

Assignment 3

Import libraries

```
In [22]: import pandas as pd
import numpy as np

import nltk
import html
nltk.download('stopwords')
nltk.download('wordnet')
from nltk.corpus import stopwords
from nltk.tokenize import RegexpTokenizer
from nltk.stem import PorterStemmer, WordNetLemmatizer
import re
import unicodedata
from bs4 import BeautifulSoup
from sklearn.feature_extraction.text import CountVecorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from collections import Counter
import seaborn as sns
import matplotlib.pyplot as plt
from wordcloud import WordCloud
from sklearn.preprocessing import LabelEncoder
from statistics import mean
from sklearn.model_selection import train_test_split, StratifiedFold
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.model_selection import GridSearchCV

from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn import svm
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.multiclass import OneVsRestClassifier
import warnings
warnings.filterwarnings('ignore')

[nltk_data] Downloading package stopwords to
[nltk_data] /Users/angelabau/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] /Users/angelabau/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

1. Data cleaning

- All HTML tags and attributes are removed. I imported BeautifulSoup package, an HTML-parsing libraries, and used 'html.parser' as an implementation of HTML-parsing.
- Html character codes are replaced with an ASCII equivalent.
- All URLs and @ are removed.
- All characters in the text are in lowercase.
- All stop words are removed. I imported stopwords package from nltk.corpus which contains all the stopwords.
- If a tweet is empty after pre-processing, it should be preserved as such.
- Lemmaization is the process of converting a word to its base form. Wordnet is an large, freely and publicly available lexical database for the English language aiming to establish structured semantic relationships between words. In order to lemmatize, you need to download Wordnet, create an instance of the WordNetLemmatizer() and call the lemmatize() function on a single word.

(a) Clean sentiment data

```
In [23]: df_sentiment = pd.read_csv('sentiment_analysis.csv')
df_sentiment

Out[23]:
```

	ID	text	label
0	7.680980e+17	Josh Jenkins is looking forward to TAB Breeder...	1
1	7.680980e+17	RT @MianUsmanIwade: Congratulations Pakistan o...	1
2	7.680980e+17	RT @PEPalerts: This September, @EYsma is tak...	1
3	7.680980e+17	RT @david_gabris: Newly painted walls, thanks...	1
4	7.680980e+17	RT @CedricFischer: Excited to announce as o...	1
...
550386	8.046170e+17	@goddesses_o i can't stop watching her...mm.M...	0
550387	8.046180e+17	Poor old Tom Odell doesn't look like he would...	0
550388	8.046180e+17	Antsmasher i smashed 7 ants in this awesome...	1
550389	8.046180e+17	@LizHudson @KymWyllie @Evasmless @meanBok @L...	1
550390	8.046190e+17	Bixbeat Mixtape Vol2 is here with great arts...	1

550391 rows x 3 columns

```
In [3]: tokenizer = RegexpTokenizer(r'\w+')
lemmatizer = WordNetLemmatizer()
def clean_data_1(text):
    # convert 'RT @' to 'R' and Remove the word after @
    rt = re.compile('RT @').sub('R', str(text))
    text = re.sub('RT @', str(text))
    at = re.compile('@(?=\w+)(\w+)').sub('at', str(text))
    text = re.sub('at', str(text))

    #remove html tags and attributes
    text = BeautifulSoup(text, 'html.parser').get_text()

    #replace ASCII
    text = html.unescape(text)
    text = unicodedata.normalize('NFKD', text).encode('ascii', 'ignore').decode('utf-8')

    #remove URLs
    text = re.sub(r'\w+:\/\/[2]([d]w-|+)(\.[d]w-|+)*(?:[?]/[^\s/]+)*', '', text)

    #lowercase
    text = text.lower()

    #remove stop words
    text = [word for word in tokenizer.tokenize(text) if word not in stopwords.words('english')]

    #lemmatize words
    text = [lemmatizer.lemmatize(word, pos='v') for word in text]
    text = [lemmatizer.lemmatize(word, pos='n') for word in text]
    return text
```

