

```

In [1]: 1 from collections import defaultdict
2 import numpy as np
3 import pandas as pd
4
5 import matplotlib_venn as venn
6 from matplotlib_venn import venn2, venn2_circles, venn3, venn3_circles
7 import matplotlib.pyplot as plt
8 %matplotlib inline
9
10 from sklearn.model_selection import train_test_split
11 import nltk
12 from nltk.corpus import stopwords, wordnet
13 from nltk import pos_tag
14 from nltk.stem import WordNetLemmatizer
15 from nltk.tokenize import regexp_tokenize, word_tokenize, RegexpTokenizer
16
17 nltk.download('stopwords')
18 nltk.download('averaged_perceptron_tagger')
19 nltk.download('wordnet')
20
21 %load_ext autoreload
22 %autoreload 2

```

executed in 1.59s, finished 18:50:04 2021-06-07

```

[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\elena\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] C:\Users\elena\AppData\Roaming\nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\elena\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!

```

```

In [2]: 1 # These functions are taken from one of the NPL Lectures at Flatiron School
2
3 sw = stopwords.words('english')
4 def get_wordnet_pos(treebank_tag):
5     ...
6     Translate nltk POS to wordnet tags
7     ...
8     if treebank_tag.startswith('J'):
9         return wordnet.ADJ
10    elif treebank_tag.startswith('V'):
11        return wordnet.VERB
12    elif treebank_tag.startswith('N'):
13        return wordnet.NOUN
14    elif treebank_tag.startswith('R'):
15        return wordnet.ADV
16    else:
17        return wordnet.NOUN
18
19 def doc_preparer(doc, stop_words=sw):
20     ...
21
22     :param doc: a document from the satire corpus
23     :return: a document string with words which have been
24             lemmatized,
25             parsed for stopwords,
26             made lowercase,
27             and stripped of punctuation and numbers.
28     ...
29
30     regex_token = RegexpTokenizer(r"([a-zA-Z]+(?:'[a-z]+)?)")
31     doc = regex_token.tokenize(doc)
32     doc = [word.lower() for word in doc]
33     doc = [word for word in doc if word not in sw]
34     print(doc)
35     doc = pos_tag(doc)
36     doc = [(word[0], get_wordnet_pos(word[1])) for word in doc]
37     lemmatizer = WordNetLemmatizer()
38     doc = [lemmatizer.lemmatize(word[0], word[1]) for word in doc]
39     return ' '.join(doc)

```

executed in 62ms, finished 18:50:04 2021-06-07

```

In [3]: 1 # This part of the code is specific to the venn diagram demo of corpus word split
2 corpus = pd.read_csv('satire_nosatire.csv')
3
4 X = corpus.body
5 y = corpus.target
6
7 X_train, X_test, y_train, y_test = train_test_split(X,
8                                                     y,
9                                                     random_state=42,
10                                                    test_size=0.25)
11
12 token_docs = [doc_preparer(doc, sw) for doc in X_train]
13
14 X_t, X_val, y_t, y_val = train_test_split(token_docs, y_train,
15                                           test_size=0.25, random_state=42)
16
17 df_X_t=pd.DataFrame(X_t)
18 df_y_t=pd.DataFrame(y_t)
19 df_train=pd.concat([df_X_t, df_y_t], axis=1, join="inner")
20 df_train.rename(columns={0: 'phrase'})
21
22 list_satire_phrases=[]
23 df_train_satire=df_train.loc[df_train.target==1]
24 L = [''.join(df_train_satire[x].astype(str)) for x in df_train_satire]
25 df_one_phrase = pd.DataFrame([L], columns=df_train_satire.columns)
26 df_one_phrase=df_one_phrase.rename(columns={0: 'long_phrase'})
27 satire_words=list(df_one_phrase.long_phrase)
28 list_of_satire_words=satire_words[0].split()
29 set_of_satire_words=set(list_of_satire_words)
30
31 list_news_phrases=[]
32 df_train_news=df_train.loc[df_train.target==0]
33 L = [''.join(df_train_news[x].astype(str)) for x in df_train_news]
34 df_one_phrase_news = pd.DataFrame([L], columns=df_train_news.columns)
35 df_one_phrase_news=df_one_phrase_news.rename(columns={0: 'long_phrase'})
36 news_words=list(df_one_phrase_news.long_phrase)
37 list_of_news_words=news_words[0].split()
38 set_of_news_words=set(list_of_news_words)

```

executed in 13.5s, finished 18:50:17 2021-06-07

['perpetually', 'offended', 'social', 'justice', 'warrior', 'keen', 'activist', 'socialist', 'matters', 'ranging', 'lgbtq', 'p', 'affairs', 'feminism', 'borders', 'soviet', 'ideology', 'marxism', 'ideal', 'communist', 'state', 'hugh', 'mungus', 'state', 'perpetual', 'offence', 'wake', 'morning', 'first', 'thing', 'hear', 'words', 'good', 'morning', 'shout', 'roommate', 'form', 'racist', 'linguistic', 'capitalist', 'imperialist', 'sexist', 'white', 'male', 'created', 'offensive', 'greeting', 'assumes', 'good', 'morning', 'western', 'capitalist', 'bourgeois', 'society', 'people', 'africa', 'good', 'morning', 'shit', 'morning', 'live', 'corrugated', 'iron', 'shack', 'walk', 'five', 'hours', 'fill', 'bucket', 'full', 'dirty', 'muddy', 'water', 'fucking', 'drink', 'mr', 'mungus', 'attends', 'berkeley', 'college', 'california', 'outraged', 'demos', 'society', 'biology', 'science', 'offensive', 'well', 'concept', 'offended', 'offensive', 'outraged', 'offended', 'offence', 'offensive', 'hateful', 'perpetual', 'state', 'concept', 'offended', 'offending', 'offended', 'offensive', 'stance', 'state', 'offence', 'offending', 'offended', 'doubt', 'mind', 'creation', 'capitalist', 'racist', 'sexist', 'system', 'imprisons', 'people', 'offended', 'suspect', 'created', 'offence', 'first', 'place', 'everything', 'offensive', 'perpetually', 'state', 'offence', 'right', 'offended', 'concept', 'offended', 'offence', 'offensive', 'manner']

