

Load X

Generate text vectors using the comments. Each comment is also tagged with the listing_id it is associated with.

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import nltk
        4 import re
        5 import matplotlib.pyplot as plt
        6 %matplotlib inline
```

```
In [3]: 1 reviews = pd.read_csv('data/reviews.csv', usecols=['listing_id', 'id', 'comm
```

```
In [4]: 1 reviews.head()
```

```
Out[4]:
```

	listing_id	id	comments
0	2454	406718547	The host canceled this reservation 23 days bef...
1	2539	55688172	Great host
2	2539	97474898	Nice room for the price. Great neighborhood. J...
3	2539	105340344	Very nice apt. New remodeled.
4	2539	133131670	Great place to stay for a while. John is a gre...

```
In [5]: 1 # Remove NAs
        2
        3 reviews.dropna(inplace=True)
```

```
In [6]: 1 reviews.shape
```

```
Out[6]: (1114779, 3)
```

Alternative X: English only reviews

Many reviews are in non-English languages which may impact model performance. I will filter those out using langdetect .

```
In [98]: 1 def detect_lang(text):
        2     from langdetect import detect
        3     try:
        4         return detect(text)
        5     except:
        6         return "unknown"
```

```
In [99]: 1 reviews['lang'] = reviews.comments.apply(detect_lang)
```

```
In [104]: 1 reviews[reviews.lang=='en'].to_csv('reviews_eng.csv')
```

Load y

Load latest price for each listing.

```
In [7]: 1 prices = pd.read_csv('data/cleaned_listings.csv', usecols=['id', 'price'])
```

```
In [8]: 1 prices.columns = ['listing_id', 'price']
```

```
In [9]: 1 prices.head()
```

```
Out[9]:
```

	listing_id	price
0	2454	137.0
1	2539	149.0
2	2595	225.0
3	3330	70.0
4	3647	150.0

```
In [39]: 1 # prices['price'] = prices.price.apply(lambda x: np.NaN if x == 'NaN' else f
```

```
In [10]: 1 prices.shape
```

```
Out[10]: (50220, 2)
```

```
In [11]: 1 reviews_price = pd.merge(reviews, prices, how='left', on='listing_id').dropn
```

```
In [12]: 1 reviews_price.head()
```

```
Out[12]:
```

	listing_id	id	comments	price
0	2454	406718547	The host canceled this reservation 23 days bef...	137.0
1	2539	55688172	Great host	149.0
2	2539	97474898	Nice room for the price. Great neighborhood. J...	149.0
3	2539	105340344	Very nice apt. New remodeled.	149.0
4	2539	133131670	Great place to stay for a while. John is a gre...	149.0

```
In [13]: 1 reviews_price.shape
```

```
Out[13]: (1114749, 4)
```

```
In [14]: 1 reviews_price.to_csv('reviews_price.csv', index=False)
```

Load Alternative y

With only english reviews

```
In [109]: 1 reviews_eng = pd.read_csv('reviews_eng.csv')
```

```
In [111]: 1 reviews_eng.dtypes
```

```
Out[111]: id          int64  
listing_id    int64  
comments      object  
lang          object  
dtype: object
```

```
In [114]: 1 prices = pd.read_csv('price.csv')  
2 prices.columns = ['listing_id', 'price']
```

```
In [115]: 1 reviews_price_alt = pd.merge(reviews_eng, prices, how='left', on='listing_id')
```

```
In [123]: 1 reviews_price_alt.drop(['lang'], axis=1).to_csv("reviews_price_alt.csv", ind
```

To save memory, do not run above code. Instead, go straight to **Vectorizing text** and load X or alternative X as you like.

Vectorizing text

Need to turn each review into an array with equal length.

First try vectorizing by word countTFIDF with skip-gram

Consider using custom tokenizer due to memory constraints - default tokenizers use a lot of memory, maybe because they are unable to filter words. However, this flexibility could lead to slower code and/or inability to filter punctuations. Custom tokenizer not optimized for speed. Code is left here for future tweaks.

