

▼ Phase 4 Recommendation Modeling

Authors: Fennec Nighintgale, Matthew Lipman

In this notebook we will be using the CRISP-DM process to walk through using Faiss K-means algorithm and naive Bayes to predict sentiment on space travel related tweets and YouTube comments.

```
1 # !pip install faiss-gpu
2 # !pip install vaderSentiment
```

```
1 # from google.colab import drive
2 # drive.mount('/content/drive')
```

```
1 import faiss
2 import math
3 import pandas as pd
4 import numpy as np
5 from datetime import datetime
```

```
1 import scipy
2 from scipy.linalg import norm
3 from scipy.sparse import csr_matrix, coo_matrix, hstack, vstack
4 from imblearn.combine import SMOTEENN
```

```
1 from textblob import TextBlob
2 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
```

```
1 from sklearn.decomposition import TruncatedSVD
2 from sklearn.naive_bayes import GaussianNB, BernoulliNB, ComplementNB
3 from sklearn.preprocessing import MaxAbsScaler
4
5 from sklearn.feature_extraction.text import TfidfVectorizer, \
6                                     CountVectorizer
7 from sklearn.metrics import silhouette_samples, confusion_matrix, \
8                             accuracy_score, recall_score, precision_score, \
9                             roc_curve, auc, roc_auc_score, confusion_matrix, \
10                            classification_report, f1_score, r2_score
11 from sklearn.model_selection import GridSearchCV, cross_val_score, \
12                                train_test_split, RandomizedSearchCV
```

```
1 from mlxtend.plotting import plot_confusion_matrix
```

```

1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 from matplotlib import cm

```

▼ Functions

```

1 def get_sentiment_an(df):
2     """
3     Given a pandas dataframe, this will shuffle it & get sentiment
4     analysis values from TextBlob & VaderSentiment & assign columns inplace.
5     Will return completed dataframe & dataframe columns we'll be using for
6     modeling as a coo_matrix.
7         Columns returned:
8     'polarity' : mean of text sentiment (TextBlob)
9     'neg': negative sentiment value (VaderSentiment)
10    'pos': positive sentiment value (VaderSentiment)
11    'neu': neutral sentiment value (VaderSentiment)
12    """
13    # shuffling our dataframe so data is no longer sorted
14    df = df.sample(frac = 1)
15    df = df.reset_index().rename(columns={'index':'id'})
16    # getting rid of all unnessecary columns
17    df = df[['favorite_count', 'repost_count', 'text', 'id']]
18    # feature generation using polarization
19    df['polarity'] = df['text'].apply(lambda x:
20                                     TextBlob(x).sentiment.polarity)
21    # feature generation using intensity analyzer
22    analyzer = SentimentIntensityAnalyzer()
23    df['pol'] = ''
24    for i in range(len(df)):
25        sentence = df.at[i, 'text']
26        df.at[i, 'pol'] = analyzer.polarity_scores(sentence)
27    for i in range(len(df)):
28        for h in ['neg', 'neu', 'pos']:
29            df.at[i, h] = float(df['pol'][i][h])
30    #turning our other variables into a matrix to combine our features
31    maxab = MaxAbsScaler()
32    drop = ['text', 'pol', 'id']
33    hstack = maxab.fit_transform(coo_matrix(df.drop(labels = drop, axis=1)))
34    return df, hstack

```

```

1 def word_transform(df):
2     """
3     Pass in a pandas dataframe and this will return a coo_matrix
4     featuring the TFIDF vectorization of your text, with a maximum
5     of 100 features of single words and 100 features of bigrams.
6     """
7     sing = TfidfVectorizer(max_features=100).fit_transform(df['text'])

```

```

8  bi = TfidfVectorizer(ngram_range = (2, 2),
9                      max_features=100).fit_transform(df['text'])
10 vect = coo_matrix(np.append(bi.toarray(), sing.toarray(), axis=1))
11 return vect

```

```

1 def svd_pca(vect, haystack):
2     """
3     Pass in your coo_matrix of your word vectors & the scaled
4     features from your dataframe this will perform Principal Component
5     Analysis in order to keep your most relevant 100 features, calculated
6     in chunks to avoid RAM issues. It will return a numpy
7     array containing all of the data for your model.
8     """
9     svd = TruncatedSVD(n_components=95, n_iter=250000,
10                      random_state=0, algorithm='arpack').fit(vect)
11     # splitting it so we never keep it all in RAM at once as an array
12     splits = []
13     shape = vect.shape[1]
14     tbs = vect
15     # since you cant directly index a matrix the easiest way is train_test
16     # chose 50000 since my computer can handle it
17     while shape > 50000:
18         shape -= 50000
19         # first, run through the whole matrix
20         split1, split2 = train_test_split(tbs, random_state=0,
21                                         train_size=shape)
22         #all of our leftover rows (not in the 50k) are in split2
23         tbs = split2
24         # put the chunk of our matrix into a list to iterate over later
25         splits.append(split1)
26     # append the remainder (whatever was left < 50k) to the end of our list
27     splits.append(tbs)
28     # start a new matrix by fit transforming our first chunk
29     init = coo_matrix(svd.fit_transform(splits[0]))
30     # start a loop that will fit transform and stack matrices together
31     for item in splits[1:]:
32         # fit transform will return it as a dense array
33         init = vstack([init, coo_matrix(svd.fit_transform(item))])
34     # combine our text matrix & our repost/favorite matrix
35     return np.append(haystack.toarray(),
36                    init.toarray(), axis=1).astype(np.float32)

