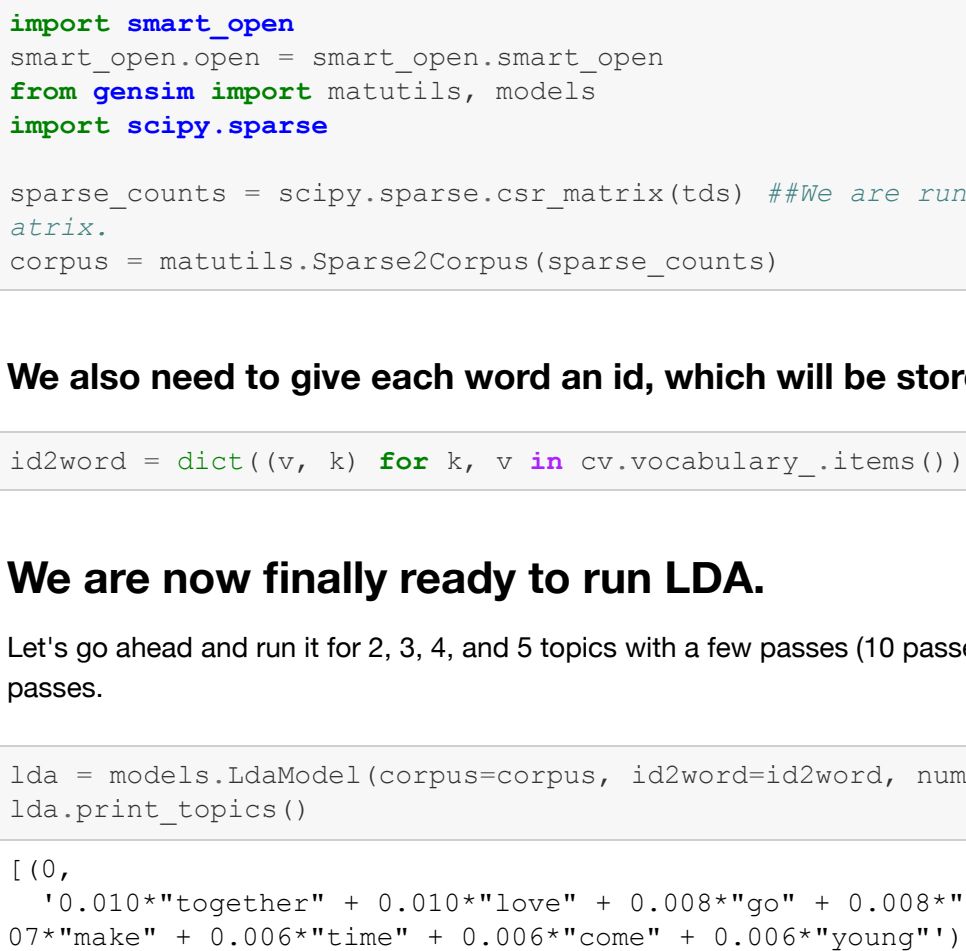
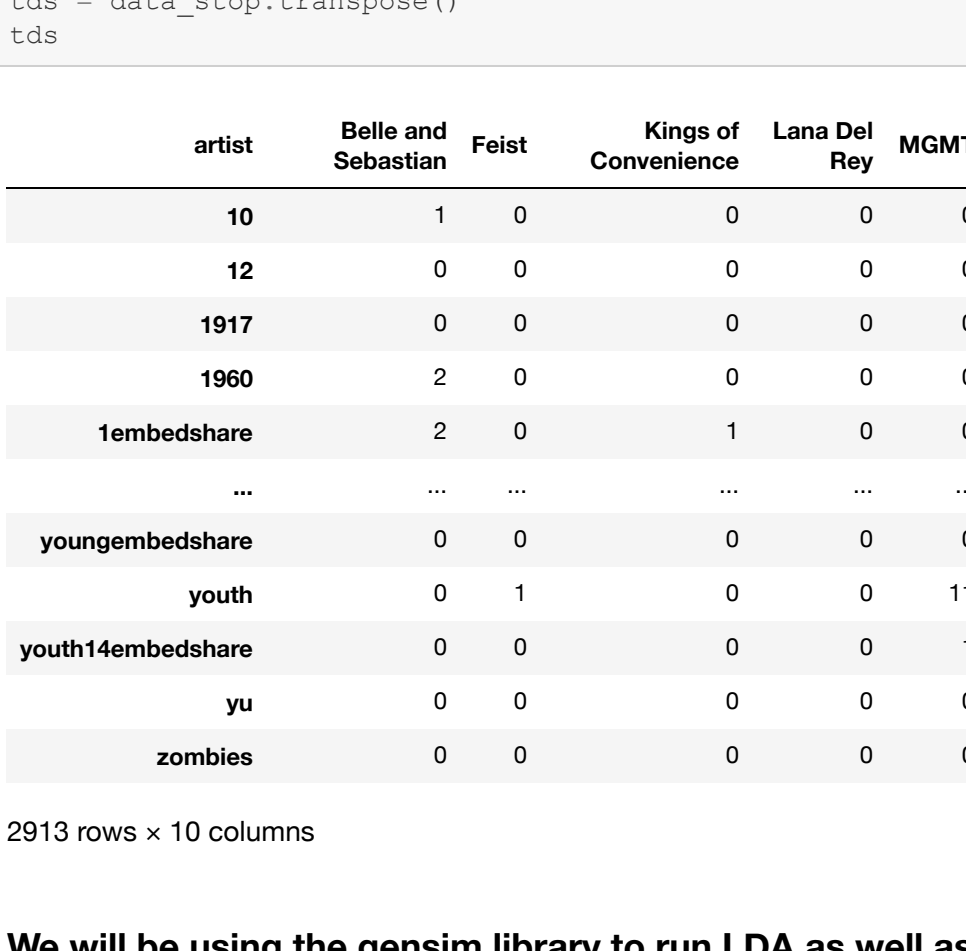
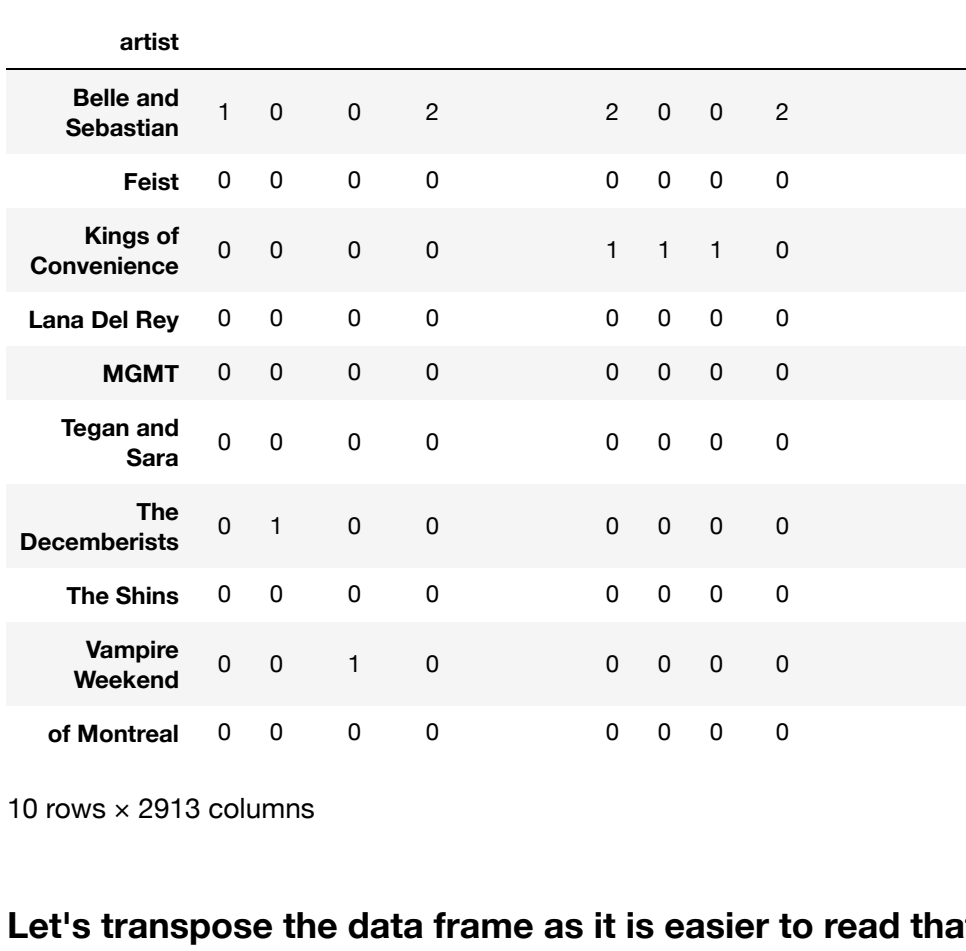
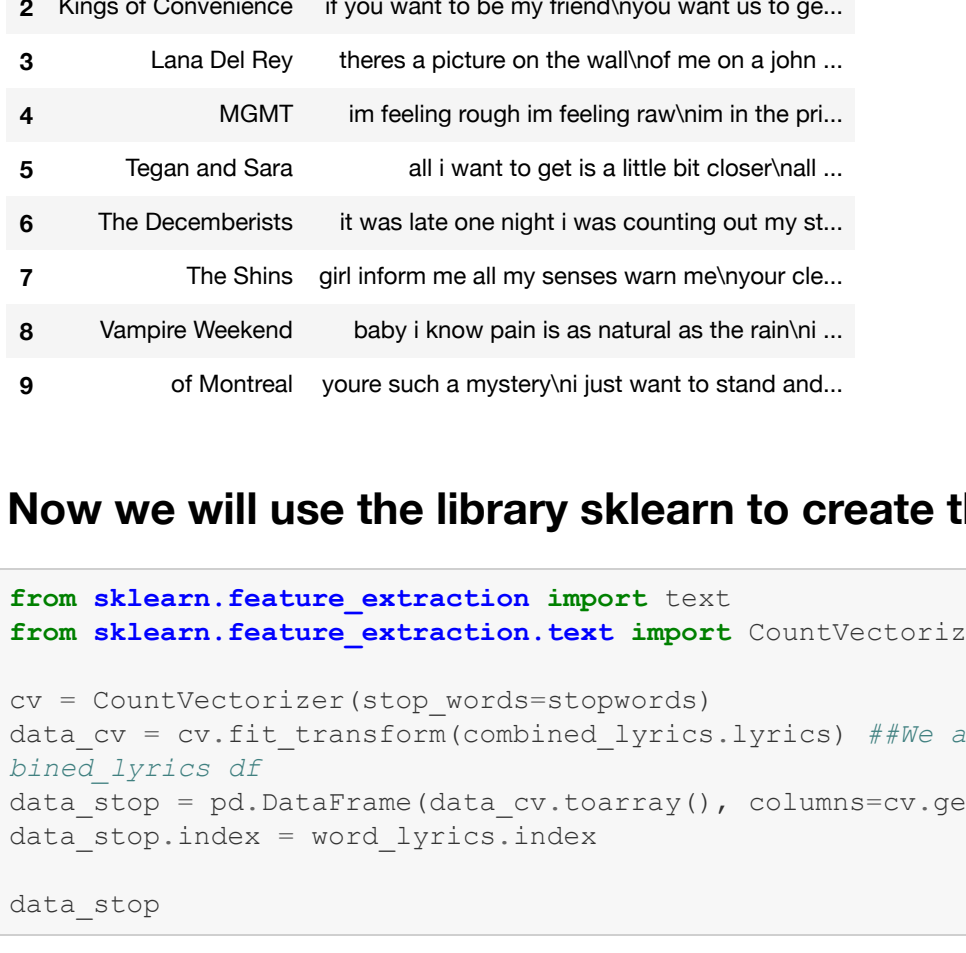
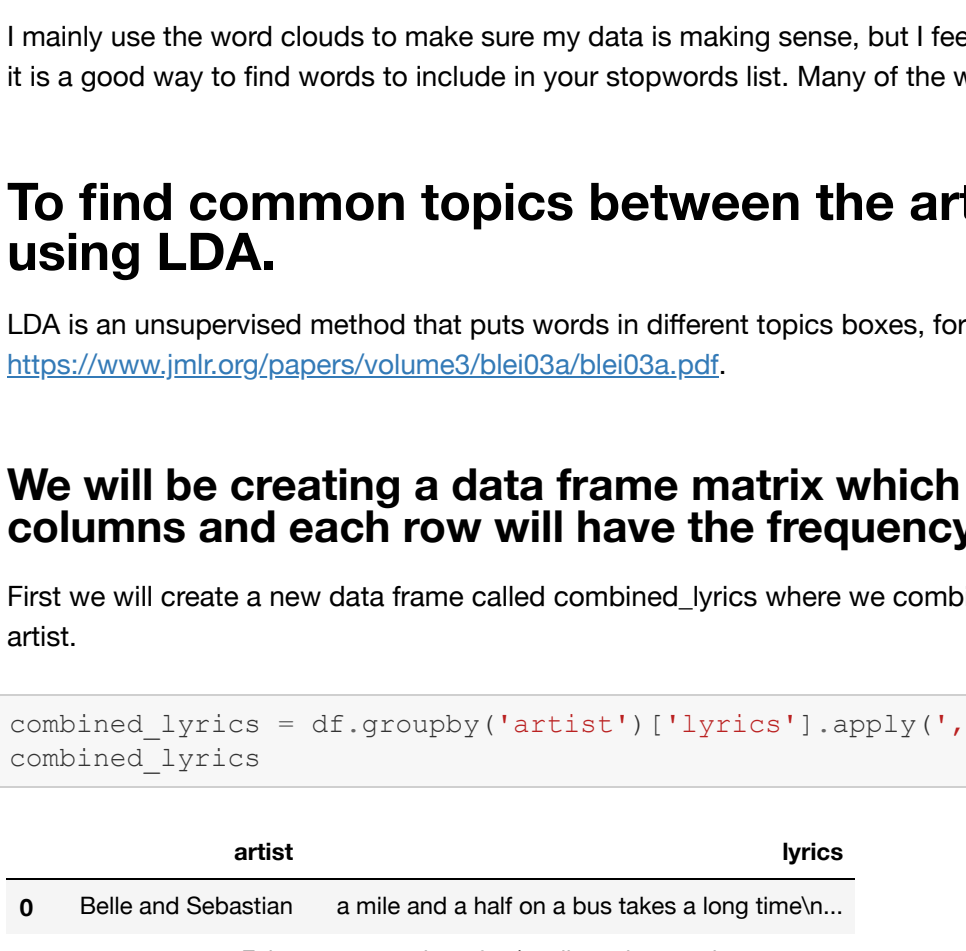
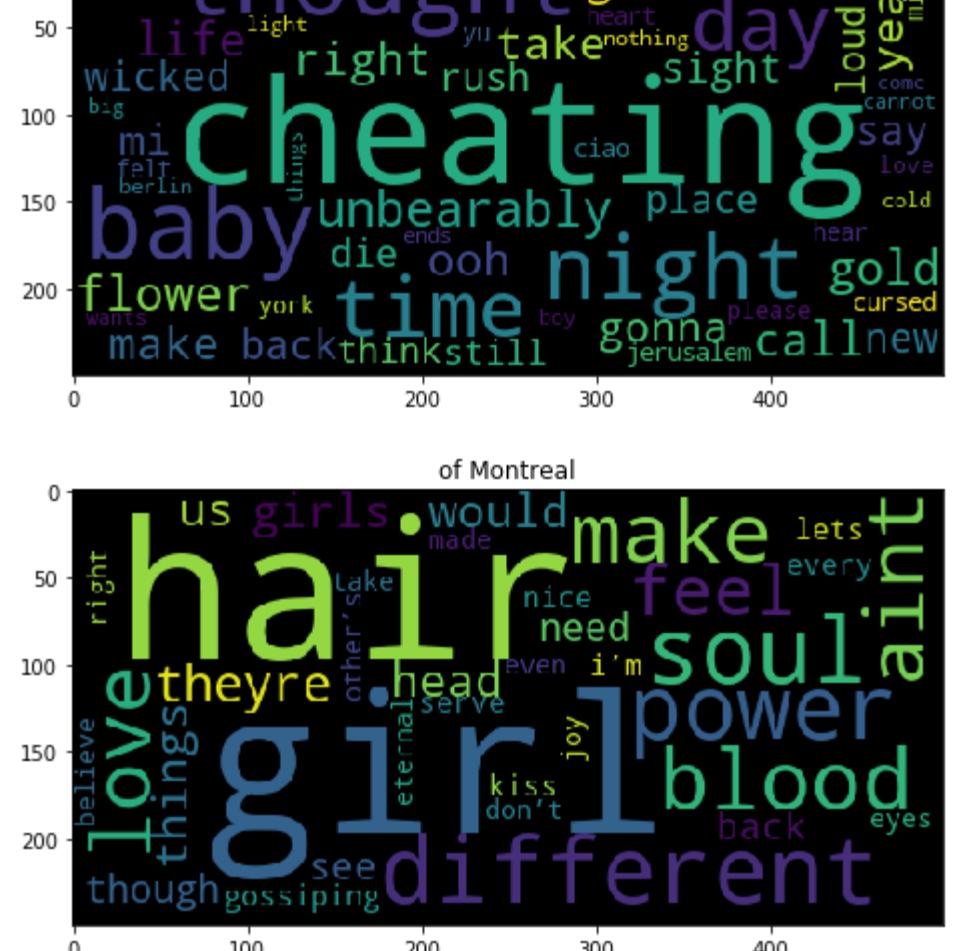
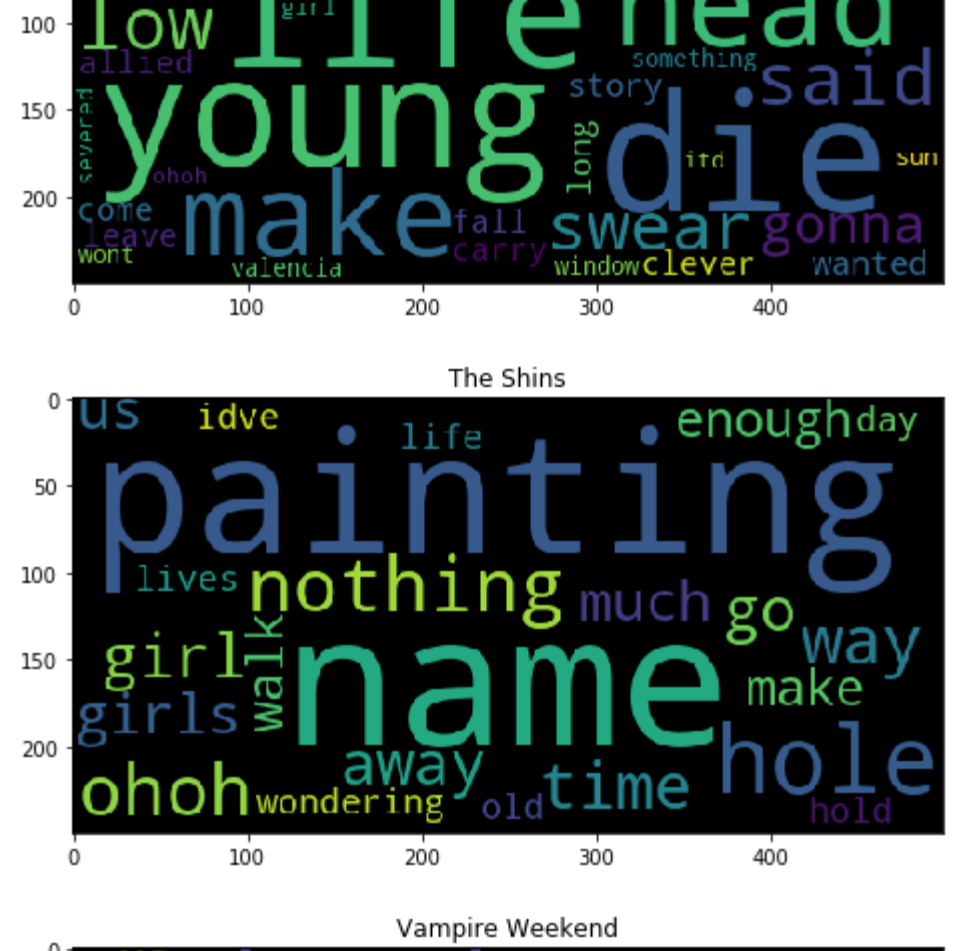
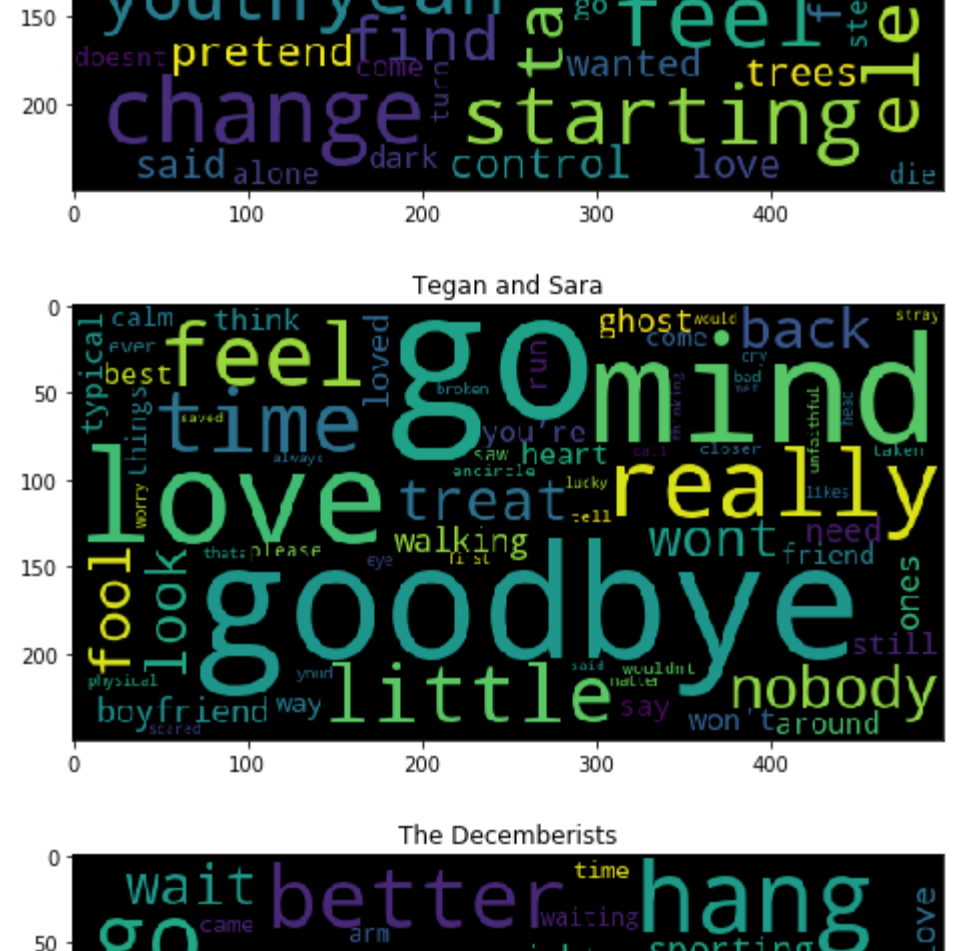
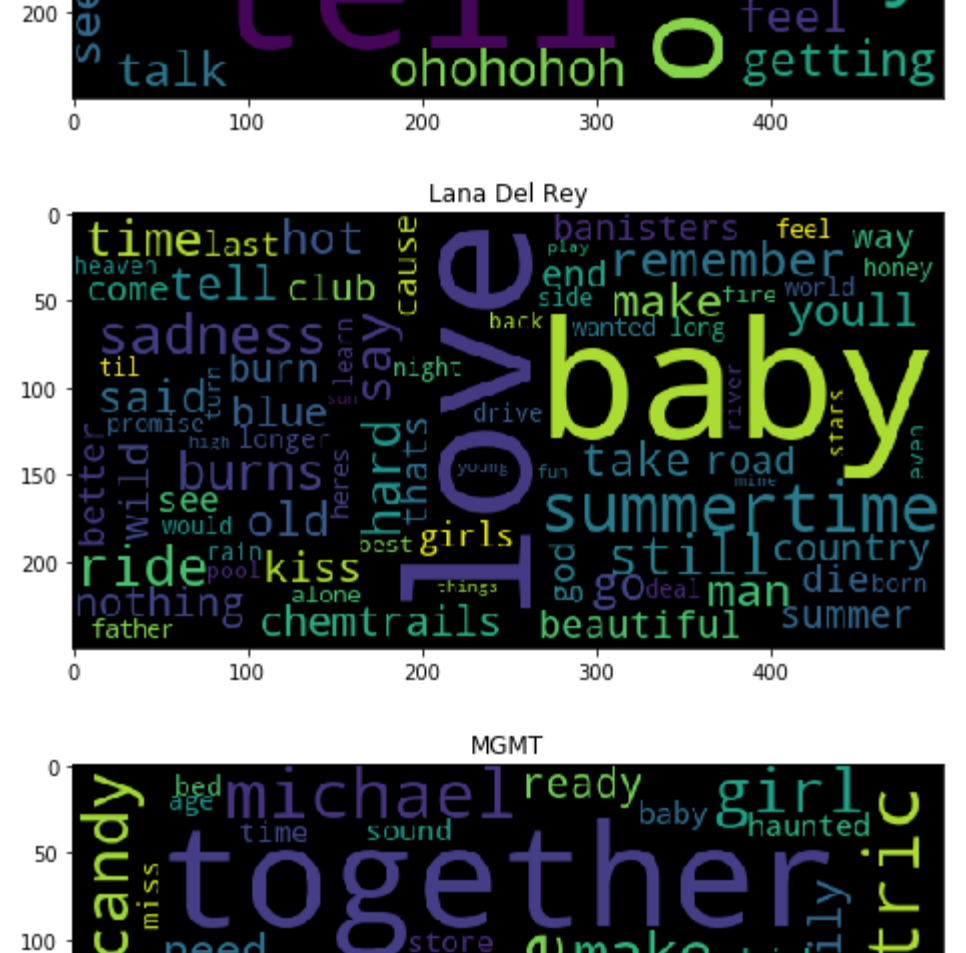
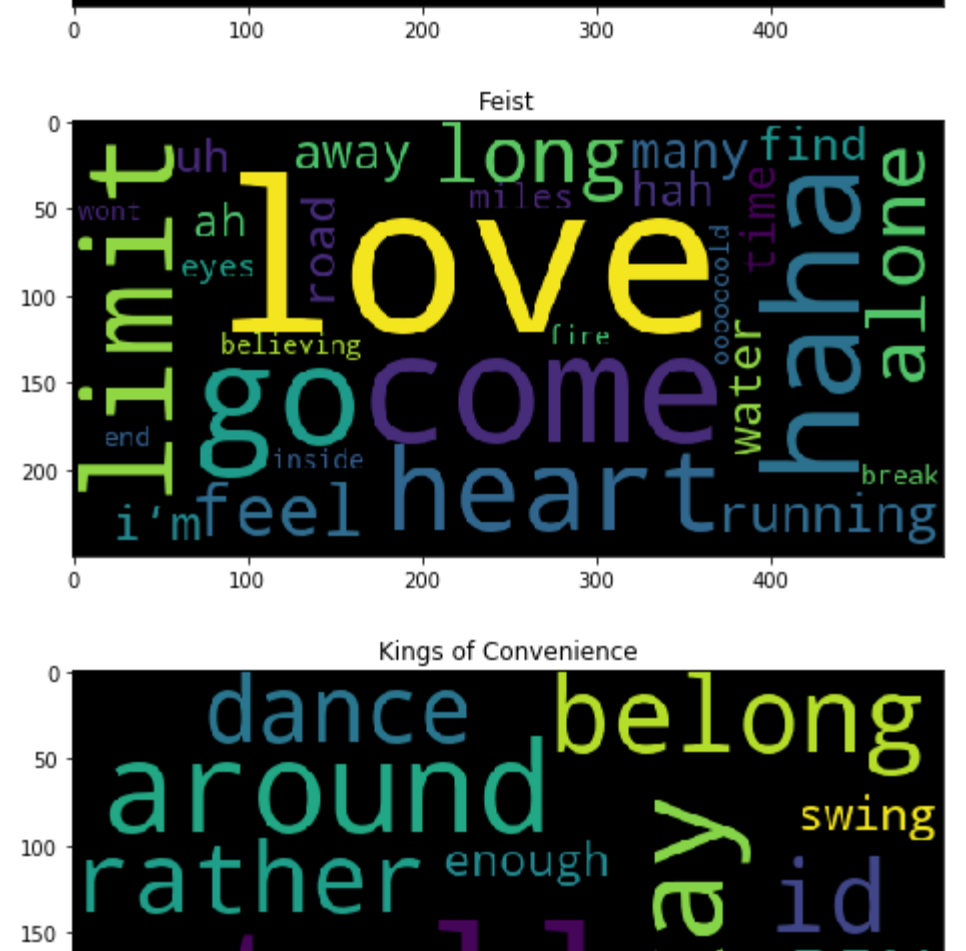



```
In [24]: import matplotlib.pyplot as plt
