

ZHVI dataset comes from [https://www.kaggle.com/datasets/paultimothymooney/zillow-house-price-data?select=Sale\\_Prices\\_City.csv](https://www.kaggle.com/datasets/paultimothymooney/zillow-house-price-data?select=Sale_Prices_City.csv)

Unemployment rate dataset comes from <https://www.kaggle.com/datasets/axeltorbenson/unemployment-data-19482021>

Inflation Rate(CPI) Dataset <https://www.kaggle.com/datasets/varpit94/us-inflation-data-updated-till-may-2021>

Interest rate dataset <https://www.kaggle.com/datasets/raoofiali/us-interest-rate-weekly>

GDP Growth Rate dataset <https://www.kaggle.com/datasets/rajkumarpandey02/economy-of-the-united-states>

```
1 #!pip install ydata-profiling
2 #!pip install tensorflow
3
4 import pandas as pd
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import statsmodels.api as sm
8 import kagglehub
9 import math
10 import os
11 import warnings
12
13 #from ydata_profiling import ProfileReport
14 from sklearn.model_selection import train_test_split
15 from sklearn.linear_model import Ridge
16 from sklearn.linear_model import Lasso
17 from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error, mean_abs
18
19 from sklearn.ensemble import RandomForestRegressor
20 from sklearn.preprocessing import PolynomialFeatures
21 from sklearn.preprocessing import StandardScaler
22 #from tensorflow.keras.models import Sequential
23 #from tensorflow.keras.layers import Dense
24 from IPython.display import clear_output, display, HTML
25
26 warnings.filterwarnings("ignore")
27 clear_output()
```



Adding Housing Data

```

1 # Download housing data
2 path = kagglehub.dataset_download("paultimothymooney/zillow-house-price-data")
3
4 print("Files in the dataset:")
5 for root, dirs, files in os.walk(path):
6     for file in files:
7         print(os.path.join(root, file))

```



Files in the dataset:

```

/kaggle/input/zillow-house-price-data/Sale_Prices_State.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_2bedroom.csv
/kaggle/input/zillow-house-price-data/DaysOnZillow_State.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_Sfr.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_SingleFamilyResidence.csv
/kaggle/input/zillow-house-price-data/State_Zri_SingleFamilyResidenceRental.csv
/kaggle/input/zillow-house-price-data/City_Zri_SingleFamilyResidenceRental.csv
/kaggle/input/zillow-house-price-data/City_Zri_AllHomesPlusMultifamily.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_3bedroom.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_1Bedroom.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_5BedroomOrMore.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_Studio.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_3Bedroom.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_4Bedroom.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_5BedroomOrMore.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_AllHomes.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_2bedroom.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_Condominum.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_AllHomes.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_4bedroom.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_3bedroom.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_1bedroom.csv
/kaggle/input/zillow-house-price-data/Sale_Prices_City.csv
/kaggle/input/zillow-house-price-data/DaysOnZillow_City.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_3Bedroom.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_Sfr.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_AllHomes.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_SingleFamilyResidence.csv
/kaggle/input/zillow-house-price-data/State_Zri_AllHomesPlusMultifamily.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_2Bedroom.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_5BedroomOrMore.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_Condominum.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_4Bedroom.csv
/kaggle/input/zillow-house-price-data/City_MedianRentalPrice_Studio.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_4bedroom.csv
/kaggle/input/zillow-house-price-data/State_Zhvi_1bedroom.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_5BedroomOrMore.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_2Bedroom.csv
/kaggle/input/zillow-house-price-data/City_Zhvi_AllHomes.csv
/kaggle/input/zillow-house-price-data/State_MedianRentalPrice_1Bedroom.csv

```



```

1 csv_path = os.path.join(path, "City_Zhvi_AllHomes.csv")
2 df = pd.read_csv(csv_path)
3 print(df.head())

```

```

⇒ Unnamed: 0  RegionID  SizeRank  RegionName  RegionType  StateName  State  \
0           0      6181         0    New York      City      NY      NY
1           1     12447         1  Los Angeles      City      CA      CA
2           2     39051         2    Houston      City      TX      TX
3           3     17426         3    Chicago      City      IL      IL
4           4      6915         4  San Antonio      City      TX      TX

           Metro      CountyName  1996-01-31  ...  \
0  New York-Newark-Jersey City  Queens County  196258.0  ...
1  Los Angeles-Long Beach-Anaheim  Los Angeles County  185649.0  ...
2  Houston-The Woodlands-Sugar Land  Harris County  93518.0  ...
3  Chicago-Naperville-Elgin  Cook County  130920.0  ...
4  San Antonio-New Braunfels  Bexar County  94041.0  ...

           2019-06-30  2019-07-31  2019-08-31  2019-09-30  2019-10-31  2019-11-30  \
0  659421.0  659007.0  658239.0  656925.0  655613.0  654394.0
1  712660.0  713807.0  715688.0  718245.0  721896.0  725180.0
2  186844.0  187464.0  188070.0  188496.0  189125.0  189612.0
3  248372.0  248646.0  248725.0  248483.0  248278.0  248090.0
4  182732.0  183350.0  183930.0  184846.0  185490.0  186244.0

           2019-12-31  2020-01-31  2020-02-29  2020-03-31
0  653930.0  653901.0  653565.0  652307.0
1  730358.0  735910.0  744137.0  752508.0
2  190179.0  190395.0  190938.0  191907.0
3  248029.0  248220.0  248599.0  249152.0
4  186420.0  186962.0  187129.0  187718.0

```

[5 rows x 300 columns]

```

1 # remove rows with NaN
2 df_cleaned = df.dropna()
3 print("DataFrame after removing rows with any NaN values:")
4 print(df_cleaned.head())
5 data = df_cleaned

```

```

⇒ DataFrame after removing rows with any NaN values:
   Unnamed: 0  RegionID  SizeRank  RegionName  RegionType  StateName  State  \
0           0      6181         0    New York      City      NY      NY
1           1     12447         1  Los Angeles      City      CA      CA
2           2     39051         2    Houston      City      TX      TX
3           3     17426         3    Chicago      City      IL      IL
4           4      6915         4  San Antonio      City      TX      TX

           Metro      CountyName  1996-01-31  ...  \
0  New York-Newark-Jersey City  Queens County  196258.0  ...
1  Los Angeles-Long Beach-Anaheim  Los Angeles County  185649.0  ...
2  Houston-The Woodlands-Sugar Land  Harris County  93518.0  ...
3  Chicago-Naperville-Elgin  Cook County  130920.0  ...
4  San Antonio-New Braunfels  Bexar County  94041.0  ...

```

	2019-06-30	2019-07-31	2019-08-31	2019-09-30	2019-10-31	2019-11-30	\
0	659421.0	659007.0	658239.0	656925.0	655613.0	654394.0	
1	712660.0	713807.0	715688.0	718245.0	721896.0	725180.0	
2	186844.0	187464.0	188070.0	188496.0	189125.0	189612.0	
3	248372.0	248646.0	248725.0	248483.0	248278.0	248090.0	
4	182732.0	183350.0	183930.0	184846.0	185490.0	186244.0	


	2019-12-31	2020-01-31	2020-02-29	2020-03-31
0	653930.0	653901.0	653565.0	652307.0
1	730358.0	735910.0	744137.0	752508.0
2	190179.0	190395.0	190938.0	191907.0
3	248029.0	248220.0	248599.0	249152.0
4	186420.0	186962.0	187129.0	187718.0

[5 rows x 300 columns]

```

1 # Remove location identifier since only one city has data for each month/year
2 data.drop('State',axis=1,inplace=True)
3 data.drop('CountyName',axis=1,inplace=True)
4 data.drop('SizeRank',axis=1,inplace=True)
5 data.drop('Metro',axis=1,inplace=True)
6 data.drop('Unnamed: 0',axis=1,inplace=True)
7 data.drop('RegionID',axis=1,inplace=True)
8 data.drop('RegionType',axis=1,inplace=True)
9 data.drop('StateName',axis=1,inplace=True)
10 data = data.reset_index(drop=True)
11
12 # Select single city (New York)
13 data = data[data['RegionName']=='New York']
14 data.drop('RegionName',axis=1,inplace=True)
15 print(data)

```



	1996-01-31	1996-02-29	1996-03-31	1996-04-30	1996-05-31	1996-06-30	\
0	196258.0	195693.0	195383.0	194836.0	194652.0	194520.0	

	1996-07-31	1996-08-31	1996-09-30	1996-10-31	...	2019-06-30	\
0	194447.0	194313.0	194271.0	194341.0	...	659421.0	

	2019-07-31	2019-08-31	2019-09-30	2019-10-31	2019-11-30	2019-12-31	\
0	659007.0	658239.0	656925.0	655613.0	654394.0	653930.0	

	2020-01-31	2020-02-29	2020-03-31
0	653901.0	653565.0	652307.0

[1 rows x 291 columns]

## Adding Interest Rate Data

```

1 path = kagglehub.dataset_download("raoofiali/us-interest-rate-weekly")
2

```

```

3 print("Files in the dataset:")
4 for root, dirs, files in os.walk(path):
5     for file in files:
6         print(os.path.join(root, file))
7
8 xlsx_path = os.path.join(path, "Us-Interest Rate-Weekly.xlsx")
9 ir_df = pd.read_excel(xlsx_path)
10 ir_df.drop('Unnamed: 0',axis=1,inplace=True)
11 print(ir_df.head())
12 print(ir_df.tail())

```



Files in the dataset:

/kaggle/input/us-interest-rate-weekly/Us-Interest Rate-Weekly.xlsx

	Date	Value
0	1971-08-04	5.50
1	1971-08-15	5.50
2	1971-08-16	5.75
3	1971-08-31	5.75
4	1971-09-01	5.13
	Date	Value
1678	2024-02-29	5.5
1679	2024-03-19	5.5
1680	2024-03-20	5.5
1681	2024-04-30	5.5
1682	2024-05-01	5.5

```

1 # convert date format
2 ir_df['Date'] = pd.to_datetime(ir_df['Date'])
3
4 # Filter to include only rows between January 1996 and March 2020 to match housing data
5 start_date = pd.to_datetime('1996-01-01')
6 end_date = pd.to_datetime('2020-03-31')
7 filtered_ir_df = ir_df[(ir_df['Date'] >= start_date) & (ir_df['Date'] <= end_date)]
8
9 # Resample the data to get the monthly average
10 ir_df = filtered_ir_df.resample('M', on='Date').mean().reset_index()
11
12 # create time index
13 ir_df['Year'] = ir_df['Date'].dt.year
14 ir_df['Month'] = ir_df['Date'].dt.month
15 ir_df['TimeIndex'] = (ir_df['Year'] - ir_df['Year'].min()) * 12 + (ir_df['Month'] - 1)
16 ir_df.drop('Date',axis=1,inplace=True)
17
18 print(ir_df.head())
19 print(ir_df.tail())

```



	Value	Year	Month	TimeIndex
0	5.375	1996	1	0
1	5.250	1996	2	1
2	5.250	1996	3	2
3	5.250	1996	4	3
4	5.250	1996	5	4

	Value	Year	Month	TimeIndex
286	1.750	2019	11	286
287	1.750	2019	12	287
288	1.750	2020	1	288
289	1.750	2020	2	289
290	1.125	2020	3	290

## Adding Inflation Rate Data

```

1 path = kagglehub.dataset_download("varpit94/us-inflation-data-updated-till-may-2021")
2
3 print("Files in the dataset:")
4 for root, dirs, files in os.walk(path):
5     for file in files:
6         print(os.path.join(root, file))
7
8 csv_path = os.path.join(path, "US CPI.csv")
9 cpi_df = pd.read_csv(csv_path)
10
11 print(cpi_df.head())
12 print(cpi_df.tail())

```



```

Files in the dataset:
/kaggle/input/us-inflation-data-updated-till-may-2021/US CPI.csv
   Yearmon  CPI
0  01-01-1913  9.8
1  01-02-1913  9.8
2  01-03-1913  9.8
3  01-04-1913  9.8
4  01-05-1913  9.7
   Yearmon      CPI
1298  01-03-2021  264.877
1299  01-04-2021  267.054
1300  01-05-2021  269.195
1301  01-06-2021  271.696
1302  01-07-2021  273.003


```

```

1 cpi_df['Yearmon'] = pd.to_datetime(cpi_df['Yearmon'], format='%d-%m-%Y')
2
3 start_date = pd.to_datetime('1996-01-01')
4 end_date = pd.to_datetime('2020-03-31')
5 filtered_cpi_df = cpi_df[(cpi_df['Yearmon'] >= start_date) & (cpi_df['Yearmon'] <= end_date)]
6 filtered_cpi_df = filtered_cpi_df.reset_index(drop=True)
7
8 filtered_cpi_df['Year'] = filtered_cpi_df['Yearmon'].dt.year
9 filtered_cpi_df['Month'] = filtered_cpi_df['Yearmon'].dt.month
10 filtered_cpi_df['TimeIndex'] = (filtered_cpi_df['Year'] - filtered_cpi_df['Year'].min())
11 filtered_cpi_df = filtered_cpi_df.reset_index(drop=True)
12
13 print(filtered_cpi_df)

```





	Yearmon	CPI	Year	Month	TimeIndex
0	1996-01-01	154.400	1996	1	0
1	1996-02-01	154.900	1996	2	1
2	1996-03-01	155.700	1996	3	2
3	1996-04-01	156.300	1996	4	3
4	1996-05-01	156.600	1996	5	4
..	...	...	...	...	...
286	2019-11-01	257.208	2019	11	286
287	2019-12-01	256.974	2019	12	287
288	2020-01-01	257.971	2020	1	288
289	2020-02-01	258.678	2020	2	289
290	2020-03-01	258.115	2020	3	290


[291 rows x 5 columns]

## Adding Unemployment rate data

```

1 # download unemployment rate data
2 path = kagglehub.dataset_download("axeltorbenson/unemployment-data-19482021")
3
4 print("Files in the dataset:")
5 for root, dirs, files in os.walk(path):
6     for file in files:
7         print(os.path.join(root, file))
8
9 # Load CSV file
10 csv_path = os.path.join(path, "unemployment_rate_data.csv")
11 un_df = pd.read_csv(csv_path)
12
13 print(un_df.head())
14 print(un_df.tail())

```



Files in the dataset:

/kaggle/input/unemployment-data-19482021/unemployment\_rate\_data.csv

	date	unrate	unrate_men	unrate_women	unrate_16_to_17	\
0	1/1/1948	4.0	4.2	3.5	10.8	
1	2/1/1948	4.7	4.7	4.8	15.0	
2	3/1/1948	4.5	4.5	4.4	13.2	
3	4/1/1948	4.0	4.0	4.1	9.9	
4	5/1/1948	3.4	3.3	3.4	6.4	

	unrate_18_to_19	unrate_20_to_24	unrate_25_to_34	unrate_35_to_44	\
0	9.6	6.6	3.6	2.6	
1	9.5	8.0	4.0	3.2	
2	9.3	8.6	3.5	3.2	
3	8.1	6.8	3.5	3.1	
4	7.2	6.3	2.8	2.5	

	unrate_45_to_54	unrate_55_over
0	2.7	3.6
1	3.4	4.0
2	2.9	3.5

3		2.9		3.2	
4		2.3		2.9	
	date	unrate	unrate_men	unrate_women	unrate_16_to_17 \
882	7/1/2021	5.7	5.5	5.8	12.8
883	8/1/2021	5.3	5.1	5.5	10.7
884	9/1/2021	4.6	4.6	4.5	9.2
885	10/1/2021	4.3	4.2	4.4	8.6
886	11/1/2021	3.9	3.9	3.9	9.7
	unrate_18_to_19	unrate_20_to_24	unrate_25_to_34	unrate_35_to_44 \	
882	9.9	9.5	6.3	4.8	
883	11.0	9.1	5.8	4.4	
884	12.6	7.7	5.0	3.8	
885	12.7	6.8	4.5	3.6	
886	11.0	6.6	3.8	3.6	
	unrate_45_to_54	unrate_55_over			
882	4.0	4.6			
883	4.2	4.1			
884	3.7	3.3			
885	3.5	3.3			
886	2.8	3.1			

```

1 # select same range of dates of housing data and only the overall unemployment rate
2 un_df = un_df.iloc[576:576+291][['unrate', 'date']]
3 un_df = un_df.reset_index(drop=True)
4
5 # Convert the date column to get specific year and month feature
6 un_df['date'] = pd.to_datetime(un_df['date'])
7 un_df['Year'] = un_df['date'].dt.year
8 un_df['Month'] = un_df['date'].dt.month
9 un_df['TimeIndex'] = (un_df['Year'] - un_df['Year'].min()) * 12 + (un_df['Month'] - un_df
10 un_df.drop('date', axis=1, inplace=True)