```

```

1 def vect_and_stck(dataframe, first_run):
2     """
3     Pass in a DataFrame & this function will preprocess and return your
4     dataframe & all of the columns necessary to model as an array
5     The first step is to get sentiment analysis scores,
6     which is only necessary the first time you run this, otherwise
7     they can be turned to False to save time,
8     sentiment analysis is followed by TDIDF vectorization on

```

```

9  single words & bigrams, and lastly it ends with PCA to reduce
10 our dimensions to a reasonable amount. Make drop_clust True to
11 return a copy of your dataframe with no cluster column before
12 rerunning your model on it.
13 """
14 if first_run:
15     # getting sentiment analysis columns the first time
16     dataframe, haystack = get_sentiment_an(dataframe)
17 else:
18     # keeping only of the columns we need
19     clmn = ['neg', 'neu', 'pos', 'polarity',
20            'repost_count', 'favorite_count']
21     haystack = coo_matrix(dataframe[clmn])
22 # word vectorizing
23 vect = word_transform(dataframe)
24 # return our matrix with PCA done to reduce
25 # dimensionality to 100 features as was recommended
26 # for PCA in the documentation
27 return svd_pca(vect, haystack), dataframe

1 def plt_elbow(matrix, df, display):
2     """
3     Plots your clusters from k 1-15 and returns the best one as an int,
4     if display is true it will display the elbow plot, if not it will
5     just return your value.
6     """
7     max = 15 if display else 5
8     K = range(1,max+1)
9     # start numpy arrays to store results in
10    inrt = np.zeros(max)
11    diff = np.zeros(max)
12    diff2 = np.zeros(max)
13    diff3 = np.zeros(max)
14    for k in K:
15        kmeans, predictions, not_k = clusterizer(matrix, df, False, k=k)
16        inrt[k - 1] = kmeans.obj[-1]
17        # first difference
18        if k > 1:
19            diff[k - 1] = inrt[k - 1] - inrt[k - 2]
20        # second difference
21        elif k > 2:
22            diff2[k - 1] = diff[k - 1] - diff[k - 2]
23        # third difference
24        elif k > 3:
25            diff3[k - 1] = diff2[k - 1] - diff2[k - 2]
26    # use differences & numpy argmin to determine best cluster
27    elbow = np.argmin(diff3[3:]) + 3
28    if display:
29        print(f'Elbow {str(elbow)}')
30        plt.plot(K, inertias, 'b*-')
31        plt.plot(K[elbow-1], inertias[elbow-1], marker='o',

```

```

32         markersize=12, markeredgewidth=2, markeredgewidth='r')
33     plt.ylabel('Inertia')
34     plt.xlabel('K')
35     plt.show()
36     return int(elbow)

1 def chunked_samples(dataframe, stack, fit_predict, drop):
2     """
3     Given your dataframe, feature array/matrix & predictions
4     Computes Silhouette Scores using sklearn and chunking
5     the data into sets of 75000 and returns a deep copy of
6     your dataframe with the silhouette scores in the column
7     'sil' Scores will not be perfect as this function chunks them, and
8     therefore cannot get complete pairwise distances
9     but without chunking silhouette scores cannot be run on
10    any computer or service we can find and it will be pretty close.
11    Use drop = True to calculate the z-score of your silhouette scores
12    and get rid of the bottom 10%, assuming they're likely to be outliers.
13    """
14    start = 0
15    clust = 75000
16    # since this is about the most my computer can handle we'll
17    # stick to a cap of 75,000
18    for i in range(math.ceil(stack.shape[0]/75000)):
19        # gathering our chunks
20        df = dataframe[start:start+clust].reset_index(drop=True).copy()
21        fp = np.ravel(fit_predict[start:start+clust][0:])
22        # making sure our dtypes are correct
23        if type(stack) != np.ndarray:
24            hstk = stack.toarray()[start:start+clust][0:]
25        else:
26            hstk = stack[start:start+clust][0:]
27        # getting our silhouette comments
28        df['sil'] = silhouette_samples(hstk, fp)
29        # we want to return our dataframe at the end so
30        # on our first iteration we'll start with a chunk
31        # and append to it every future iteration
32        head = df if i == 0 else head.append(df)
33        if drop:
34            # if drop = True, this will calculate and drop the
35            # lowest ~10th percentile of our silhouette scores
36            # assuming they're likely to be outliers
37            mean = head['sil'].mean()
38            std = head['sil'].std()
39            for j in head.index:
40                head.at[j, 'silstd'] = (head.at[j, 'sil'] - mean)/std
41            head = head.loc[head['silstd'] > -1.28].drop('silstd', 1)
42            head.reset_index(drop=True, inplace=True)
43            start += 75000
44    return head

