data Frame column names

```
In [4]: df.columns
```

```
Out[4]: Index(['Date', 'Mortgage Average', 'Unemployment Rate', 'Housing Inventory',
      'Population Growth', 'Inflation', 'Permit-Issuing Places',
      'Median Household Income', 'Price'],
      dtype='object')
```

Drop Columns which is not required for the ML models

```
In [5]: df = df.drop('Date',axis = 1)
```

Fill "NA" values with forward fill

```
In [6]: df.isnull().sum()
```

```
Out[6]: Mortgage Average      0
Unemployment Rate          0
Housing Inventory          0
Population Growth          69
Inflation                  69
Permit-Issuing Places      0
Median Household Income    70
Price                     0
dtype: int64
```

```
In [7]: df = df.fillna(method='ffill')
```

Check NA values in the Data frame

```
In [8]: df.isnull().sum()
```

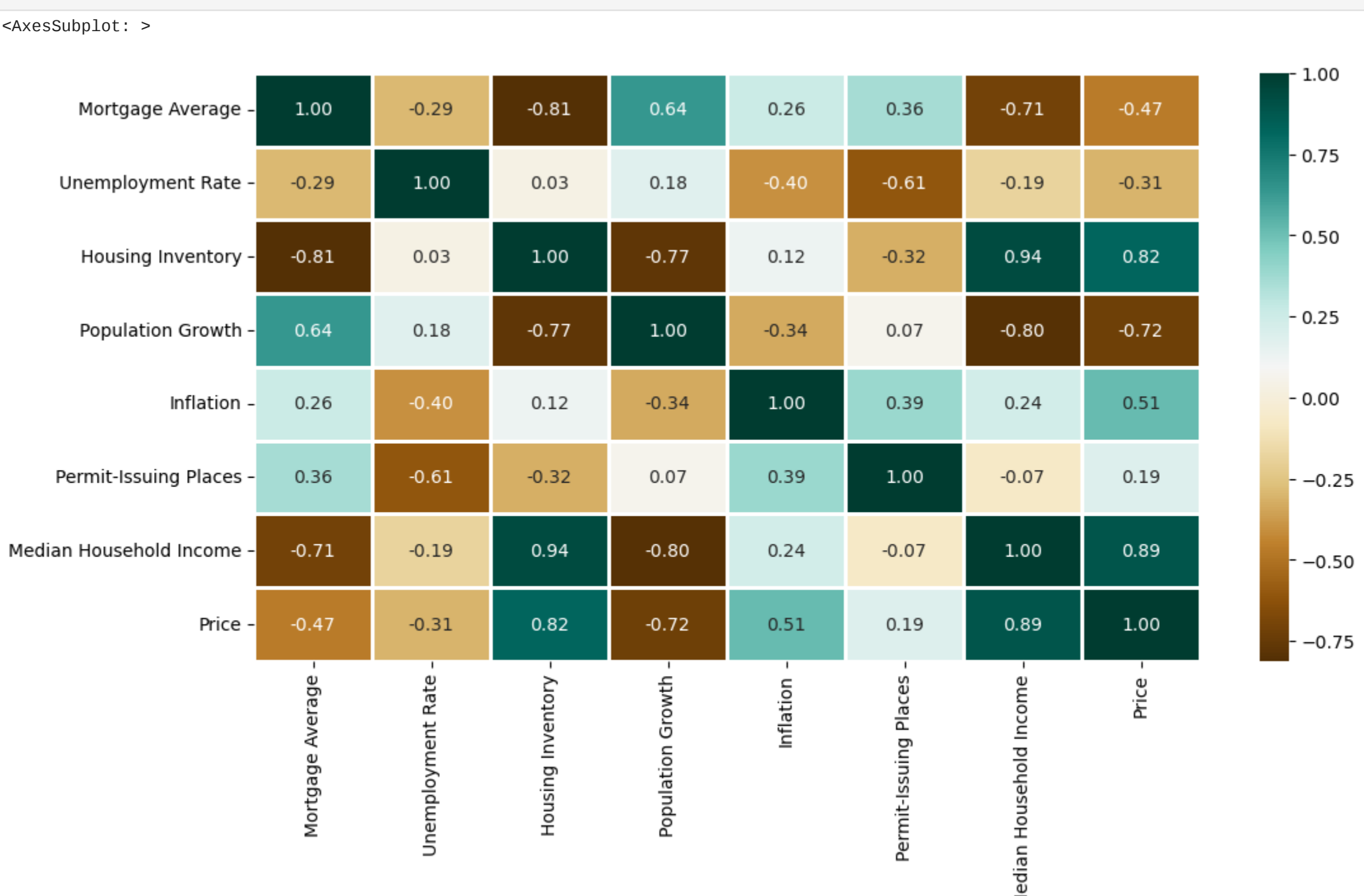
```
Out[8]: Mortgage Average      0
Unemployment Rate          0
Housing Inventory          0
Population Growth          0
Inflation                  0
Permit-Issuing Places      0
Median Household Income    0
Price                     0
dtype: int64
```