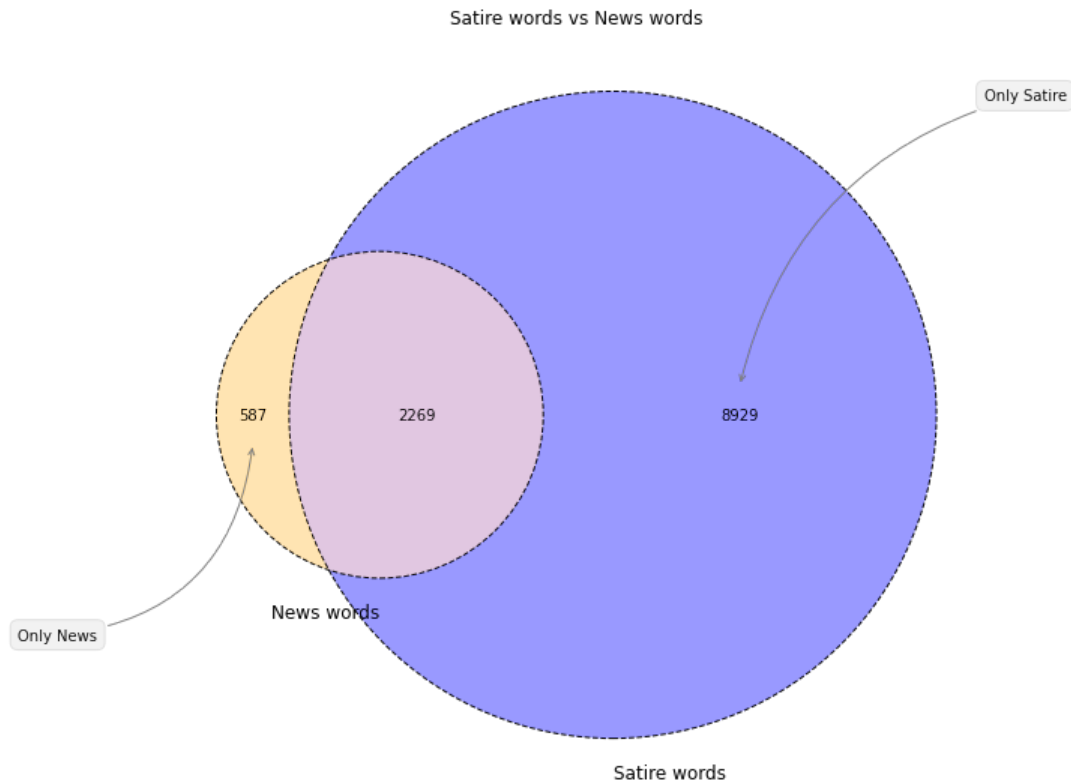
['ana', 'luz', 'sister', 'law', 'ronald', 'blanco', 'looked', 'grimly', 'neighbours', 'murdered', 'honduran', 'man', 'was', 'hed', 'away', 'rills', 'blood', 'left', 'bullet', 'ridden', 'body', 'lain', 'outside', 'house', 'troubled', 'barrio', 'outskirts', 'tegucigalpa', 'one', 'many', 'scenes', 'witnessed', 'year', 'assignment', 'honduras', 'thousands', 'people', 'sought', 'escape', 'violence', 'poverty', 'joining', 'migrant', 'caravan', 'hope', 'making', 'safety', 'across', 'mexico', 'u', 'border', 'problems', 'small', 'central', 'american', 'country', 'grabbed', 'international', 'attention', 'u', 'president', 'donald', 'trump', 'cracked', 'illegal', 'immigration', 'honduras', 'years', 'one', 'world's', 'murderous', 'countries', 'though', 'official', 'data', 'show', 'homicide', 'rate', 'fallen', 'sharply', 'continues', 'highly', 'challenge']

```

In [4]: 1 plt.figure(figsize=(10, 10))
2
3
4 sets=[set_of_news_words, set_of_satire_words]
5 labels=('News words', 'Satire words')
6
7 v=venn2(sets, set_labels = labels, set_colors=("orange", "blue"))
8
9 v.get_patch_by_id('10').set_alpha(0.3)
10
11
12 venn2_circles(subsets=sets,
13               linestyle="dashed", linewidth=1)
14
15 plt.annotate('Only News',
16             xy=v.get_label_by_id('10').get_position() - np.array([0, 0.05]), xytext=(-130,-130),
17             ha='center', textcoords='offset points', bbox=dict(boxstyle='round, pad=0.5', fc='gray', alpha=0.1),
18             arrowprops=dict(arrowstyle='->', connectionstyle='arc3,rad=0.4',color='gray'))
19
20 plt.annotate('Only Satire',
21             xy=v.get_label_by_id('01').get_position() - np.array([0, -0.05]), xytext=(190,190),
22             ha='center', textcoords='offset points', bbox=dict(boxstyle='round, pad=0.5', fc='gray', alpha=0.1),
23             arrowprops=dict(arrowstyle='->', connectionstyle='arc3,rad=0.3',color='gray'))
24
25
26 plt.title('Satire words vs News words')
27 plt.show()

```

executed in 223ms, finished 18:50:18 2021-06-07



In [ ]:	1	
---------	---	--