```

## Adding GDP Growth %

```

1 # Download data
2 path = kagglehub.dataset_download("rajkumarpandey02/economy-of-the-united-states")
3
4 print("Path to dataset files:", path)
5
6 print("Files in the dataset:")
7 for root, dirs, files in os.walk(path):
8     for file in files:
9         print(os.path.join(root, file))
10
11 csv_path = os.path.join(path, "Economy of the United States.csv")
12 gdp_df = pd.read_csv(csv_path)
13
14 print(gdp_df.head())
15 print(gdp_df.tail())

```







Path to dataset files: /kaggle/input/economy-of-the-united-states

Files in the dataset:

/kaggle/input/economy-of-the-united-states/Economy of the United States.csv

Unnamed: 0	Year	GDP (in Bil. US\$PPP)	GDP per capita (in US\$ PPP)	\
0	0 1980	2857.3	12552.9	
1	1 1981	3207.0	13948.7	
2	2 1982	3343.8	14405.0	
3	3 1983	3634.0	15513.7	
4	4 1984	4037.7	17086.4	

	GDP (in Bil. US\$nominal)	GDP per capita (in US\$ nominal)	\
0	2857.3	12552.9	
1	3207.0	13948.7	
2	3343.8	14405.0	
3	3634.0	15513.7	
4	4037.7	17086.4	

	GDP growth (real)	Inflation rate (in Percent)	Unemployment (in Percent)	\
0	-0.30%	13.50%	7.20%	
1	2.50%	10.40%	7.60%	
2	-1.80%	6.20%	9.70%	
3	4.60%	3.20%	9.60%	
4	7.20%	4.40%	7.50%	

	Government debt (in % of GDP)
0	NaN
1	NaN
2	NaN
3	NaN
4	NaN

Unnamed: 0	Year	GDP (in Bil. US\$PPP)	GDP per capita (in US\$ PPP)	\
43	43 2023	26185.2	78421.9	
44	44 2024	27057.2	80779.3	
45	45 2025	28045.3	83463.2	
46	46 2026	29165.5	86521.2	
47	47 2027	30281.5	89546.4	

	GDP (in Bil. US\$nominal)	GDP per capita (in US\$ nominal)	\
43	26185.2	78421.9	
44	27057.2	80779.3	
45	28045.3	83463.2	
46	29165.5	86521.2	
47	30281.5	89546.4	

	GDP growth (real)	Inflation rate (in Percent)	Unemployment (in Percent)	\
43	1.00%	3.50%	4.60%	
44	1.20%	2.20%	5.40%	
45	1.80%	2.00%	5.40%	
46	2.10%	2.00%	4.90%	
47	1.90%	2.00%	4.70%	

	Government debt (in % of GDP)
43	122.90%
44	126.00%
45	129.40%

46	132.20%
47	134.90%

```

1 gdp_df = gdp_df[gdp_df['Year'] >= 1996]
2 gdp_df = gdp_df[gdp_df['Year'] <= 2020]
3 gdp_df = gdp_df.reset_index(drop=True)
4 gdp_df = gdp_df[['Year', 'GDP growth (real)']]
5
6 gdp_df['GDP growth (real)'] = gdp_df['GDP growth (real)'].str.replace('%', '')
7 gdp_df['GDP Growth'] = pd.to_numeric(gdp_df['GDP growth (real)'])
8 gdp_df.drop('GDP growth (real)',axis=1,inplace=True)
9
10 # add instance for each month
11 gdp_df = gdp_df.loc[gdp_df.index.repeat(12)].reset_index(drop=True)
12 gdp_df['Month'] = (gdp_df.groupby('Year').cumcount() % 12) + 1
13 gdp_df = gdp_df.iloc[:-9]
14
15 print(gdp_df.head())
16 print(gdp_df.tail())

```

```

↩
   Year  GDP Growth  Month
0  1996          3.8      1
1  1996          3.8      2
2  1996          3.8      3
3  1996          3.8      4
4  1996          3.8      5
   Year  GDP Growth  Month
286  2019          2.3     11
287  2019          2.3     12
288  2020         -3.4      1
289  2020         -3.4      2
290  2020         -3.4      3

```

```

1 # reshape data to have rows correspond to each time, with features being the time, price
2 reshaped_data = []
3
4 # Loop through each column to get feature dates
5 for column in data.columns:
6     year, month, day = map(int, column.split('-'))
7
8     # Loop through each row to get price for the current date
9     for index, row in data.iterrows():
10         zhvi = row[column]
11
12         reshaped_data.append({
13             'ZHVI': zhvi,
14             'Year': year,
15             'Month': month,
16             'Year-Month': f'{year}-{month}'
17         })
18

```

```

19 reshaped_df = pd.DataFrame(reshaped_data)
20
21 # Add a time index
22 reshaped_df['TimeIndex'] = (reshaped_df['Year'] - reshaped_df['Year'].min()) * 12 + (res
23
24 # Sort data by month/year
25 full_df = reshaped_df.sort_values(by=['Year', 'Month']).reset_index(drop=True)
26 full_df['Unemployment Rate'] = un_df['unrate']
27 full_df['CPI'] = filtered_cpi_df['CPI']
28 full_df['Interest Rate'] = ir_df['Value']
29 full_df['GDP Growth'] = gdp_df['GDP Growth']
30 print("Reshaped DataFrame:")
31 print(full_df)

```



Reshaped DataFrame:

	ZHVI	Year	Month	Year-Month	TimeIndex	Unemployment Rate	CPI \
0	196258.0	1996	1	1996-1	0	6.3	154.400
1	195693.0	1996	2	1996-2	1	6.0	154.900
2	195383.0	1996	3	1996-3	2	5.8	155.700
3	194836.0	1996	4	1996-4	3	5.4	156.300
4	194652.0	1996	5	1996-5	4	5.4	156.600
..	...	...	...	...	...	...	...
286	654394.0	2019	11	2019-11	286	3.3	257.208
287	653930.0	2019	12	2019-12	287	3.4	256.974
288	653901.0	2020	1	2020-1	288	4.0	257.971
289	653565.0	2020	2	2020-2	289	3.8	258.678
290	652307.0	2020	3	2020-3	290	4.5	258.115

	Interest Rate	GDP Growth
0	5.375	3.8
1	5.250	3.8
2	5.250	3.8
3	5.250	3.8
4	5.250	3.8
..	...	...
286	1.750	2.3
287	1.750	2.3
288	1.750	-3.4
289	1.750	-3.4
290	1.125	-3.4

[291 rows x 9 columns]



```

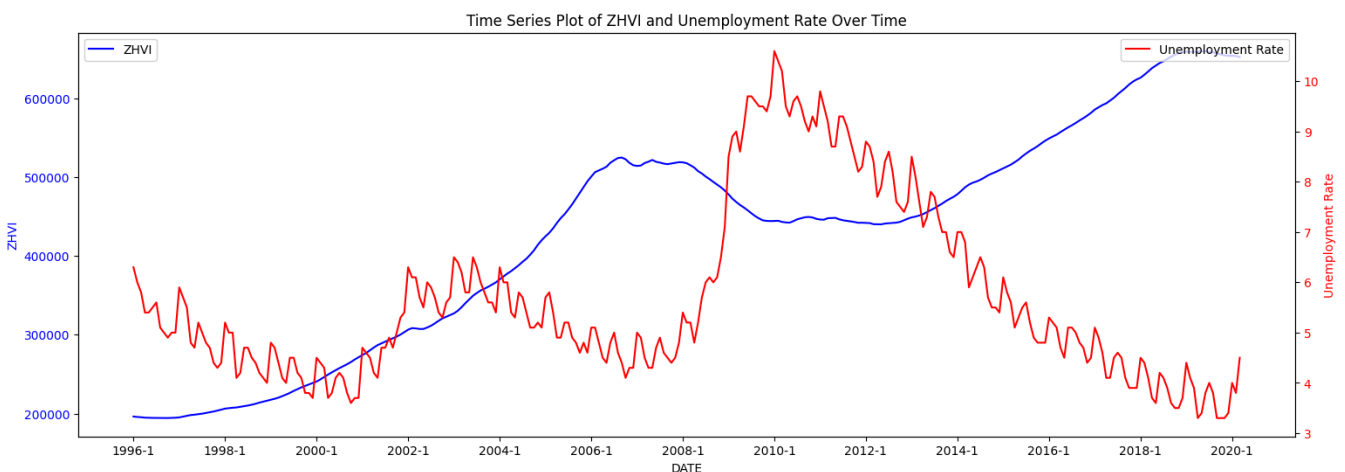
1 # Create figure and primary axis
2 fig, ax1 = plt.subplots(figsize=(18, 6))
3
4 # Plot the first ZHVI dataset
5 ax1.plot(full_df['Year-Month'], full_df['ZHVI'], color='blue', label='ZHVI')
6 ax1.set_xlabel('DATE')
7 ax1.set_ylabel('ZHVI', color='blue')
8 ax1.tick_params(axis='y', labelcolor='blue')
9
10 # Create a second axis sharing the same x-axis

```

```

11 ax2 = ax1.twinx()
12
13 # Plot the Unemployment Rate data
14 ax2.plot(un_df['TimeIndex'], un_df['unrate'], color='red', label='Unemployment Rate')
15 ax2.set_ylabel('Unemployment Rate', color='red')
16 ax2.tick_params(axis='y', labelcolor='red')
17
18 x_ticks = np.arange(0, 290, 24)
19 ax1.set_xticks(x_ticks)
20
21 plt.title('Time Series Plot of ZHVI and Unemployment Rate Over Time')
22
23 # legend
24 ax1.legend(loc='upper left')
25 ax2.legend(loc='upper right')
26
27 plt.show()

```



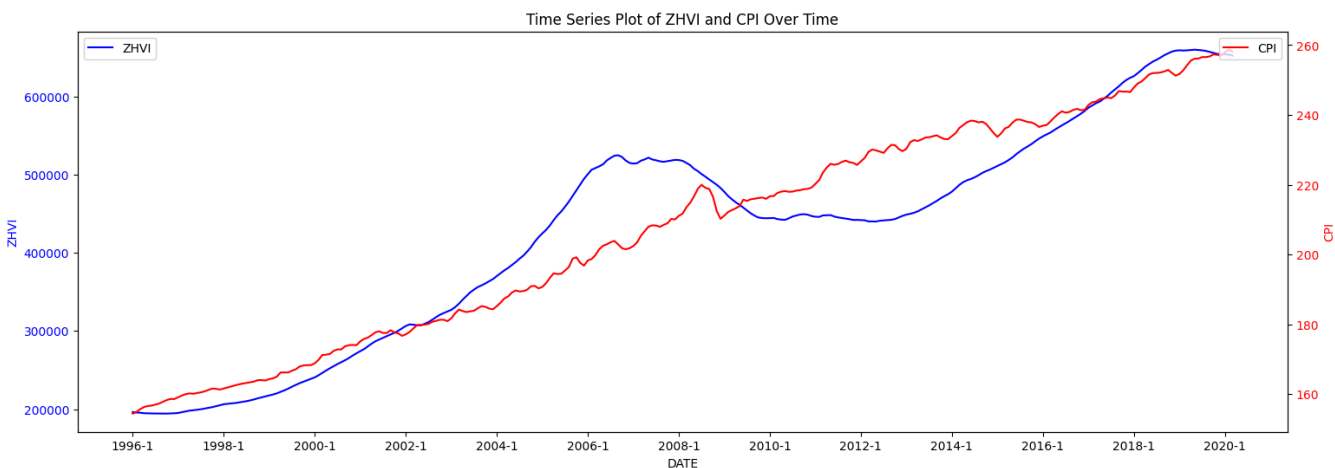
```

1 # Create figure and primary axis
2 fig, ax1 = plt.subplots(figsize=(18, 6))
3

```

```
4 # Plot the first ZHVI dataset
5 ax1.plot(full_df['Year-Month'], full_df['ZHVI'], color='blue', label='ZHVI')
6 ax1.set_xlabel('DATE')
7 ax1.set_ylabel('ZHVI', color='blue')
8 ax1.tick_params(axis='y', labelcolor='blue')
9
10 # Create a second axis sharing the same x-axis
11 ax2 = ax1.twinx()
12
13 # Plot the Unemployment Rate data
14 ax2.plot(filtered_cpi_df['TimeIndex'], filtered_cpi_df['CPI'], color='red', label='CPI')
15 ax2.set_ylabel('CPI', color='red')
16 ax2.tick_params(axis='y', labelcolor='red')
17
18 x_ticks = np.arange(0, 290, 24)
19 ax1.set_xticks(x_ticks)
20
21 plt.title('Time Series Plot of ZHVI and CPI Over Time')
22
23 # legend
24 ax1.legend(loc='upper left')
25 ax2.legend(loc='upper right')
26
27 plt.show()
```





```

1 # Create figure and primary axis
2 fig, ax1 = plt.subplots(figsize=(18, 6))
3
4 # Plot the first ZHVI dataset
5 ax1.plot(full_df['Year-Month'], full_df['ZHVI'], color='blue', label='ZHVI')
6 ax1.set_xlabel('DATE')
7 ax1.set_ylabel('ZVHI', color='blue')
8 ax1.tick_params(axis='y', labelcolor='blue')
9
10 # Create a second axis sharing the same x-axis
11 ax2 = ax1.twinx()
12
13 # Plot the Unemployment Rate data
14 ax2.plot(ir_df['TimeIndex'], ir_df['Value'], color='red', label='Interest Rate')
15 ax2.set_ylabel('Interest Rate', color='red')
16 ax2.tick_params(axis='y', labelcolor='red')
17
18 x_ticks = np.arange(0, 290, 24)
19 ax1.set_xticks(x_ticks)
20

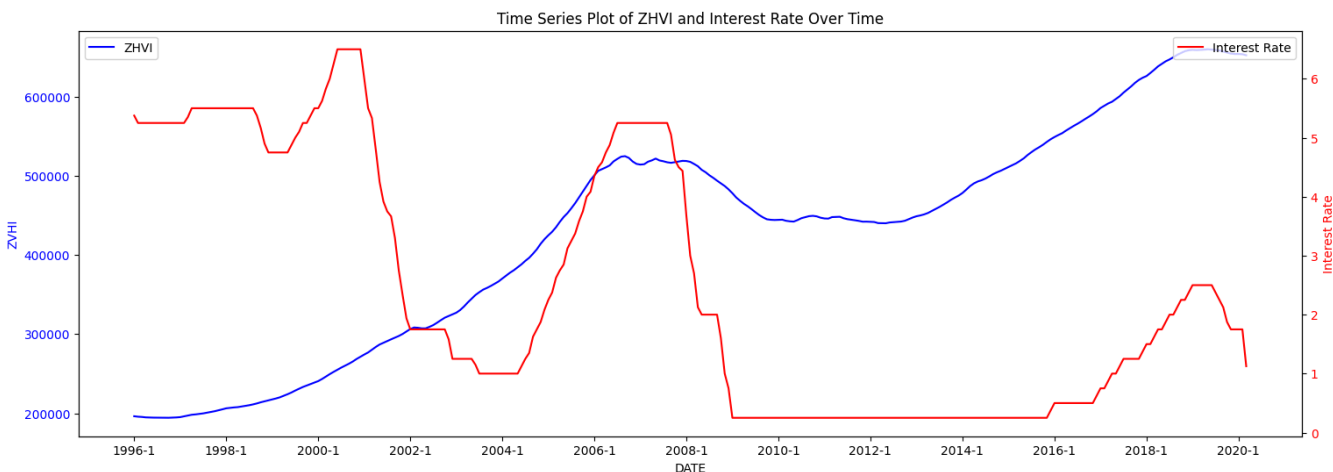
```



```

21 plt.title('Time Series Plot of ZHVI and Interest Rate Over Time')
22
23 # legend
24 ax1.legend(loc='upper left')
25 ax2.legend(loc='upper right')
26
27 plt.show()

