

```
In [16]: 1 #General Libraries for
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import matplotlib
7 from matplotlib.pyplot import figure
8 from sklearn.metrics import mean_absolute_error
9 import statsmodels.tsa.statespace.sarimax
10 from statsmodels.tsa.stattools import adfuller
11 import itertools
12
13 import warnings
14 warnings.filterwarnings('ignore')
15 warnings.warn('DelftStack')
16 warnings.warn('Do not show this message')
17 print("No Warning Shown")
```

No Warning Shown

```
In [2]: 1 df = pd.read_csv('zillow_data.csv')
```

Step 2: Data Preprocessing

```
In [4]: 1 def get_datetimes(df):
2
3     """
4     Takes a dataframe:
5     returns only those column names that can be converted into datetime
6     as datetime objects.
7     NOTE number of returned columns may not match total number of column
8     """
9
10    return pd.to_datetime(df.columns.values[7:], format='%Y-%m')
```

```
In [5]: 1 get_datetimes(df)
```

```
Out[5]: DatetimeIndex(['1996-04-01', '1996-05-01', '1996-06-01', '1996-07-01',
                        '1996-08-01', '1996-09-01', '1996-10-01', '1996-11-01',
                        '1996-12-01', '1997-01-01',
                        ...,
                        '2017-07-01', '2017-08-01', '2017-09-01', '2017-10-01',
                        '2017-11-01', '2017-12-01', '2018-01-01', '2018-02-01',
                        '2018-03-01', '2018-04-01'],
                        dtype='datetime64[ns]', length=265, freq=None)
```

Step 3: EDA and Visualization

```
In [6]: 1 font = {'family' : 'normal',
2           'weight' : 'bold',
3           'size' : 11}
4
5 matplotlib.rc('font', **font)
6
7 # NOTE: if you visualizations are too cluttered to read, try calling 'p
```

Step 4: Reshape from Wide to Long Format

```
In [7]: 1 def melt_total_data(df):
2         """
3         Takes the zillow_data dataset in wide form or a subset of the zillo
4         Returns a long-form datetime dataframe
5         with the datetime column names as the index and the values as the '
6
7         If more than one row is passes in the wide-form dataset, the values
8         will be the mean of the values from the datetime columns in all of
9         """
10
11         melted = pd.melt(df, id_vars=['RegionName', 'RegionID', 'SizeRank',
12         melted['time'] = pd.to_datetime(melted['time'], infer_datetime_form
13         melted = melted.dropna(subset=['value'])
14         return melted.groupby('time').aggregate({'value': 'mean'})
```

```
In [8]: 1 def melt_data(df):
2         """
3         Takes the zillow_data dataset in wide form or a subset of the zillo
4         Returns a long-form datetime dataframe
5         with the datetime column names as the index and the values as the '
6
7         If more than one row is passes in the wide-form dataset, the values
8         will be the mean of the values from the datetime columns in all of
9         """
10
11         melted = pd.melt(df, id_vars=['RegionName', 'RegionID', 'SizeRank',
12         melted['time'] = pd.to_datetime(melted['time'], infer_datetime_form
13         melted = melted.dropna(subset=['value'])
14         return melted
```

