

Assignment 2 Answer Example

Instructor: Dr. Emily Zhu

Due date: April 4th, 2021

Please note that the answers/methods are NOT unique.
The provided codes are for reference ONLY.

1. The God, The Bad and The Ugly: Ratings for Austin Airbnb Properties

We now aim to identify the good, the so-so and the bad properties by studying six(6) review scores for accuracy, cleanliness, checkin, communication, location and value, of those listings with available rating records, in *reviews.csv* (in-sample/training data). By deriving the evaluating criteria from these existing listings, please classify 28 ratings in *Reviews4Cluster.h5* (out-of-sample/testing set) into 3 clusters of good, medium and bad listings, respectively.

- Use hierarchical clustering to classify properties regarding the given six review scores. What problem(s) do you have when running Python? Please try to draw a random sample of 150 properties and run hierarchical clustering on the smaller sample.
- Use K-means to cluster properties regarding the given six review scores. Compare the outcomes from hierarchical and K-Means clustering. Are the clustering results the same under different clustering methods?
- Will the clustering above be the same if replacing missing values with modes of the corresponding review scores?
Hint: You may make a table to illustrate the comparisons.
- Do you think it reasonable to classify those ratings into 3 categories?
Hint: You may want to check Within-Cluster-Sum-of-Squares (WCSS) in K-Means.

Answer. We list the answers to the first two questions in a table for the purpose of comparison:

	Clustering			
ID*	kmean**		hierarchy	
	delete NAN	replace NAN	delete NAN	replace NAN
1	1	1	0	0
2	2	1	0	0
3	1	2	0	0
4	0	0	0	0
5	1	1	0	0
6	1	1	1	1
7	0	1	2	2
8	2	1	0	0
9	1	2	0	0
10	0	0	0	0
11	1	1	0	0
12	1	1	1	1
13	0	1	2	2
14	2	2	0	0
15	2	0	0	0
16	2	1	0	0
17	2	1	1	1
18	2	1	2	2
19	0	0	0	0
20	1	2	0	0
21	1	2	1	1
22	0	2	2	2
23	0	0	0	0
24	0	0	1	1
25	0	0	2	2
26	1	1	0	0
27	0	1	2	2
28	0	1	1	1

*The ID refers to the new review ratings in *Reviews4Cluster.h5*.

** K-Means is not a deterministic algorithm, in this case, you may end up with different prediction results.

```

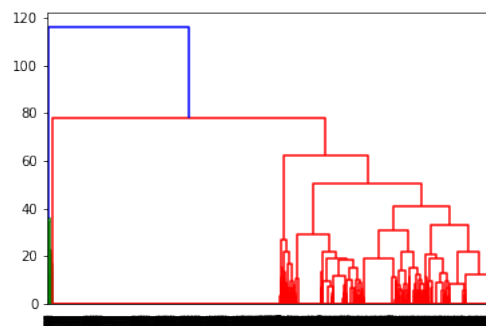
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 from sklearn.cluster import AgglomerativeClustering
6
7 #---Historical Review Ratings-----
8 df_a = pd.read_csv('reviews.csv', usecols = ['review_scores_accuracy',
9         'review_scores_cleanliness', 'review_scores_checkin', '
10         review_scores_communication', 'review_scores_location', '
11         review_scores_value'])
12 df_b = df_a.dropna(axis=0) #drop NaN for the 1st question
13 df_b = df_a.fillna(df_a.mode()) # replace missing values with the modes
14         for the 2nd question
15
16 #---Retrieve New Ratings for Prediction-----
    
```

```

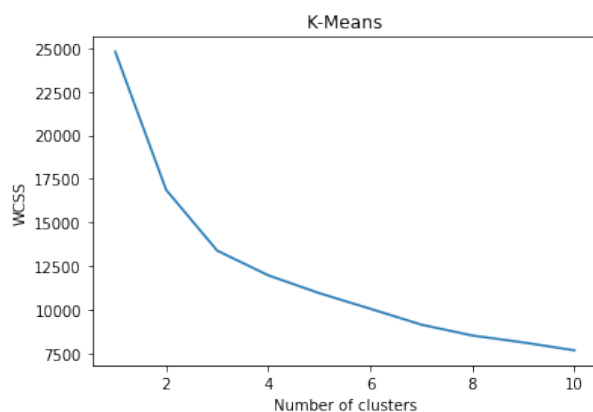
13 data_store = pd.HDFStore('Reviews4Cluster.h5') # Retrieve data using key
14 df_new = data_store['df_review_cluster']
15 data_store.close()
16
17 #----Apply K-Means Clustering-----
18 kmeans = KMeans(n_clusters = 3)
19 # fit kmeans object to data
20 kmeans.fit(df_b)
21 # print location of clusters learned by kmeans object
22 print(kmeans.cluster_centers_)
23 # save new clusters for chart
24 y_km = kmeans.fit_predict(df_new)
25
26 #----Apply Hierarchy Clustering-----
27 import scipy.cluster.hierarchy as sch
28 from sklearn.cluster import AgglomerativeClustering
29 #You should plot dendrogram first***
30 dendrogram = sch.dendrogram(sch.linkage(df_b, method='ward'))
31 hc = AgglomerativeClustering(n_clusters=3, affinity = 'euclidean',
32                               linkage = 'ward')
33 y_hc = hc.fit_predict(df_new)
34
35 #-----***You May Need this: Expand the limit for recursion-----
36 import sys
37 sys.setrecursionlimit(10000)
38 #DO REMEMBER to set back to default after running Hierarchy Clustering
39 sys.setrecursionlimit(1000)