```



```

1 # Create figure and primary axis
2 fig, ax1 = plt.subplots(figsize=(18, 6))
3
4 # Plot the first ZHVI dataset
5 ax1.plot(full_df['Year-Month'], full_df['ZHVI'], color='blue', label='ZHVI')
6 ax1.set_xlabel('DATE')
7 ax1.set_ylabel('ZHVI', color='blue')
8 ax1.tick_params(axis='y', labelcolor='blue')
9
10 # Create a second axis sharing the same x-axis
11 ax2 = ax1.twinx()
12
13 # Plot the Unemployment Rate data

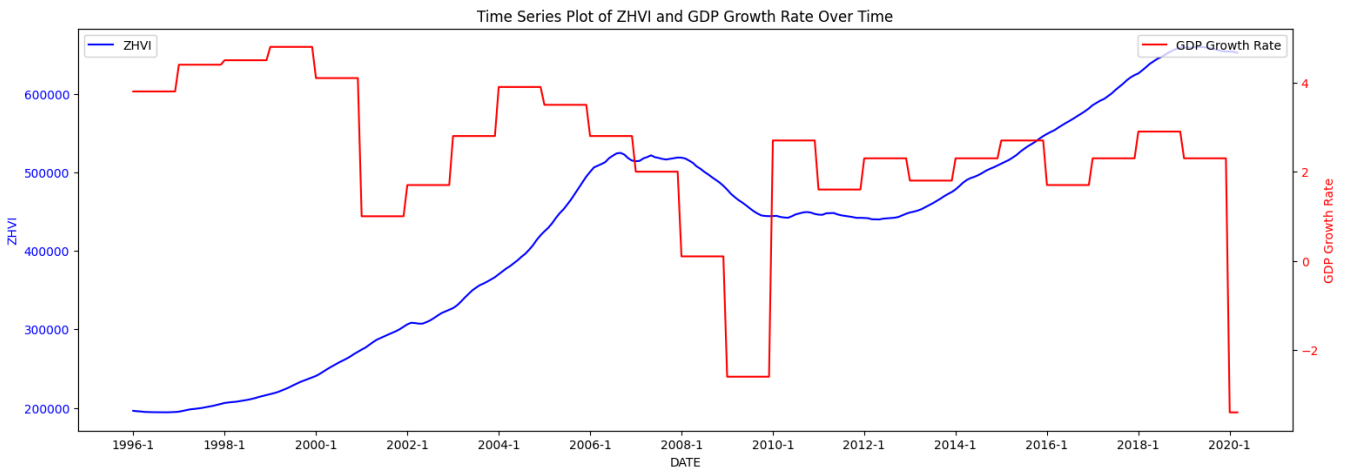
```



```

14 ax2.plot(ir_df['TimeIndex'], gdp_df['GDP Growth'], color='red', label='GDP Growth Rate')
15 ax2.set_ylabel('GDP Growth Rate', color='red')
16 ax2.tick_params(axis='y', labelcolor='red')
17
18 x_ticks = np.arange(0, 290, 24)
19 ax1.set_xticks(x_ticks)
20
21 plt.title('Time Series Plot of ZHVI and GDP Growth Rate Over Time')
22
23 # legend
24 ax1.legend(loc='upper left')
25 ax2.legend(loc='upper right')
26
27 plt.show()

```



## OLS Model

```

1 # Split data into training and test
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month']

```



```
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rat
7 y_train = train['ZHVI']
8
9 # Prediction test
10 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate',
11
12 # add polynomial features and scale
13 scaler = StandardScaler()
14 poly = PolynomialFeatures(degree=2)
15
16 X_train_poly = poly.fit_transform(X_train)
17 X_test_poly = poly.transform(X_test)
18
19 X_train_scaled = scaler.fit_transform(X_train_poly)
20 X_test_scaled = scaler.transform(X_test_poly)
21
22 # add constant
23 X_train_scaled = sm.add_constant(X_train_scaled)
24 X_test_scaled = sm.add_constant(X_test_scaled)
25
26 # Fit OLS model
27 model = sm.OLS(y_train, X_train_scaled)
28 results = model.fit()
29
30 predictions = results.predict(X_test_scaled)
31 test['Predicted_ZHVI'] = predictions
32
33 y_test = test['ZHVI']
34 y_pred = test['Predicted_ZHVI']
35 OLS_pred = test['Predicted_ZHVI']
36
37 # model evaluation
38 rmse = math.sqrt(mean_squared_error(y_test, y_pred))
39 print(f"OLS Root Mean Squared Error (RMSE): {rmse}")
40 mape = mean_absolute_percentage_error(y_test, y_pred)
41 print("OLS Mean Absolute Percentage Error(MAPE):", mape)
42 MAE = mean_absolute_error(y_test, y_pred)
43 print("OLS Mean Absolute Error(MAE):", MAE)
44 r2 = r2_score(y_pred, y_test)
45 print(f"OLS R-squared(R^2): {r2}")
46
47 # Plot truth vs prediction
48 plt.figure(figsize=(18, 6))
49 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
50 plt.plot(test['Year-Month'], test['Predicted_ZHVI'], color='blue', label='Predicted ZHVI')
51 plt.xlabel('Year')
52 plt.ylabel('Price')
53 x_ticks = np.arange(0, 80, 6)
54 plt.xticks(x_ticks)
```



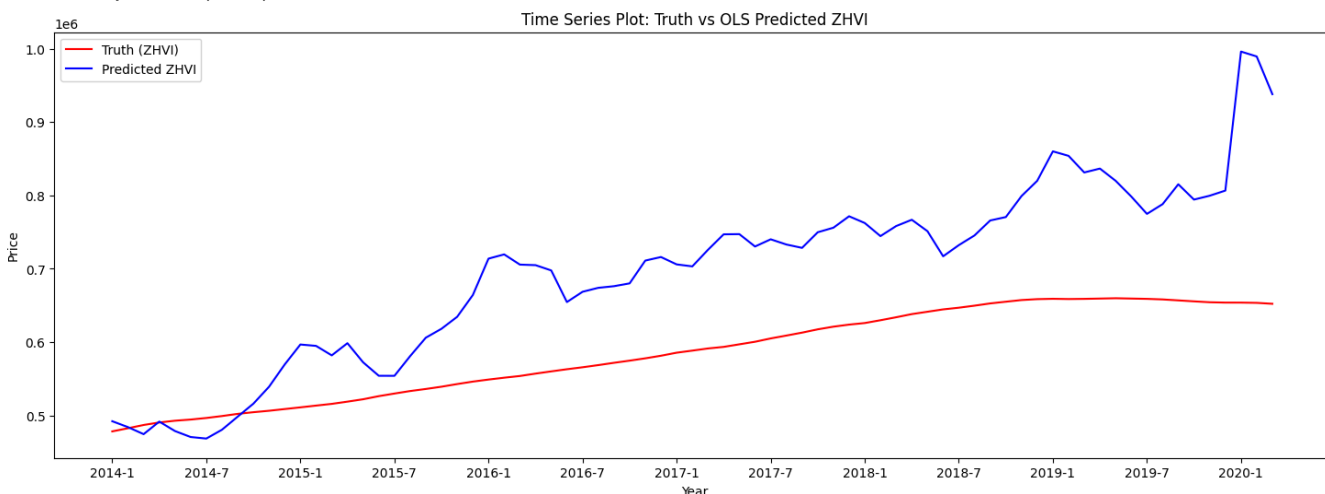
```

55 plt.title('Time Series Plot: Truth vs OLS Predicted ZHVI')
56 plt.legend(loc='upper left')
57 plt.show()

```



OLS Root Mean Squared Error (RMSE): 130532.09645357582  
 OLS Mean Absolute Percentage Error(MAPE): 0.1853761687014548  
 OLS Mean Absolute Error(MAE): 112483.27602848076  
 OLS R-squared( $R^2$ ): -0.14393633055157418



```

1 # Get feature names
2 feature_names = ['const'] + list(poly.get_feature_names_out(X_train.columns))
3
4 # Create dataframe to store coefficients and feature names
5 coefficients_df = pd.DataFrame({
6     'Feature': feature_names,
7     'Coefficient': results.params
8 })
9
10 # sort features by coeff magnitude
11 coefficients_df['Absolute_Coefficient'] = np.abs(coefficients_df['Coefficient'])
12 coefficients_df = coefficients_df.sort_values(by='Absolute_Coefficient', ascending=False)
13

```

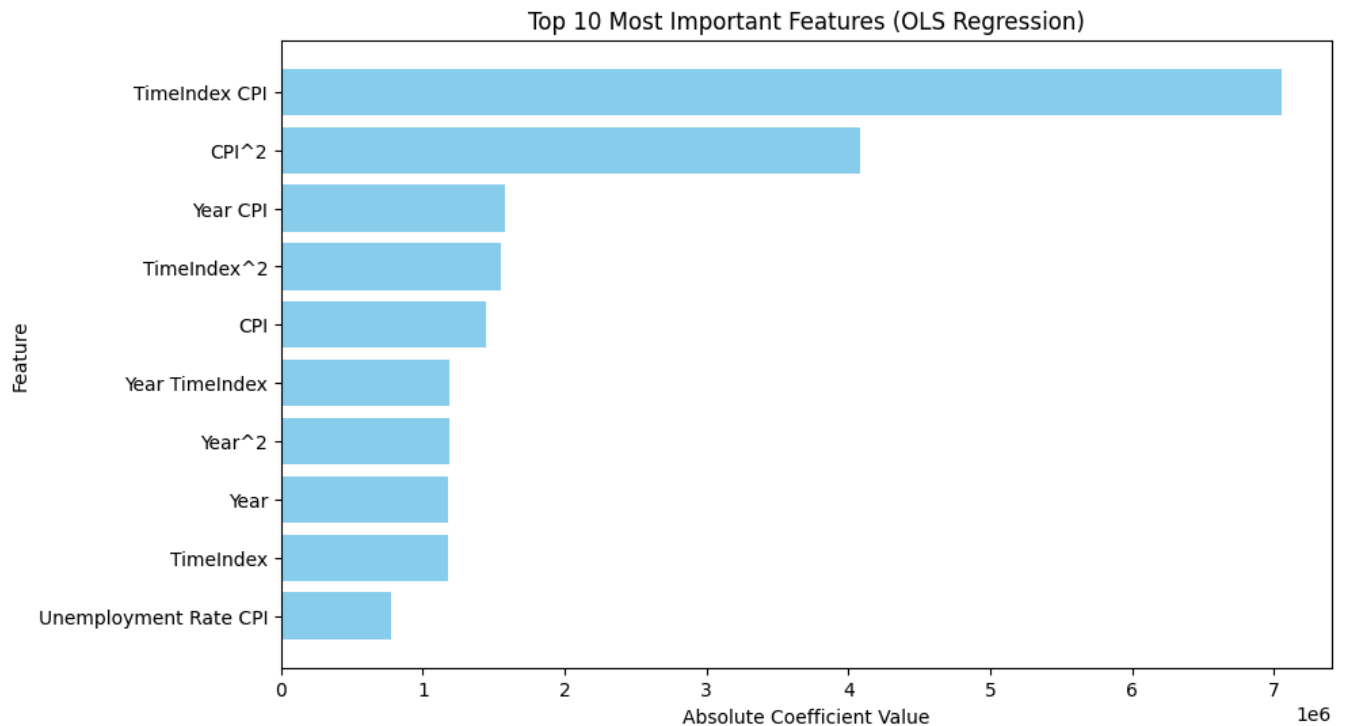


```
14 print("Sorted Top 10 OLS Regression Coefficients:")
15 print(coefficients_df[0:10])
16
17 # plot coeff
18 plt.figure(figsize=(10, 6))
19 plt.barh(coefficients_df['Feature'][:10], coefficients_df['Absolute_Coefficient'][:10],
20 plt.xlabel('Absolute Coefficient Value')
21 plt.ylabel('Feature')
22 plt.title('Top 10 Most Important Features (OLS Regression)')
23 plt.gca().invert_yaxis()
24 plt.show()
```



Sorted Top 10 OLS Regression Coefficients:

	Feature	Coefficient	Absolute_Coefficient
x24	TimeIndex CPI	-7.053306e+06	7.053306e+06
x31	CPI^2	4.083916e+06	4.083916e+06
x13	Year CPI	-1.580288e+06	1.580288e+06
x22	TimeIndex^2	1.550302e+06	1.550302e+06
x6	CPI	-1.441312e+06	1.441312e+06
x11	Year TimeIndex	1.185220e+06	1.185220e+06
x9	Year^2	1.183376e+06	1.183376e+06
x2	Year	1.181150e+06	1.181150e+06
x4	TimeIndex	1.181129e+06	1.181129e+06
x28	Unemployment Rate CPI	7.756565e+05	7.756565e+05



## Lasso Regression Model

```

1 # Split data into training and test
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month'
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate
7 y_train = train['ZHVI']

```

```
8
9 # Prepare test data for prediction
10 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI','Interest Rate'],
11
12 # add polynomial features and scale
13 scaler = StandardScaler()
14 poly = PolynomialFeatures(degree=2)
15 X_train_poly = poly.fit_transform(X_train)
16 X_test_poly = poly.transform(X_test)
17
18 X_train_scaled = scaler.fit_transform(X_train_poly)
19 X_test_scaled = scaler.transform(X_test_poly)
20
21 # Fit Lasso regression model
22 alphas = [0.001,0.005,0.01, 0.05, 0.1, 0.5, 1, 2, 3, 5, 10,25,50,75,100,150]
23 results = []
24 lowest_alpha = alphas[0]
25 lowest_mape = float('inf')
26
27 for alpha in alphas:
28     lasso_model = Lasso(alpha=alpha)
29     lasso_model.fit(X_train_scaled, y_train)
30
31 # Predict
32 predictions = lasso_model.predict(X_test_scaled)
33 test['Predicted_ZHVI'] = predictions
34
35 # Model evaluation
36 y_test = test['ZHVI']
37 y_pred = test['Predicted_ZHVI']
38 lasso_pred = test['Predicted_ZHVI']
39
40 mape = mean_absolute_percentage_error(y_test, y_pred)
41 results.append(mape)
42 if mape < lowest_mape:
43     lowest_mape = mape
44     lowest_alpha = alpha
45
46 print("Lowest MAPE:", lowest_mape)
47 print("Lowest Alpha:", lowest_alpha)
48
49 # Plot hyperparameter tuning
50 plt.plot(alphas, results, marker='o')
51 plt.xscale('log')
52
53 plt.xlabel('Alpha (log scale)')
54 plt.ylabel('MAPE')
55 plt.title('Lasso Regression Alpha Tuning - MAPE vs Alpha (Logarithmic X-axis)')
56
57 plt.show()
58
```