          from wordcloud import WordCloud

def word_cloud(data):
    c, r = data.shape
    for i in range(c):
        s = data["word_counter"].iloc[i]
        wordcloud = WordCloud(width = 500, height = 250).generate_from_frequencies(s)

        plt.figure(figsize=(8,4))
        plt.imshow(wordcloud)
        plt.title(data.index[i])

wordcloud = word_cloud(word_lyrics)
wordcloud
```



```
id4 = models.LdaModel(corpus=corpus, id2word=id2word, num_topics=3, passes=10)
lda.print_topics(
    {0,
      "0.011\"baby\" + 0.010\"tell\" + 0.006\"see\" + 0.005\"around\" + 0.006\"poor\" + 0.005\"love\" + 0.005\"b\"
      eat\" + 0.005\"tonight\" + 0.003\"jump\" + 0.005\"line\""},
    1),
    {0.015\"love\" + 0.013\"go\" + 0.010\"cheating\" + 0.009\"time\" + 0.007\"mind\" + 0.007\"come\" + 0.006
      \"result\" + 0.006\"back\" + 0.006\"goodbye\" + 0.005\"feel\""},
    2),
    {0,
```

```
lda = models.LdaModel(corpus=corpus, id2word=id2word, num_topics=4, passes=10)
lda.print_topics()

[('0.0171**love + 0.011**go + 0.007**time + 0.007**feel + 0.006**come + 0.005**say + 0.005**bac
  k + 0.005**baby + 0.005**make + 0.005**mind'),
 ('0.013**cheating + 0.010**boy + 0.008**tell + 0.006**thought + 0.006**time + 0.006**dav + 0.0
```

```
["0.024*die" + 0.021*life" + 0.020*young" + 0.011*head" + 0.011*make" + 0.010*hang" + 0.008*g  
o" + 0.007*better" + 0.005*said" + 0.005*low"]],  
(3,  