```

The dendrogram may look like



The number of clusters is 3, which is reasonable as Within-Cluster-Sum-of-Squares (WCSS) in K-Means shows below



```

1 wcss = []
2 for i in range(1, 11):
3     kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state =
4         42)
5     kmeans.fit(df_b)
6     wcss.append(kmeans.inertia_)
7 plt.plot(range(1, 11), wcss)
    
```

No Traffic Jams In the Extra Mile: Weather, Traffic & Planning

Next, we study the traffic volumes, how weather affects traffic and temperature prediction based on the file *Metro_Interstate_Traffic_AMPeak_cleaned.csv*, an extract from the original dataset¹.

2. How Does Weather Affect Traffic Volumes?

We focus on traffic volumes in rush hours and take 8AM for instance. In other words, the sample of our interest consists of observations at 8AM ONLY (one observation per day from 2015 to 2018). We would like to examine how the following factors affect the *traffic volumes*.

- *weekend*, including Saturdays, Sundays and holidays in the column *holiday*;
- *temperature* as shown in the column *temp*;
- *amount of clouds* as shown in the column *clouds_all*;
- *spot snowfall* as shown in the column *snow_1h*.

Please compare the following two models, M1 & M2, in terms of statistical significance (R^2 , adjusted R^2) and predicting power (out-of-sample RMSE). You can divide the sample into 80% training and 20% testing sets. Which model would you recommend and why?

(Hint: We do NOT need to consider the time series in this case, because we do NOT study any association with time.)

traffic volume = $b_0 + b_1(\text{weekend}) + b_2(\text{clouds}) + b_3(\text{temperature}) + b_4(\text{snowfall})$ (M1);

traffic volume = $b_0 + b_1(\text{weekend}) + b_2(\text{clouds}) + b_3(\text{snowfall})$ (M2).

Answer. We prefer M2 as it is more statistically significant (higher R square and Adjusted R square) and stronger prediction power (lower RMSE).

	M1	M2
R square	.756	.738
Adjusted R	.755	.737
RMSE	830.04	839.83

Please note that the answers are NOT unique, as the test and training sets are split at random.