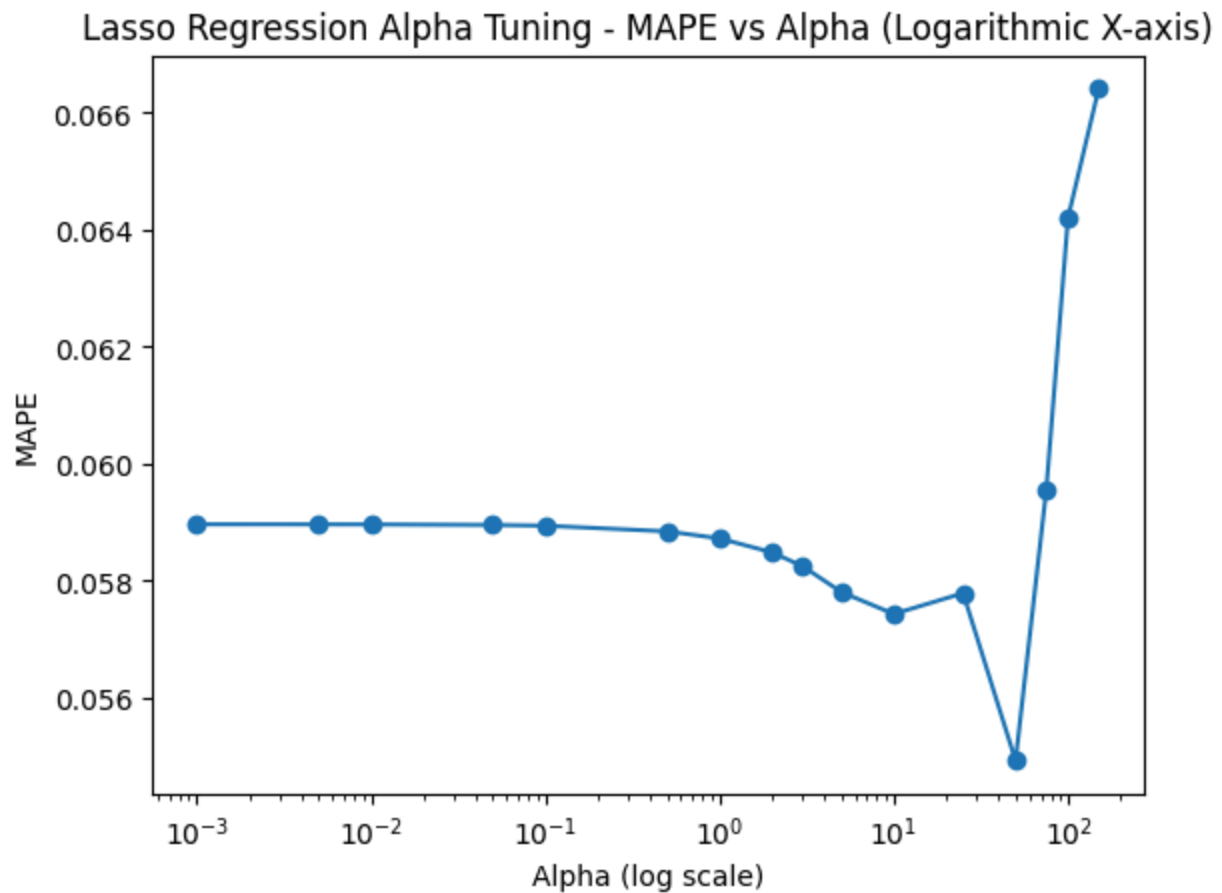
```
59 alpha = lowest_alpha
60
61 lasso_model = Lasso(alpha=alpha)
62 lasso_model.fit(X_train_scaled, y_train)
63
64 # Predict
65 predictions = lasso_model.predict(X_test_scaled)
66 test['Predicted_ZHVI'] = predictions
67
68 # Model evaluation
69 y_test = test['ZHVI']
70 y_pred = test['Predicted_ZHVI']
71 lasso_pred = test['Predicted_ZHVI']
72
73 lasso_rmse = math.sqrt(mean_squared_error(y_test, y_pred))
74 print(f"\n\nLasso Root Mean Squared Error (RMSE): {lasso_rmse}")
75 lasso_mape = mean_absolute_percentage_error(y_test, y_pred)
76 print("Lasso Mean Absolute Percentage Error (MAPE):", lasso_mape)
77 lasso_MAE = mean_absolute_error(y_test, y_pred)
78 print("Mean Absolute Error (MAE):", lasso_MAE)
79 lasso_r2 = r2_score(y_pred, y_test)
80 print(f"R-squared(R^2): {lasso_r2}")
81
82 # Plot truth vs prediction
83 plt.figure(figsize=(18, 6))
84 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
85 plt.plot(test['Year-Month'], test['Predicted_ZHVI'], color='blue', label='Predicted ZHVI')
86 plt.xlabel('Date')
87 plt.ylabel('ZHVI')
88 x_ticks = np.arange(0, 80, 6)
89 plt.xticks(x_ticks)
90 plt.title('Time Series Plot: Truth vs Lasso Regression Predicted ZHVI')
91 plt.legend(loc='upper left')
92 plt.show()
```





Lowest MAPE: 0.05492470799176251

Lowest Alpha: 50

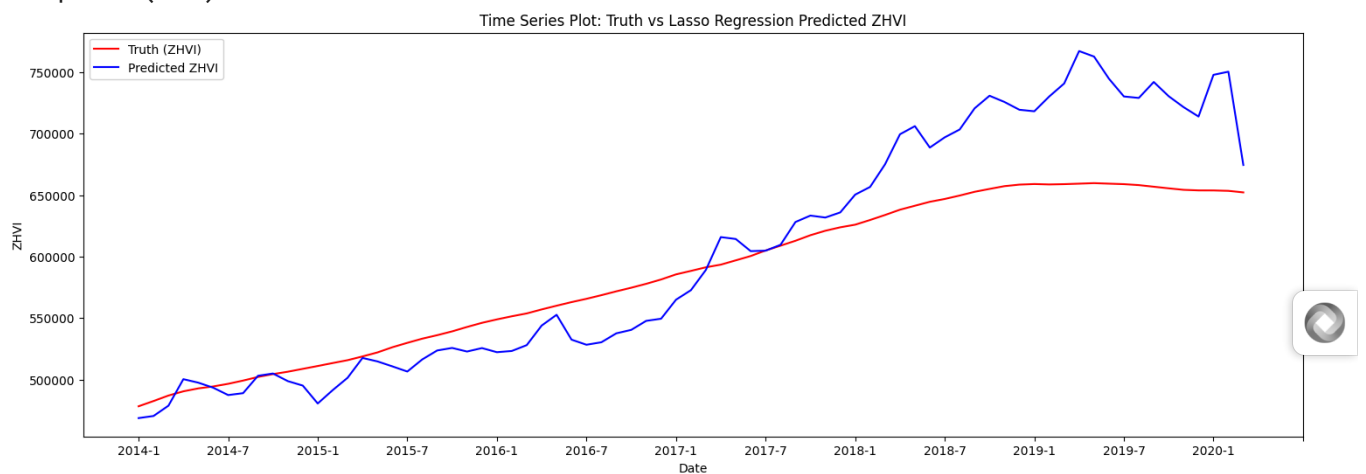


Lasso Root Mean Squared Error (RMSE): 44618.016149679155

Lasso Mean Absolute Percentage Error (MAPE): 0.05492470799176251

Mean Absolute Error (MAE): 34030.374729689334

R-squared( $R^2$ ): 0.7853309788737721



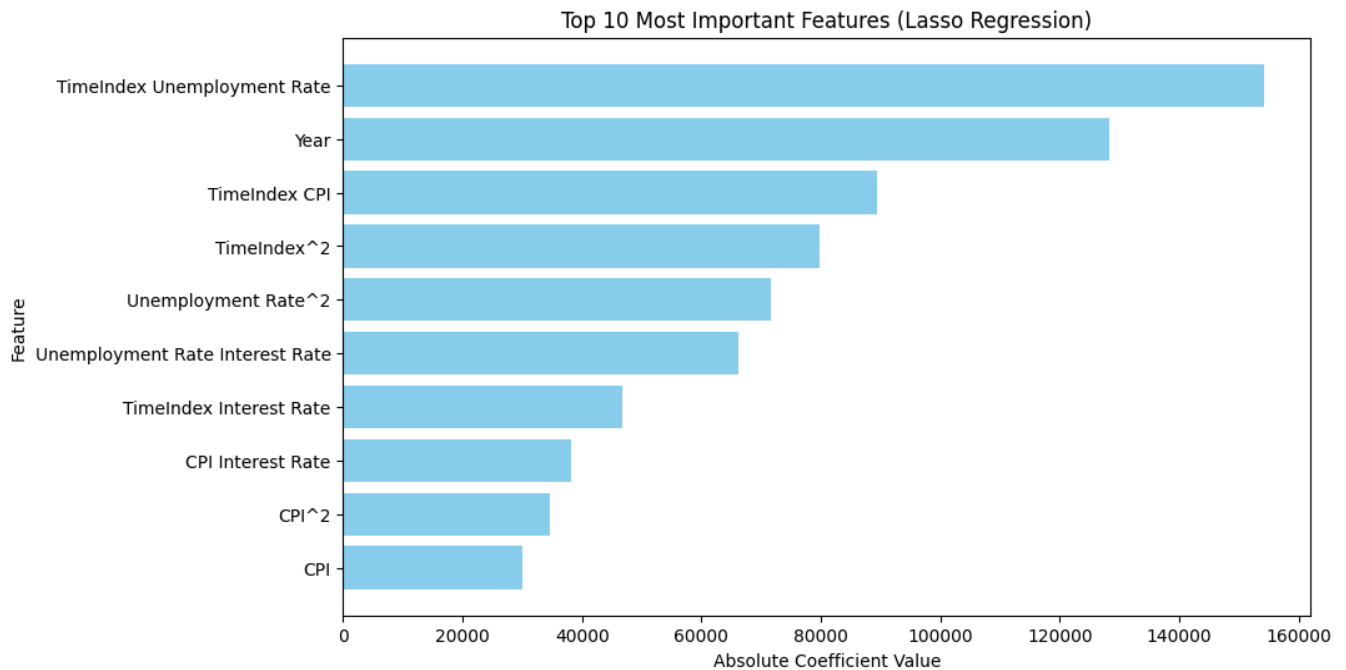
```
1 # Get feature names
2 feature_names = poly.get_feature_names_out(X_train.columns)
3
4 # Create dataframe to store coefficients and feature names
5 coefficients_df = pd.DataFrame({
6     'Feature': feature_names,
7     'Coefficient': lasso_model.coef_
8 })
9
10 # sort features by coeff magnitude
11 coefficients_df['Absolute_Coefficient'] = np.abs(coefficients_df['Coefficient'])
12 coefficients_df = coefficients_df.sort_values(by='Absolute_Coefficient', ascending=False)
13
14 print("Sorted Top 10 Lasso Regression Coefficients:")
15 print(coefficients_df[0:10])
16
17 # plot coeff
18 plt.figure(figsize=(10, 6))
19 plt.barh(coefficients_df['Feature'][0:10], coefficients_df['Absolute_Coefficient'][0:10],
20 plt.xlabel('Absolute Coefficient Value')
21 plt.ylabel('Feature')
22 plt.title('Top 10 Most Important Features (Lasso Regression)')
23 plt.gca().invert_yaxis()
24 plt.show()
```





Sorted Top 10 Lasso Regression Coefficients:

	Feature	Coefficient	Absolute_Coefficient
22	TimeIndex Unemployment Rate	-154296.036934	154296.036934
1	Year	128222.423511	128222.423511
23	TimeIndex CPI	89468.745570	89468.745570
21	TimeIndex^2	-79800.859453	79800.859453
26	Unemployment Rate^2	71651.217699	71651.217699
28	Unemployment Rate Interest Rate	66275.131194	66275.131194
24	TimeIndex Interest Rate	46722.278427	46722.278427
31	CPI Interest Rate	-38094.302067	38094.302067
30	CPI^2	34600.125584	34600.125584
5	CPI	29987.267681	29987.267681



## Ridge Regression Model

```

1 # Split data into training and test
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month'
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate
7 y_train = train['ZHVI']

```

```
8
9 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI','Interest Rate'],
10
11 # add polynomial features and scale
12 scaler = StandardScaler()
13 poly = PolynomialFeatures(degree=2)
14 X_train_poly = poly.fit_transform(X_train)
15 X_test_poly = poly.transform(X_test)
16
17 X_train_scaled = scaler.fit_transform(X_train_poly)
18 X_test_scaled = scaler.transform(X_test_poly)
19
20 # Fit Ridge regression model
21 alphas = [0.001,0.005,0.01, 0.05, 0.075, 0.1, 0.25, 0.35, 0.5, 1, 2, 3, 5, 10]
22
23 lowest_alpha = alphas[0]
24 lowest_mape = float('inf')
25 results = []
26 for alpha in alphas:
27     ridge_model = Ridge(alpha=alpha)
28     ridge_model.fit(X_train_scaled, y_train)
29
30     # Predict
31     predictions = ridge_model.predict(X_test_scaled)
32     test['Predicted_ZHVI'] = predictions
33
34     # Model evaluation
35     y_test = test['ZHVI']
36     y_pred = test['Predicted_ZHVI']
37     ridge_pred = test['Predicted_ZHVI']
38
39     mape = mean_absolute_percentage_error(y_test, y_pred)
40     results.append(mape)
41     if mape < lowest_mape:
42         lowest_mape = mape
43         lowest_alpha = alpha
44
45 print("Lowest MAPE:", lowest_mape)
46 print("Lowest Alpha:", lowest_alpha)
47
48 # Plot hyperparameter tuning
49 plt.plot(alphas, results, marker='o')
50 plt.xscale('log')
51
52 plt.xlabel('Alpha (log scale)')
53 plt.ylabel('MAPE')
54 plt.title('Ridge Regression Alpha Tuning - MAPE vs Alpha (Logarithmic X-axis)')
55
56 plt.show()
57
58 alpha = lowest_alpha
```