Check the correlation the data frame and also plot a heat map for visualization

```
In [9]: df.corr()
```

	Mortgage Average	Unemployment Rate	Housing Inventory	Population Growth	Inflation	Permit-Issuing Places	Median Household Income	Price
Mortgage Average	1.000000	-0.289183	-0.811049	0.635879	0.257359	0.359120	-0.705010	-0.469061
Unemployment Rate	-0.289183	1.000000	0.031282	0.177335	-0.402189	-0.606142	-0.185951	-0.307772
Housing Inventory	-0.811049	0.031282	1.000000	-0.771436	0.123763	-0.316246	0.936284	0.821443
Population Growth	0.635879	0.177335	-0.771436	1.000000	-0.336552	0.065970	-0.798108	-0.724927
Inflation	0.257359	-0.402189	0.123763	-0.336552	1.000000	0.386807	0.235315	0.513766
Permit-Issuing Places	0.359120	-0.606142	-0.316246	0.065970	0.386807	1.000000	-0.074142	0.188586
Median Household Income	-0.705010	-0.185951	0.936284	-0.798108	0.235315	-0.074142	1.000000	0.887793
Price	-0.469061	-0.307772	0.821443	-0.724927	0.513766	0.188586	0.887793	1.000000

```
In [10]: plt.figure(figsize=(12, 6))
sns.heatmap(df.corr(),
            cmap = 'BrBG',
            fmt = '.2f',
            linewidths = 2,
            annot = True)
```



Machine learning models which can predict the home price

```
In [11]: from sklearn.model_selection import train_test_split

X = df.drop(['Price'], axis=1)
Y = df['Price']

# Split the training set into
# training and validation set
X_train, X_test, y_train, y_test = train_test_split(
    X, Y, train_size=0.7, test_size=0.3, random_state=2)
```