["0.037*together" + 0.011*change" + 0.010*feel" + 0.009*starting" + 0.008*yeah" + 0.007*electr  
ic" + 0.006*michael" + 0.006*youth" + 0.006*take" + 0.005*girl"]])  
  
lda = models.LdaModel(corpus=corpus, id2word=id2word, num_topics=5, passes=10)
```

```

(10)
    0.0011*“fire” + 0.0103*“die” + 0.0087*“boy” + 0.0087*“young” + 0.0087*“make” + 0.0077*“love” + 0.0077*“be”
    + 0.0067*“girl” + 0.0050*“go” + 0.0050*“see””,
    (1)
    0.0097*“name” + 0.0077*“hole” + 0.0077*“painting” + 0.0077*“nothing” + 0.0067*“chose” + 0.0057*“way” + 0.
    0057*“time” + 0.0047*“way” + 0.0047*“girl” + 0.0047*“away””,
    (2)
    0.0027*“go” + 0.0177*“goodbye” + 0.0167*“mind” + 0.0157*“love” + 0.0127*“really” + 0.0117*“time” + 0.011
    7*“feeling” + 0.0097*“tell” + 0.0097*“body” + 0.0097*“treat””,
    (3)
    0.0013*“cheating” + 0.0137*“baby” + 0.0107*“love” + 0.0087*“tell” + 0.0077*“time” + 0.0077*“way” + 0.006
    7*“summer” + 0.0067*“go” + 0.0067*“night” + 0.0067*“thought””,
    (4)