```

```

1 def quick_silhouette(dataframe, stack, predict):
2     """
3     Pass in your dataframe, your feature array/matrix,
4     & your predictions, and this will use chunked_samples
5     to get your score & try to quickly form a matplotlib visual
6     Scores won't be 100% accurate as you cannot get pairwise distances
7     with chunking, but they use a sample size of 75,000 so they should
8     be close, and the decrease in accuracy is well worth the 6-7 hours
9     in speed it saves
10    """
11    #first we calculate our silhouette scores
12    dataframe = chunked_samples(dataframe, stack, predict, False)
13    # convert our predictions and their scores to numpy arrays
14    predictions = dataframe.cluster.to_numpy()
15    silhouette_vals = dataframe.sil.to_numpy()
16    # get our labels
17    labels = np.unique(predictions)
18    # start plotting
19    y_ax_lower, y_ax_upper = 0, 0
20    yticks = []
21    # go through our labels and assign our predictions & values to them
22    for i, c in enumerate(labels):
23        c_silhouette_vals = silhouette_vals[predictions == c]
24        # sort to keeo our clusters together
25        c_silhouette_vals.sort()
26        # change plot size & ticks accordingly
27        y_ax_upper += len(c_silhouette_vals)
28        yticks.append((y_ax_lower + y_ax_upper) / 2.)
29        # get colors & plot!
30        color = cm.jet(float(i) / labels.shape[0])
31        plt.barh(range(y_ax_lower, y_ax_upper), c_silhouette_vals,
32                height=1.0, color=color)
33        y_ax_lower += len(c_silhouette_vals)
34    # find our average and visualize it
35    plt.axvline(np.mean(silhouette_vals), color="red", linestyle="--")
36    # fix up our labels and ticks
37    plt.yticks(yticks, labels + 1)
38    plt.ylabel('Cluster')
39    plt.xlabel('Silhouette Coefficient')
40    plt.show()

1 def get_imb_clstr(dataframe):
2     """
3     Pass in your dataframe and get the cluster with the most values in it
4     returned as an interger
5     """
6     clstr = dataframe['cluster'].value_counts()
7     # largest cluster gets returned at the first index
8     clstr = int(clstr.to_frame().reset_index()['index'][0])
9     return clstr

```

```

1 def reset_clusters(new_predict, dataframe):
2     """
3     Pass in your new predictions and your dataframe
4     and it will renumber all of your clusters.
5     Be sure to use this before making visuals, if you don't
6     clusters that have been reassigned will just appear to be missing.
7     Edits your dataframe inplace. This is intentionally slower in
8     the hopes of not crashing my ram by iterating over each line one by one
9     """
10    # get a list of our starting clusters
11    _clusters = [int(x) for x in np.unique(new_predict)]
12    # generate a list of new clusters
13    _range = list(range(len(_clusters)))
14    to_rename = {}
15    # generate a dictionary with key value pairs of new : old clusters
16    for num in _range:
17        to_rename[_clusters[num]] = num + 1
18    # go through the dataframe and edit our clusters
19    for i in range(len(dataframe)):
20        dataframe.at[i, 'cluster'] = to_rename[dataframe.at[i, 'cluster']]

```

```

1 def add_new_clst(df, dataframe, k, num_clust):
2     """
3     Pass in your new predictions and your dataframe
4     and it will reassign the cluster you modeled to its new
5     cluster assignments. Edits dataframe inplace.
6     """
7     # when generating new clusters after your first modeling
8     # you'll need to assign them to new unique cluster numbers
9     for i in df.index:
10        # use the id's we assigned earlier to match up our clusters
11        _id = df.at[i, 'id']
12        index = dataframe.loc[dataframe['id'] == _id].index[0]
13        # make sure we always get unique clusters, and that numbers being
14        # reassigned to the same cluster always turn out the same
15        n_cluster = (df.at[i, 'cluster'] + 2) * 10 + num_clust*2
16        # reassign our clusters
17        dataframe.at[index, 'cluster'] = n_cluster

```

```

1 def clusterizer(matrix, dataframe, calc_k, k=3, display_elbow=False):
2     """
3     Pass in your feature array/matrix and your dataframe
4     If you would like to have k be calculated for you using the
5     elbow method with a maximum number of 5 clusters, do so by
6     setting calc_k = True, otherwise, k will be 3
7     """
8     if calc_k:
9         if display_elbow:
10            k = plt.elbow(matrix, dataframe, True)