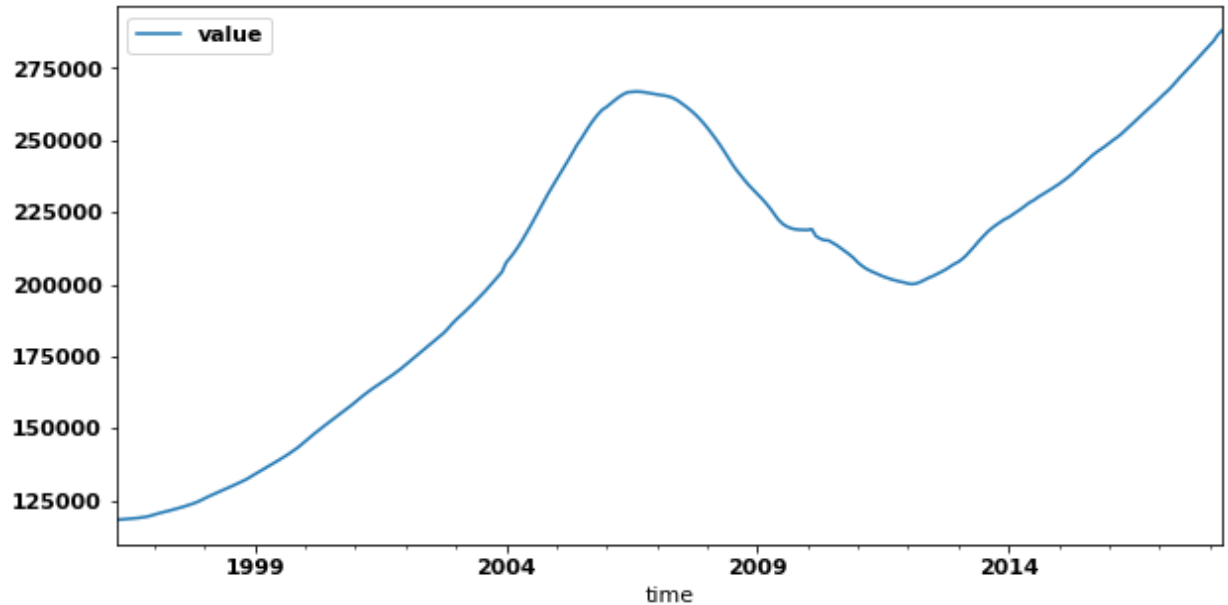
Step 5: ARIMA Modeling

```
In [9]: 1 #Create Initial Dataframe for Model
2
3 total_df = melt_total_data(df)
```

```
In [10]: 1 #Plot Model to Determine potential patterns
2
3 total_df.plot(figsize=(10,5))
```

Out[10]: <AxesSubplot:xlabel='time'>

findfont: Font family ['normal'] not found. Falling back to DejaVu Sans.
findfont: Font family ['normal'] not found. Falling back to DejaVu Sans.



```
In [11]: 1 #Conducting Dickey Fuller Test to Determine Stationarity
2
3 def adtest(frame):
4     dickeyfuller = adfuller(frame)
5     ad_results = pd.DataFrame(dickeyfuller[0:4], index=['Test Statistic',
6                                                         'p-value',
7                                                         '# Lags Used',
8                                                         'Number of Observations Used'])
9
10     return ad_results.loc['p-value']
11
12 adtest(total_df)
```

Out[11]: 0 0.339082
Name: p-value, dtype: float64

```

In [12]: 1 #Due to Dickey Fuller test returning a P Value greater than 0.05, diffe
2
3 good_diff_values = []
4
5 for n in range(1, 60, 1):
6     total_df_diff = total_df.diff(n).dropna()
7     p_val = adtest(total_df_diff)
8     if p_val[0] < 0.051:
9         good_diff_values.append([n, p_val[0]])
10
11 differencing_vals = pd.DataFrame(good_diff_values).sort_values(by=0, as
12
13 differencing_vals.head()

```

Out[12]:

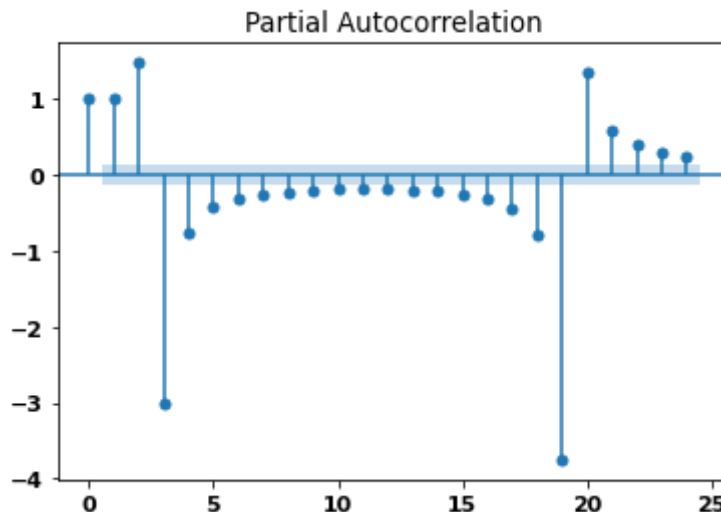
	0	1
0	19	0.017197
1	20	0.009884
2	21	0.021668
3	22	0.015678
4	23	0.022229

```

In [13]: 1 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
2
3 for item in differencing_vals[0]:
4     plot_pacf(total_df.diff(item).dropna());
5     # plot_pacf(total_df.diff(19).dropna());