Use Linear Regression

```
In [12]: from sklearn.linear_model import LinearRegression

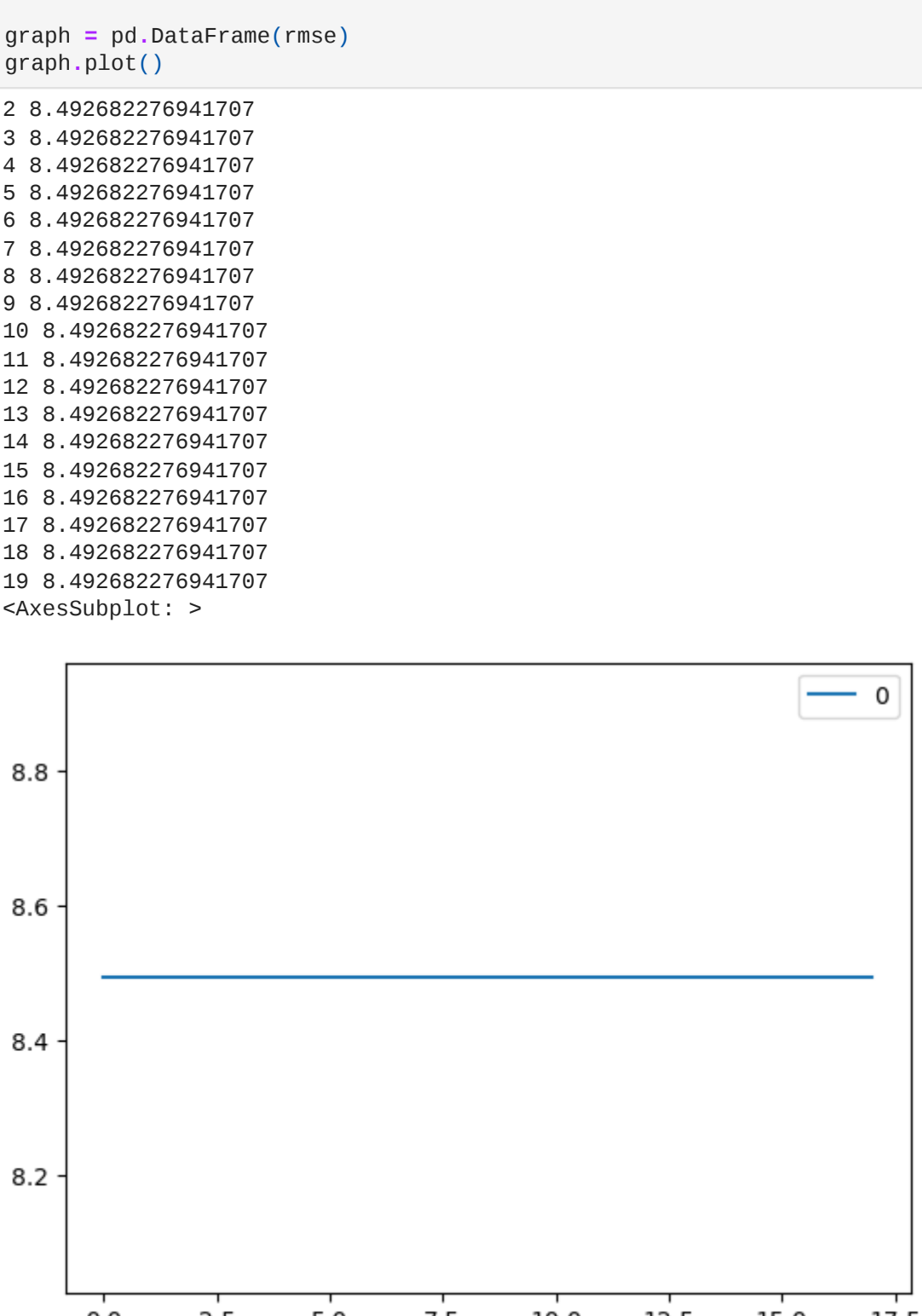
lr_model = LinearRegression()
lr_model.fit(X_train, y_train)
y_pred = lr_model.predict(X_test)
print("By using Linear Regression we got:")
print("Training data accuracy",lr_model.score(X_train, y_train))
print("Testing data accuracy",lr_model.score(X_test, y_test))

By using Linear Regression we got:
Training data accuracy 0.9762721807921758
Testing data accuracy 0.9626654707815229
```

```
In [13]: from sklearn.metrics import mean_squared_error
from math import sqrt
rmse = []

for k in range(2, 20):
    lr_model = LinearRegression()
    lr_model.fit(X_train, y_train)
    y_pred = lr_model.predict(X_test)

    error = sqrt(mean_squared_error(y_test, y_pred))
    rmse.append(error)
    print(k, error)
```