```

```

10     k = plt_elbow(matrix, dataframe, True)
11     else:
12         k = plt_elbow(matrix, dataframe, False)
13     # making sure our dtypes always line up
14     if type(matrix) != np.ndarray:
15         matrix = matrix.toarray()
16     # we dont want our clusters to be imbalanced, so we want them
17     # to be roughly equal to the number of points/ number of clusters
18     shape = matrix.shape[0]
19     # start our k-means
20     kmeans = faiss.Kmeans(d = matrix.shape[1], k = k, nredo = 100,
21                           update_index = True, seed = 42, verbose=True,
22                           max_points_per_centroid = math.ceil(shape/k),
23                           niter = 20)
24     # train our k means
25     kmeans.train(matrix)
26     # get our predictions
27     predict = kmeans.index.search(matrix, 1)[1]
28     # assign them inplace to our dataframe
29     dataframe['cluster'] = predict
30     return kmeans, predict, k

1 def get_slice(df, cluster):
2     '''
3     Pass in your dataframe and return just the values predicted to be in the
4     largest cluster.
5     '''
6     return df.loc[df['cluster'] == cluster].reset_index(drop=True)

1 def showconfusionmatrix(y_t, y_hat_t, title):
2     """
3     Plots confusion matrix for provided y_train, y_hat_train
4     OR y_test, y_hat_test using mlxtend
5     """
6     # assign parameters to dictionary so line wont exceed 78 char
7     p = {conf_mat:confusion_matrix(y_t, y_hat_t), colorbar:True,
8          show_absolute:True, show_normed:True, cmap:'jet'}
9     # gather and plot confusion matrix
10    fig, ax = plot_confusion_matrix(**p)
11    # set title and axis names
12    ax.set_title(f'{title} Model')
13    ax.set_ylabel('Actual Data')
14    ax.set_xlabel('Predicted Data')
15    plt.show()

1 def printscores(y_test, y_hat_test, X_test, y_train,
2                 y_hat_train, X_train, model):
3     """Uses sklearn.metrics to compute testing and training scores
4     Results include: Accuracy, Precision, Recall, F1, Residual Counts,
5     Residuals, CV Accuracy With Standard Deviation & R^2"""

```



```

6 print('_____')
7 print('Training Accuracy Score: ' +
8       str(round(accuracy_score(y_train, y_hat_train)* 100, 4)))
9 train_precision_scores = []
10 for item in precision_score(y_test, y_hat_test, average=None):
11     train_precision_scores.append(round(item, 2))
12 print('Training Precision Scores: ' + str(train_precision_scores))
13 train_recall_scores = []
14 for item in recall_score(y_test, y_hat_test, average=None):
15     train_recall_scores.append(round(item, 2))
16 print('Training Recall Scores: ' + str(train_recall_scores))
17 train_f1_scores = []
18 for item in f1_score(y_test, y_hat_test, average=None):
19     train_f1_scores.append(round(item, 2))
20 print('Training F1 Scores: ' + str(train_f1_scores))
21 print('_____')
22 trainresiduals = np.abs(y_train - y_hat_train)
23 print('Training residual counts:')
24 print(str(pd.Series(trainresiduals).value_counts())[:-29])
25 print('_____')
26 print('Testing Accuracy Score: ' +
27       str(round(accuracy_score(y_test, y_hat_test) * 100, 4)))
28 test_precision_scores = []
29 for item in precision_score(y_test, y_hat_test, average=None):
30     test_precision_scores.append(round(item, 2))
31 print('Testing Precision Scores: ' + str(test_precision_scores))
32 test_recall_scores = []
33 for item in recall_score(y_test, y_hat_test, average=None):
34     test_recall_scores.append(round(item, 2))
35 print('Testing Recall Scores: ' + str(test_recall_scores))
36 test_f1_scores = []
37 for item in f1_score(y_test, y_hat_test, average=None):
38     test_f1_scores.append(round(item, 2))
39 print('Testing F1 Scores: ' + str(test_f1_scores))
40 print('_____')
41 testresiduals = np.abs(y_test - y_hat_test)
42 print('Testing residual counts:')
43 print(str(pd.Series(testresiduals).value_counts())[:-29])
44 print('_____')
45 r2_scoretst = r2_score(y_test, y_hat_test)
46 print(f'Testing R2 Score: {round(r2_scoretst, 4)}')
47 r2_scoretrn = r2_score(y_train, y_hat_train)
48 print(f'Training R2 Score: {round(r2_scoretrn, 4)}')
49 print('_____')
50 print('Cross validated model accuracy:')
51 scores = cross_val_score(model, X_test, y_test, cv=10)
52 mean = round(scores.mean(), 4)
53 std = round(scores.std())
54 print(f'Testing: {mean} with a standard deviation of {std}')
55 scores = cross_val_score(model, X_train, y_train, cv=10)
56 mean = round(scores.mean(), 4)
57 std = round(scores.std())