```
In [14]: 1 # from nltk import download  
2 # from nltk.corpus import stopwords  
3 # download('stopwords')  
4 # stop_words = set(stopwords.words('english'))
```

```
[nltk_data] Downloading package stopwords to /home/twang/nltk_data...  
[nltk_data]   Unzipping corpora/stopwords.zip.
```

```
In [22]: 1 # def Tokenizer(str_input):
2 #     from re import sub
3 #     from nltk.stem.porter import PorterStemmer
4 #     from string import punctuation
5
6 #     # Remove punctuations
7 #     str_input = str_input.translate(None, punctuation)
8
9 #     # Remove any word containing a number (of which there are many)
10 #     words = sub(u'(?ui)\b[a-zA-Z0-9]*[0-9]+[a-zA-Z0-9]*\b', " ", str_inp
11
12 #     porter_stemmer=PorterStemmer()
13 #     words = [porter_stemmer.stem(word) for word in words if not word in st
14 #     return words
```

```
In [43]: 1 import pandas as pd
2 import numpy as np
3 reviews_price = pd.read_csv('reviews_price.csv')
```

```
In [44]: 1 reviews_price.head()
```

```
Out[44]:
```

	listing_id	id	comments	price
0	2515	198	Stephanie was a wonderful host! Her apartment ...	59.0
1	21456	29826	We had a delightful stay at Dana's fantastical...	140.0
2	21456	30680	Dana's place is charming, and very well-locate...	140.0
3	21456	32640	great stay, i would recommend her anytime...	140.0
4	21456	34234	Dana is a warm and welcoming host. We enjoyed...	140.0

```
In [15]: 1 from sklearn.decomposition import TruncatedSVD
2 from sklearn.pipeline import make_pipeline
3
4 from sklearn.feature_extraction.text import CountVectorizer
5 from sklearn.model_selection import train_test_split
6
7 X_train, X_test, y_train, y_test = train_test_split(reviews_price[['listing_
8                                                         random_state = 45)
```

```
In [16]: 1 X_train.to_csv('X_train.csv')
2 X_test.to_csv('X_test.csv')
3 pd.DataFrame(y_train, columns = ['price']).to_csv('y_train.csv')
4 pd.DataFrame(y_test, columns = ['price']).to_csv('y_test.csv')
```

```
In [17]: 1 del reviews_price
```

```

In [63]: 1 class count_preprocessing:
2
3     def __init__(self, svd_components = 50):
4         # import pandas as pd
5         # import numpy as np
6         from sklearn.decomposition import TruncatedSVD
7         from sklearn.feature_extraction.text import CountVectorizer
8         from sklearn.pipeline import make_pipeline
9
10        # Make instances
11        ## Used a token pattern to use only words without numbers
12        self.count = CountVectorizer(stop_words='english',\
13                                     lowercase=True,\
14                                     token_pattern=u'(?ui)\b[a-zA-Z]*[a-zA-
15                                     min_df=0.025)
16        # self.count = CountVectorizer(analyzer=Tokenizer, min_df = 0.1)
17        self.svd = TruncatedSVD(n_components=svd_components)
18        self.pipeline = make_pipeline(self.count, self.svd)
19
20    def fit_transform(self, text_data):
21        """
22        Fits and transforms text data pandas series into matrix of svd_compo
23        Input:
24            text_data: column or series of text data
25        Output:
26            A matrix with same number of rows as input and svd_components co
27        """
28        return(self.pipeline.fit_transform(text_data))
29
30    def transform(self, text_data):
31        """
32        Transforms new data into the right format
33        Input:
34            text_data: column or series of text data (for testing most likel
35        Output:
36            A matrix with same number of rows as input and svd_components co
37        """
38        return(self.pipeline.transform(text_data))

```

```

In [65]: 1 %%time
2
3 count_prep = count_preprocessing()
4
5 count_train = count_prep.fit_transform(X_train.comments)
6 count_test = count_prep.transform(X_test.comments)

```

CPU times: user 1min 4s, sys: 3.35 s, total: 1min 7s
 Wall time: 44.8 s