```

I don't feel I have gathered anything that gives a commonality to each topic. Let's try something different.

Let's create a filter that we only consider the nouns from the lyrics.

First we will need to create a function that uses the `nltk` library to classify all the words in our lyrics and then just get those that are labeled as nouns. We will create a new data frame called nouns which is a copy of our combined lyrics data frame with an added column which contains the nouns contained in the lyrics column.

```
def nouns(text):
    """Given a string of text, tokenize the text and pull out only the nouns."""
    is_noun = lambda pos: pos[2:] == "NN"
    tokenized = word_tokenize(text)
    all_nouns = [word for (word, pos) in pos_tag(tokenized) if is_noun(pos)]
    return ' '.join(all_nouns)

data_nouns = combined_lyrics.copy()
data_nouns['nouns'] = data_nouns['lyrics'].apply(lambda x: nouns(x))
data_nouns
```

```
Out[79]:
```

artist	lyrics	nouns
--------	--------	-------

1	Faist	one two three four five me that you love me ...	tell nights you teenage hours do not nothing ...
2	Kings of Convenience	if you want to be my friend/y you want us to ge...	friend please observation in demeanor friend ...
3	Lana Del Rey	there's pictures on the wall/no in a p... ..	there picture wall john beer he oklahoma ...
4	MGMT	im feeling rough im feeling raw/im in the j... ..	im in im prime life feels music many moments w...
5	Tegan and Sara	i all want to get a little bit closer/all ...	i i closer breath bit rush doors wind whir...
6	The Decemberists	it was late one night i was counting out my st...	right stitches side read time brings thim...
7	The Shins	girl info me all my senses warm mellow/ye co...	girl senses eyes backwards shifts eye lips ...
8	Vampire Weekend	baby i know pain is as natural as the rain/... ..	baby pain rain rain california baby int...

```
In [80]: from sklearn.feature_extraction import Text
         from sklearn.feature_extraction.text import CountVectorizer

cvm = CountVectorizer(stop_words=stopwords)
data_cvm = cvm.fit_transform(data_rooms.rooms)

data_stopn = pd.DataFrame(data_cvm.toarray(), columns=cvm.get_feature_names())
data_stopn.index = word_lyrics.index

data_stopn
```

	Belle and Sebastian	1	0	0	0	0	0	1	0	2	...	0	2	
	Feist	0	0	0	0	0	1	0	0	0	0	...	0	0
	Kings of Convenience	0	0	0	1	0	0	0	0	0	0	...	0	0
	Lana Del Rey	0	0	0	0	0	0	0	0	0	0	...	0	0
	MGMT	0	0	0	0	0	0	0	0	0	0	...	0	0
	Tegan and Sara	0	0	0	0	0	0	1	0	0	0	...	0	0

DecemBERTes										
The Shirts	0	0	0	1	0	0	0	0	0	0
Wampse Weekend	0	0	0	0	0	0	0	4	0	0
of Montreal	0	0	0	0	1	0	0	0	0	1

10 rows x 1008 columns

```
In [81]: corpusn = matutils.SparseCorpus(scipy.sparse.csr_matrix(data_stopn.transpose()))
id2wordn = dict((v, k) for k, v in cvn.vocabulary.items())
```

Again let's run lda for the nouns of the lyrics with 2, 3, and 4 topics.

```
In [82]: ldan = models.LdaModel(corpus=corpusnum, num_topics=2, id2word=id2words, passes=10)
ldan.print_topics()

Out[82]: [0,
          '0.017*+time+ 0.016*+life+ 0.013*+love+ 0.011*+baby+ 0.011*+mind+ 0.010*+heart+ 0.008
          *+night+ 0.008*+human+ 0.007*+nobody+ 0.007*+go+',
          '0.012*+boy+ 0.010*+girl+ 0.008*+time+ 0.007*+day+ 0.007*+way+ 0.006*+nothing+ 0.006*+n
          ame+ 0.006*+jump+ 0.005*+life+ 0.005*+heart*']
```

```
[Out]: [0, ('0.021*yearh + 0.011*youth + 0.011*fzeli' + 0.009*girl' + 0.008*family + 0.008*kid' + 0.008*candy' + 0.007*control' + 0.007*treats' + 0.007*storen'),  
      (1,  
       '0.021*life' + 0.011*time' + 0.011*heave' + 0.010*girl' + 0.010*way' + 0.008*nothing' + 0.008*nights' + 0.007*wan' + 0.007*things' + 0.007*power'),  
      (2,  
       '0.015*time' + 0.015*love' + 0.012*boy' + 0.011*mind' + 0.010*heart' + 0.009*way' + 0.008*baby' + 0.009*go' + 0.008*goodbye' + 0.007*numbertime')]
```

```
In [84]: ldan = models.LdaModel(corpus=corpus_n, num_topics=4, id2word=id2word_n, passae=10)
```

```
Out[4]: [{0,
  0.015*happy + 0.015*love + 0.004*summertime + 0.013*time + 0.011*sadness + 0.010*road +
  0.010*happy + 0.009*happy + 0.009*nothing + 0.008*time*},
 {0,
  0.019*love + 0.013*joy + 0.011*time + 0.009*girl + 0.008*head + 0.007*war + 0.007*da
  y + 0.007*name + 0.006*time + 0.006*nothing*},
 {0,
  0.024*mind + 0.022*time + 0.017*joy + 0.016*goodbye + 0.016*want + 0.014*love + 0.013*ti
  ght + 0.012*heart + 0.011*nothing + 0.011*happy*},
 {0,
  0.016*hair + 0.014*power + 0.012*girl + 0.012*blood + 0.010*soul + 0.009*head + 0.009
  *things + 0.009*girls + 0.008*love + 0.006*joy*}]
```

Now let's consider nouns and adjectives only.

```
data_nouns_adj
```

```
Out[185]:
```

	artist	lyrics	nouns	nouns_adj
0	Beke and Sebastian	a mile and a half on a bus takes a long time	mile half bus time odour prison food time	mile half bus time odour prison food time
1	Feist	one two three four/tell me that you love me	tell nights youth teenage hopes door nothing ...	tell nights youth teenage hopes door nothing ...
2	Kings of Convenience	if you want to be my friend/you want us to go	friend please observation in demand friend m	friend please observation in demand friend
3	Lana Del Rey	there's a picture on the wall/no me on a	there's picture will john beer hell oklahoma	there's picture will john beer hell oklahoma