```

```

58     print(f'Training: {mean} with a standard deviation of {std}')
59     print('_____')

1 def printreports(y_test, y_hat_test, y_train, y_hat_train):
2     """
3     Provided all input & predicted y values, prints classification
4     report for training and testing data right next to each other"""
5     print('_____')
6     print('          Testing Report:')
7     report1 = classification_report(y_test, y_hat_test)
8     print(report1)
9     print('_____')
10    print('          Training Report:')
11    report2 = classification_report(y_train, y_hat_train)
12    print(report2)
13    print('_____')

```

▼ Unsupervised Model Setup

First, we begin by opening up our data and doing some preprocessing, our data has already been cleaned and tokenized and lemmatized, so here we mostly focus on feature generation and dimensionality reduction

```

1 # open data, low_memory = False because we have some mixed dtype columns,
2 # they aren't one's we use anyways so not worried about it
3 df = pd.read_csv("Documents/flatiron/hospitality project/data/final_clean_6_word.csv",
4                 low_memory=False)

1 # tokenize and preprocess our data, turning it into a matrix
2 matrix, df = vect_and_stck(df, True)
3 # taking an initial count to see how much we lose by the end
4 start = [len(df), datetime.now()]

```

▼ First Unsupervised Model

▼ Model Evaluation

Looking at our initial model, our clusters aren't particularly well defined, and even with multiple attempts at what k should be, they're not showing any substantial improvement

```

1 # starting a baseline model for our initial model

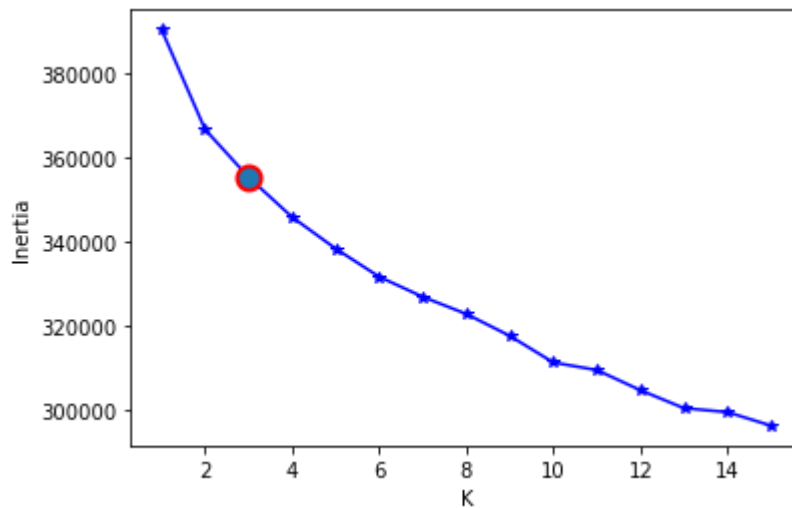
```

```

2 first_model = df.copy()
3 first_means, predict, k = clusterizer(matrix, first_model,
4                                     True, display_elbow=True)

```

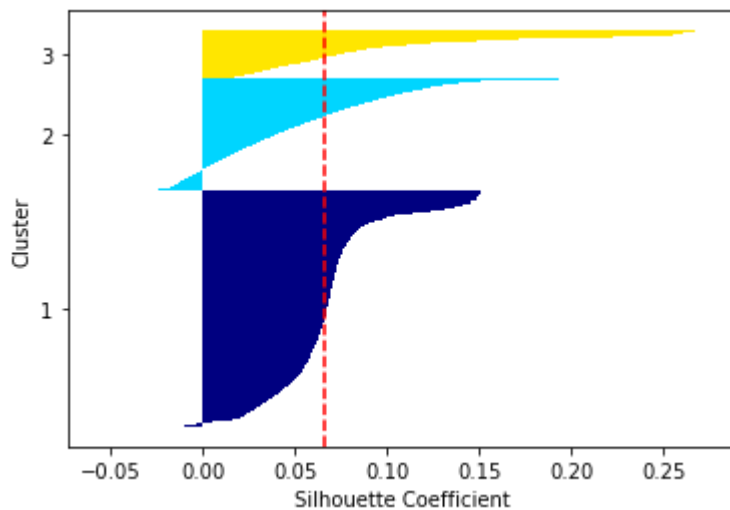
Elbow 3



```

1 # initial silhouette scores very low
2 quick_silhouette(first_model, matrix, predict)

```



▼ Unsupervised Model Refinement

```

1 # start our model & make our initial predictions
2 second_means, predict, k = clusterizer(matrix, df, False)

```

```

1 num_clust = 2
2 cl = get_imb_clstr(df)
3 # if our largest cluster is too large, we break it down into more clusters
4 while len(get_slice(df, cl)) >= len(df)/num_clust-1 and num_clust <= 9:
5     # dropping outliers and shuffling our data