```
In [76]: 1 count_prep.count.get_feature_names()
```

```
Out[76]: ['access',  
          'accommodating',  
          'airbnb',  
          'amazing',  
          'amenities',  
          'apartment',  
          'appartement',  
          'area',  
          'arrival',  
          'arrived',  
          'available',  
          'away',  
          'awesome',  
          'bars',  
          'bathroom',  
          'beautiful',  
          'bed',  
          'bedroom',  
          'best',  
          ...]
```

After using only high-frequency words (by setting `min_df` to 0.05) we can see that the features make a lot of sense. One down side is that there is no good way to stem the words. For that we will have to use the custom tokenizer which may be much slower.

```
In [77]: 1 pd.DataFrame(count_train, index=X_train.index).to_csv('count_train.csv', ind
```

```
In [78]: 1 pd.DataFrame(count_test, index=X_test.index).to_csv('count_test.csv', index=
```

Slightly more sophisticated vectorization: TFIDF with skip-gram

```

In [18]: 1 class tfidf_preprocessing:
2
3     def __init__(self, svd_components = 50):
4         # import pandas as pd
5         # import numpy as np
6         from sklearn.decomposition import TruncatedSVD
7         from sklearn.feature_extraction.text import TfidfVectorizer
8         from sklearn.pipeline import make_pipeline
9
10        # Make instances
11        ## Used a token pattern to use only words without numbers
12        self.tfidf = TfidfVectorizer(stop_words='english', lowercase=True, \
13                                     token_pattern=u'(?ui)\b[a-zA-Z]*[a-zA-
14                                     max_df=0.9, \
15                                     min_df=0.025, \
16                                     ngram_range=(1,2))
17        self.svd = TruncatedSVD(n_components=svd_components)
18        self.pipeline = make_pipeline(self.tfidf, self.svd)
19
20    def fit_transform(self, text_data):
21        """
22        Fits and transforms text data pandas series into matrix of svd_compo
23        Input:
24            text_data: column or series of text data
25        Output:
26            A matrix with same number of rows as input and svd_components co
27        """
28        return(self.pipeline.fit_transform(text_data))
29
30    def transform(self, text_data):
31        """
32        Transforms new data into the right format
33        Input:
34            text_data: column or series of text data (for testing most likel
35        Output:
36            A matrix with same number of rows as input and svd_components co
37        """
38        return(self.pipeline.transform(text_data))

```

```

In [19]: 1 %%time
2
3 tfidf_prep = tfidf_preprocessing()
4
5 tfidf_train = tfidf_prep.fit_transform(X_train.comments)
6 tfidf_test = tfidf_prep.transform(X_test.comments)

```

CPU times: user 1min 37s, sys: 4.25 s, total: 1min 41s
 Wall time: 1min 18s

```
In [20]: 1 tfidf_prep.tfidf.get_feature_names()
```

```
Out[20]: ['access',  
          'accommodating',  
          'airbnb',  
          'amazing',  
          'amenities',  
          'apartment',  
          'apartment clean',  
          'appartement',  
          'area',  
          'arrival',  
          'arrived',  
          'available',  
          'away',  
          'awesome',  
          'bars',  
          'bathroom',  
          'beautiful',  
          'bed',  
          'bedroom',  
          ...]
```

```
In [21]: 1 %%time  
2 pd.DataFrame(tfidf_train, index=X_train.listing_id).to_csv('tfidf_train.csv')  
3 pd.DataFrame(tfidf_test, index=X_test.listing_id).to_csv('tfidf_test.csv', i
```

CPU times: user 50.5 s, sys: 1.04 s, total: 51.5 s

Wall time: 51.7 s

Topic Modeling

Best accomplished using tfidf vectors - more in-context meaning is preserved. Will explore a few clustering techniques.

