

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

```
In [9]: 1 data = pd.read_csv('data/listings.csv')

/home/twang/anaconda3/lib/python3.6/site-packages/IPython/core/interactiveshell
1.py:3020: DtypeWarning: Columns (61,87,88) have mixed types. Specify dtype opt
ion on import or set low_memory=False.
    interactivity=interactivity, compiler=compiler, result=result)
```

```
In [29]: 1 # Concat text columns
        2 text = data[['name', 'summary', 'description']].astype('U').apply(lambda x:
```

```
In [31]: 1 # Clean text, get rid of punctuations and numbers, and Lower letter case
        2 from nltk.corpus import stopwords
        3 stop_words = set(stopwords.words('english'))
        4
        5 def Clean_text(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 punctuati
        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 st
        18    return " ".join(word for word in words if not word in stop_words)
```

```
In [32]: 1 text=text.apply(Clean_text)
```

```
In [41]: 1 # Vectorize using tfidf
        2
        3 from sklearn.decomposition import TruncatedSVD
        4 from sklearn.feature_extraction.text import TfidfVectorizer
        5 from sklearn.pipeline import make_pipeline
        6
        7 tfidf = TfidfVectorizer(max_df=0.9, \
        8                       min_df=0.01, \
        9                       ngram_range=(1,3))
        10 svd = TruncatedSVD(n_components=100)
        11 pipeline = make_pipeline(tfidf, svd)
        12
```

```
In [43]: 1 text_vec = pipeline.fit_transform(text)
```

```
In [109]: 1 text_vec_df = pd.DataFrame(text_vec, index=data.id)
```

```
In [49]: 1 # Now, lets select dimensions most correlated with
2
3 prices = data.price.apply(lambda x: np.NaN if x == 'NaN' else float(str(x).l
```

```
In [125]: 1 corrs = []
2
3 for column in text_vec_df:
4     corrs.append((column, np.corrcoef(text_vec_df[column], prices)[0,1]))
```

```
In [126]: 1 corrs
```

```
Out[126]: [(0, -0.06648791206665025),
(1, -0.028831301490426873),
(2, 0.1376116399887893),
(3, -0.14064679796341056),
(4, -0.015045949884558632),
(5, -0.017338075788816746),
(6, -0.15018636924743195),
(7, -0.03251137010051627),
(8, -0.03402721371523929),
(9, 0.002502574276354029),
(10, -0.017603960019613424),
(11, -0.012825034604486103),
(12, 0.015011891979576192),
(13, -0.008502979562095274),
(14, -0.021421223808995887),
(15, 0.08916933204700668),
(16, 0.03600523544519698),
(17, 0.05221572799303739),
(18, 0.07873021129446113),
(19, 0.011067431110016746)]
```

```
In [127]: 1 # Take the top 10 most correlated columns by absolute value
2 best_cols = [item[0] for item in sorted([(item[0], abs(item[1])) for item in
```

```
In [128]: 1 best_cols
```

```
Out[128]: [6, 3, 2, 15, 18, 0, 24, 17, 20, 27]
```

Quick analysis of the most useful dimensions of svd

NOTE: Check with people knowledgeable with the matter. I am not sure what is the correct interpretation of weights in a factor.

```
In [102]: 1 vocab = [item[0] for item in sorted(tfidf.vocabulary_.items(), key = operato
```

```
In [131]: 1 # The best column has over 0.1 correlation
          2 voc_weight_6 = sorted(zip(vocab, svd.components_[6,:]), key = lambda x:x[1],
```

```
In [135]: 1 voc_weight_6[-20:]
```

```
Out[135]: [('center', -0.07272840672472593),
            ('high', -0.07314626498261155),
            ('brownstone', -0.07491977130806952),
            ('studio', -0.0783623473510409),
            ('beautiful', -0.08270145726799635),
            ('luxury', -0.08282724873363825),
            ('garden', -0.08313558721756222),
            ('new york city', -0.08584393989391555),
            ('york city', -0.08585334043852712),
            ('views', -0.09038845236272569),
            ('building', -0.09170300233121975),
            ('modern', -0.09791023397410266),
            ('city', -0.09839316504879354),
            ('prospect park', -0.10030319935963242),
            ('prospect', -0.12129553554707352),
            ('loft', -0.12883929977176176),
            ('new york', -0.15294873544812715),
            ('york', -0.1537952360737561),
            ('new', -0.20245435856556063),
            ('brooklyn', -0.23794178814603853)]
```

```
In [141]: 1 # Top 20 words by weight
          2 vocab_weights = []
          3
          4 for col in best_cols:
          5     vocab_weights.append(sorted(zip(vocab, svd.components_[col,:]), key = la
```

Hopefully from the loadings in these best factors we can gain insight into what makes an apartment pricey.

```
In [142]: 1 # Select best columns
          2 text_vec_df = text_vec_df[best_cols]
```

```
In [143]: 1 text_vec_df.columns = ['descVec_'+str(x) for x in best_cols]
```

```
In [162]: 1 # Join to cleaned dataset
          2 cleaned_data = pd.read_csv('data/cleaned_with_nlp_listings_2.csv')
```

```
In [146]: 1 cleaned_with_desc = pd.merge(cleaned_data, text_vec_df.reset_index(), on='id
```

