



CASE STUDY - topic modeling and feature engineering

[Feature engineering](https://en.wikipedia.org/wiki/Feature_engineering) (https://en.wikipedia.org/wiki/Feature_engineering) is the process of using domain knowledge of your data to create features that can be leveraged by machine learning models, because sometimes it is used in a context where features are transformed for machine learning, but the inclusion of domain knowledge is not implied.

It is unfortunately common that for large datasets engineered features are not easy to create. When there are many features generally only a small number play an important role. Furthermore, domain insight is even more difficult to fold into the model when there are hundreds or thousands of features to keep in mind. However, there is a middle ground--locked up in language. In this case study we will use topic modeling to gather insight from text. Ideally, the result of these types of experiments would be shared with domain experts that are relevant when it comes to your business opportunity.

```
In [1]: %%capture
pip install -U pip
```

```
In [2]: %%capture
pip install pyLDAvis
```

```
In [3]: %%capture
pip install ../data/en_core_web_sm-2.3.1.tar.gz --user
```

```
In [4]: ##IMPORTANT: Please restart the Kernel after running the above 3 cells
```

```
In [1]: import os
import re
import sys
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.utils import shuffle
from sklearn.datasets import load_files
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.decomposition import LatentDirichletAllocation
from sklearn.pipeline import Pipeline
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
from string import punctuation, printable
from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS

try:
    import pyLDAvis
    import pyLDAvis.sklearn
except:
    raise Exception("'pip install pyldavis' before running this notebook")

pyLDAvis.enable_notebook()
plt.style.use('seaborn')
%matplotlib inline

DATA_DIR = os.path.join("../", "data")
```

Synopsis

Goal: AAVAIL has recently enabled comments on the core streaming service. The data science team knows that this will be an incredibly important source of data going to inform customer retention, product quality, product market fit and more. Comments are going live next week and being the diligent data scientist that you are your plan a pipeline that will consume the comments and create visualizations that can be used to communicate with domain experts.

Outline

1. EDA - summary tables, use tSNE to visualize the data
2. Create a transformation pipelines for NMF and LDA
3. Use Idaviz and wordclouds to get insight into the clusters

Data

Even before receiving the first comment, we want to start building our Pipeline using a proxy dataset. In this study Case we will work with a dataset publicly available dataset of

- [Here \(http://www.nltk.org/nltk_data\)](http://www.nltk.org/nltk_data) is the web page that references all the public dataset that NLTK provide. In this Study Case we will work with the 'Sentiment Polarity D' dataset has already been downloaded and is in the data folder of the working directory)
- For more examples of applications with these data see [NLTK's book chapter that uses these data \(https://www.nltk.org/book/ch06.html\)](https://www.nltk.org/book/ch06.html)

```
In [2]: filename = os.path.join(DATA_DIR, 'movie_reviews.csv')
df = pd.read_csv(filename)
X = df['review'].tolist()
print(X[4])
```

```
b"kolya is one of the richest films i've seen in some time . \nzdenek sverak plays a confirmed old bachelor ( who's like
as a czech cellist increasingly impacted by the five-year old boy that he's taking care of . \nthough it ends rather abr
ed to spend more time with these characters-- the acting , writing , and production values are as high as , if not highe
\nthis father-and-son delight-- sverak also wrote the script , while his son , jan , directed-- won a golden globe for b
e days after i saw it , walked away an oscar . \nin czech and russian , with english subtitles . \n"
```

QUESTION 1

The main focus of this exercise is to enable visualization of topics, but these topics can be used as additional features for prediction tasks. The goal of this case study is to ensure natural language processing pipelines and topic modeling tools.

There are many ways to process tokens (words, dates, emojis etc). NLTK is often used to pre-process text data before the tokens are vectorized. Generally, the tokens are modified (https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html). The next code block provides a lemmatization function that makes use of the library `spacy` to install it and download the English language reference material as follows. Stopwords are words that are very common or otherwise irrelevant we use a default list here, but it is a pipeline that needs to be customized for the subject area.