```
In [1]: 1 # Load Data  
2 import numpy as np  
3 import pandas as pd  
4  
5 tfidf_train = pd.read_csv('tfidf_train.csv', index_col=0)  
6 tfidf_test = pd.read_csv('tfidf_test.csv', index_col=0)
```

```
In [2]: 1 from sklearn.cluster import KMeans
```

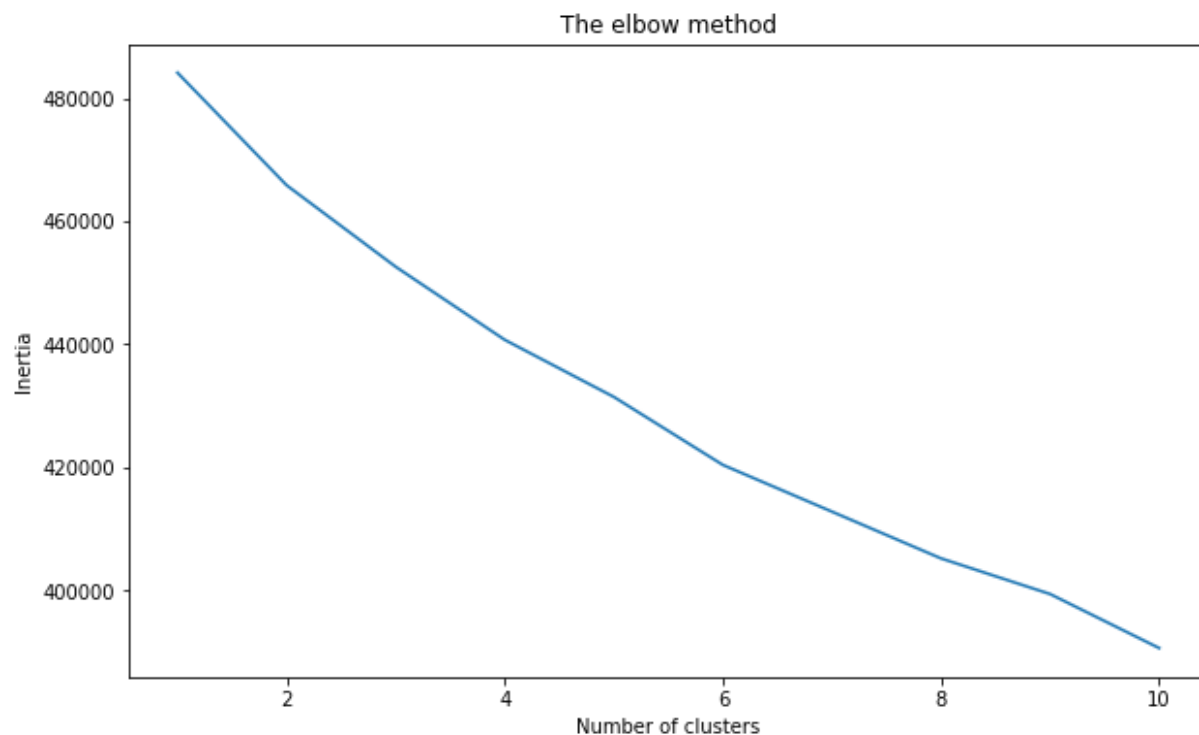


```
In [11]: 1 def cluster(array, random_state, n_clusters=4):
2         '''
3         Fits and predicts k-means clustering on "array"
4
5         Parameters
6         -----
7         array: A numpy array
8         random_state: Random seed, e.g. check_random_state(0)
9         n_clusters: The number of clusters. Default: 4
10
11        Returns
12        -----
13        A tuple (sklearn.KMeans, np.ndarray)
14        '''
15        from sklearn.cluster import KMeans
16        model = KMeans(n_clusters=n_clusters, random_state = random_state)
17        model.fit(array)
18        clusters = model.predict(array)
19
20        return model, clusters
21
22        #here we return fitted (model) and predicted (clusters) arrays as a tuple
```