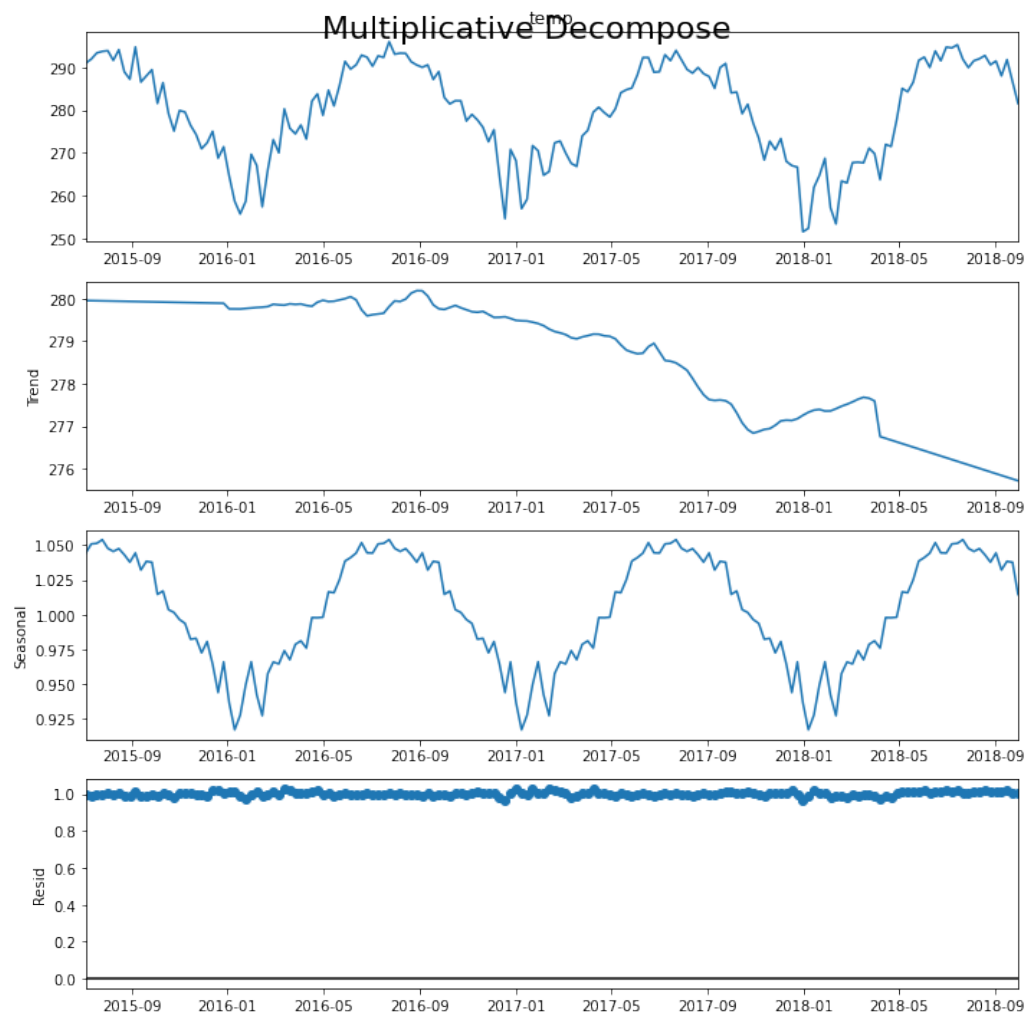
¹Data source and descriptions of the data can be found at <https://archive.ics.uci.edu/ml/datasets/Metro+Interstate+Traffic+Volume>