If you prefer to use NLTK then you could use a simple lemmatizer like the WordLemmatizer.

```
In [3]: import spacy
STOPLIST = ENGLISH_STOP_WORDS
STOPLIST = set(list(STOPLIST) + ["foo", "film", "movie", "make"])

if not 'nlp' in locals():
    print("Loading English Module...")
    nlp = spacy.load('en_core_web_sm')

def lemmatize_document(doc, stop_words=None):
    """
    takes a list of strings where each string is a document
    returns a processed list of strings
    """

    if not stop_words:
        stop_words = set([])

    ## ensure working with string
    doc = str(doc)
    doc = doc.replace('\n', '')
    doc = doc.replace('\t', '')

    # First remove punctuation form string
    if sys.version_info.major == 3:
        PUNCT_DICT = {ord(punc): None for punc in punctuation}
        doc = doc.translate(PUNCT_DICT)

    # remove unicode
    clean_doc = "".join([char for char in doc if char in printable])

    # Run the doc through spaCy
    doc = nlp(clean_doc)

    # Lemmatize and lower text
    tokens = [re.sub(r"\W+", "", token.lemma_.lower()) for token in doc]
    tokens = [t for t in tokens if len(t) > 1]

    return ' '.join(w for w in tokens if w not in stop_words)

## example usage
corpus = ["You can fool some of the people all of the time, and all of the people some of the time, but you can not fool the time". -- Abraham Lincoln']
processed = [lemmatize_document(doc, STOPLIST) for doc in corpus]
print(processed[0])
processed = [lemmatize_document(doc, None) for doc in corpus]
print("\n"+processed[0])
```

Loading English Module...

pron fool people time people time pron fool people time abraham lincoln

pron can fool some of the people all of the time and all of the people some of the time but pron can not fool all of the

```
In [8]: ## YOUR CODE HERE

## Preprocess all the reviews of the corpus with the lemmatize_document() function to create a list of cleaned reviews.

## Applying the lemmatize_document() function to all the documents of the corpus takes several minutes.
## In order to save you some time we preprocessed the texts with the line of code commented bellow and saved
## the processed documents in a .txt file. You can either re-preprocess the text by uncommenting the lines above
## or you can directly read the processed_text.txt file as shown bellow.

# from tqdm import tqdm
# tqdm.pandas()
# processed = df.progress_apply(lambda x : lemmatize_document(x['review'], STOPLIST), axis=1).tolist()

processed = []
with open(os.path.join(DATA_DIR, 'processed_text.txt'), 'r') as f :
    for line in f:
        processed.append(line)

print("processing done.")
```

processing done.

QUESTION 2

Use the CountVectorizer from sklearn to vectorize the documents.

Additional resources:

- [scikit-learn CountVectorizer](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html) (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)
- [scikit-learn working with text](https://scikit-learn.org/stable/tutorial/text_analytics/working_with_text_data.html) (https://scikit-learn.org/stable/tutorial/text_analytics/working_with_text_data.html)

Because this is an exercise in visualization set the `max_features` to something like 500. In the context of supervised learning it is reasonable to grid-search to optimize this |

```
In [9]: ## YOUR CODE HERE

max_features = 500

# Create a CountVectorizer object
tf_vectorizer = CountVectorizer(max_df=0.95, min_df=2,
                                max_features=max_features,
                                stop_words='english')

# Fit and transform this object to the processed reviews
tf = tf_vectorizer.fit_transform(processed)
print("ready")

ready
```

QUESTION 3

Fit a LDA model to the corpus. For example, you could use something like the following.

```
n_topics = 10
lda_model = LatentDirichletAllocation(n_components=n_topics, max_iter=5,
                                      learning_method='online',
                                      learning_offset=50.,
                                      random_state=0)

lda_model.fit(tf)
```

- [scikit-learn's LDA](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.LatentDirichletAllocation.html) (https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.LatentDirichletAllocation.html)
- [scikit-learn's user guide for LDA](https://scikit-learn.org/stable/modules/decomposition.html#latentdirichletallocation) (https://scikit-learn.org/stable/modules/decomposition.html#latentdirichletallocation)

```
In [10]: ## YOUR CODE HERE
n_topics = 10

# Create an LDA object
lda_model = LatentDirichletAllocation(n_components=n_topics, max_iter=5,
                                      learning_method='online',
                                      learning_offset=50.,
                                      random_state=0)

# Fit the model to the bag of word we created earlier
lda_model.fit(tf)
```

```
Out[10]: LatentDirichletAllocation(batch_size=128, doc_topic_prior=None,
                                   evaluate_every=-1, learning_decay=0.7,
                                   learning_method='online', learning_offset=50.0,
                                   max_doc_update_iter=100, max_iter=5,
                                   mean_change_tol=0.001, n_components=10, n_jobs=None,
                                   perp_tol=0.1, random_state=0, topic_word_prior=None,
                                   total_samples=1000000.0, verbose=0)
```