```
In [12]: 1 def plot_inertia(array, start=1, end=10):
2         '''
3         Increase the number of clusters from "start" to "end" (inclusive).
4         Finds the inertia of k-means clustering for different k.
5         Plots inertia as a function of the number of clusters.
6
7
8         Parameters
9         -----
10        array: A numpy array.
11        start: An int. Default: 1
12        end: An int. Default: 10
13
14        Returns
15        -----
16        A matplotlib.Axes instance.
17        '''
18        from sklearn.utils import check_random_state
19
20        inertia = []
21        for i in range(start, end+1):
22            model, clusters = cluster(array, check_random_state(0), i)
23            inertia.append(model.inertia_)
24
25        x_axis = list(range(start, end+1))
26
27        fig, ax = plt.subplots(figsize=(10,6))
28
29        ax.set_title('The elbow method')
30        ax.set_ylabel('Inertia')
31        ax.set_xlabel('Number of clusters')
32        plt.plot(x_axis, inertia)
33
34
35        return ax
```

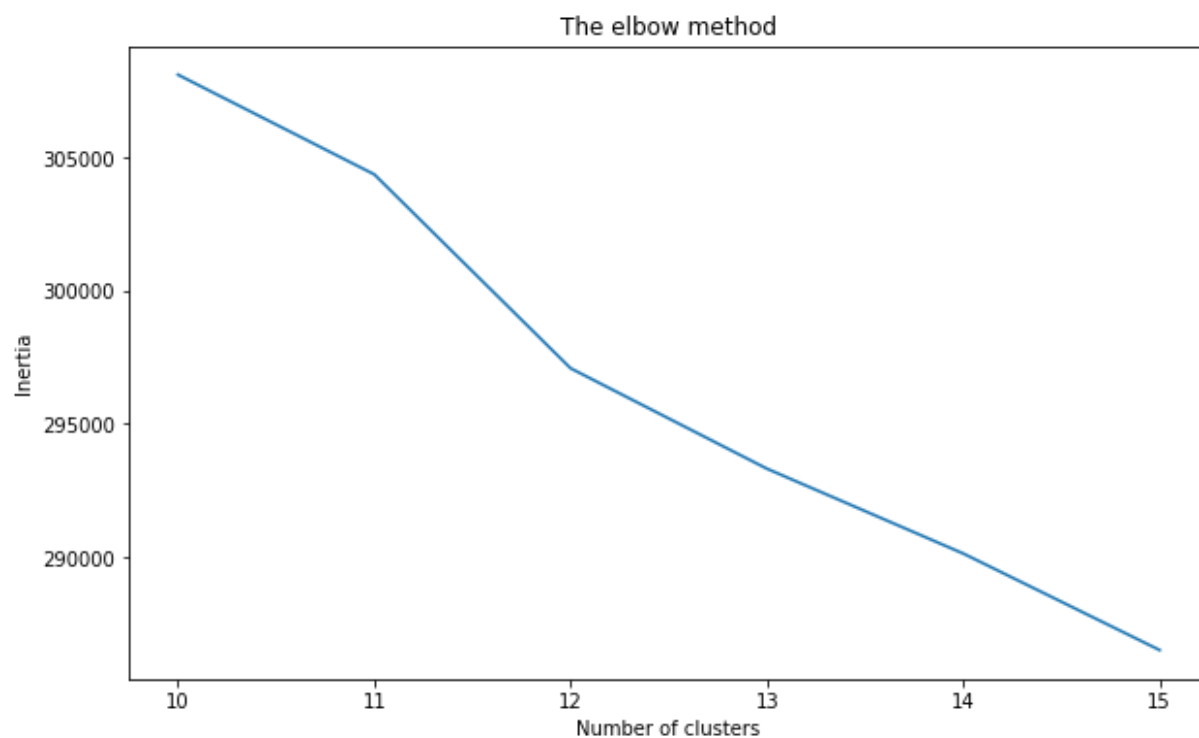
```
In [53]: 1 # Using elbow method to determine how many topics to have  
2 plot_inertia(tfidf_train)
```

Out[53]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9d0125cd30>



```
In [86]: 1 plot_inertia(tfidf_train, start=10, end=15)
```

Out[86]: <matplotlib.axes._subplots.AxesSubplot at 0x7f74249713c8>