```
In [4]: %time
df_sentiment['text']=list(map(clean_data_1, df_sentiment['text']))
CPU times: user 13min 1s, sys: 3min 18s, total: 16min 20s
Wall time: 17min 33s
```

```
In [5]: df_sentiment

Out[5]:
```

	ID	text	label
0	7.680980e+17	[josh,jenkins,look,forward,tab,breeder,c...]	1
1	7.680980e+17	[congratulation,pakistan,become,not,teasem...]	1
2	7.680980e+17	[september,take,maine,mendozas,surprise,t...]	1
3	7.680980e+17	[newly,paint,wall,thank,million,custodial...]	1
4	7.680980e+17	[excite,announce,july,2017,feschotte,lab...]	1
...
550386	8.046170e+17	[stop,watch,mm]	0
550387	8.046180e+17	[poor,old,tom,odell,look,like,would,kno...]	0
550388	8.046180e+17	[antsmasher,smash,7,ant,awesome,game,hfl...]	1
550389	8.046180e+17	[morning,girl,wonderful,friday]	1
550390	8.046190e+17	[bixbeat,mixtape,vol,2,great,artiste,joi...]	1

550391 rows x 3 columns

(b) Clean elections data

For Canadian election data, I added one more cleaning which is to remove b/b/ at the beginning of the tweets and the \n in the tweets.

```
In [6]: df_election = pd.read_csv('Canadian_elections_2019.csv')
df_election

Out[6]:
```

	sentiment	negative_reason	text
0	negative	Women Reproductive right and Racism	b"/RosieBarton So instead of your suggestion, ...
1	positive		b"/AllWomenSpacewalk it's really@SpaceStatio...
2	negative	Economy	b"/Brantford It's going to cost YOU \$94 BILLIO...
3	positive		b"/Canada #CanadaElection2019 #CanadaVotes #n...
4	negative	Economy	b"/Canada #taxpayers are sick &mp; tired of h...
...
2128	negative	Scandal	b"/You know he was doing a good enough job smea...
2129	negative	Scandal	b"/You missed the comment. Deflecting the issue...
2130	positive		b"/Your daily reminder/n@theJagmeetSingh endor...
2131	negative	Women Reproductive right and Racism	b"/Yup...not going to reopen the abortion debat...
2132	positive		b"/Zing/n#NDP #eln43 https://t.co/xo6Qe4y/3p

2133 rows x 3 columns

```
In [9]: def clean_data_2(text):
    #remove unknown characters b and \n
    text = text.replace('\n','')
    text = re.sub(r'([b|B])', '', text)

    #remove emoji
    emoji = re.compile(r'[\x]([a-z\d]{2})')
    text = re.sub(emoji, '', text)

    # Remove the word after @
    text = re.compile('RT @').sub('R', str(text))
    at = re.compile('@(?=\w+)(\w+)').sub('at', str(text))
    text = re.sub('at', str(text))

    #remove html tags and attributes
    text = BeautifulSoup(text, 'html.parser').get_text()

    #replace ASCII
    text = html.unescape(text)
    text = unicodedata.normalize('NFKD', text).encode('ascii', 'ignore').decode('utf-8')

    #remove URLs
    text = re.sub(r'\w+:\/\/[2]([d]w-|+)(\.[d]w-|+)*(?:[?]/[^\s/]+)*', '', text)

    #lowercase
    text = text.lower()

    #remove stop words
    text = [word for word in tokenizer.tokenize(text) if word not in stopwords.words('english')]

    #lemmatize words
    text = [lemmatizer.lemmatize(word, pos='v') for word in text]
    text = [lemmatizer.lemmatize(word, pos='n') for word in text]
    return text
```

```
In [10]: df_election['text']=list(map(clean_data_2, df_election['text']))