```

```

5 # dropping outliers and shuffling our data
6 df = chunked_samples(df, matrix, predict, True)
7 df = df.sample(frac=1).reset_index(drop=True)
8 # processing our data without outliers
9 matrix, df = vect_and_stck(df, False)
10 # re-modeling our data with less outliers
11 btrr_mn, btrr_prd, k = clusterizer(matrix, df, False, 3, False)
12 # slicing our dataframe to get disproportionate cluster
13 cl = get_imb_clstr(df)
14 cut_matrix, cut_df = vect_and_stck(get_slice(df, cl), False)
15 # plugging our new k value into a new model optimized for our cluster
16 clst_mn, clst_prd, k = clusterizer(cut_matrix, cut_df, True, 3, False)
17 # updating our original dataframe with our new clusters clusters
18 add_new_clst(cut_df, df, k, num_clust)
19 # updating our predictions
20 predict = df['cluster'].to_numpy()
21 # updating our largest cluster
22 cl = get_imb_clstr(df)
23 # keeping a count of clusters
24 num_clust += k-1

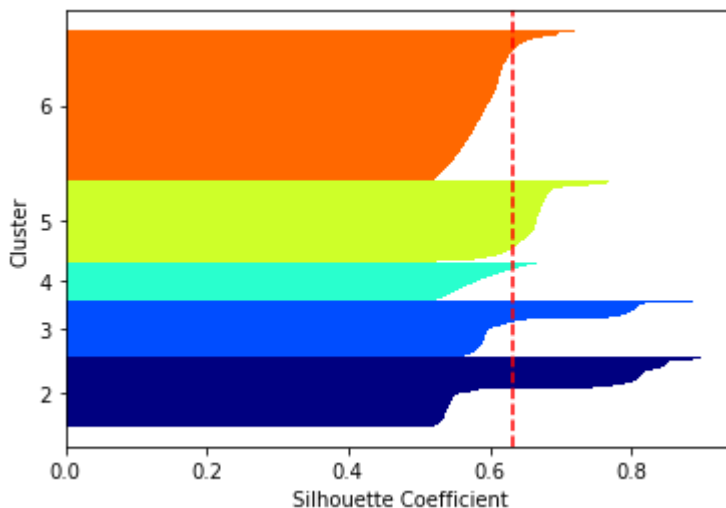
```

▼ Model Evaluation

```

1 # grab our final clusters & plot our silhouettes
2 predict = df['cluster'].to_numpy()
3 reset_clusters(predict, df)
4 quick_silhouette(df, matrix, predict)

```



```

1 # calculating how long our modeling took & how much data was lost
2 end = [len(df), datetime.now()]
3 time = end[1]-start[1]
4 print(f' {round(end[0]/start[0], 4)}% Unclassifiable')
5 print(f' Model took {round(time.total_seconds()/60, 2)} min to complete')

```

```
0.5903% Unclassifiable
Model took 126.75 min to complete
```

Our model showed substantial improvement after iterating through & reclustering, losing minimal amounts of data. With our mean silhouette scores looking good and in a reasonable range, we move on to taking a look at visualizing our clusters and discerning the differences between them, which you can find in our notebook /wordclouds.ipynb

▼ Bayesian Modeling

We chose bayesian modeling because it is said to work especially well for natural language processing, 2 of the methods available from Sklearn were not compatible with our data, but the other 2 are tested & refined below.

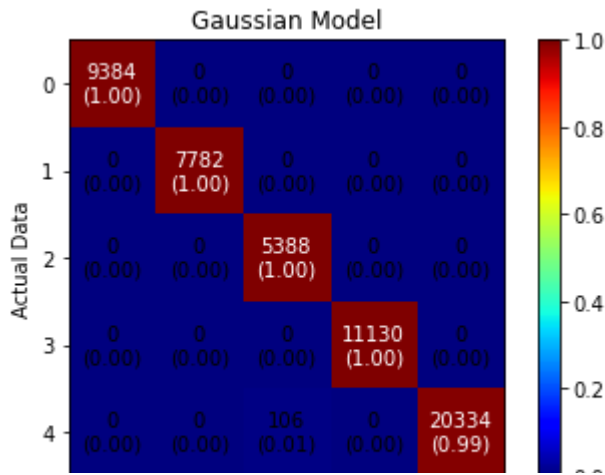
```
1 # since we have a bit of a class imbalance that we want to fix
2 # we're going to do a combination of upsampling and unsampling
3 # to even them out without overdoing it
4 smote = SMOTEENN()
5
6 X_train, X_test, y_train, y_test = train_test_split(matrix,
7                                                     predict,
8                                                     random_state=42)
9 # we're only applying this to our training data, so we don't leak
10 # anything over into our test that could cause overfitting
11 X_train, y_train = smote.fit_resample(X_train, y_train)
```