```
In [23]: 1 # Don't really see any pattern - will go with a middle value
2 from sklearn.cluster import KMeans
3 kmeans = KMeans(n_clusters=6)
4 kmeans.fit(tfidf_train)
```

```
Out[23]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
n_clusters=6, n_init=10, n_jobs=None, precompute_distances='auto',
random_state=None, tol=0.0001, verbose=0)
```

```
In [52]: 1 topic_labels = pd.concat([pd.Series(kmeans.labels_), X_train.listing_id.reset_index(drop=True)])
```

```
In [54]: 1 topic_labels.columns = ['kmeans_label', 'listing_id', 'id']
```

```
In [69]: 1 topic_labels['comments'] = X_train.comments.reset_index(drop=True)
```

```
In [70]: 1 topic_labels.head()
```

```
Out[70]:
```

	kmeans_label	listing_id	id	comments
0	5	4572225	39668754	Allie's apartment was very clean and located i...
1	3	1542279	239649632	Appartement de 4 chambres Airbnb, avec cuisine...
2	5	15820512	262350020	Guillermo's apartment is in a great neighbourh...
3	5	5061165	85828237	Javier and Vicky were very welcoming and accom...
4	2	21214554	347637143	The place was amazing and the location was gre...

```
In [56]: 1 topic_labels.kmeans_label.value_counts()
```

```
Out[56]: 2    413989
5    183479
1    146937
4     70183
3     62920
0     14291
Name: kmeans_label, dtype: int64
```

```
In [71]: 1 topic_labels_test = pd.concat([pd.Series(kmeans.predict(tfidf_test)), X_test
2 topic_labels_test.columns = ['kmeans_label', 'listing_id', 'id']
3 topic_labels_test['comments'] = X_test.comments.reset_index(drop=True)]
```

In [72]:

```
1 topic_labels_test.head()
```

Out[72]:

	kmeans_label	listing_id	id	comments
0	2	4218034	113949495	We loved our stay at the Urban Jungle! The dup...
1	4	20817201	314215888	Is a good place just to sleep and go. Perfect ...
2	2	4152752	185785747	Niko's place is in a centralized area near all...
3	1	20990607	211627258	Angie is a great host, accommodating, very swe...
4	2	1759462	66002562	The place is lovely, and the location is amazi...

In [63]:

```
1 topic_labels_test.kmeans_label.value_counts()
```

Out[63]:

```
2 103757
5  45943
1  36419
4  17570
3  15678
0   3583
Name: kmeans_label, dtype: int64
```

Will sample a few reviews and see if the topics make sense.

In [64]:

```
1 from random import sample
2 sample_train = sample(range(X_train.shape[0]), 60)
3 sample_test = sample(range(X_test.shape[0]), 30)
```

In [74]:

```
1 topic_labels.loc[sample_train].sort_values('kmeans_label')
```

Out[74]:

	kmeans_label	listing_id	id	comments
15909	0	4426877	35444267	The host canceled this reservation 2 days befo...
787634	1	2197401	297255596	Great place, they were very accommodating.
153334	1	182649	345105753	Great location in the heart of Williamsburg. ...
620250	1	654612	44301099	Kathleen was so accommodating and quick to res...
153402	1	26171844	304954339	great location, a/c was crucial. communicatio...
195609	1	11437634	198769425	Great location near subway, clean apartment an...
525122	1	27195276	351350781	Very stylish and fun apartment! Madia is a gre...
679742	1	19675743	334903072	lovely host, friendly, helpful, just a very ni...
169375	1	8304809	214293652	Exactly as advertised! Clean, stylish and com...
156095	1	7983640	341295284	Erica's place is fantastic and a great deal fo...
114399	1	8126282	152644804	Great apt in the middle of the village \nLoved...
102272	1	3534012	47461665	Great place. perfect location. I agree with al...