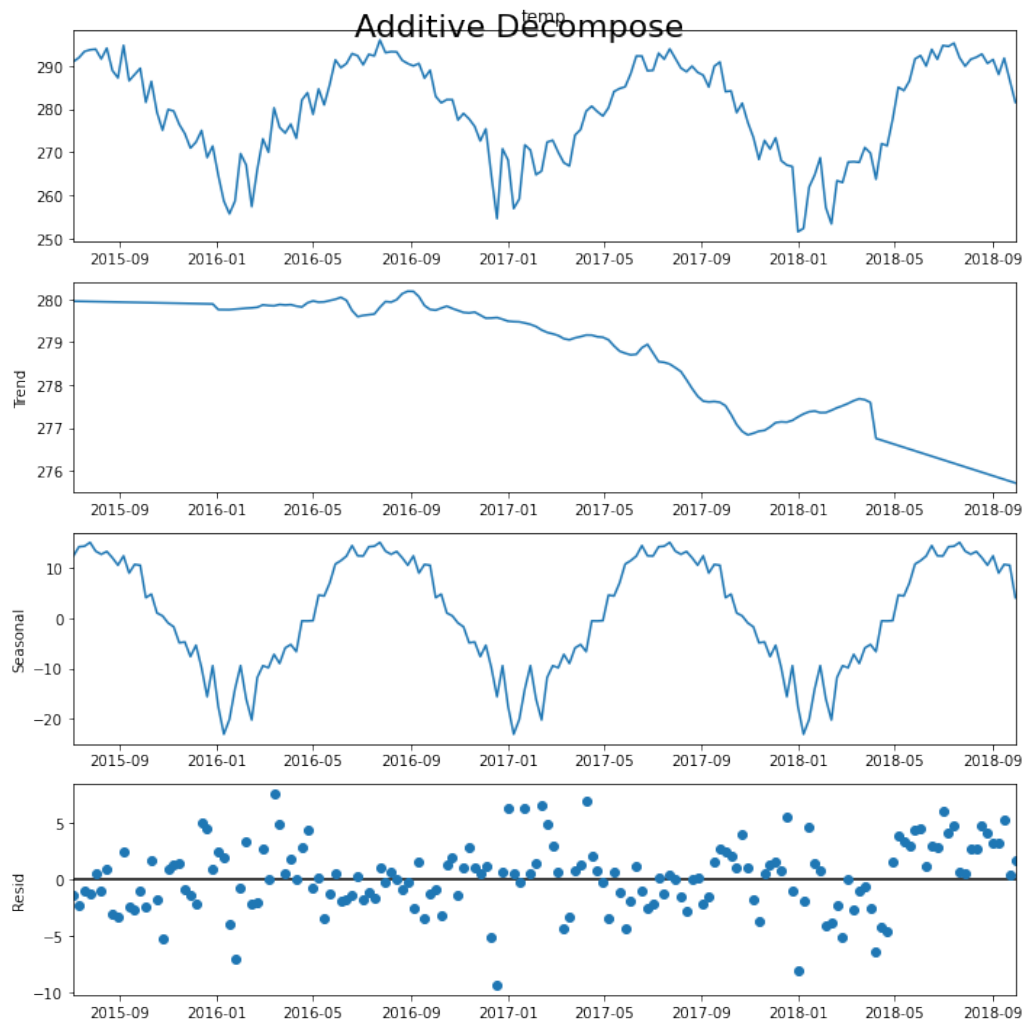
```
1 import pandas as pd
2 import numpy as np
3
4 df_clean = pd.read_csv('Metro_Interstate_Traffic_cleaned.csv',
5                        parse_dates = True, encoding = "utf_8_sig")
6
7 #----Regression Model 1-----
8 import statsmodels.formula.api as smf
9 # create a fitted model
10 lm1 = smf.ols(formula = 'traffic_volume ~ weekend + temp + clouds_all +
11                    rain_1h + snow_1h', data = df_slt).fit()
12 lm1.summary()
13
14 #----Regression Model 2-----
15 from sklearn.linear_model import LinearRegression
16 from sklearn import metrics
17 from sklearn.model_selection import train_test_split
18 # Out-of-Sample Evaluation
19 X = df_slt[['weekend', 'temp', 'clouds_all', 'rain_1h', 'snow_1h']]
20 y = df_slt.traffic_volume
21
22 # Split data
23 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =
24                                                    0.2, random_state=1)
25
26 # Instantiate model
27 lm2 = LinearRegression()
28
29 # Fit Model
30 lm2.fit(X_train, y_train)
31
32 # Predict
33 y_pred = lm2.predict(X_test)
34
35 # RMSE
36 print(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
```

3. Can We Predict Average Weekly Temperature?

- Resample the variable temperature *Temp* into weekly basis; and
- decompose the time series *Temp* and which decomposition (additive or multiplicative) is appropriate?
 - cross validate RMSE on the later 30% of the time series.
 - predict average weekly temperature for the next year starting on Oct. 7, 2018.

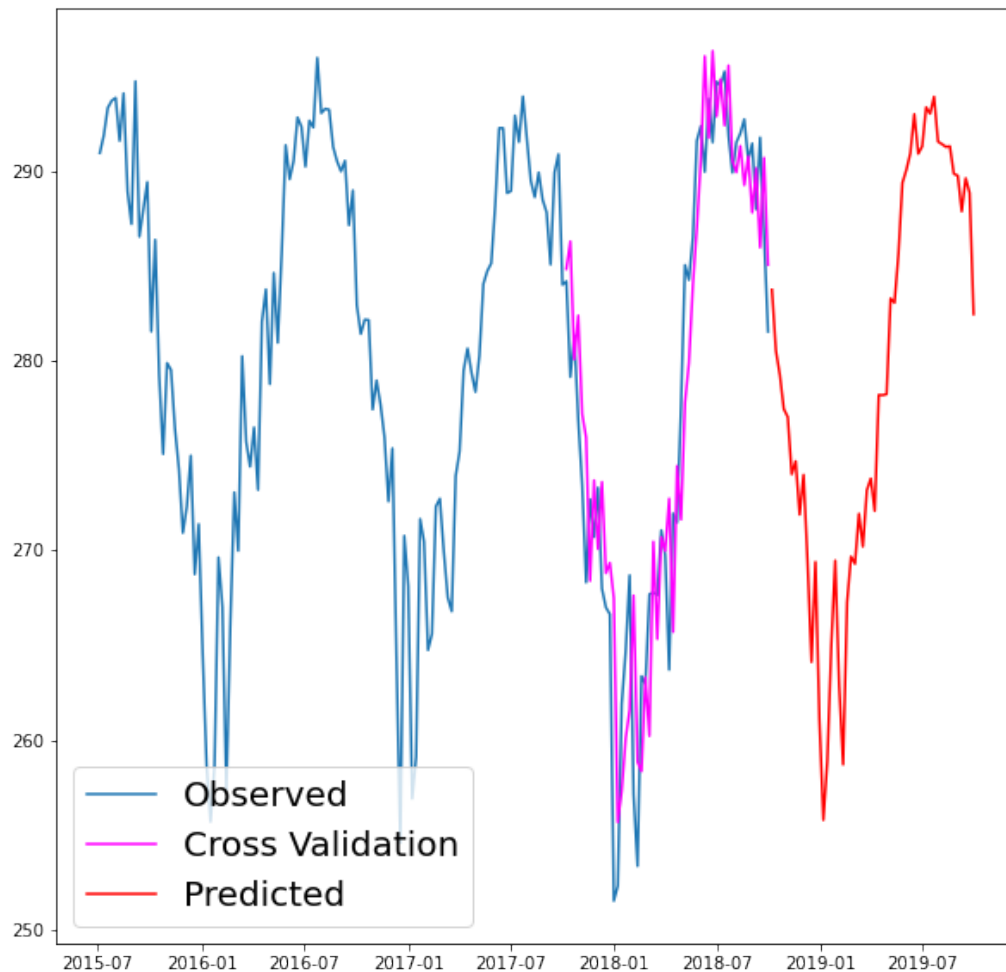
Answer.