QUESTION 4

Visualize the corpus using [pyldavis](https://github.com/bmabey/pyLDavis) (https://github.com/bmabey/pyLDavis).

```
pyLDavis.sklearn.prepare(lda_model,tf, tf_vectorizer, R=20)

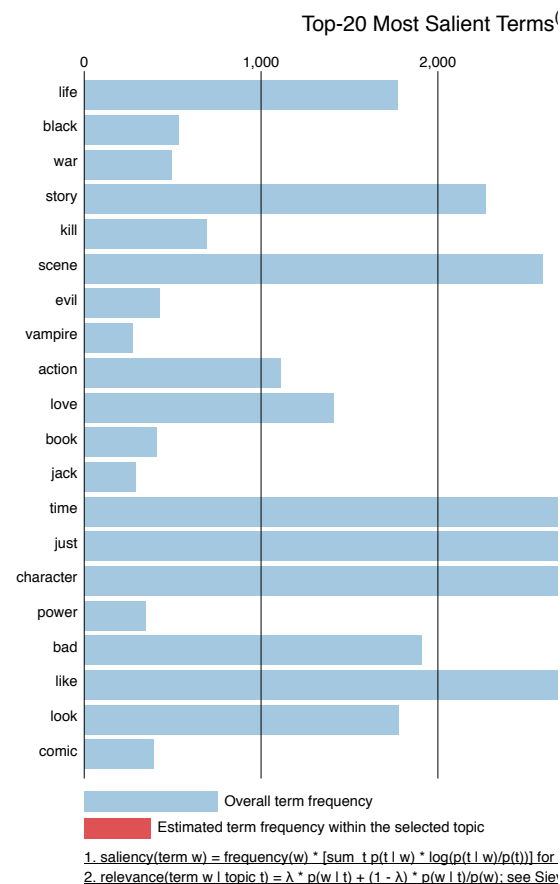
PyLDavis documentation (https://pyldavis.readthedocs.io/en/latest)
PyLDavis demos (https://pyldavis.readthedocs.io/en/latest/readme.html#video-demos)
```

Out[11]:

Selected Topic: 0	Previous Topic	Next Topic	Clear Topic
-------------------	----------------	------------	-------------

Slide to adjust relevance metric:(2)

$\lambda = 1$



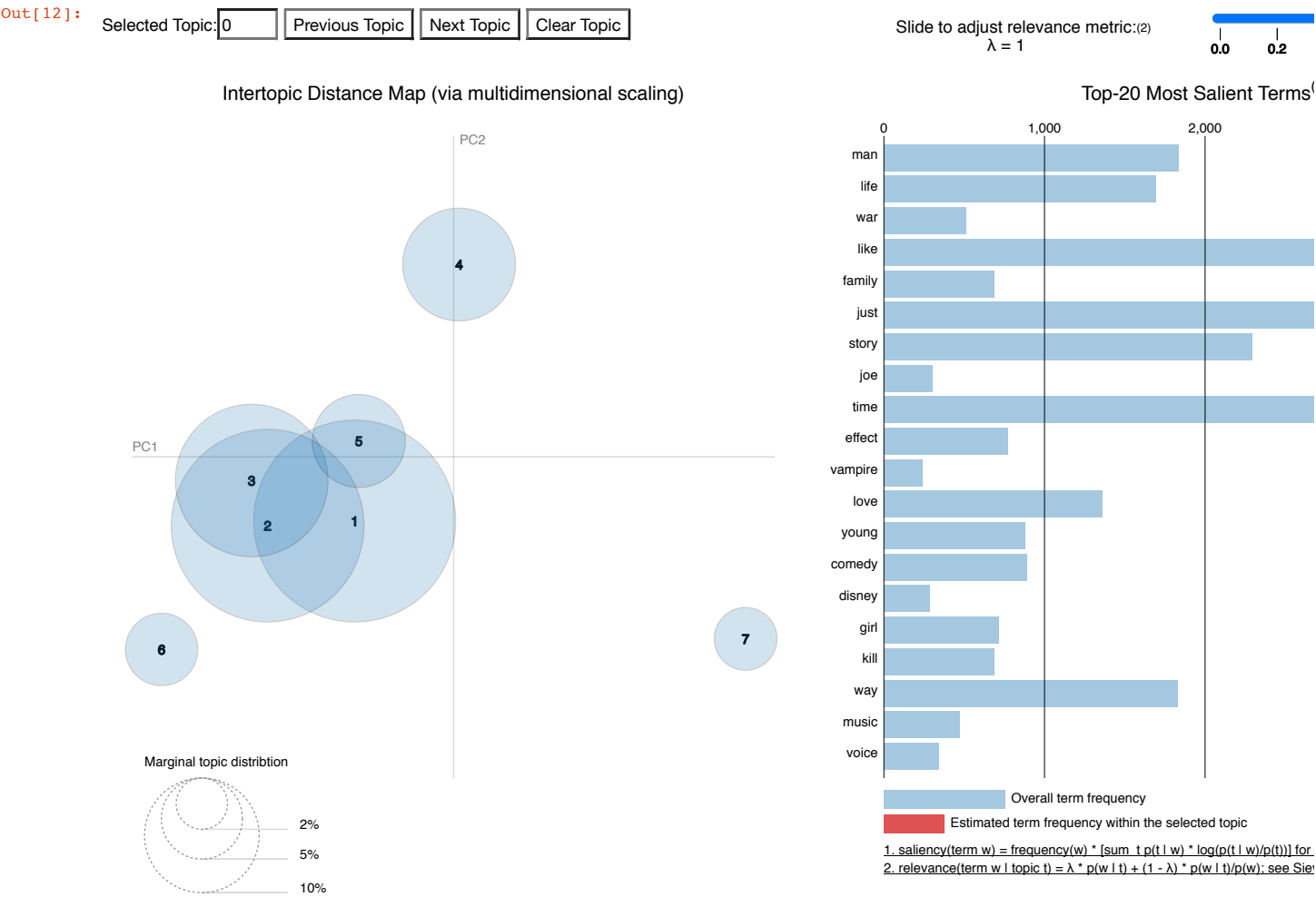
Try different numbers of clusters until there is decent separation in the visualization