```

findfont: Font family ['normal'] not found. Falling back to DejaVu Sans.



```
In [14]: 1 #Splitting Train and Test Data
          2
          3 X_train = total_df.head(265-12)
          4 y_train = total_df.tail(12)
```

```
In [17]: 1 #Iterate through potential pdq values for total dataset
          2
          3 # Define the p, d and q parameters to take any value between 0 and 2
          4 p = d = q = range(0,4,1)
          5
          6 # Generate all different combinations of p, q and q triplets
          7 pdq = list(itertools.product(p, d, q))
          8
          9 # Generate all different combinations of seasonal p, q and q triplets (
         10
         11 pdqs = [(x[0], x[1], x[2], 12) for x in list(itertools.product(p, d, q))]
```

```
In [18]: 1 # #Ascertain Best Combos
          2
          3 # from sklearn.metrics import mean_absolute_error
          4
          5 # import statsmodels
          6
          7 # best_iteration = []
          8
          9 # for combo in pdq:
         10 #     for seasonal_combo in pdqs:
         11 #         final_model = statsmodels.tsa.statespace.sarimax.SARIMAX(X_tr
         12 #                                                                 orde
         13 #                                                                 seas
         14 #                                                                 enfo
         15 #                                                                 enfo
         16
         17 #         final_model_fit = final_model.fit()
         18
         19 #         preds = final_model_fit.get_forecast(steps = 12)
         20
         21 #         preds = preds.summary_frame()[['mean']]
         22
         23 #         mae = mean_absolute_error(preds, y_train)
         24
         25 #         comp_error = mae / int(y_train.tail(1)['value'])
         26
         27 #         best_iteration.append([combo, seasonal_combo, mae, comp_error])
         28
         29 #         print([combo, seasonal_combo, mae, comp_error])
```

```

In [19]: 1 # #Ascertain Best Combination:
          2
          3 # best_iteration = pd.DataFrame(best_iteration)
          4
          5 # best_iteration.columns = ['pdq', 'PDQS', 'MAE', 'MAE / M12']
          6
          7 # best_iteration = best_iteration.sort_values(by='MAE', ascending = Fal
          8
          9 # best_iteration['MAE'] = round(best_iteration['MAE'])
         10
         11 # #Isolate best features
         12
         13 # best_pdq = list(best_iteration[['pdq']].tail(1)['pdq'])[0]
         14 # best_PDQS = list(best_iteration[['PDQS']].tail(1)['PDQS'])[0]
         15
         16 # #Saving best features below:
         17
         18 # best_iteration.tail(1)
         19
         20 # #      pdq PDQS      MAE MAE / M12
         21 # # 2119      (2, 0, 1)   (0, 1, 3, 12)   88.0      0.000307
         22
         23 best_pdq = (2,0,1)
         24 best_PDQS = (0,1,3,12)

```

```

In [20]: 1 #Fit Best Model
          2
          3 import statsmodels
          4
          5 final_model = statsmodels.tsa.statespace.sarimax.SARIMAX(X_train,
          6                                                         order = best_pdq,
          7                                                         seasonal_order = best_PDQS,
          8                                                         enforce_stationarity = False,
          9                                                         enforce_invertibility = False)
         10
         11 fit_model = final_model.fit()