Because the residuals of additive decomposition are randomly scattered, it is an appropriate decomposition.

RMSE of cross validation is 4.965452802567684 degrees.



```

1 df1 = df_full[['temp', 'date']]
2 df1.set_index('date', inplace = True)
3 df1.index = df1.index.astype('datetime64[ns]')
4 df = df1.resample('W', label='right', closed = 'right').mean()
5
6 ##-----Decompose-----
7 from matplotlib import pyplot as plt
8 from statsmodels.tsa.seasonal import seasonal_decompose
9 # Multiplicative Decomposition
10 result_mul = seasonal_decompose(df['temp'], model= 'multiplicative',
11     extrapolate_trend= 'freq')
12 #Setting extrapolate_trend='freq' takes care of any missing values in
13     the trend and residuals at the beginning of the series.
14 # Additive Decomposition
15 result_add = seasonal_decompose(df['temp'], model= 'additive',
16     extrapolate_trend= 'freq')
17 # Plot
18 plt.rcParams.update({'figure.figsize': (10,10)})
19 result_mul.plot().suptitle('Multiplicative Decompose', fontsize=22)
20 result_add.plot().suptitle('Additive Decompose', fontsize=22)
21 plt.show()
22
23 ##----Predict trend, seasonality and residual separately-----
24 from statsmodels.tsa.ar_model import AutoReg, ar_select_order
25 from sklearn.metrics import mean_squared_error
26 series = df['temp']
    
```



```
24 historic = series.iloc[:int(len(series) * 0.7)]
25 test = series.iloc[int(len(series) * 0.7):]
26 historic = historic.to_list()
27 prediction = []
28
29 for i in range(len(test)):
30     sel = ar_select_order(historic, 13, old_names=False)
31     sel.ar_lags
32     model_fit = sel.model.fit()
33     pred = model_fit.predict(start=len(historic), end=len(historic),
34                             dynamic=False)
35     predictions.append(pred[0])
36     historic.append(test[i])
37 cross_val = pd.Series(predictions, index=test.index)
38
39 test_score = mean_squared_error(test, cross_val, squared = False)
40 print(test_score)
41
42 ##----Prediction for the next year----
43 #existing data
44 historic = series
45 #prediction for the next 52 weeks (1 year)
46 n_pred = 52
47 #the future weeks starts on 2018-09-30, because the last existing data
48 #ends on 2018-10-07
49 date_pred = pd.date_range("2018-10-07", periods=n_pred, freq="w")
50
51 sel = ar_select_order(historic, 13, glob = True, seasonal = True,
52                       old_names=False)
53 sel.ar_lags
54 model_fit = sel.model.fit()
55 pred = model_fit.predict(start=len(historic), end=len(historic) + n_pred
56                           - 1, dynamic=False)
57 predictions = pd.Series(pred.values, index=date_pred)
58
59 # plot results
60 plt.plot(series, label='Observed')
61 plt.plot(cross_val, color='magenta', label='Cross Validation')
62 plt.plot(predictions, color='red', label='Predicted')
63 plt.legend(fontsize = 20)
64 plt.show()
```