In [11]: df_election
```

```
Out[11]:
```

	sentiment	negative_reason	text
0	negative	Women Reproductive right and Racism	[instead,suggestion,agree,canadian,woman,...]
1	positive		[allwomanspacewalk,real,etobiconorth,city...]
2	negative	Economy	[brantford,go,cast,94,billion,next,4,ye...]
3	positive		[canada,canadaelection2019,canadavotes,eln...]
4	negative	Economy	[canada,taxpayer,sick,tire,hard,earn,don...]
...
2128	negative	Scandal	[know,good,enough,job,smear,campaign,rig...]
2129	negative	Scandal	[miss,comment,deflect,issue,answer,best...]
2130	positive		[daily,reminder,endorse,strategic,vote,el...]
2131	negative	Women Reproductive right and Racism	[yup,go,reopen,abortion,debate,eln43,sc...]
2132	positive		[zing,ndp,eln43]

2133 rows x 3 columns

2. Exploratory data analysis

There are three political parties in consideration which are Liberal, Conservative and NDP. To determine the political affiliation, check whether the tweet contains any of these identifiers related to certain political parties or candidates. The identifiers are listed as below.

(a) Create keywords lists for political parties

For each political party, there is a list that contain the relevant keywords which can reflect political preference. The keywords are determined by leader, slogan, hashtag.

Relevant Information for Different Political Parties:

- Liberal**
 - Leader: Justin Trudeau
 - Twitter Account: @JustinTrudeau
 - Slogan: 'chooseforward'
 - Hashtag: #lpc
 - Identifiers: "Justin Trudeau"; "Justin"; "Choose forward"; "Proven Leadership for a strong Canada"; "Safer Canada/Stronger Economy"; "Protect our Economy"; "liberal"; "#LAvScan"; "Trudeau must go"; "TrudeauMustGo"; "#LiberalsMustGo"; "#TrudeauBlackface"; "#NotAsAdvised"; "#LiberalLiesAndDeception";
- Conservative**
 - Leader: Andrew Scheer
 - Twitter Account: @AndrewScheer
 - Slogan: 'It's time for you to get ahead'
 - Hashtag: #cpc
 - Identifiers: "Andrew Scheer"; "Andrew"; "Scheer"; "#Scheer4PM"; "#ScheerMajority"; "It's time for you to get ahead"; "conservative"; "Real Change (Now)"; "#ScheerLies"; "#ScheerStupidity"; "#ScheerDelusion"; "#Scheer2019"
- New Democratic**
 - Leader: Jagmeet Singh
 - Twitter Account: @thelagmeetSingh
 - Slogan: 'in it for you'
 - Hashtag: #ndp
 - Identifiers: "Jagmeet Singh"; "Jagmeet"; "Singh"; "We fight for you"; "In it for you"; "Ready for Change"; "#SinghIsSinging"; "Jaggersnaut"; "#UpRiSiSingh"
- None**
 - Includes Bloc Quebecois and Green parties

(b) Return only one party or return none

Since there may be cases that multiple parties are mentioned, I need to compare which party is mentioned most frequently and return only one political party. Also, there may be cases where tweet exhibits no political affiliation and thus return None to show no political interest.

```
In [12]: def party(text):
    #define the words that show political inclination
    l_ls = ('trudeau', 'justin', 'justintrudeau',
            'liberal', 'liberals', 'red',
            'safecanada', 'strongereconomy', 'proteceforward',
            'teamtrudeau', 'lpc', 'voteliberal', 'chooseforward')

    c_ls = ('andrew scheer', 'conservatives', 'cpc', 'andrew',
            'scheer', 'scheer2019', 'scheermajority', 'scheer4pm',
            'antismscheer', 'conservative', 'voteconservative', 'getahead',
            'it is time for you to get ahead')

    n_ls = ('new democratic party', 'ndp', 'tpndp', 'thejagmeet',
            'jagmeet', 'singh', 'jagmeet Singh', 'democratic', 'democrats',
            'in it for you', 'singhupswing', 'uprisingh', 'we fight for you', 'jaggersnaut',
            'uprisingsingh')

    mention = []