```

```
In [21]: 1 fit_model.summary()
```

Out[21]: SARIMAX Results

Dep. Variable:	value	No. Observations:	253
Model:	SARIMAX(2, 0, 1)x(0, 1, [1, 2, 3], 12)	Log Likelihood	-1519.133
Date:	Mon, 25 Oct 2021	AIC	3052.266
Time:	15:43:01	BIC	3075.458
Sample:	04-01-1996	HQIC	3061.648
	- 04-01-2017		
Covariance Type:	opg		

	coef	std err	z	P> z	[0.025	0.975]
ar.L1	1.9791	0.043	46.081	0.000	1.895	2.063
ar.L2	-0.9789	0.043	-22.730	0.000	-1.063	-0.895
ma.L1	-0.4010	0.108	-3.724	0.000	-0.612	-0.190
ma.S.L12	0.4827	0.124	3.897	0.000	0.240	0.725
ma.S.L24	-1.0618	0.105	-10.159	0.000	-1.267	-0.857
ma.S.L36	0.3717	0.090	4.135	0.000	0.196	0.548
sigma2	1.382e+05	2.15e+04	6.416	0.000	9.6e+04	1.8e+05

Ljung-Box (L1) (Q):	0.70	Jarque-Bera (JB):	1287.79
Prob(Q):	0.40	Prob(JB):	0.00
Heteroskedasticity (H):	0.49	Skew:	0.13
Prob(H) (two-sided):	0.00	Kurtosis:	15.34

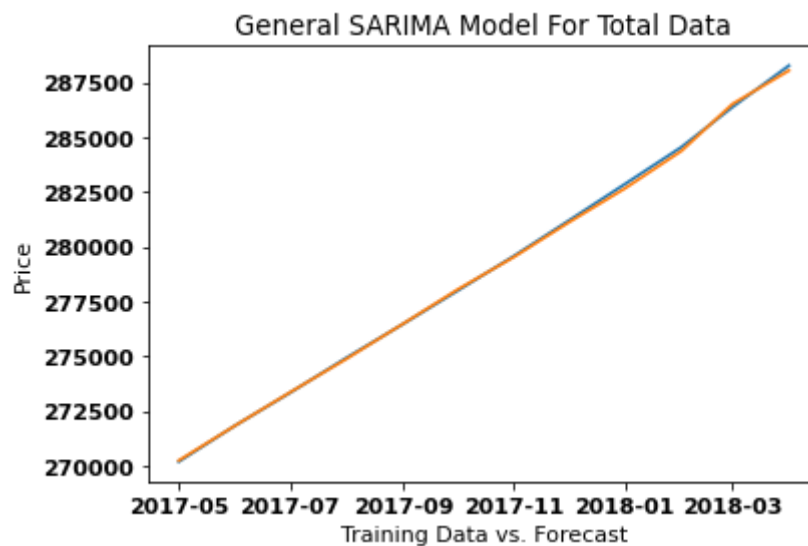
Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).

[2] Covariance matrix is singular or near-singular, with condition number 6.7e+14. Standard errors may be unstable.

```
In [74]: 1 #Get Preds and Map Against y_train
2
3 preds = fit_model.get_forecast(steps=12).summary_frame()['mean']
4
5 plt.plot(preds)
6 plt.plot(y_train)
7
8 print(f'Mean Absolute Error: {mean_absolute_error(preds, y_train)}')
9
10 plt.title('General SARIMA Model For Total Data')
11 plt.xlabel('Training Data vs. Forecast')
12 plt.ylabel('Price');
13
```

Mean Absolute Error: 88.33125573327318




```
In [24]: 1 #Creating dataframe to assess future values of all regions
2
3 final_df = melt_data(df)
4
5 final_df['time'] = pd.to_datetime(final_df['time'])
6
7 final_df.set_index('time', inplace=True)
8
9 final_df
```

Out[24]:

	RegionName	RegionID	SizeRank	City	State	Metro	CountyName	value
time								
1996-04-01	60657	84654	1	Chicago	IL	Chicago	Cook	334200.0
1996-04-01	75070	90668	2	McKinney	TX	Dallas-Fort Worth	Collin	235700.0
1996-04-01	77494	91982	3	Katy	TX	Houston	Harris	210400.0
1996-04-01	60614	84616	4	Chicago	IL	Chicago	Cook	498100.0
1996-04-01	79936	93144	5	El Paso	TX	El Paso	El Paso	77300.0
...
2018-04-01	1338	58333	14719	Ashfield	MA	Greenfield Town	Franklin	209300.0
2018-04-01	3293	59107	14720	Woodstock	NH	Claremont	Grafton	225800.0
2018-04-01	40404	75672	14721	Berea	KY	Richmond	Madison	133400.0
2018-04-01	81225	93733	14722	Mount Crested Butte	CO	NaN	Gunnison	664400.0
2018-04-01	89155	95851	14723	Mesquite	NV	Las Vegas	Clark	357200.0

3744704 rows x 8 columns

```
In [25]: 1 #Isolating all regions
2
3 regions = pd.DataFrame(final_df['RegionID'].value_counts())
4
5 #Isolating Regions with complete data sets
6
7 complete_regions = regions[regions['RegionID'] == 265]
8
9 #Isolating regions that have experienced the most growth in the last ye
10
11 complete_regions = complete_regions.reset_index()
12
13 complete_regions.columns = ['RegionID', 'DataPoints']
14
15 annual_growth = complete_regions[['RegionID']]
```

```
In [26]: 1 #Create Columns to Compare Annual growth
2
3 annual_growth['Current Year'] = pd.to_datetime('2018-04-01')
4 annual_growth['Previous Year'] = pd.to_datetime('2017-04-01')
```

```
In [27]: 1 merge_table = final_df.reset_index()  
2  
3 merge_table['time'] = pd.to_datetime(merge_table['time'])  
4  
5 merge_table
```

Out[27]:

	time	RegionName	RegionID	SizeRank	City	State	Metro	CountyName	v
0	1996-04-01	60657	84654	1	Chicago	IL	Chicago	Cook	3342
1	1996-04-01	75070	90668	2	McKinney	TX	Dallas-Fort Worth	Collin	2357
2	1996-04-01	77494	91982	3	Katy	TX	Houston	Harris	2104
3	1996-04-01	60614	84616	4	Chicago	IL	Chicago	Cook	4981
4	1996-04-01	79936	93144	5	El Paso	TX	El Paso	El Paso	773
...
3744699	2018-04-01	1338	58333	14719	Ashfield	MA	Greenfield Town	Franklin	2093
3744700	2018-04-01	3293	59107	14720	Woodstock	NH	Claremont	Grafton	2258
3744701	2018-04-01	40404	75672	14721	Berea	KY	Richmond	Madison	1334
3744702	2018-04-01	81225	93733	14722	Mount Crested Butte	CO	NaN	Gunnison	6644
3744703	2018-04-01	89155	95851	14723	Mesquite	NV	Las Vegas	Clark	3572

3744704 rows x 9 columns

```

In [28]: 1  # #Joining DataFrames to get current and previous values
2
3  growth_table = annual_growth.merge(merge_table,
4                                     how='inner', left_on=['Current Year', 'RegionID'],
5
6  growth_table = growth_table[['Current Year', 'Previous Year', 'RegionID', 'Current Value', 'Previous Value', 'Growth']]
7
8  growth_table.columns = ['Current Year', 'Previous Year', 'RegionID', 'Current Value', 'Previous Value', 'Growth']
9
10 growth_table = growth_table.merge(merge_table, how='inner',
11                                   left_on=['Previous Year', 'RegionID'],
12
13 growth_table = growth_table[['RegionID', 'Previous Year', 'Current Year', 'Current Value', 'Previous Value', 'Growth']]
14
15 growth_table.columns = ['RegionID', 'Previous Year', 'Current Year', 'Current Value', 'Previous Value', 'Growth']
16
17 growth_table['Growth'] = growth_table['Current Value'] / growth_table['Previous Value']
18
19 #Isolating Top 20 Regions that Grew between 2017 and 2018
20
21 top_regions = growth_table.sort_values(by='Growth', ascending = False).
22
23 top_regions