In [12]:

YOUR CODE HERE

```
n_topics = 7
lda_model = LatentDirichletAllocation(n_components=n_topics, max_iter=5,
                                     learning_method='online',
                                     learning_offset=50.,
                                     random_state=0)

lda_model.fit(tf)
lda_transformed = lda_model.transform(tf)
pyLDAvis.sklearn.prepare(lda_model, tf, tf_vectorizer, R=20)
```



The visualization here can help determine a reasonable number of number of clusters and it can serve as a communication tool. If the goal was to find topics that are associated with folks in marketing to refine the clustering. There are a couple of parameters that can be used to modify the clustering and visualization. The discovery of feature engineering.

QUESTION 6

If you were to use the topics from this model to inform clustering or supervised learning you would first need to be able to extract and represent them as a matrix. Along the same report with tabular descriptions of the data then you will need to be able to extract topic representations. Here is a starter function

```
In [13]: def get_top_words(model, feature_names, n_top_words):
        """
        Get the top words defining the different topics of the LDA model
        INPUT : the LDA model, the names of the features of the bag of word (these are the actual words in the vocabulary)
        and the number of top words.
        RETURN : A dictionary where the keys are the topic's ID and the values are the lists of the n_top_words top words.

        """
        top_words = {}
        for topic_idx, topic in enumerate(model.components_):
            _top_words = [feature_names[i] for i in topic.argsort()[::-n_top_words - 1:-1]]
            top_words[str(topic_idx)] = _top_words
        return top_words
```

Use the function to print the top k words for each topic

```
In [14]: ## YOUR CODE HERE

## set n_top_words
top_words = 15
## get the vectorizer's feature names
tf_feature_names = np.array(tf_vectorizer.get_feature_names())
## get the top words for each topic
top_words = get_top_words(lda_model, tf_feature_names, top_words)
all_top_words = np.array(list(set().union(*[v for v in top_words.values()])))

## print the topics and the top words of each topic
for key, vals in top_words.items():
    print(key, " ".join(vals))
print("total words: %s"%len(all_top_words))

0 joe like computer just music rock deep brother effect look time cop really say song
1 war vampire man life ryan battle george time save love way kill beautiful scene world
2 family disney story jackie voice kid king little character like mr year time child good
3 life girl performance mother character daughter father child boy young come man love batman story
4 character good story like scene time action play just work man plot great alien year
5 character good comedy just like play funny time big come thing work laugh scene role
6 like bad just good character time know scene think really say play look thing way
total words: 66
```

QUESTION (EXTRA CREDIT) 7

If you used `transform` on your original tokens you should have a `2000 x k` array where `k` is the number of topics you choose. Create a PCA or tSNE visualization that p-dimensional space then uses colors to indicate which documents belong to a topic (e.g. probability > 0.5).