```
In [75]: 1 topic_labels_test.loc[sample_test].sort_values('kmeans_label')
```

```
Out[75]:
```

	kmeans_label	listing_id	id	comments
32170	1	7478570	86722635	Great location! Excellent hosts!
86274	1	9783	24941	he's just great!! very polite, very gentle, ve...
31625	1	4455094	150372823	Great place, cozy.
71277	1	21056653	281342900	Great spot
115652	1	23306594	305665377	Great apartment, very well located, all the a...
83730	1	10043483	215594176	We had a great time in this cozy, clean, well-...
163075	1	1102858	355979746	Great space, great location. Eric was very com...
195946	1	9227929	297915139	Great location, prompt and detailed responses.
46016	2	10413576	322612581	Mahmood was incredibly accommodating. It was a...
185902	2	1728437	33955915	My hosts Cheryl and Caroline were great! Enjoy...
15179	2	14935001	198796833	Odi was a great host. She was very helpful in ...
14276	2	21306457	336619584	Worth it!

Honestly I think this clustering works because it separates out foreign language reviews (cluster 2) as well as a cluster for negative reviews (cluster 1). I think this encodes some information and would add to the model.

Make output, group by listing_id for further modeling use.

Will give each listing a vector of length 6 for proportion of each topic.

```
In [79]: 1 topics = pd.concat([topic_labels, topic_labels_test], axis=0)
```

```
In [80]: 1 topics.kmeans_label.value_counts()
```

```
Out[80]: 2    517746
5    229422
1    183356
4     87753
3     78598
0     17874
Name: kmeans_label, dtype: int64
```

```
In [83]: 1 topics = pd.get_dummies(topics, columns = ['kmeans_label'], prefix='topic')
```

```
In [89]: 1 def prop(x):
2         return sum(x)/len(x)
```

```
In [100]: 1 topics_prop = topics.groupby('listing_id').agg({'topic_0': prop, 'topic_1':
```

```
In [101]: 1 topics_prop.head()
```

```
Out[101]:
```

	topic_0	topic_1	topic_2	topic_3	topic_4	topic_5
listing_id						
2454	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000
2539	0.000000	0.222222	0.777778	0.000000	0.000000	0.000000
2595	0.000000	0.232558	0.302326	0.093023	0.046512	0.325581
3330	0.051282	0.102564	0.435897	0.051282	0.051282	0.307692
3831	0.004329	0.134199	0.571429	0.121212	0.086580	0.082251

```
In [103]: 1 topics_prop.to_csv('kmeans_topics.csv', index=True)
```

DBSCAN kernel died (maybe too many rows) so not gonna use this technique.

```
In [61]: 1 # from sklearn.cluster import DBSCAN
2 # dbsc = DBSCAN(eps = 0.5, min_samples = 1000)
3 # dbsc.fit(tfidf_train)
```

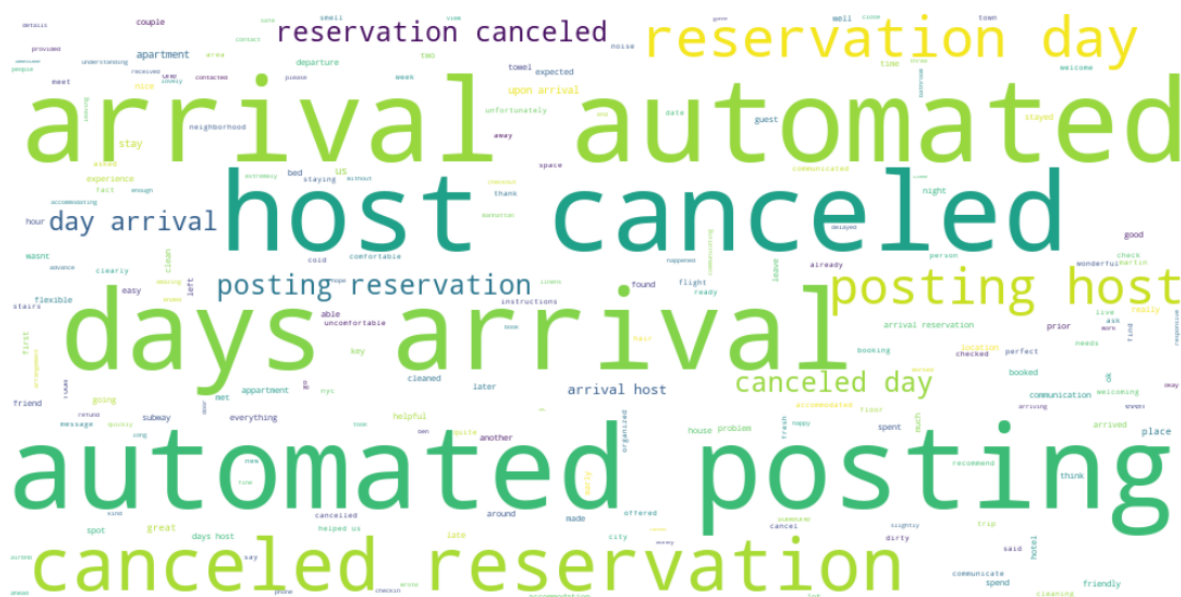
Analysis of topics