```

Out[28]:

	RegionID	Previous Year	Current Year	Previous Value	Current Value	Growth
10656	60610	2017-04-01	2018-04-01	123800.0	186700.0	0.508078
6050	66014	2017-04-01	2018-04-01	36800.0	52900.0	0.437500
12044	60607	2017-04-01	2018-04-01	131200.0	188300.0	0.435213
12183	60561	2017-04-01	2018-04-01	141500.0	202000.0	0.427562
8174	73219	2017-04-01	2018-04-01	67200.0	95500.0	0.421131
...
2214	72702	2017-04-01	2018-04-01	87000.0	110800.0	0.273563
6359	98652	2017-04-01	2018-04-01	58500.0	74500.0	0.273504
6740	63882	2017-04-01	2018-04-01	90400.0	115000.0	0.272124
12891	73115	2017-04-01	2018-04-01	86400.0	109900.0	0.271991
10698	97528	2017-04-01	2018-04-01	1772600.0	2254100.0	0.271635

100 rows × 6 columns

```
In [29]: 1 #Creating function to get forecast by region
2
3 region_melt = melt_data(df)
4
5 def get_2019_forecast(Region_ID):
6
7     X_train_region = region_melt[region_melt['RegionID'] == Region_ID]
8
9     X_train_region['time'] = pd.to_datetime(X_train_region['time'])
10
11     X_train_region.set_index('time', inplace = True)
12
13     X_train_region = X_train_region[['value']]
14
15     regional_model = statsmodels.tsa.statespace.sarimax.SARIMAX(X_train
16                                                                order = best_pdq,
17                                                                seasonal_order = best_PDQS,
18                                                                enforce_stationarity = False
19                                                                enforce_invertibility = False
20
21     regional_model_fit = regional_model.fit()
22
23     regional_preds = regional_model_fit.get_forecast(12)
24
25     regional_preds_df = regional_preds.summary_frame()[['mean']]
26
27     return [Region_ID, int(regional_preds_df[regional_preds_df.index ==
```

```
In [30]: 1 forecasted_values = []
2
3 for item in list(top_regions['RegionID']):
4     forecasted_values.append(get_2019_forecast(item))
```

```

In [31]: 1 #Isolating top regions by forecast
2
3 top_future_regions = pd.DataFrame(forecasted_values)
4
5 top_future_regions.columns = ['RegionID', '2019 value']
6
7 top_100_forecasted_growth = top_future_regions.merge(growth_table, how=
8
9 top_100_forecasted_growth['Future Growth'] = (top_100_forecasted_growth
10 top_100_forecasted_growth
11
12 #Isolating Top 10 Growers
13
14 top_10_growers_df = top_100_forecasted_growth.sort_values('Future Growth
15
16 top_10_growers = top_100_forecasted_growth.sort_values('Future Growth',
17
18 top_10_growers = list(top_10_growers['RegionID'])
19
20 top_10_growers