Use Random forest regression

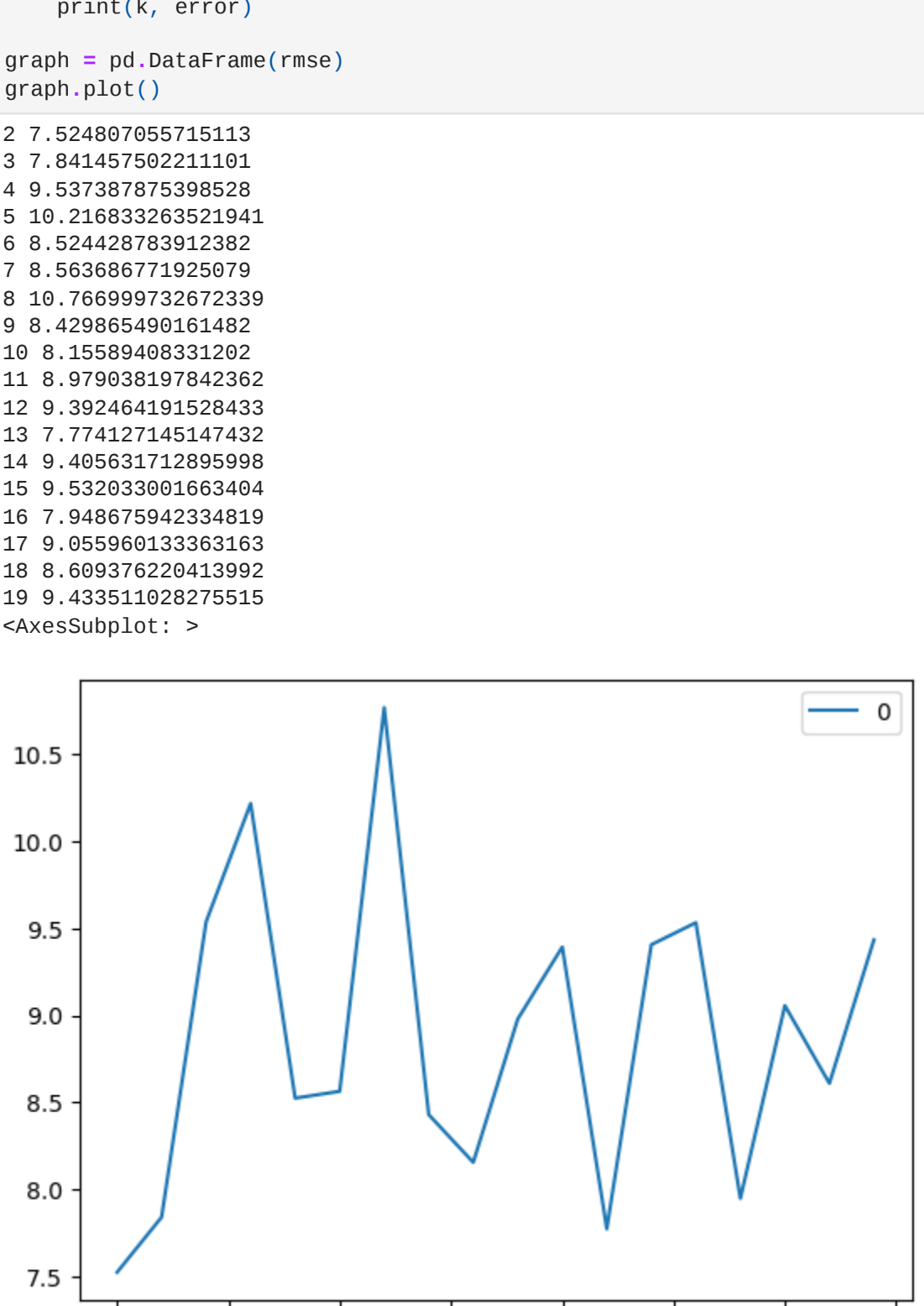
```
In [14]: from sklearn.ensemble import RandomForestRegressor
rfm=RandomForestRegressor()
rfm.fit(X_train,y_train)
y_pred = rfm.predict(X_test)
print("By using Random Forest Regression we got:")
print("Training data accuracy",rfm.score(X_train,y_train))
print("Testing data accuracy",rfm.score(X_test,y_test))

By using Random Forest Regression we got:
Training data accuracy 0.9964916398758433
Testing data accuracy 0.9656225151047019
```

```
In [15]: from sklearn.metrics import mean_squared_error
from math import sqrt
rmse = []

for k in range(2, 20):
    rfm=RandomForestRegressor()
    rfm.fit(X_train,y_train)
    y_pred = rfm.predict(X_test)

    error = sqrt(mean_squared_error(y_test, y_pred))
    rmse.append(error)
    print(k, error)
```



Final Insight from the Data

With the above chart we can say the following features are effecting the home price in US

1. House inventory is directly effect the price of houses.
2. Inflation also effect the price of houses.
3. Median household income also effect the price of houses.

Linear regression work fine with our Data with high accuracy

```
In [ ]:
```