In [147]: 1 cleaned_with_desc.isnull().sum()

Out[147]: id 0
 host_response_rate 0
 host_is_superhost 0
 host_listings_count 0
 host_total_listings_count 0
 host_has_profile_pic 0
 host_identity_verified 0
 is_location_exact 0
 accommodates 0
 bathrooms 0
 bedrooms 0
 beds 0
 price 0
 security_deposit 0
 cleaning_fee 0
 guests_included 0
 extra_people 0
 minimum_nights 0
 maximum_nights 0

In [151]: 1 %%time
 2 # Make imputer for the missing values using KNN
 3
 4 from fancyimpute import KNN
 5
 6 cleaned_with_desc_filled = KNN(k=5).fit_transform(cleaned_with_desc)

Imputing row 1/39926 with 10 missing, elapsed time: 524.312
 Imputing row 101/39926 with 0 missing, elapsed time: 524.319
 Imputing row 201/39926 with 0 missing, elapsed time: 524.324
 Imputing row 301/39926 with 0 missing, elapsed time: 524.328
 Imputing row 401/39926 with 0 missing, elapsed time: 524.334
 Imputing row 501/39926 with 0 missing, elapsed time: 524.343
 Imputing row 601/39926 with 0 missing, elapsed time: 524.350
 Imputing row 701/39926 with 0 missing, elapsed time: 524.354
 Imputing row 801/39926 with 0 missing, elapsed time: 524.359
 Imputing row 901/39926 with 0 missing, elapsed time: 524.363
 Imputing row 1001/39926 with 0 missing, elapsed time: 524.365
 Imputing row 1101/39926 with 0 missing, elapsed time: 524.372
 Imputing row 1201/39926 with 0 missing, elapsed time: 524.374
 Imputing row 1301/39926 with 0 missing, elapsed time: 524.383
 Imputing row 1401/39926 with 0 missing, elapsed time: 524.387
 Imputing row 1501/39926 with 0 missing, elapsed time: 524.391
 Imputing row 1601/39926 with 0 missing, elapsed time: 524.398
 Imputing row 1701/39926 with 0 missing, elapsed time: 524.404
 Imputing row 1801/39926 with 10 missing, elapsed time: 524.406
 Imputing row 1901/39926 with 0 missing, elapsed time: 524.412

In [159]: 1 cleaned_with_desc_filled = pd.DataFrame(cleaned_with_desc_filled, columns=cl

In [160]: 1 cleaned_with_desc_filled.to_csv('cleaned_with_desc_filled.csv')

In [161]: 1 del cleaned_with_desc, data

Model building

```
In [175]: 1 # Make a pipeline to train with random search tuning
2
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import cross_val_score
5 from sklearn import metrics
6 from sklearn.model_selection import RandomizedSearchCV
7 from sklearn.preprocessing import StandardScaler
8 import xgboost as xgb
```

Baseline: see performance without appended columns

```
In [164]: 1 cleaned_data.price.head()
```

```
Out[164]: 0    137.0
1    149.0
2    225.0
3     70.0
4     89.0
Name: price, dtype: float64
```

```
In [ ]: 1 cleaned_data.set_index('id', inplace=True)
```

```
In [176]: 1 scaler_baseline = StandardScaler()
```

```
In [166]: 1 y_baseline = cleaned_data.price
2 X_baseline = cleaned_data.loc[:, cleaned_data.columns != 'price']
```

```
In [177]: 1 fit_baseline = xgb.XGBRegressor(objective='reg:linear', random_state = 23)
2 cv_baseline = cross_val_score(fit_baseline, scaler_baseline.fit_transform(X_
```

```
/home/twang/anaconda3/lib/python3.6/site-packages/sklearn/preprocessing/data.p
y:625: DataConversionWarning: Data with input dtype int64, float64 were all con
verted to float64 by StandardScaler.
```

```
    return self.partial_fit(X, y)
```

```
/home/twang/anaconda3/lib/python3.6/site-packages/sklearn/base.py:462: DataConv
ersionWarning: Data with input dtype int64, float64 were all converted to float
64 by StandardScaler.
```

```
    return self.fit(X, **fit_params).transform(X)
```

```
In [178]: 1 np.sqrt(-cv_baseline.mean())
```

```
Out[178]: 171.53075105041998
```