```

Out[31]: [99877, 73177, 66015, 60706, 91259, 60642, 92306, 399593, 66828, 70395]

```

In [51]: 1 avg_growth = np.average(top_100_forecasted_growth['Future Growth'])

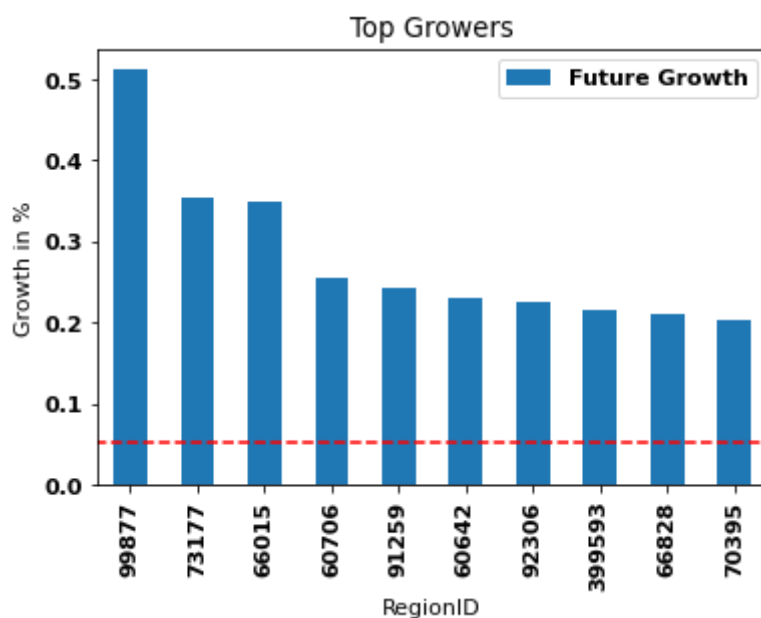
```

```

In [75]: 1 top_10_growers_df.plot(x='RegionID', y='Future Growth', kind='bar')
2 plt.hlines(avg_growth, xmin=-1, xmax=50, linestyle='dashed', color='re
3 plt.title('Top Growers')
4 plt.ylabel('Growth in %')

```

Out[75]: Text(0, 0.5, 'Growth in %')



```
In [32]: 1 def region_acc_compute(ID, choice):
2
3     region_df = region_melt[region_melt['RegionID'] == ID]
4
5     region_df['time'] = pd.to_datetime(region_df['time'])
6
7     region_df.set_index('time', inplace = True)
8
9     region_df = region_df[['value']]
10
11     X_train = region_df.head(265-12)
12     y_train = region_df.tail(12)
13
14     regional_model = statsmodels.tsa.statespace.sarimax.SARIMAX(X_train
15                                                                order = best_pdq,
16                                                                seasonal_order = best_PDQS,
17                                                                enforce_stationarity = False
18                                                                enforce_invertibility = False)
19
20     fit_model = regional_model.fit()
21     preds = fit_model.get_forecast(steps = 12).summary_frame()['mean']
22
23     mae = mean_absolute_error(preds, y_train)
24     comp_error = round(float(mae / preds[-1:]), 2)
25
26     if choice == 'stats':
27
28         return [ID, mae, comp_error]
29
30     else:
31
32         return preds
```

```
In [34]: 1 region_acc_df = []
2
3     for item in top_10_growers:
4         region_acc_df.append(region_acc_compute(item, 'stats'))
5
```

```

In [42]: 1 #Top 10 Most Accurate
2
3 region_acc_df = pd.DataFrame(region_acc_df)
4
5 region_acc_df.columns = ['RegionID', 'MAE', 'Comp Error']
6
7 most_accurate_growers = region_acc_df.sort_values(by='Comp Error', asce
8
9 most_accurate_growers
10
11 #Add Descriptors to Region
12
13 most_accurate_growers = most_accurate_growers.merge(region_melt, how='i
14
15 most_accurate_growers = most_accurate_growers[['RegionID', 'City', 'Sta
16
17 most_accurate_growers
18
19 most_accurate_growers.to_excel('acc_growers.xlsx')

```

```

In [81]: 1 comp_error_avg = np.average(region_acc_df['Comp Error'])