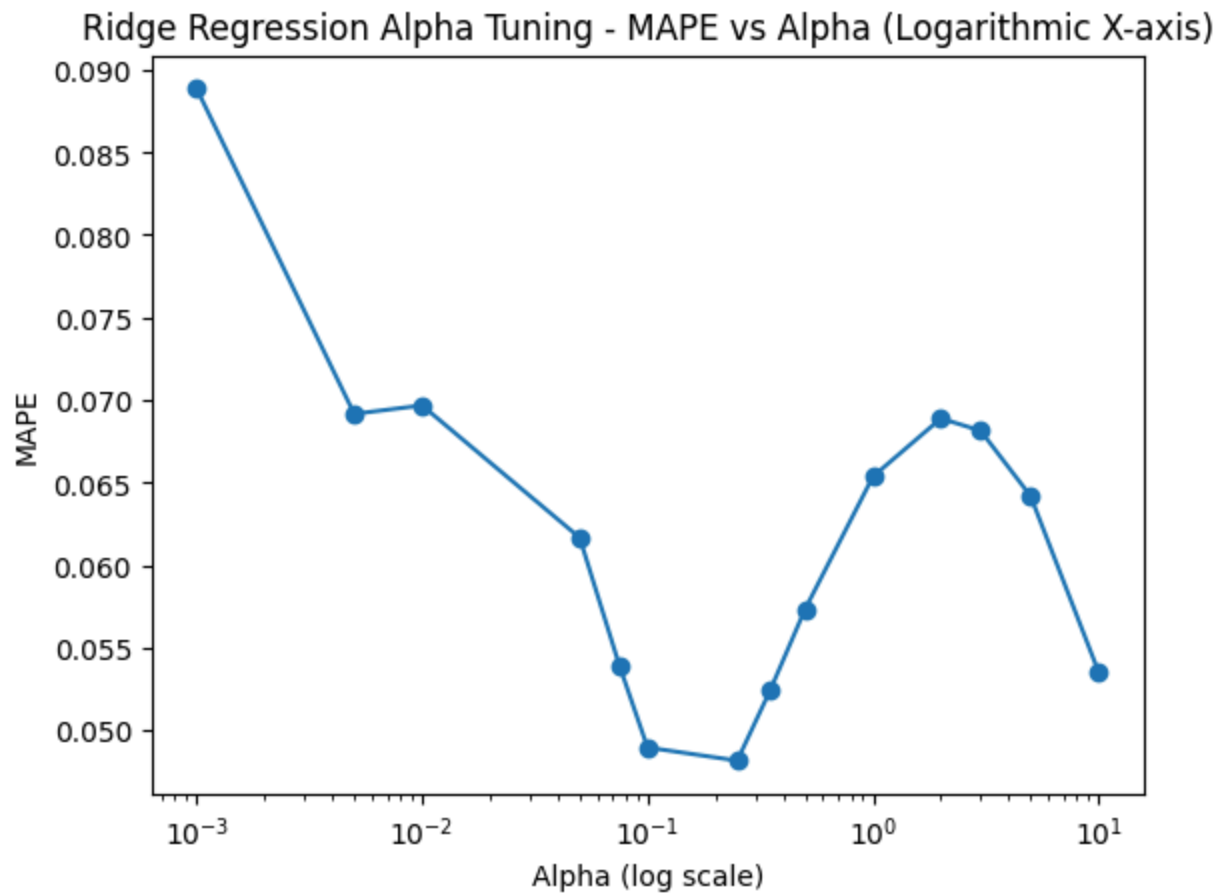
```
59 ridge_model = Ridge(alpha=alpha)
60 ridge_model.fit(X_train_scaled, y_train)
61
62 # Predict
63 predictions = ridge_model.predict(X_test_scaled)
64 test['Predicted_ZHVI'] = predictions
65
66 # Model evaluation
67 y_test = test['ZHVI']
68 y_pred = test['Predicted_ZHVI']
69 ridge_pred = test['Predicted_ZHVI']
70
71 ridge_rmse = math.sqrt(mean_squared_error(y_test, y_pred))
72 print(f"\n\nRR Root Mean Squared Error (RMSE): {ridge_rmse}")
73 ridge_mape = mean_absolute_percentage_error(y_test, y_pred)
74 print("RR Mean Absolute Percentage Error(MAPE):", ridge_mape)
75 ridge_MAE = mean_absolute_error(y_test, y_pred)
76 print("RR Mean Absolute Error(MAE):", ridge_MAE)
77 ridge_r2 = r2_score(y_pred, y_test)
78 print(f"R-squared(R^2): {ridge_r2}")
79
80 # Plot truth vs prediction
81 plt.figure(figsize=(18, 6))
82 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
83 plt.plot(test['Year-Month'], test['Predicted_ZHVI'], color='blue', label='Predicted ZHVI')
84 plt.xlabel('Year')
85 plt.ylabel('ZHVI')
86 x_ticks = np.arange(0, 80, 6)
87 plt.xticks(x_ticks)
88 plt.title('Time Series Plot: Truth vs Ridge Regression Predicted ZHVI')
89 plt.legend(loc='upper left')
90 plt.show()
```





Lowest MAPE: 0.04815805234530082

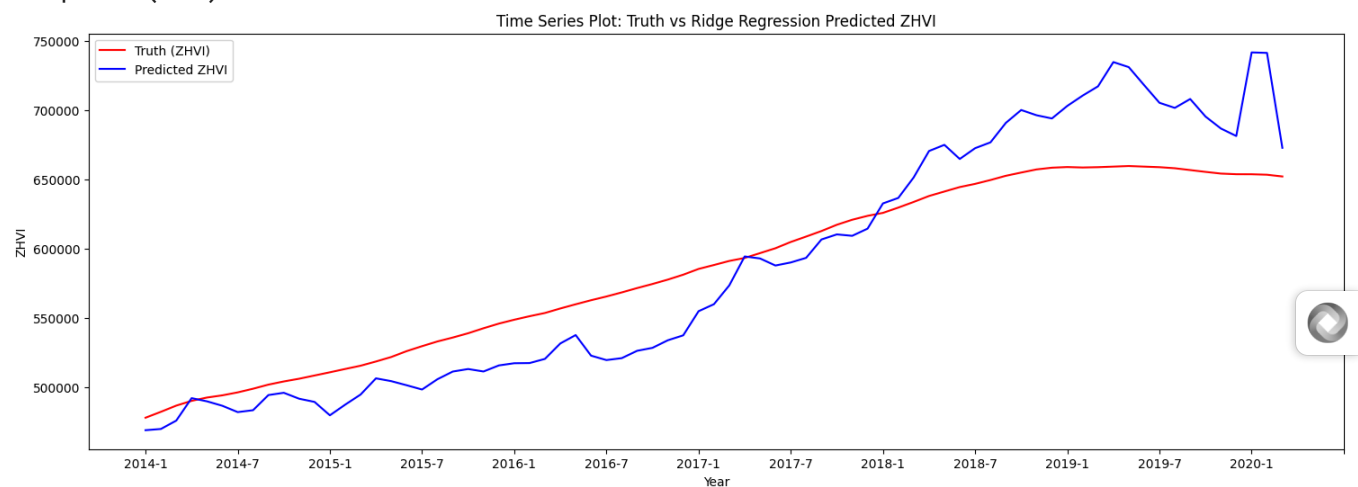
Lowest Alpha: 0.25



RR Root Mean Squared Error (RMSE): 34586.997710405834

RR Mean Absolute Percentage Error(MAPE): 0.04815805234530082

RR Mean Absolute Error(MAE): 28891.172716798137

R-squared(R<sup>2</sup>): 0.846871147352243

```
1 # Get feature names
2 feature_names = poly.get_feature_names_out(X_train.columns)
3
4 # Create dataframe to store coefficients and feature names
5 coefficients_df = pd.DataFrame({
6     'Feature': feature_names,
7     'Coefficient': ridge_model.coef_
8 })
9
10 # sort features by coeff magnitude
11 coefficients_df['Absolute_Coefficient'] = np.abs(coefficients_df['Coefficient'])
12 coefficients_df = coefficients_df.sort_values(by='Absolute_Coefficient', ascending=False)
13
14 print("Sorted Top 10 Ridge Regression Coefficients:")
15 print(coefficients_df[0:10])
16
17 # plot coeff
18 plt.figure(figsize=(10, 6))
19 plt.barh(coefficients_df['Feature'][0:10], coefficients_df['Absolute_Coefficient'][0:10], c
20 plt.xlabel('Absolute Coefficient Value')
21 plt.ylabel('Feature')
22 plt.title('Top 10 Most Important Features (Ridge Regression)')
23 plt.gca().invert_yaxis()
24 plt.show()
```

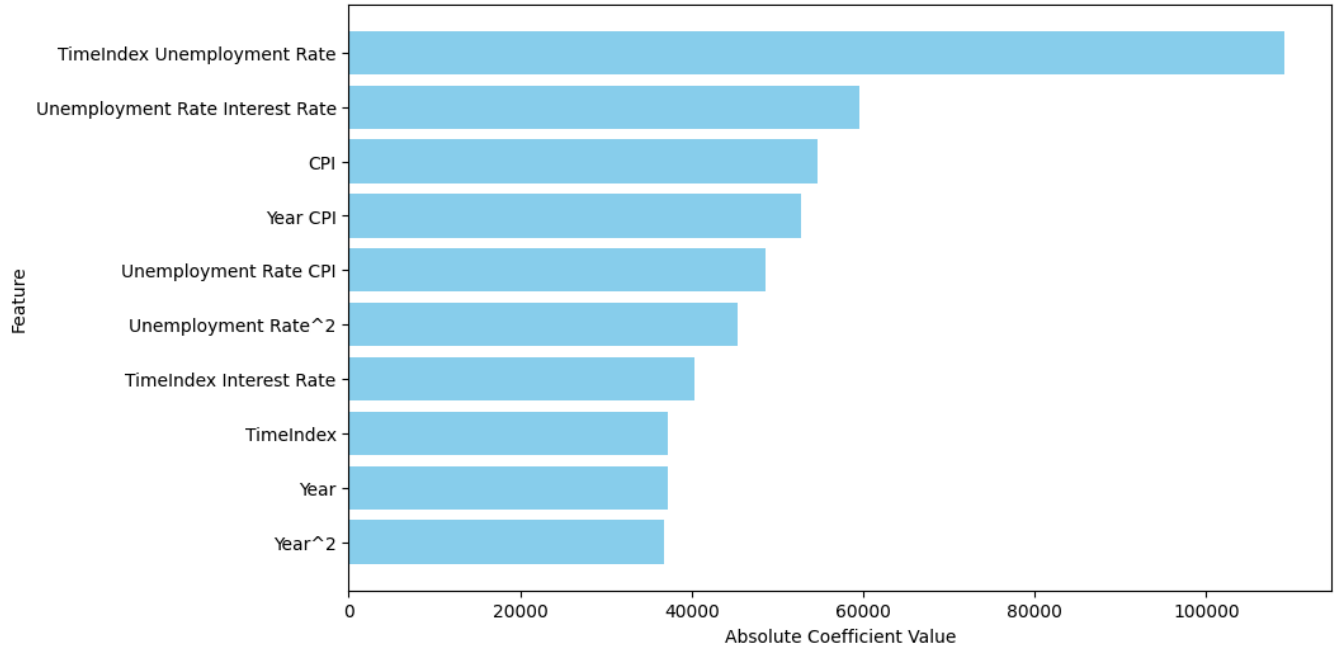




Sorted Top 10 Ridge Regression Coefficients:

	Feature	Coefficient	Absolute_Coefficient
22	TimeIndex Unemployment Rate	-109178.576829	109178.576829
28	Unemployment Rate Interest Rate	59548.087563	59548.087563
5	CPI	54636.320560	54636.320560
12	Year CPI	52733.852309	52733.852309
27	Unemployment Rate CPI	-48590.347941	48590.347941
26	Unemployment Rate^2	45311.048503	45311.048503
24	TimeIndex Interest Rate	40266.286996	40266.286996
3	TimeIndex	37172.836936	37172.836936
1	Year	37166.437453	37166.437453
8	Year^2	36844.519700	36844.519700

Top 10 Most Important Features (Ridge Regression)



```

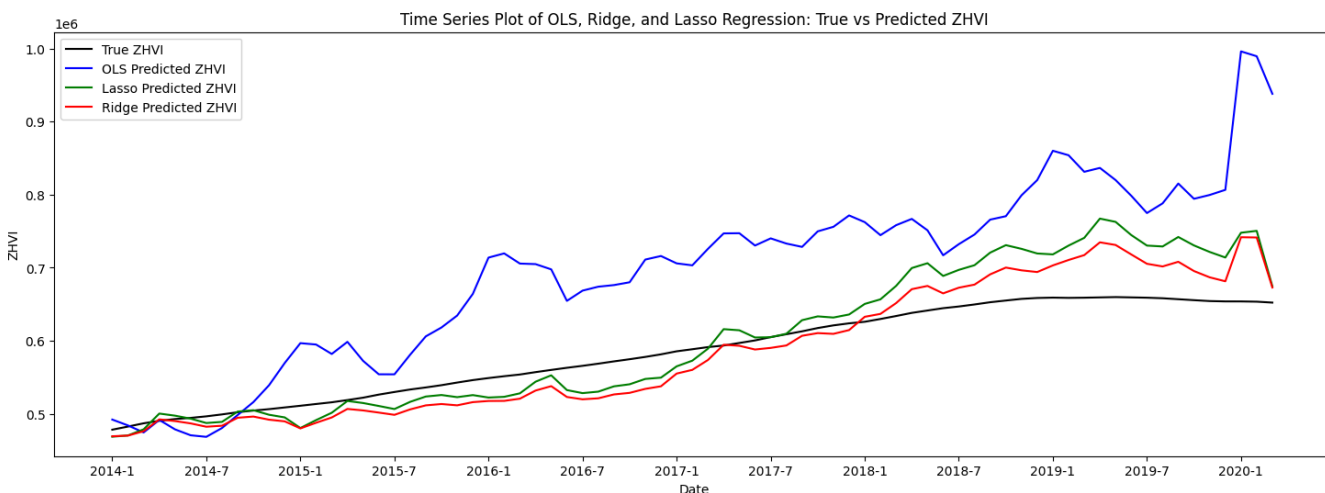
1 plt.figure(figsize=(18, 6))
2 plt.plot(test['Year-Month'], test['ZHVI'], color='black', label='True ZHVI')
3 plt.plot(test['Year-Month'], OLS_pred, color='blue', label='OLS Predicted ZHVI')
4 plt.plot(test['Year-Month'], lasso_pred, color='green', label='Lasso Predicted ZHVI')
5 plt.plot(test['Year-Month'], ridge_pred, color='red', label='Ridge Predicted ZHVI')
6
7 plt.xlabel('Date')
8 plt.ylabel('ZHVI')
9 x_ticks = np.arange(0, 80, 6)
10 plt.xticks(x_ticks)

```

```

11 plt.title('Time Series Plot of OLS, Ridge, and Lasso Regression: True vs Predicted ZHVI')
12
13 plt.legend(loc='upper left')
14
15 plt.show()

```



## milestone 2

```

1 !pip install keras-tuner
2 !pip install -q streamlit
3 !npm install localtunnel
4
5 from tensorflow.keras.models import Sequential
6 from tensorflow.keras.layers import Dense, BatchNormalization, Dropout, Activation
7 from tensorflow.keras.optimizers import Adam
8 from tensorflow.keras.callbacks import EarlyStopping
9 from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error, mean_abs
10 from sklearn.ensemble import RandomForestRegressor
11 from sklearn.model_selection import GridSearchCV
12 from sklearn.preprocessing import StandardScaler, PolynomialFeatures

```



```
13 from sklearn.model_selection import RandomizedSearchCV, TimeSeriesSplit
14 from sklearn.linear_model import LinearRegression
15 from tabulate import tabulate
16
17 import tensorflow as tf
18 import keras_tuner as kt
19 import math
20 import matplotlib.pyplot as plt
21 import numpy as np
22 import random
23 import xgboost as xgb
24 import statsmodels.api as sm
25 import joblib
26 import xgboost as xgb
27
28 clear_output()
```