    #identify the political party based on lists of keywords
    for word in text:
        if word in l_ls:
            mention.append('liberal')
        if word in c_ls:
            mention.append('conservatives')
        if word in n_ls:
            mention.append('new democratic party')

    # If there are multiple parties mentioned, return the party mentioned most frequent
    if len(mention) > 1:
        most_frequent = max(set(mention), key=lambda x: mention.count(x))
        mention = [most_frequent]

    # If there is no party mentioned in the list, return None
    elif len(mention) == 0:
        mention.append('None')

    #return value rather than list
    return mention[0]
```

```
In [13]: # add a column to df_election to show political affiliation
df_election['political_party'] = ''
df_election
```

```
Out[13]:
```

	sentiment	negative_reason	text	political_party
0	negative	Women Reproductive right and Racism	[instead,suggestion,agree,canadian,woman,...]	
1	positive		[allwomanspacewalk,real,etobiconorth,city...]	
2	negative	Economy	[brantford,go,cast,94,billion,next,4,ye...]	
3	positive		[canada,canadaelection2019,canadavotes,eln...]	
4	negative	Economy	[canada,taxpayer,sick,tire,hard,earn,don...]	
...
2128	negative	Scandal	[know,good,enough,job,smear,campaign,rig...]	
2129	negative	Scandal	[miss,comment,deflect,issue,answer,best...]	
2130	positive		[daily,reminder,endorse,strategic,vote,el...]	
2131	negative	Women Reproductive right and Racism	[yup,go,reopen,abortion,debate,eln43,sc...]	
2132	positive		[zing,ndp,eln43]	

2133 rows x 4 columns

```
In [16]: #call the function to determine party affiliation
df_election['political_party']=list(map(party, df_election['text']))
df_election

Out[16]:
```

	sentiment	negative_reason	text	political_party
0	negative	Women Reproductive right and Racism	[instead,suggestion,agree,canadian,woman,...]	conservatives
1	positive		[allwomanspacewalk,real,etobiconorth,city...]	None
2	negative	Economy	[brantford,go,cast,94,billion,next,4,ye...]	liberal
3	positive		[canada,canadaelection2019,canadavotes,eln...]	None
4	negative	Economy	[canada,taxpayer,sick,tire,hard,earn,don...]	None
...
2128	negative	Scandal	[know,good,enough,job,smear,campaign,rig...]	None
2129	negative	Scandal	[miss,comment,deflect,issue,answer,best...]	None
2130	positive		[daily,reminder,endorse,strategic,vote,el...]	liberal
2131	negative	Women Reproductive right and Racism	[yup,go,reopen,abortion,debate,eln43,sc...]	None
2132	positive		[zing,ndp,eln43]	new democratic party

2133 rows x 4 columns

(c) Present graphical figures for political parties

By presenting the political parties in histogram, I surprisingly found that the dominating political party is "None". It may be that the keywords are not enough to classify political parties or people are more interested in other parties. Other than "None", liberal party and conservative party are two main portion in tweets, while NDP occupy only a little portion.

```
In [20]: # Histogram of political parties
plt.figure(figsize=(10,5))
ax = sns.countplot(x='political_party', data=df_election)
ax.set_title('The distribution of the political affiliations of the tweets')
ax.set_xlabel('Political Parties')
ax.set_ylabel('Numbers')
plt.show()
```

(d) Present graphical figures for sentiment

By presenting the sentiment in histogram, I found that positive sentiment is about double more than the negative sentiment.

```
In [20]: # Plot distribution of tweets with sentiment values for each party
plt.figure(figsize=(5,4))
ax = sns.countplot(x='label', data=df_sentiment)
ax.set_title('The distribution of the sentiment')
ax.set_xlabel('Positive')
ax.set_ylabel('Numbers')
plt.show()
```

(e) Present distribution of negative reason for election

To visualize some other aspects for election, I draw a piechart to show the distribution of negative reasons. From the piechart shown, it is obvious that "scandal", "tell lies" and "others" are three dominating categories in ten negative reasons. On the contrary, "healthcare" and "marjuana" are two least important factors. Therefore, people care more about leader's characteristics.