```

```

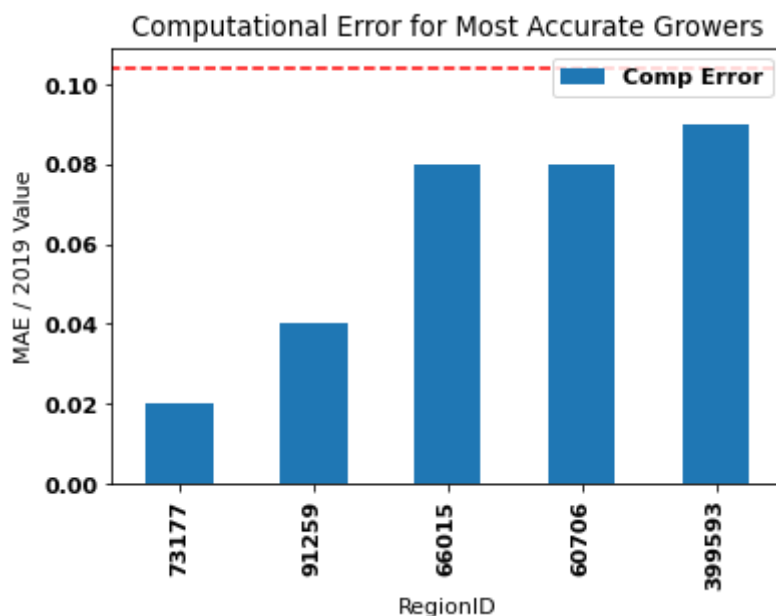
In [85]: 1 most_accurate_growers.plot('RegionID', 'Comp Error', kind='bar')
2
3 plt.hlines(comp_error_avg, xmin=-1, xmax=50, linestyle='dashed', color
4 plt.ylabel('MAE / 2019 Value')
5 plt.title('Computational Error for Most Accurate Growers')
6

```

```

Out[85]: Text(0.5, 1.0, 'Computational Error for Most Accurate Growers')

```




```
In [37]: 1 #Create function to plot each forecast
2
3 def plot_top_growers(region):
4
5     test_preds = region_acc_compute(region, 'preds')
6
7     sample_region_df = region_melt[region_melt['RegionID'] == region]
8
9     sample_region_df['time'] = pd.to_datetime(sample_region_df['time'])
10
11     sample_region_df.set_index('time', inplace = True)
12
13     sample_region_df = sample_region_df['value']
14
15     y_train_sample_df = sample_region_df.tail(12)
16
17     figure(figsize=(12, 6), dpi=80)
18
19     plt.plot(sample_region_df['2015:'], color='black')
20     plt.plot(test_preds, color='blue')
21     plt.plot(y_train_sample_df, color='green')
22
23     #     county = most_accurate_growers[most_accurate_growers['RegionID'] == region]
24     #     city = most_accurate_growers[most_accurate_growers['RegionID'] == region]
25     #     state = most_accurate_growers[most_accurate_growers['RegionID'] == region]
26
27     location = f'{region}'
28
29     plt.title(f'Forecast For {location}')
30
31     test_preds.to_csv(f'{region}_preds.csv')
```

```
In [38]: 1 #plot top growers to determine accuracy
          2
          3 for item in list(most_accurate_growers['RegionID']):
          4     plot_top_growers(item)
```

```
/Users/angelogayanelo/opt/anaconda3/lib/python3.8/site-packages/statsmodels/tsa/base/tsa_model.py:524: ValueWarning: No frequency information was provided, so inferred frequency MS will be used.
  warnings.warn('No frequency information was')
/Users/angelogayanelo/opt/anaconda3/lib/python3.8/site-packages/statsmodels/tsa/base/tsa_model.py:524: ValueWarning: No frequency information was provided, so inferred frequency MS will be used.
  warnings.warn('No frequency information was')
/Users/angelogayanelo/opt/anaconda3/lib/python3.8/site-packages/statsmodels/tsa/base/tsa_model.py:524: ValueWarning: No frequency information was provided, so inferred frequency MS will be used.
  warnings.warn('No frequency information was')
/Users/angelogayanelo/opt/anaconda3/lib/python3.8/site-packages/statsmodels/tsa/base/tsa_model.py:524: ValueWarning: No frequency information was provided, so inferred frequency MS will be used.
  warnings.warn('No frequency information was')
/Users/angelogayanelo/opt/anaconda3/lib/python3.8/site-packages/statsmodels/base/model.py:566: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals
  warnings.warn("Maximum Likelihood optimization failed to converge. Check mle_retvals", ConvergenceWarning)
```

Step 6: Interpreting Results

```
In [ ]: 1
```

```
In [ ]: 1
```