## DNN Model

```
1 # Hyperparameter tuning
2
3 random.seed(158)
4
5 # Preprocess data
6 scaler = StandardScaler()
7 X_train_scaled = scaler.fit_transform(X_train)
8 X_test_scaled = scaler.transform(X_test)
9
10 # Define the model building function for Keras Tuner
11 def build_model(hp):
12
13     model = Sequential()
14
15     # Tune the number of units in the first Dense layer
16     model.add(Dense(
17         units=hp.Int('units_1', min_value=64, max_value=256, step=64),
18         activation='relu',
19         input_shape=(X_train_scaled.shape[1],)
20     ))
21     model.add(Dropout(
22         rate=hp.Float('dropout_1', min_value=0.0, max_value=0.5, step=0.1)
23     ))
24
25     # Tune the number of hidden layers (1-3)
26     for i in range(hp.Int('num_layers', 1, 3)):
27         model.add(Dense(
28             units=hp.Int(f'units_{i+2}', min_value=32, max_value=128, step=32),
29             activation='relu')
30         )
31     model.add(Dropout(
```





```
32         rate=hp.Float(f'dropout_{i+2}', min_value=0.0, max_value=0.5, step=0.1)
33     ))
34
35     model.add(Dense(1))
36
37     # Tune the learning rate
38     learning_rate = hp.Choice('learning_rate', values=[1e-2, 1e-3, 1e-4])
39
40     optimizer = Adam(learning_rate=learning_rate)
41     model.compile(optimizer=optimizer, loss='mse', metrics=['mae'])
42
43     return model
44
45 # Set up the tuner
46 tuner = kt.RandomSearch(
47     build_model,
48     objective='val_loss',
49     max_trials=20,
50     executions_per_trial=2,
51     directory='keras_tuner',
52     project_name='zhvi_prediction'
53 )
54
55 # Early stopping callback
56 early_stopping = EarlyStopping(
57     monitor='val_loss',
58     patience=10,
59     restore_best_weights=True
60 )
61
62 # Perform hyperparameter search
63 tuner.search(
64     X_train_scaled, y_train,
65     epochs=100,
66     batch_size=32,
67     validation_split=0.2,
68     callbacks=[early_stopping],
69     verbose=1
70 )
71
72 # Get the best hyperparameters
73 best_hps = tuner.get_best_hyperparameters(num_trials=1)[0]
74
75 print(f"""
76 The hyperparameter search is complete. The optimal hyperparameters are:
77 - Units in first layer: {best_hps.get('units_1')}
78 - Number of hidden layers: {best_hps.get('num_layers')}
79 - Learning rate: {best_hps.get('learning_rate')}
80 - Dropout rate (first layer): {best_hps.get('dropout_1')}
81 """)
82
```



```
83 # Build the model with the best hyperparameters
84 best_model = tuner.hypermodel.build(best_hps)
85
86 # Train the best model
87 history = best_model.fit(
88     X_train_scaled, y_train,
89     epochs=100,
90     batch_size=32,
91     validation_split=0.2,
92     callbacks=[early_stopping],
93     verbose=1
94 )
95
96 # Plot training history
97 plt.figure(figsize=(12, 6))
98 plt.plot(history.history['loss'], label='Training Loss')
99 plt.plot(history.history['val_loss'], label='Validation Loss')
100 plt.xlabel('Epoch')
101 plt.ylabel('Loss (MSE)')
102 plt.title('Neural Network Training History')
103 plt.legend()
104 plt.show()
105
106 # Predictions
107 nn_predictions = best_model.predict(X_test_scaled).flatten()
108 test['NN_Predicted_ZHVI'] = nn_predictions
109
110 # Model evaluation
111 nn_rmse = math.sqrt(mean_squared_error(y_test, nn_predictions))
112 print(f"Neural Network Root Mean Squared Error (RMSE): {nn_rmse}")
113 nn_mape = mean_absolute_percentage_error(y_test, nn_predictions)
114 print("Neural Network Mean Absolute Percentage Error (MAPE):", nn_mape)
115 nn_mae = mean_absolute_error(y_test, nn_predictions)
116 print("Neural Network Mean Absolute Error (MAE):", nn_mae)
117 nn_r2 = r2_score(y_test, nn_predictions)
118 print(f"Neural Network R-squared (R^2): {nn_r2}")
119
120 # Plot truth vs prediction
121 plt.figure(figsize=(18, 6))
122 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
123 plt.plot(test['Year-Month'], test['NN_Predicted_ZHVI'], color='purple', label='Neural Network Predicted ZHVI (Optimized)')
124 plt.xlabel('Date')
125 plt.ylabel('ZHVI')
126 x_ticks = np.arange(0, 80, 6)
127 plt.xticks(x_ticks)
128 plt.title('Time Series Plot: Truth vs Neural Network Predicted ZHVI (Optimized)')
129 plt.legend(loc='upper left')
130 plt.show()
```



➡ Trial 20 Complete [00h 00m 14s]  
val\_loss: 3773630720.0

Best val\_loss So Far: 2747978496.0  
Total elapsed time: 00h 11m 01s

The hyperparameter search is complete. The optimal hyperparameters are:

- Units in first layer: 192
- Number of hidden layers: 3
- Learning rate: 0.001
- Dropout rate (first layer): 0.4

Epoch 1/100

6/6 ————— 2s 58ms/step - loss: 142135230464.0000 - mae: 357076.5938 - val

Epoch 2/100

6/6 ————— 0s 17ms/step - loss: 137808904192.0000 - mae: 351127.2188 - val

Epoch 3/100

6/6 ————— 0s 17ms/step - loss: 137643327488.0000 - mae: 350684.6562 - val

Epoch 4/100

6/6 ————— 0s 17ms/step - loss: 143591079936.0000 - mae: 359643.5312 - val

Epoch 5/100

6/6 ————— 0s 17ms/step - loss: 132160069632.0000 - mae: 343093.5000 - val

Epoch 6/100

6/6 ————— 0s 17ms/step - loss: 137535651840.0000 - mae: 351015.3125 - val

Epoch 7/100

6/6 ————— 0s 20ms/step - loss: 136340348928.0000 - mae: 349860.6562 - val

Epoch 8/100

6/6 ————— 0s 18ms/step - loss: 136219549696.0000 - mae: 350163.6875 - val

Epoch 9/100

6/6 ————— 0s 25ms/step - loss: 134784786432.0000 - mae: 346724.8438 - val

Epoch 10/100

6/6 ————— 0s 17ms/step - loss: 140556959744.0000 - mae: 354887.2188 - val

Epoch 11/100

6/6 ————— 0s 18ms/step - loss: 137576038400.0000 - mae: 350448.6250 - val

Epoch 12/100

6/6 ————— 0s 18ms/step - loss: 137234776064.0000 - mae: 350552.9375 - val

Epoch 13/100

6/6 ————— 0s 18ms/step - loss: 139266768896.0000 - mae: 353964.1875 - val

Epoch 14/100

6/6 ————— 0s 18ms/step - loss: 131869065216.0000 - mae: 341946.0938 - val

Epoch 15/100

6/6 ————— 0s 21ms/step - loss: 137546448896.0000 - mae: 351033.6250 - val

Epoch 16/100

6/6 ————— 0s 18ms/step - loss: 131390136320.0000 - mae: 342966.9375 - val

Epoch 17/100

6/6 ————— 0s 17ms/step - loss: 142992736256.0000 - mae: 359089.5000 - val

Epoch 18/100

6/6 ————— 0s 19ms/step - loss: 135150436352.0000 - mae: 345449.7500 - val

Epoch 19/100

6/6 ————— 0s 18ms/step - loss: 128360128512.0000 - mae: 337479.7500 - val

Epoch 20/100

6/6 ————— 0s 18ms/step - loss: 132962910208.0000 - mae: 343252.2188 - val

Epoch 21/100

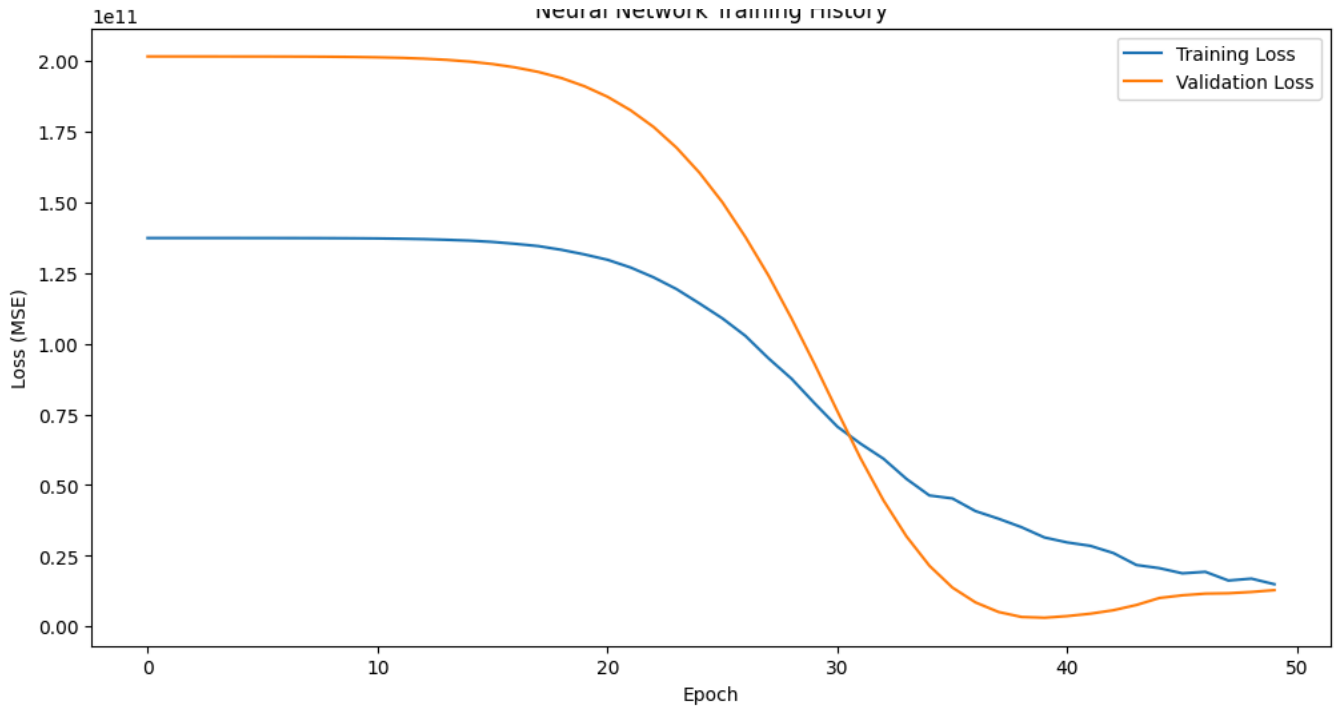
6/6 ————— 0s 25ms/step - loss: 135287021568.0000 - mae: 346663.2500 - val

Epoch 22/100

6/6 ————— 0s 20ms/step - loss: 121984532480.0000 - mae: 327748.7188 - val

```
Epoch 23/100
6/6 ————— 0s 18ms/step - loss: 116732141568.0000 - mae: 316579.8750 - val_
Epoch 24/100
6/6 ————— 0s 18ms/step - loss: 112045711360.0000 - mae: 310469.1250 - val_
Epoch 25/100
6/6 ————— 0s 18ms/step - loss: 113209040896.0000 - mae: 310353.0938 - val_
Epoch 26/100
6/6 ————— 0s 17ms/step - loss: 110794670080.0000 - mae: 303401.6250 - val_
Epoch 27/100
6/6 ————— 0s 18ms/step - loss: 101304696832.0000 - mae: 285922.2188 - val_
Epoch 28/100
6/6 ————— 0s 25ms/step - loss: 91559772160.0000 - mae: 265344.7500 - val_
Epoch 29/100
6/6 ————— 0s 18ms/step - loss: 84936499200.0000 - mae: 248896.7031 - val_
Epoch 30/100
6/6 ————— 0s 17ms/step - loss: 76313583616.0000 - mae: 228155.7656 - val_
Epoch 31/100
6/6 ————— 0s 25ms/step - loss: 69197414400.0000 - mae: 220884.8906 - val_
Epoch 32/100
6/6 ————— 0s 17ms/step - loss: 71095009280.0000 - mae: 228538.0312 - val_
Epoch 33/100
6/6 ————— 0s 17ms/step - loss: 62709309440.0000 - mae: 215566.1250 - val_
Epoch 34/100
6/6 ————— 0s 17ms/step - loss: 50817744896.0000 - mae: 194482.9688 - val_
Epoch 35/100
6/6 ————— 0s 21ms/step - loss: 47198117888.0000 - mae: 192297.9062 - val_
Epoch 36/100
6/6 ————— 0s 18ms/step - loss: 45227626496.0000 - mae: 191217.2031 - val_
Epoch 37/100
6/6 ————— 0s 18ms/step - loss: 40039145472.0000 - mae: 180178.0469 - val_
Epoch 38/100
6/6 ————— 0s 25ms/step - loss: 38438428672.0000 - mae: 176614.7188 - val_
Epoch 39/100
6/6 ————— 0s 17ms/step - loss: 35704451072.0000 - mae: 169358.5781 - val_
Epoch 40/100
6/6 ————— 0s 30ms/step - loss: 33121267712.0000 - mae: 165117.3594 - val_
Epoch 41/100
6/6 ————— 0s 25ms/step - loss: 31072069632.0000 - mae: 159930.7500 - val_
Epoch 42/100
6/6 ————— 0s 31ms/step - loss: 29801713664.0000 - mae: 152902.3125 - val_
Epoch 43/100
6/6 ————— 0s 23ms/step - loss: 26341781504.0000 - mae: 145040.1875 - val_
Epoch 44/100
6/6 ————— 0s 28ms/step - loss: 21740562432.0000 - mae: 130391.8125 - val_
Epoch 45/100
6/6 ————— 0s 32ms/step - loss: 22103099392.0000 - mae: 129428.6250 - val_
Epoch 46/100
6/6 ————— 0s 34ms/step - loss: 17934903296.0000 - mae: 117617.8984 - val_
Epoch 47/100
6/6 ————— 0s 31ms/step - loss: 18871126016.0000 - mae: 119005.3906 - val_
Epoch 48/100
6/6 ————— 0s 31ms/step - loss: 15961183232.0000 - mae: 107405.6875 - val_
Epoch 49/100
6/6 ————— 0s 19ms/step - loss: 16721813504.0000 - mae: 110311.5703 - val_
Epoch 50/100
6/6 ————— 0s 18ms/step - loss: 14423368704.0000 - mae: 101693.5312 - val_
```

Manual Network Training History



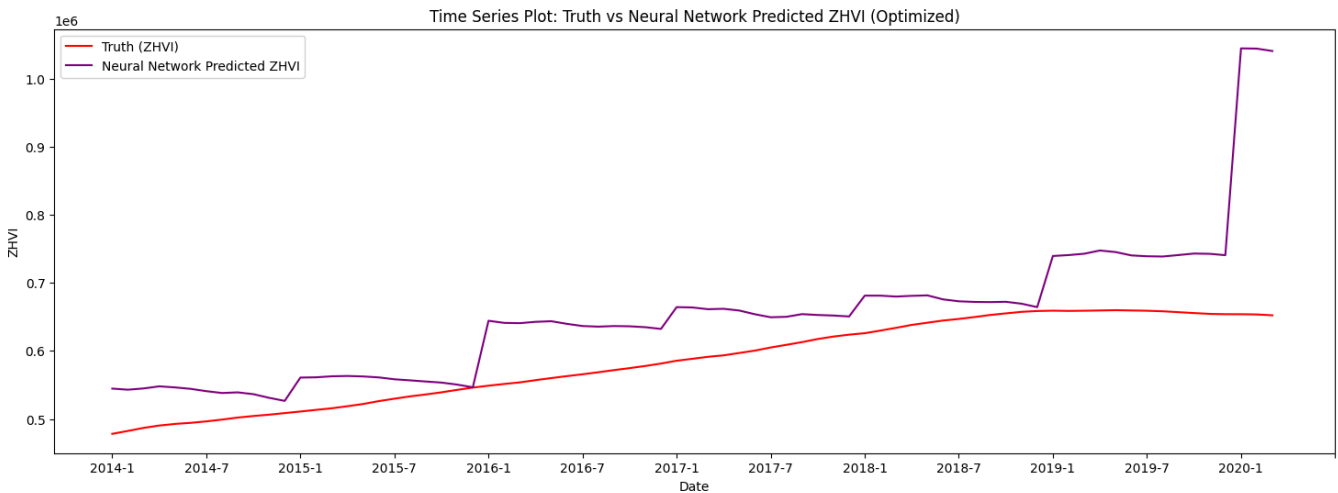
3/3 ———— 0s 37ms/step

Neural Network Root Mean Squared Error (RMSE): 96556.13473193542

Neural Network Mean Absolute Percentage Error (MAPE): 0.11076739748861884

Neural Network Mean Absolute Error (MAE): 66117.98833333333

Neural Network R-squared ( $R^2$ ): -1.5689345558410244





```
1 # Split data into training and test
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month'
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI','Interest Rate
7 y_train = train['ZHVI']
8
9 # Prepare test data for prediction
10 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI','Interest Rate',
11
12 # set seed for easy reproducibility of different configurations
13 random.seed(158)
14
15 # preprocess data
16 scaler = StandardScaler()
17 X_train_scaled = scaler.fit_transform(X_train)
18 X_test_scaled = scaler.transform(X_test)
19
20 # Build neural network
21 model = Sequential([
22     Dense(128, activation='relu', input_shape=(X_train_scaled.shape[1],)),
23     Dropout(0.2),
24     Dense(64, activation='relu'),
25     Dropout(0.2),
26     Dense(32, activation='relu'),
27     Dropout(0.2),
28     Dense(1)
29 ])
30
31 optimizer = Adam(learning_rate=0.001)
32 model.compile(optimizer=optimizer, loss='mse', metrics=['mae'])
33
34 # Train the model
35 history = model.fit(
36     X_train_scaled, y_train,
37     epochs=100,
38     batch_size=32,
39     validation_split=0.2,
40     verbose=1
41 )
42
43 clear_output()
44
45 # Predictions
46 nn_predictions = model.predict(X_test_scaled).flatten()
47 test['NN_Predicted_ZHVI'] = nn_predictions
48
49 # Model evaluation
50 nn_rmse = math.sqrt(mean_squared_error(y_test, nn_predictions))
51 print(f"Neural Network Root Mean Squared Error (RMSE): {nn_rmse}")
```



```
52 nn_mape = mean_absolute_percentage_error(y_test, nn_predictions)
53 print("Neural Network Mean Absolute Percentage Error (MAPE):", nn_mape)
54 nn_mae = mean_absolute_error(y_test, nn_predictions)
55 print("Neural Network Mean Absolute Error (MAE):", nn_mae)
56 nn_r2 = r2_score(y_test, nn_predictions)
57 print(f"Neural Network R-squared (R^2): {nn_r2}")
58
59 # Plot truth vs prediction
60 plt.figure(figsize=(18, 6))
61 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
62 plt.plot(test['Year-Month'], test['NN_Predicted_ZHVI'], color='purple', label='Neural Ne
63 plt.xlabel('Date')
64 plt.ylabel('ZHVI')
65 x_ticks = np.arange(0, 80, 6)
66 plt.xticks(x_ticks)
67 plt.title('Time Series Plot: Truth vs Neural Network Predicted ZHVI')
68 plt.legend(loc='upper left')
69 plt.show()
```