4	MGMT	in hearing room/in hearing hallway	in prime time last night	in prime time last night	in prime time last night
5			models will	models will	models will
6	Teagan and Sara	all it want to get is a little bit closer	I bit I close breath bit nush doors wind high...	I bit I close breath bit nush doors wind high...	I bit I close breath bit nush doors wind high...
7	The Decembertists	it was late one night I was counting out my st...	night I stitches side road I time breaths	night I stitches side road I time breaths	night I stitches side road I time breaths
8	The Shins	girl info me all my senses waim mairgnoe eye	girl senses eyes backwards signa spit eye	girl senses eyes backwards signa spit eye	girl senses eyes backwards signa spit eye
9	The Vampire Weekend	baby I know pain is as natural as the surviv...	baby I pain rain rain rain rain baby I sirt...	baby I pain rain rain rain rain baby I sirt...	baby I pain rain rain rain rain baby I sirt...
9	of Montreal	you're such a mystery/in just want to stand i...	you're mystery/ i sea ocean hair I knot blues	you're mystery/ i sea ocean hair I knot blues	you're mystery/ i sea ocean hair I knot blues

```
In [86]: cvna = CountVectorizer(stop_words=stopwords, max_df=.8)
data_cvna = cvna.fit_transform(data_points_max_norm_adj)
data_dtmsa = pd.DataFrame(data_cvna.toarray()), columns=cvna.get_feature_names()
data_dtmsa.index = data_points_adj['artist']
data_dtmsa

Out[86]:
```

	aberration	abrades	accusations	act	acts	admit	advantage	affair	affection	afraid	...	yesterday	york	you2me2bedshu
Belle and Sebastian	1	0	0	0	0	0	0	1	0	2	...	0	2	

[illegible]

```
of Montreal      0 0          0 0 1 0          0 0          0 0 ...      1 0
```

```
10 rows x 1597 columns
```

```
In [87]: corpusna = matutils.Sparse2Corpus(scipy.sparse.csr_matrix(data_dmtna.transpose()))  
         id2wordna = dict((v, k) for k, v in corpusna.vocabulary.items())
```

```
In [88]: ldana = models.LdaModel(corpus=corpusna, num_topics=2, id2word=id2wordna, passes=10)  
         ldana.print_topics()
```

```
Out[88]: [(0,  
           [0,
```

```
(  
    "0.015*mindr + 0.012*goe + 0.011*goodyear + 0.009*moody + 0.007*ones + 0.007*bahe + 0.007  
    *yeah" + 0.006*nothing" = 0.006*foo", 0.006*byfriend"]])  
  
In [92]: idna = models.IdaModel(corpus=corpus, num_topics=3, id2word=id2words, passes=10)  
         idna.print_topics()  
  
Out[92]: [0,0,  
    "0.019*baby" + 0.009*summertime + 0.008*nothing" + 0.007*xan + 0.007*sadness + 0.007*girl  
    + 0.006*daily + 0.006*yearyear + 0.006*men + 0.006*midnight",  
    (  
    "0.017*hair + 0.011*bahar + 0.010*power" + 0.010*init + 0.009*girl" + 0.009*blood" + 0.008  
    *goe + 0.008*goodyear + 0.008*moody + 0.007*ones + 0.007*bahe + 0.007*yeah" + 0.006*nothing"  
    = 0.006*foo", 0.006*byfriend"]])
```

```
[0.015 "good", 0.012 "mind", 0.012 "have", 0.008 "nobody" + 0.010 "go" + 0.010 "goodbye" + 0.008 "day" + 0.007 "name" + 0.007 "beat" + 0.006 "nothing")]
```

```
In [90]: ldms = models.LdaModel(corpus=corpusna, num_topics=4, id2word=id2wordna, passes=10)
ldms.print_topics()
```

```
Out[90]: [(0,
            "0.020*head + 0.009*hang" + 0.008*gon" + 0.008*story" + 0.008*nobody" + 0.008*dance" + 0.007*something + 0.007>window" + 0.007*valencia" + 0.006*girl"),
            (1,
            "0.013*baby" + 0.010*day" + 0.010*wan" + 0.008*haha" + 0.008*girl" + 0.008*yeah" + 0.008*moo"
            na + 0.008*limit" + 0.008*youth" + 0.007*nothing"),
            (2,
```

```
{0.0007: "cause", 0.0007: "and", 0.0007: "burns", 0.0007: "and", 1},
{0.0181: "boy", 0.0151: "mind", 0.0142: "goodbye", 0.0142: "go", + 0.0093: "nobody", + 0.0093: "jump" + 0.008
"girl" + 0.008: "head" + 0.008: "heart" + 0.008: "parcys"]}
```

Eventually I had failed to determine what each topic means, let us understand what topic each artist would belong to if we were to use the nouns and adjective method with four topics.

```
In [91]: corpus_transformed = Idnaa[corpusna]
list(zip([a for (a,b) in corpus_transformed, data_stop.index()])
```

```
Out[91]: [{3, 'Belle and Sebastian'},
```

(2, 'Lana Del Rey'),
(1, 'MGMT'),
(3, 'Tegan and Sara'),
(0, 'The Decemberists'),
(2, 'The Shins'),
(1, 'Vampire Weekend'),
(3, 'of Montreal').

I have failed to find what the topics mean for all the different methods I have used:

1. Using words that excluded stop words (and others I included)
2. Including stop words

In other words I have not found a method where the topics made sense on what they were representing.

Conclusion: Artists that belong to the indie/pop genre are very similar to each other.

I initially set it out for which indie/pop artists were similar to Kings of Convenience. However, my conclusion is that most of the artists from that genre are too similar to try to find differences among them. I reach this conclusion because they had very similar ratings for polarity and subjectivity. Also, because using LDA (with three different methods) I could not find any clear topics that were shared by some artists and excluded by others. My next experiment will be trying this notebook with artists from different genres.

Limitations

As mentioned in the conclusion some limitations include that artists were too similar to each other as they were all from the same genre. Also, I might get better results if I include more artists and more songs for each artist.

References

https://www.youtube.com/watch?v=xvnsFTUeOmc&ab_channel=PyOhiq