Lift: with new columns, see how much lift is generated

```
In [171]: 1 cleaned_with_desc_filled.set_index('id', inplace=True)
```

```
In [179]: 1 scaler_lift = StandardScaler()

In [172]: 1 y_lift = cleaned_with_desc_filled.price
          2 X_lift = cleaned_with_desc_filled.loc[:, cleaned_with_desc_filled.columns !=

In [181]: 1 fit_lift = xgb.XGBRegressor(objective='reg:linear', random_state = 23)
          2 cv_lift = cross_val_score(fit_lift, scaler_lift.fit_transform(X_lift), y_lift

In [182]: 1 np.sqrt(-cv_lift.mean())

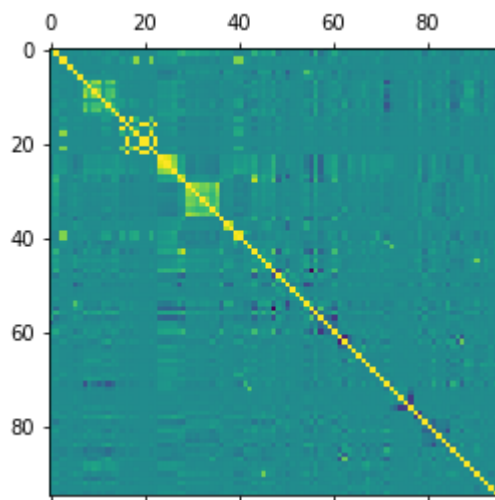
Out[182]: 172.6613213053378
```

Conclusion

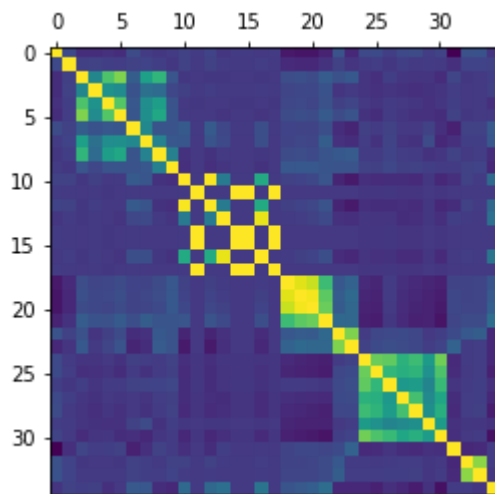
It does seem that the new columns did not add predictive power to the model. This is surprising, but understandable, as these new vectors did not really have a strong correlation with price (in the 10% range). It seems we just don't have enough information in these vectors to really explain variations in price.

I also looked at collinearity and found that we can further process the data to remove some of the super highly collinear columns potentially, and this might merit further investigation. Perhaps these could be result of mistakes in data cleaning.

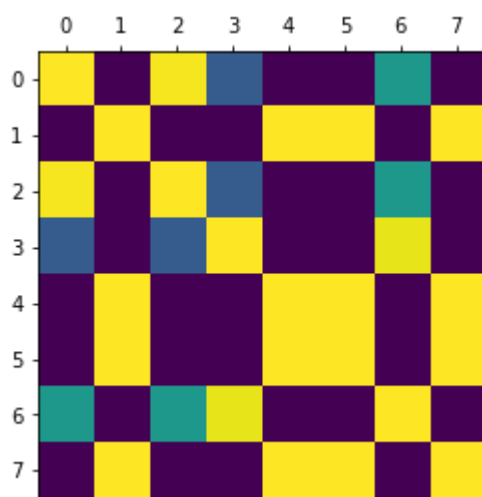
```
In [184]: 1 plt.matshow(X_lift.corr())
          2 plt.show()
```



```
In [186]: 1 plt.matshow(X_lift.iloc[:, 5:40].corr())
          2 plt.show()
```



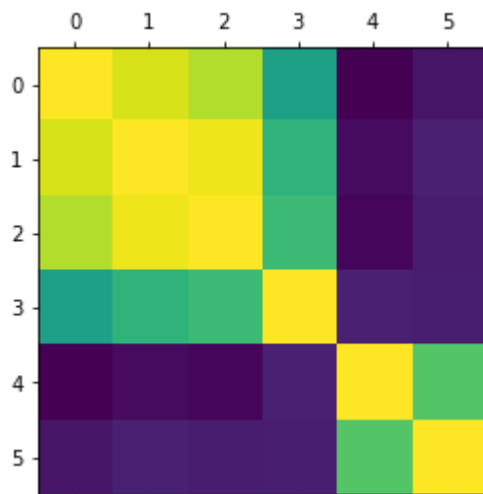
```
In [188]: 1 plt.matshow(X_lift.iloc[:, 15:23].corr())
          2 plt.show()
```



```
In [189]: 1 X_lift.iloc[:, 15:23].columns
```

```
Out[189]: Index(['minimum_nights', 'maximum_nights', 'minimum_minimum_nights',
                 'maximum_minimum_nights', 'minimum_maximum_nights',
                 'maximum_maximum_nights', 'minimum_nights_avg_ntm',
                 'maximum_nights_avg_ntm'],
                 dtype='object')
```

```
In [192]: 1 plt.matshow(X_lift.iloc[:, 23:29].corr())  
2 plt.show()
```



```
In [193]: 1 X_lift.iloc[:, 23:29].columns
```

```
Out[193]: Index(['availability_30', 'availability_60', 'availability_90',  
                'availability_365', 'number_of_reviews', 'number_of_reviews_ltm'],  
              dtype='object')
```

```
In [ ]: 1
```