3/3 0s 43ms/step

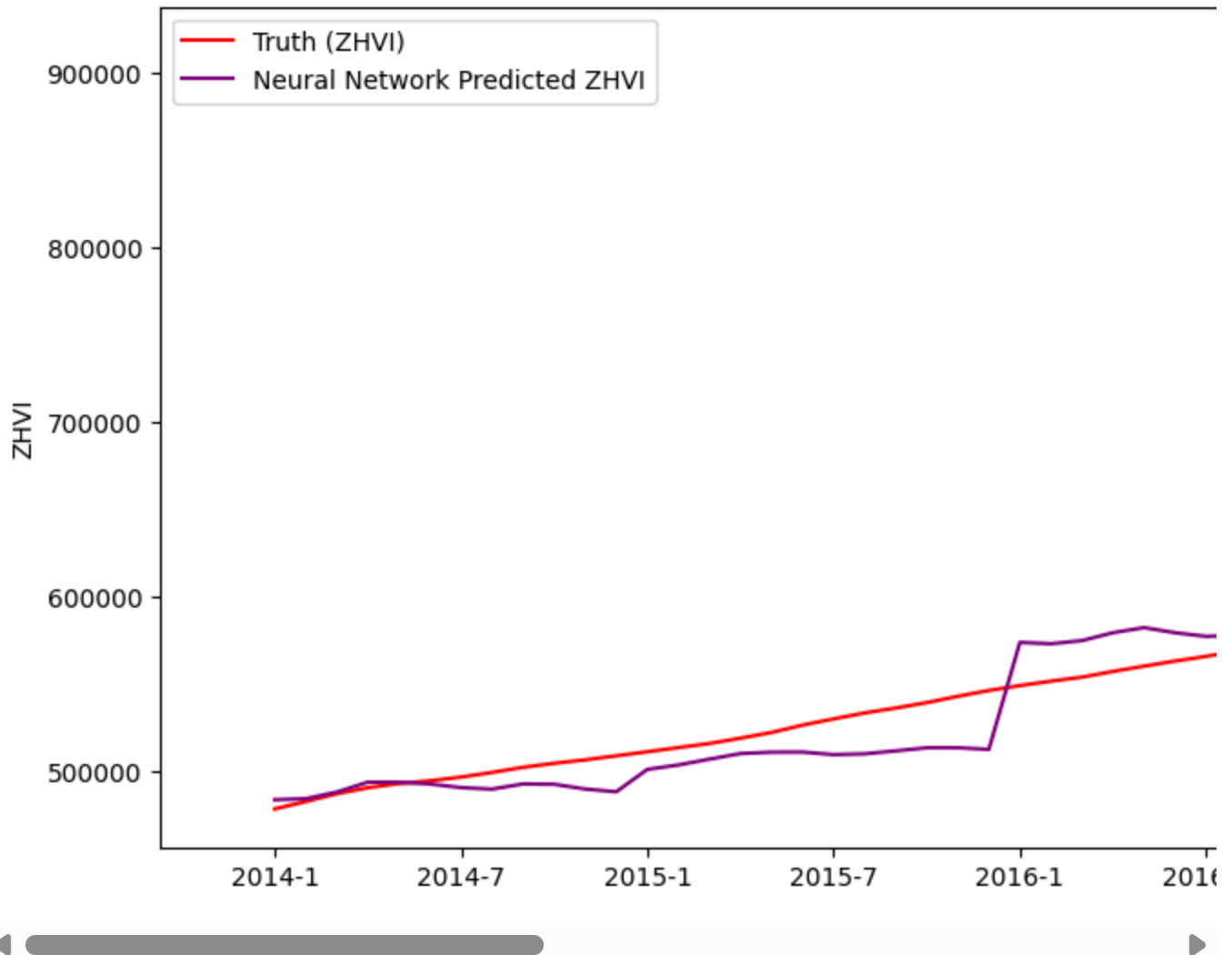
Neural Network Root Mean Squared Error (RMSE): 55288.09761836565

Neural Network Mean Absolute Percentage Error (MAPE): 0.041211447799650795

Neural Network Mean Absolute Error (MAE): 25382.914583333335

Neural Network R-squared ( $R^2$ ): 0.15771980293931132

Time Series Pl



## Random Forest Model

```

1 # Split data into training and test (same as before)
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'] < 12))]
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month'] < 12))]
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate']]
7 y_train = train['ZHVI']
8 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate']]
9 y_test = test['ZHVI']
10
11 # add polynomial features and scale
12 scaler = StandardScaler()
13 poly = PolynomialFeatures(degree=2)

```

```
14 X_train_poly = poly.fit_transform(X_train)
15 X_test_poly = poly.transform(X_test)
16
17 # Get feature names from polynomial features
18 feature_names = poly.get_feature_names_out(input_features=X_train.columns)
19
20 X_train_scaled = scaler.fit_transform(X_train_poly)
21 X_test_scaled = scaler.transform(X_test_poly)
22
23 # Hyperparameter tuning with GridSearchCV
24 param_grid = {
25     'n_estimators': [100, 200, 300],
26     'max_depth': [None, 10, 20, 30],
27     'min_samples_split': [2, 5, 10],
28     'min_samples_leaf': [1, 2, 4]
29 }
30
31 rf = RandomForestRegressor(random_state=42)
32 grid_search = GridSearchCV(estimator=rf, param_grid=param_grid,
33                             cv=5, n_jobs=-1, verbose=2, scoring='neg_mean_squared_error')
34 grid_search.fit(X_train_scaled, y_train)
35
36 # Best model
37 best_rf = grid_search.best_estimator_
38 print(f"Best Random Forest parameters: {grid_search.best_params_}")
39
40 # Predictions
41 rf_predictions = best_rf.predict(X_test_scaled)
42 test['RF_Predicted_ZHVI'] = rf_predictions
43
44 # Model evaluation
45 rf_rmse = math.sqrt(mean_squared_error(y_test, rf_predictions))
46 print(f"Random Forest Root Mean Squared Error (RMSE): {rf_rmse}")
47 rf_mape = mean_absolute_percentage_error(y_test, rf_predictions)
48 print("Random Forest Mean Absolute Percentage Error (MAPE):", rf_mape)
49 rf_mae = mean_absolute_error(y_test, rf_predictions)
50 print("Random Forest Mean Absolute Error (MAE):", rf_mae)
51 rf_r2 = r2_score(y_test, rf_predictions)
52 print(f"Random Forest R-squared (R^2): {rf_r2}")
53
54 # Feature importance
55 feature_importance = pd.DataFrame({
56     'Feature': feature_names,
57     'Importance': best_rf.feature_importances_
58 }).sort_values('Importance', ascending=False)
59
60 print("\nRandom Forest Feature Importance:")
61 print(feature_importance)
62
63 # Plot feature importance
64 plt.figure(figsize=(10, 6))
```



```
65 plt.barh(feature_importance['Feature'], feature_importance['Importance'], color='skyblue')
66 plt.xlabel('Importance Score')
67 plt.ylabel('Feature')
68 plt.title('Random Forest Feature Importance')
69 plt.gca().invert_yaxis()
70 plt.show()
71
72 # Plot truth vs prediction
73 plt.figure(figsize=(18, 6))
74 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
75 plt.plot(test['Year-Month'], test['RF_Predicted_ZHVI'], color='green', label='Random For
76 plt.xlabel('Date')
77 plt.ylabel('ZHVI')
78 x_ticks = np.arange(0, 80, 6)
79 plt.xticks(x_ticks)
80 plt.title('Time Series Plot: Truth vs Random Forest Predicted ZHVI')
81 plt.legend(loc='upper left')
82 plt.show()
83
```





Fitting 5 folds for each of 108 candidates, totalling 540 fits

Best Random Forest parameters: {'max\_depth': None, 'min\_samples\_leaf': 1, 'min\_samples\_s

Random Forest Root Mean Squared Error (RMSE): 110934.39674566274

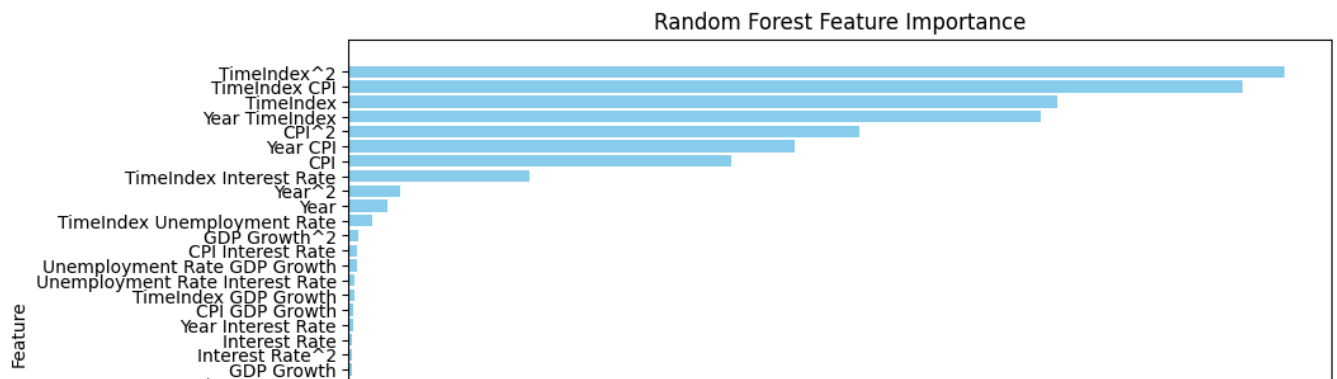
Random Forest Mean Absolute Percentage Error (MAPE): 0.17001389114529938

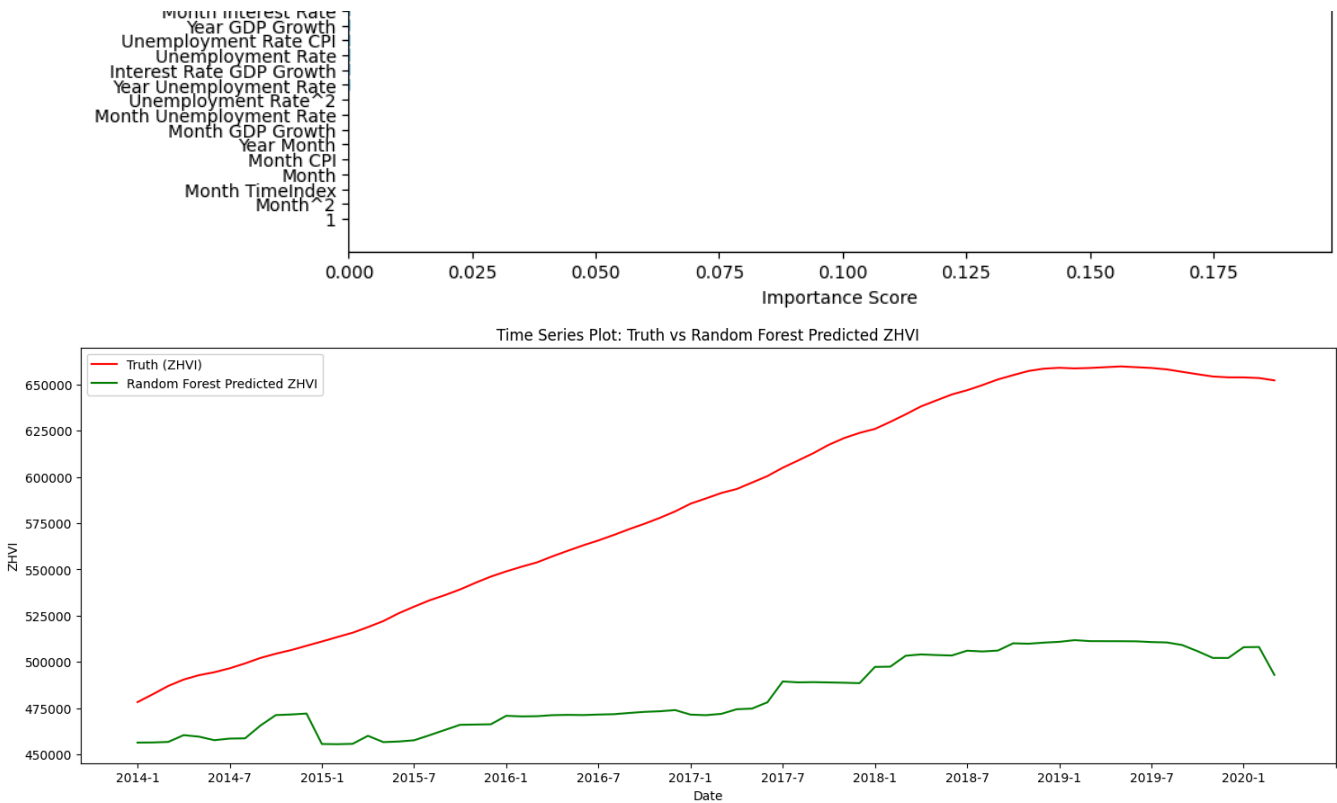
Random Forest Mean Absolute Error (MAE): 102800.76702222222

Random Forest R-squared ( $R^2$ ): -2.3909840626295815

Random Forest Feature Importance:

	Feature	Importance
21	TimeIndex^2	0.189539
23	TimeIndex CPI	0.181033
3	TimeIndex	0.143556
10	Year TimeIndex	0.140208
30	CPI^2	0.103491
12	Year CPI	0.090341
5	CPI	0.077447
24	TimeIndex Interest Rate	0.036600
8	Year^2	0.010326
1	Year	0.007912
22	TimeIndex Unemployment Rate	0.004745
35	GDP Growth^2	0.001947
31	CPI Interest Rate	0.001729
29	Unemployment Rate GDP Growth	0.001691
28	Unemployment Rate Interest Rate	0.001140
25	TimeIndex GDP Growth	0.001053
32	CPI GDP Growth	0.000835
13	Year Interest Rate	0.000798
6	Interest Rate	0.000763
33	Interest Rate^2	0.000742
7	GDP Growth	0.000715
19	Month Interest Rate	0.000534
14	Year GDP Growth	0.000463
27	Unemployment Rate CPI	0.000446
4	Unemployment Rate	0.000351
34	Interest Rate GDP Growth	0.000339
11	Year Unemployment Rate	0.000333
26	Unemployment Rate^2	0.000277
17	Month Unemployment Rate	0.000152
20	Month GDP Growth	0.000141
9	Year Month	0.000074
18	Month CPI	0.000072
2	Month	0.000072
16	Month TimeIndex	0.000070
15	Month^2	0.000066
0	1	0.000000