```
In [19]: # There are ten negative reasons
df_election['negative_reason'].unique()

Out[19]: array(['Scandal', 'Tell lies', 'Privilege', 'Climate Problem', 'Healthcare', 'Separation', 'Healthcare and Marijuana', dtype=object])

In [20]: # Obtain negative reasons
neg_reasons = df_election['negative_reason'].tolist()
# Remove None
neg_reasons = [x for x in neg_reasons if str(x) != 'nan']
neg_count = Counter(neg_reasons)

In [25]: # draw the pie chart
fig, ax = plt.subplots(figsize=(10, 10))
ax.pie(neg_count.values(), labels=neg_count.keys(), rotatelabels=True, labeldistance=0)
ax.axis('equal')
ax.set_title('The distribution of negative reasons')
plt.legend(loc='upper left')
fig.tight_layout()
plt.show()
```

Bonus: Wordcloud on Positive Sentiment for Liberal Party

I used word cloud to show the most frequent words that appear in positive sentiments on liberal party, which is the winner party in 2019 election. As shown in the wordcloud, the most popular word is "chooseforward", "liberal", and "trudeau". However, "eln43" and "cdnpoli" are not appropriate and cannot provide any useful information, which should be removed during the following prediction.

```
In [26]: # Transform list of string into a master string to pass into wordcloud
lib_positive = lib(lib['sentiment'] == 'positive')
lib_positive['text']=lib_positive['text'].map(lambda x: ' '.join(map(str, x)))
print('The best cross-validation score: %f' % (max(lib_positive['text']).split(' ').count(' '), lib_positive['text'].sum()))

In [29]: # Generate the wordcloud
wordcloud = WordCloud(width=2000, height=1000, max_font_size=200, collocations=False,
plt.figure(figsize=(12,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.title('Positive Tweets on Liberal Party')
plt.show()
```

```
In [31]: # Transform generic tweets and set max features to be 1000
vec_wf = CountVecorizer(max_features=1000, max_df=0.8)
array_wf = vec_wf.fit_transform(df['text'])

In [32]: # Split generic tweets randomly into train data and test data
Xtf_train, Xtf_test, ywf_train, ywf_test = train_test_split(array_wf, df['label'], test_size=0.2, random_state=42)
```

4. Model Implementation

(a) Choose Optimal Model by Tuning Hyper-parameters

Among seven models, logistic regression / k-NN / Naive Bayes / SVM / decision trees / Random Forest / XGBoost, I found that Logistic Regression is the best model in terms of accuracy in predictions, which is 94.29%.

```
In [35]: #Best hyperparameters are recorded in this function for WF features
def get_model_accuracy_wf(model, X_test, y_test):
    dic = {'lr': 'Logistic Regression',
           'knn': 'K Nearest Neighbors',
           'nb': 'Naive Bayes',
           'svm': 'Support Vector Machine',
           'dt': 'Decision Tree',
           'rf': 'Random Forest',
           'ab': 'Ada Boosting',
           'xgb': 'XGBClassifier'}

    train_acc, test_acc = [], []
    # predict using test data
    for model in dic.values():
        test_pred = clf.predict(X_test)
        test_acc = accuracy_score(y_test, test_pred)
    print('Model: ', dic[model])
    print('Test accuracy: %f' % format(test_acc))
```