▼ Testing models

There are our baseline models for making predictions off of our data

```
1 # the baseline model is actually incredibly close to our results
2 # while tuning it we discovered a few variables that did have considerable
3 # effect, so here is one to help showcase the difference
4 gau_model = GaussianNB(var_smoothing=1)
5 y_hat_test = gau_model.fit(X_train, y_train).predict(X_test)
6 y_hat_train = gau_model.predict(X_train)

1 showconfusionmatrix(y_test, y_hat_test, 'Gaussian')
2 # certainly too good to be true
```



```
1 print('          Gaussian Model')
2 print('          _____')
3 printscores(y_test, y_hat_test, X_test, y_train, y_hat_train, X_train, gau_model)
```

Gaussian Model

Training Accuracy Score: 99.9006
 Training Precision Scores: [1.0, 1.0, 0.98, 1.0, 1.0]
 Training Recall Scores: [1.0, 1.0, 1.0, 1.0, 0.99]
 Training F1 Scores: [1.0, 1.0, 0.99, 1.0, 1.0]

Training residual counts:
 0 30

Testing Accuracy Score: 99.8042
 Testing Precision Scores: [1.0, 1.0, 0.98, 1.0, 1.0]
 Testing Recall Scores: [1.0, 1.0, 1.0, 1.0, 0.99]
 Testing F1 Scores: [1.0, 1.0, 0.99, 1.0, 1.0]

Testing residual counts:
 0

Testing R2 Score: 0.9966
 Training R2 Score: 0.998

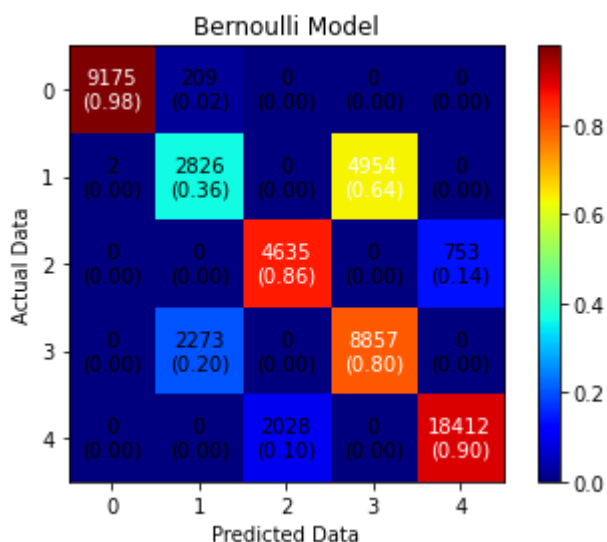
Cross validated model accuracy:
 Testing: 0.7858 with a standard deviation of 0.0
 Training: 0.999 with a standard deviation of 0.0

As you can see from our scores, this model is way too good to be true and is not producing a consistent output despite its high scores

```
1 bern_model = BernoulliNB()
2 y_hat_test = bern_model.fit(X_train, y_train).predict(X_test)
3 y_hat_train = bern_model.predict(X_train)
```

```
1 print('BernoulliNB Model Accuracy: %.3f' % bern_model.score(X_test, y_test))
```

```
1 showconfusionmatrix(y_test, y_hat_test, 'Bernoulli')
```



```
1 print('          Bernoulli Model')
2 print('          _____')
3 printreports(y_test, y_hat_test, y_train, y_hat_train)
```

Bernoulli Model					
Testing Report:					
	precision	recall	f1-score	support	
1	1.00	0.98	0.99	9384	
2	0.53	0.36	0.43	7782	
3	0.70	0.86	0.77	5388	
4	0.64	0.80	0.71	11130	
5	0.96	0.90	0.93	20440	
accuracy			0.81	54124	
macro avg	0.77	0.78	0.77	54124	
weighted avg	0.81	0.81	0.81	54124	

Training Report:					
	precision	recall	f1-score	support	
1	1.00	0.98	0.99	61186	
2	0.62	0.37	0.46	61186	
3	0.91	0.87	0.89	60348	
4	0.56	0.80	0.66	61185	
5	0.88	0.91	0.89	58953	
accuracy			0.79	302858	
macro avg	0.79	0.79	0.78	302858	
weighted avg	0.79	0.79	0.78	302858	

```

1 print('                Bernoulli Model')
2 print('                _____')
3 print(scores(y_test, y_hat_test, X_test, y_train, y_hat_train, X_train, bern_model))

```

```

Bernoulli Model
_____

Training Accuracy Score: 78.5124
Training Precision Scores: [1.0, 0.53, 0.7, 0.64, 0.96]
Training Recall Scores: [0.98, 0.36, 0.86, 0.8, 0.9]
Training F1 Scores: [0.99, 0.43, 0.77, 0.71, 0.93]

Training residual counts:
0      237781
2         6

Testing Accuracy Score: 81.1193
Testing Precision Scores: [1.0, 0.53, 0.7, 0.64, 0.96]
Testing Recall Scores: [0.98, 0.36, 0.86, 0.8, 0.9]
Testing F1 Scores: [0.99, 0.43, 0.77, 0.71, 0.93]

Testing residual counts:
0      43905
2

Testing R2 Score: 0.6812
Training R2 Score: 0.5753

Cross validated model accuracy:
Testing: 0.8185 with a standard deviation of 0.0
Training: 0.7848 with a standard deviation of 0.0

```

This model is only slightly better than random chance overall, and in some clusters, it's even worse! Certainly not the route we want to be going.