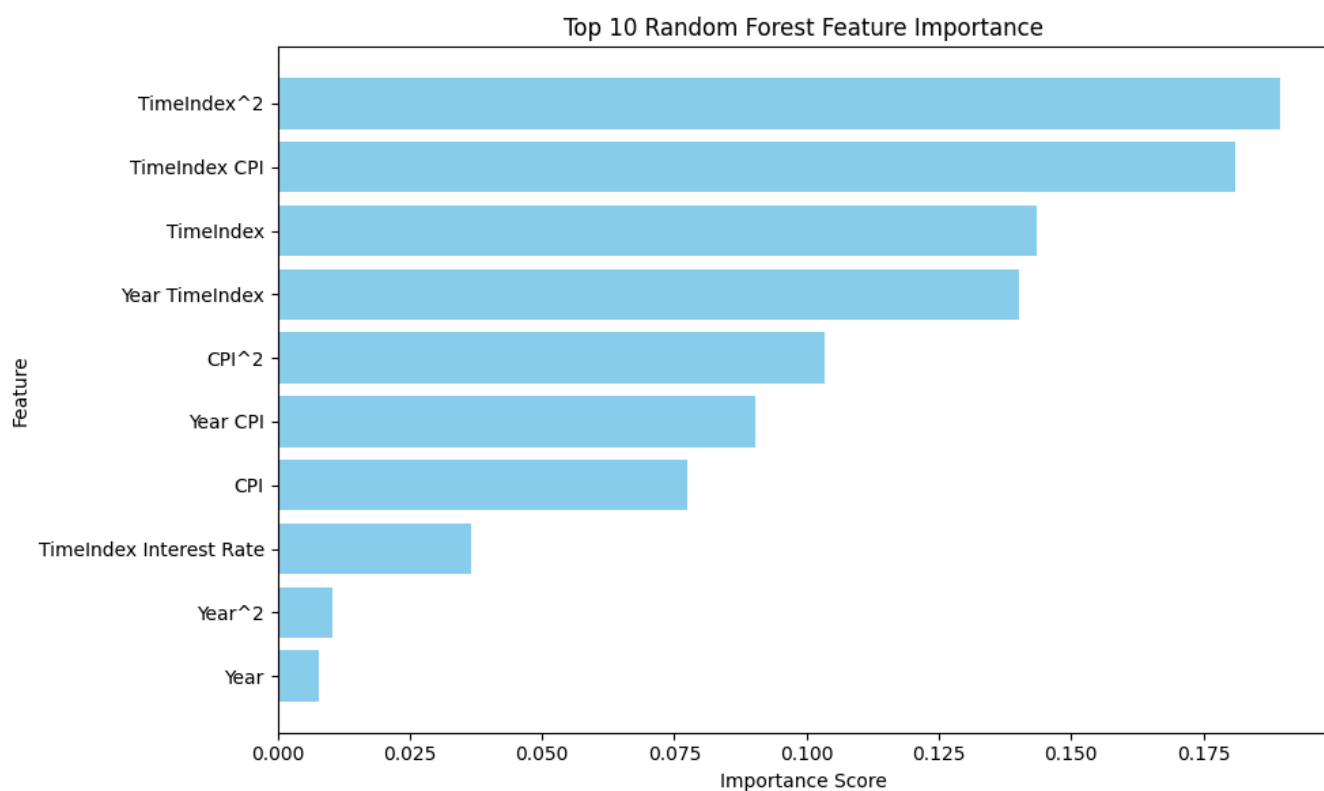




```

1 # Get top 10 features
2 top_10_features = feature_importance.head(10)
3
4 # Plot top 10 feature importance
5 plt.figure(figsize=(10, 6))
6 plt.barh(top_10_features['Feature'], top_10_features['Importance'], color='skyblue')
7 plt.xlabel('Importance Score')
8 plt.ylabel('Feature')
9 plt.title('Top 10 Random Forest Feature Importance')
10 plt.gca().invert_yaxis() # Most important at top
11 plt.tight_layout() # Prevent label cutoff
12 plt.show()

```



## XGBoost Model

```
1 # Split data into training and test (same as before)
2 train = full_df[(full_df['Year'] < 2014) | ((full_df['Year'] == 2013) & (full_df['Month'
3 test = full_df[(full_df['Year'] > 2013) | ((full_df['Year'] == 2014) & (full_df['Month'
4
5 # Define features and target
6 X_train = train[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Ra
7 y_train = train['ZHVI']
8 X_test = test[['Year', 'Month', 'TimeIndex', 'Unemployment Rate', 'CPI', 'Interest Rate
9 y_test = test['ZHVI']
10
11 # Feature Engineering Functions
12 def create_features(df, target_col='ZHVI'):
13     # Create lag features
14     for lag in [1, 2, 3, 6, 12]: # Multiple time horizons
15         df[f'{target_col}_lag_{lag}'] = df[target_col].shift(lag)
16
17     # Create rolling statistics
18     for window in [3, 6, 12]:
19         df[f'{target_col}_rolling_avg_{window}'] = df[target_col].rolling(window).mean(
20         df[f'{target_col}_rolling_std_{window}'] = df[target_col].rolling(window).std()
21
22     # Month/year indicators
23     df['month_sin'] = np.sin(2 * np.pi * df['Month']/12)
24     df['month_cos'] = np.cos(2 * np.pi * df['Month']/12)
25
26     return df
27
28 # Data Preparation
29 train = create_features(train.copy())
30 test = create_features(test.copy())
31
32 # Define features
33 base_features = ['Year', 'Month', 'TimeIndex', 'Unemployment Rate',
34                 'CPI', 'Interest Rate', 'GDP Growth', 'month_sin', 'month_cos']
35 lag_features = [col for col in train.columns if 'lag_' in col or 'rolling_' in col]
36 features = base_features + lag_features
37
38 # Handle missing values from lag features
39 X_train = train[features].dropna()
40 y_train = train.loc[X_train.index, 'ZHVI']
41 X_test = test[features].dropna()
42 y_test = test.loc[X_test.index, 'ZHVI']
43
44 # add polynomial features and scale
45 scaler = StandardScaler()
46 poly = PolynomialFeatures(degree=2)
47 X_train_poly = poly.fit_transform(X_train)
48 X_test_poly = poly.transform(X_test)
49 X_train_scaled = scaler.fit_transform(X_train_poly)
50 X_test_scaled = scaler.transform(X_test_poly)
51
```



```
52 # XGBoost with Native API for MAPE Optimization
53 def xgboost_mape_train(X_train, y_train, X_test, y_test, params):
54     # Convert to DMatrix format
55     dtrain = xgb.DMatrix(X_train, label=y_train)
56     dtest = xgb.DMatrix(X_test, label=y_test)
57
58     # Custom MAPE evaluation metric
59     def mape_eval(preds, dmatrix):
60         labels = dmatrix.get_label()
61         return 'mape', np.mean(np.abs((labels - preds) / (labels + 1e-6))) * 100
62
63     # Train model
64     model = xgb.train(
65         params,
66         dtrain,
67         num_boost_round=1000,
68         evals=[(dtrain, 'train'), (dtest, 'test')],
69         early_stopping_rounds=50,
70         feval=mape_eval,
71         verbose_eval=50
72     )
73     return model
74
75 # Parameter Tuning with scikit-learn API
76 param_grid = {
77     'max_depth': [3, 5, 7],
78     'learning_rate': [0.01, 0.05, 0.1],
79     'subsample': [0.8, 0.9, 1.0],
80     'colsample_bytree': [0.8, 0.9, 1.0],
81     'gamma': [0, 0.1, 0.2],
82     'min_child_weight': [1, 3, 5]
83 }
84
85 xgb_sklearn = xgb.XGBRegressor(
86     objective='reg:squarederror',
87     n_estimators=100,
88     random_state=42,
89     n_jobs=-1
90 )
91
92 tscv = TimeSeriesSplit(n_splits=3)
93 search = RandomizedSearchCV(
94     estimator=xgb_sklearn,
95     param_distributions=param_grid,
96     n_iter=20,
97     cv=tscv,
98     scoring='neg_mean_absolute_percentage_error',
99     verbose=1,
100     n_jobs=-1,
101     random_state=42
102 )
```





```
103
104 search.fit(X_train_scaled, y_train)
105 best_params = search.best_params_
106
107 # Final Model Training with MAPE Focus
108 final_params = {
109     **best_params,
110     'objective': 'reg:squarederror',
111     'seed': 158
112 }
113
114 # Train with native API
115 xgb_model = xgboost_mape_train(
116     X_train_scaled, y_train,
117     X_test_scaled, y_test,
118     final_params
119 )
120
121 # Evaluation
122 xgb_predictions = xgb_model.predict(xgb.DMatrix(X_test_scaled))
123 test.loc[X_test.index, 'XGBoost_Predicted_ZHVI'] = xgb_predictions
124
125 xgb_rmse = math.sqrt(mean_squared_error(y_test, xgb_predictions))
126 xgb_mape = mean_absolute_percentage_error(y_test, xgb_predictions)
127 xgb_mae = mean_absolute_error(y_test, xgb_predictions)
128 xgb_r2 = r2_score(y_test, xgb_predictions)
129
130 metrics = {
131     'RMSE': xgb_rmse,
132     'MAPE': xgb_mape,
133     'MAE': xgb_mae,
134     'R2': xgb_r2
135 }
136
137 print("\nModel Performance:")
138 for name, value in metrics.items():
139     print(f"{name}: {value:.4f}")
140
141 # Plot truth vs prediction
142 plt.figure(figsize=(18, 6))
143 plt.plot(test['Year-Month'], test['ZHVI'], color='red', label='Truth (ZHVI)')
144 plt.plot(test['Year-Month'], test['XGBoost_Predicted_ZHVI'], color='darkgreen', label='XGBoost Predicted ZHVI')
145 plt.xlabel('Date')
146 plt.ylabel('ZHVI')
147 x_ticks = np.arange(0, 80, 6)
148 plt.xticks(x_ticks)
149 plt.title('Time Series Plot: Truth vs XGBoost Predicted ZHVI')
150 plt.legend(loc='upper left')
151 plt.show()
152
```





Fitting 3 folds for each of 20 candidates, totalling 60 fits

[0]	train-rmse:97728.66420	train-mape:27.91250	test-rmse:214233.40336	test-mape:0.1449
[50]	train-rmse:1000.27290	train-mape:0.22388	test-rmse:103571.44737	test-mape:0.1449
[100]	train-rmse:172.37367	train-mape:0.03824	test-rmse:102197.62927	test-mape:0.1449
[146]	train-rmse:75.17564	train-mape:0.01729	test-rmse:102249.79864	test-mape:0.1449

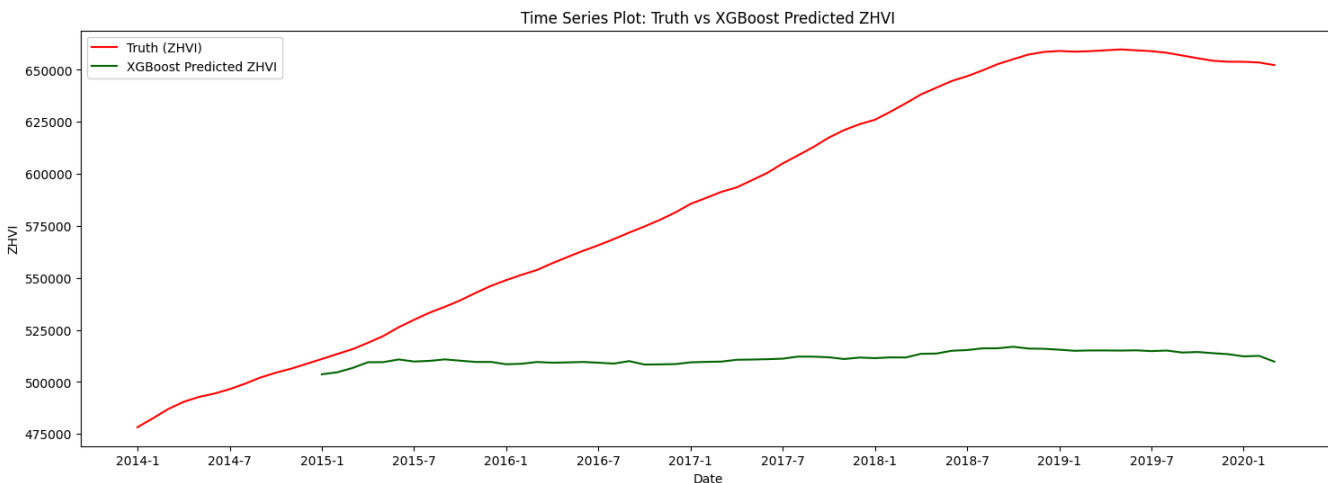
Model Performance:

RMSE: 102249.7987

MAPE: 0.1449

MAE: 90674.5417

R2: -3.2301



```
1 # Display best parameters from hyperparameter tuning
2 print("\nBest Hyperparameters from RandomizedSearchCV:")
3 for param, value in search.best_params_.items():
4     print(f"{param}: {value}")
```



Best Hyperparameters from RandomizedSearchCV:

subsample: 0.8

min\_child\_weight: 1

max\_depth: 5

learning\_rate: 0.1

```
gamma: 0.2  
colsample_bytree: 0.9
```

## Comparison of Models

```
1 # Combined Model Comparison Plot  
2 plt.figure(figsize=(18, 6))  
3 plt.plot(test['Year-Month'], test['ZHVI'], color='black', label='True ZHVI')  
4 plt.plot(test['Year-Month'], OLS_pred, color='pink', label='OLS Predicted ZHVI')  
5 plt.plot(test['Year-Month'], lasso_pred, color='green', label='Lasso Predicted ZHVI')  
6 plt.plot(test['Year-Month'], ridge_pred, color='red', label='Ridge Predicted ZHVI')  
7 plt.plot(test['Year-Month'], nn_predictions, color='purple', label='Neural Network Predicted ZHVI')  
8 plt.plot(test['Year-Month'], rf_predictions, color='orange', label='Random Forest Predicted ZHVI')  
9 plt.plot(test['Year-Month'], test['XGBoost_Predicted_ZHVI'], color='blue', label='XGBoost Predicted ZHVI')  
10  
11  
12 plt.xlabel('Date')  
13 plt.ylabel('ZHVI')  
14 x_ticks = np.arange(0, 80, 6)  
15 plt.xticks(x_ticks)  
16 plt.title('Time Series Plot: True vs All Model Predictions')  
17 plt.legend(loc='upper left')  
18 plt.show()  
19  
20 # Model Performance Comparison Table  
21 model_comparison = pd.DataFrame({  
22     'Model': ['OLS', 'Lasso', 'Ridge', 'Neural Network', 'Random Forest', 'XGBoost'],  
23     'RMSE': [rmse, lasso_rmse, ridge_rmse, nn_rmse, rf_rmse, xgb_rmse],  
24     'MAPE': [mape, lasso_mape, ridge_mape, nn_mape, rf_mape, xgb_mape],  
25     'MAE': [MAE, lasso_MAE, ridge_MAE, nn_mae, rf_mae, xgb_mae],  
26     'R2': [r2, lasso_r2, ridge_r2, nn_r2, rf_r2, xgb_r2]  
27 })  
28  
29 print("\nModel Performance Comparison:")  
30 print(model_comparison.sort_values('RMSE'))
```