```
In [105]: 1 from PIL import Image
2 from wordcloud import WordCloud, ImageColorGenerator
```

```
In [125]: 1 # Use a tokenizer to clean text
2 from nltk.corpus import stopwords
3 stop_words = set(stopwords.words('english'))
4
5 def Tokenizer(str_input):
6     from re import sub
7     # from nltk.stem.porter import PorterStemmer
8     from string import punctuation
9
10    # Remove punctuations
11    str_input = ''.join((char for char in str_input if char not in punctuation))
12
13    # Remove any word containing a number (of which there are many)
14    words = sub(u'(?ui)\b[a-zA-Z0-9]*[0-9]+[a-zA-Z0-9]*\b', "", str_input)
15
16    # porter_stemmer=PorterStemmer()
17    # words = [porter_stemmer.stem(word) for word in words if not word in stop_words]
18    return " ".join(word for word in words if not word in stop_words)
```

```
1 # Collect words from same topic into one variable for each topic
2 topic_0 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
3 topic_1 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
4 topic_2 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
5 topic_3 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
6 topic_4 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
7 topic_5 = " ".join(Tokenizer(review) for review in topics.comments[topics.to
```

```
1 # Topic 0 wordcloud
2
3 wordcloud_0 = WordCloud(background_color="white", width=1000, height=500).ge
4 wordcloud_0.to_file('img/wordcloud_0.png')
5 plt.figure(figsize=(20, 10))
6 plt.imshow(wordcloud_0, interpolation='bilinear')
7 plt.axis("off")
8 plt.show()
```




```
1 wordcloud_1 = WordCloud(background_color="white", width=1000, height=500).ge
2 wordcloud_1.to_file('img/wordcloud_1.png')
3 plt.figure(figsize=(20, 10))
4 plt.imshow(wordcloud_1, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()
```



```
1 wordcloud_2 = WordCloud(background_color="white", width=1000, height=500).ge
2 wordcloud_2.to_file('img/wordcloud_2.png')
3 plt.figure(figsize=(20, 10))
4 plt.imshow(wordcloud_2, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()
```



```
1 wordcloud_3 = WordCloud(background_color="white", width=1000, height=500).ge
2 wordcloud_3.to_file('img/wordcloud_3.png')
3 plt.figure(figsize=(20, 10))
4 plt.imshow(wordcloud_3, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()
```



```
1 wordcloud_4 = WordCloud(background_color="white", width=1000, height=500).ge
2 wordcloud_4.to_file('img/wordcloud_4.png')
3 plt.figure(figsize=(20, 10))
4 plt.imshow(wordcloud_4, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()
```



```
1 wordcloud_5 = WordCloud(background_color="white", width=1000, height=500).ge
2 wordcloud_5.to_file('img/wordcloud_5.png')
3 plt.figure(figsize=(20, 10))
4 plt.imshow(wordcloud_5, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()
```



The topic clustering mostly makes sense, although I can envision decreasing the number of clusters a bit.