▼ Final Model

Here, with a slight parameter tune, we see our final model come into realization with solid, consistent results.

```

1 final_model = GaussianNB(var_smoothing= 1e-200)
2 y_hat_test = final_model.fit(X_train, y_train).predict(X_test)
3 y_hat_train = final_model.predict(X_train)

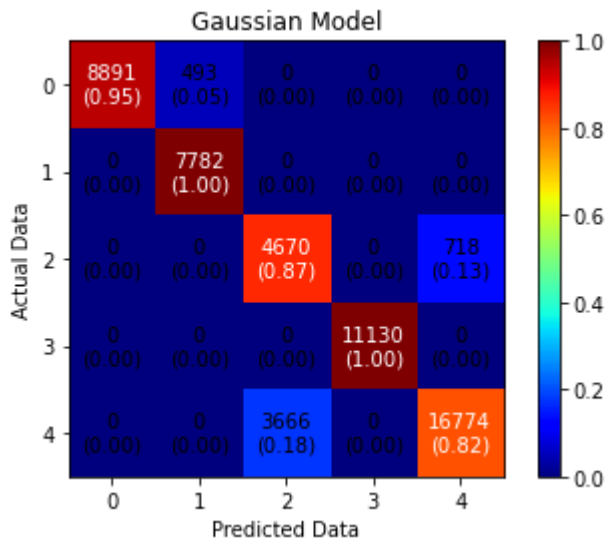
```

Evaluation

```

1 showconfusionmatrix(y_test, y_hat_test, 'Gaussian')

```

```
1 print('          Final Gaussian Model')
2 print('          _____')
3 printreports(y_test, y_hat_test, y_train, y_hat_train)
```

Final Gaussian Model

Testing Report:

	precision	recall	f1-score	support
1	1.00	0.95	0.97	9384
2	0.94	1.00	0.97	7782
3	0.56	0.87	0.68	5388
4	1.00	1.00	1.00	11130
5	0.96	0.82	0.88	20440
accuracy			0.91	54124
macro avg	0.89	0.93	0.90	54124
weighted avg	0.93	0.91	0.92	54124

Training Report:

	precision	recall	f1-score	support
1	1.00	0.95	0.97	61186
2	0.95	1.00	0.98	61186
3	0.84	0.86	0.85	60348
4	1.00	1.00	1.00	61185
5	0.86	0.83	0.84	58953
accuracy			0.93	302858
macro avg	0.93	0.93	0.93	302858
weighted avg	0.93	0.93	0.93	302858

```
1 print('          Final Gaussian Model')
```

```

2 print('_____')
3 print(scores(y_test, y_hat_test, X_test,
4             y_train, y_hat_train, X_train, final_model)

          Final Gaussian Model
          _____

Training Accuracy Score: 93.0278
Training Precision Scores: [1.0, 0.94, 0.56, 1.0, 0.96]
Training Recall Scores: [0.95, 1.0, 0.87, 1.0, 0.82]
Training F1 Scores: [0.97, 0.97, 0.68, 1.0, 0.88]

Training residual counts:
0      281742
2        1

Testing Accuracy Score: 90.9892
Testing Precision Scores: [1.0, 0.94, 0.56, 1.0, 0.96]
Testing Recall Scores: [0.95, 1.0, 0.87, 1.0, 0.82]
Testing F1 Scores: [0.97, 0.97, 0.68, 1.0, 0.88]

Testing residual counts:
0      49247
2

Testing R2 Score: 0.8572
Training R2 Score: 0.875

Cross validated model accuracy:
Testing: 0.9145 with a standard deviation of 0.0
Training: 0.93 with a standard deviation of 0.0

```

In conclusion,

Even though we failed to find **strong** sentiment between our clusters, our model itself found clear enough differences between them to make meaningful and replicable predictions and the knowledge we came away from this with has given us great inspiration to make business recommendations that will help your company move forward in the growing space hospitality industry.

Our recommendations are:

1. Leverage use of the word 'Space' Space is the most commonly found word within the discourse, with that in mind, tagging along a positive sentiment will continue to build a positive arena around space and space travel & extend your reach to the largest audience.
2. Evoke emotion and drive the conversation. Customers think with their hearts, and by utilizing words like "love", which are some of the most popular non-science words we found, you are likely to create a positive sentiment around your product. That positive emotional response adds tremendous value to your marketing capabilities.

3. Partner with the best Institutions like NASA and SpaceX not only have scientific headway, but they have built reputations. If you want to get people talking about you, working with businesses they already support can help extend your reach, as even with the variation between our clusters, the big names in them stay rather consistent.

We look forward to working with you and appreciate your time. Thank you.