```
In [36]: #Best hyperparameters are recorded in this function for WF features
def get_model_accuracy_tf(model, X_test, y_test):
    dic = {'lr': 'Logistic Regression',
           'knn': 'K Nearest Neighbors',
           'nb': 'Naive Bayes',
           'svm': 'Support Vector Machine',
           'dt': 'Decision Tree',
           'rf': 'Random Forest',
           'ab': 'Ada Boosting',
           'xgb': 'XGBClassifier'}

    train_acc, test_acc = [], []
    # predict using test data
    for model in dic.values():
        test_pred = clf.predict(X_test)
        test_acc = accuracy_score(y_test, test_pred)
    print('Model: ', dic[model])
    print('Test accuracy: %f' % format(test_acc))
```

(1) Logistic Regression

I chose a hyper-parameters for tuning, which are C. C is the inverse of regularization term. Small value of C means large regularization which penalizes more on large weights to reduce overfitting and vice versa.

Hyper-Parameter Tuning by Cross-Validation

```
In [37]: %time
parameters = {'C': [0.01, 0.1, 1, 10, 100], 'solver': ('liblinear', 'lbfgs', 'newton-cg')}
clf_lr = LogisticRegression(max_iter=1000)
clf = GridSearchCV(clf_lr, parameters, cv=5, n_jobs=-1)
clf.fit(Xwf_train, ywf_train)

print('The best cross-validation score: %f' % format(round(clf.best_score_*100, 2)))
print('The best parameter: %f' % format(round(clf.best_params_)))

The best cross-validation score: 94.29%
The best parameter:
{'C': 1, 'solver': 'liblinear'}
CPU times: user 1.81 s, sys: 316 ms, total: 2.13 s
Wall time: 1min 56s

Best Model of LR using WF
```

```
In [37]: get_model_accuracy_wf('lr', Xwf_test, ywf_test)

Model: Logistic Regression
Test accuracy: 0.9427318644847927
```

```
In [37]: get_model_accuracy_tf('lr', Xtf_test, ytf_test)

Model: Logistic Regression
Test accuracy: 0.9059642195278528
```

(2) K-NN

The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification and regression problems. KNN assumes similar things exist in close proximity.

Hyper-Parameter Tuning by Cross-Validation

```
In [37]: %time
parameters = {'n_neighbors': [5, 10, 15]}
clf_knn = KNeighborsClassifier()
clf = GridSearchCV(clf_knn, parameters, cv=5, n_jobs=-1)
clf.fit(Xwf_train, ywf_train)

print('The best cross-validation score: %f' % format(round(clf.best_score_*100, 2)))
print('The best parameter: %f' % format(round(clf.best_params_)))

The best cross-validation score: 91.3%
The best parameter:
{'n_neighbors': 5}
CPU times: user 372 ms, sys: 1.44 s, total: 1.81 s
Wall time: 2h 17min 30s

Best Model of KNN using WF
```

```
In [37]: get_model_accuracy_wf('knn', Xwf_test, ywf_test)

Model: K Nearest Neighbors
Test accuracy: 0.9150849695369372
```

```
In [37]: get_model_accuracy_tf('knn', Xtf_test, ytf_test)

Model: K Nearest Neighbors
Test accuracy: 0.5747041509708208
```

(3) XGBClassify

XGBoost is a decision-tree-based ensemble Machine Learning algorithm that uses a gradient boosting framework. In prediction problems involving unstructured data (images, text, etc.) artificial neural networks tend to outperform all other algorithms or frameworks.

Hyper-Parameter Tuning by Cross-Validation

```
In [42]: # pip install xgboost
from xgboost import XGBClassifier
parameters = {'max_depth': [6, 7, 8]}
clf_xgb = XGBClassifier()
clf = GridSearchCV(clf_xgb, parameters, cv=5, n_jobs=-1)
clf.fit(Xwf_train, ywf_train)

print('The best cross-validation score: %f' % format(round(clf.best_score_*100, 2)))
print('The best parameter: %f' % format(round(clf.best_params_)))

The best cross-validation score: 93.0%
```



```
The best parameter:
{'max_depth': 8}

Best Model of XGBClassifier using
get_model_accuracy_wf('xgb',
Model: XGBClassifier
Test accuracy: 0.930237769352
```