In [15]: `## YOUR CODE HERE`

```
def make_plot(lda_mat):

    fig = plt.figure(figsize=(15,15), facecolor='white')
    ax = fig.add_subplot(111)

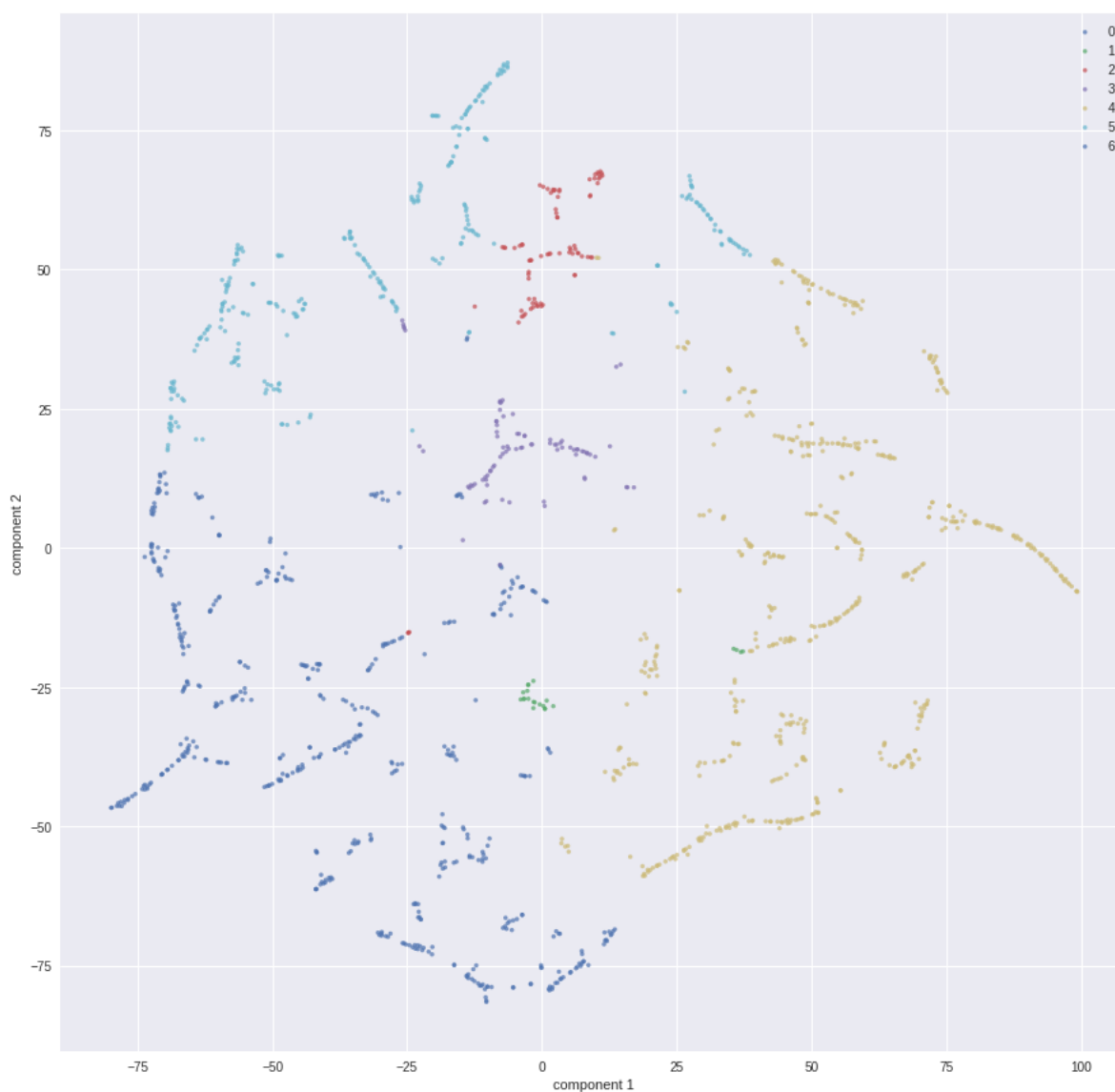
    tsne = TSNE(n_components=2, perplexity=10, init='pca')
    projected = tsne.fit_transform(lda_mat)

    #    pca = PCA(n_components=2)
    #    projected = pca.fit_transform(lda_mat)

    for class_num in np.arange(n_topics):
        topic_inds = np.where(lda_mat[:, class_num] > 0.5)[0]
        ax.scatter(projected[topic_inds, 0],
                   projected[topic_inds, 1],
                   edgecolor='none', marker='.', alpha=0.7, label=str(class_num))

    ax.set_xlabel('component 1')
    ax.set_ylabel('component 2')
    ax.legend()

make_plot(lda_transformed)
```



```
In [16]: with open(os.path.join(DATA_DIR, 'processed_text.txt'), 'w') as f:
        for item in processed:
            f.write("%s\n" % item)
```