516 FINAL PROJECT

AJ Acacio, Padmapriya Eunny, Isha Jain, Sita Sanjivini Sawhney, Riya Sirdeshmukh

In [2]:	
211 [2]1	fimporting modules
	#import praw - commented out as scraping no longer needed
	import pandas as pd
	from bs4 import BeautifulSoup
	import nitk from nitk.corpus import stopwords
	from nltk.oorpus import brown
	import conterions import conter import number and the second se
	from pprint import pprint
	import os
	import string, collections
	from string import punctuation
	from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
	from sklearn.decomposition import LatentDirichletAllocation as LDA
	from sklearn.naive_bayes import MultinomialNB
	from schearn import metrics from schev.sparse import hstack
	from sklearn.model_selection import train_test_split from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier, BaggingClassifier
	from sklearn compose import ColumnTransformer
	from sklearn.swn import LinearSVC
	from sklearn.tree import DecisionTreeClassifier from sklearn import tree
	from sklearn.linear_model import LogisticRegression, PassiveAggressiveClassifier
	from sklearn.metigis import RadiusNeighborsClassifier from sklearn.metrics import classification report, confusion matrix, ConfusionMatrixDisplay
	import statsmodels.api as sm
	import matplotlib.pyplot as plt
	from sklearn.svn import SVC
	from sklearn.preprocessing moort MinMaxScaler
	from sklearn.preprocessing import StandardScaler
	Here is the code for scraping from Reddit. Commenting it out as we froze out dataset.
In [3]:	# reddit = praw.Reddit(client_id= 'RjdDRvbOSRsZ6L=F4RFadg',
	<pre># client_secret= orwordgummykilonaolitizabiougg , # user_agent = '516_scraper_practice')</pre>
	# posts = ()
	<pre># posts = [] # cmv = reddit.subreddit('ChangeMyView')</pre>
	<pre># for post in cmv.hot(limit=100):</pre>
	<pre>/ post-depend(post.clumns=('itle', 'id', 'url', 'body', 'flair', 'sore', 'num_comments'))</pre>
	<pre># display(posts)</pre>
T. [4].	
IN [4].	#Importing the .csv that was downloaded the afternoon of May 1st
	<pre>posts['flair'] = posts['flair'].fillna("Nome")</pre>
	posts.head()
Out[4]:	title body flair score num_comments id url
	0 CMV: Student Loans Should Not be Guaranteed by I was surprised to see that this topic has not None 168 123 ug3oa7 https://www.reddit.com/r/changemyview/comments
	1 CMV: *IF* you believe isolation from lockdown edit note: \nWhile I've already awarded delta Delta(s) from OP 440 201 ufqwxc https://www.reddit.com/r/changemyview/comments
	2 CMV: US Colleges should not waste student's ti I went to a very competitive college in the US Delta(s) from OP 1892 638 ufdhhz https://www.reddit.com/r/changemyview/comments
	3 CMV: There is Nothing Wrong with Abortion What We do not treat life as something that is sacr None 15 85 ug4uhk https://www.reddit.com/r/changemyview/comments
	CMV: Western countries should invest in and su Admittedly. I do not know much about the compl Delta(s) from OP 1270 140 uf9cva https://www.reddit.com//changemview/comments
In [5]:	
	#Cleaning the text, creating tokenized versions def preproces_text(text):
	text = text.lower()
	text = **.join(text.split)) # Remove extra spaces, tabs, and new lines
	return text
	<pre>posts["title_cleaned"] = posts["title"].apply(preprocess_text)</pre>
	<pre>posts[boay_cleaned] = posts[boay_].apply(preprocess_text)</pre>
	fremoving stoppords
	sw = hitk.corpus.scopwords.words(english) + [Canv, people, us, head, give, word, dont, better, think, 'shouldt,'better,'life','someone','nothing','icrease','good','least',
	'top','like','best','new','much','without','also','day','bad','one', 'sto'''wynd''aret 'lass' traally' 'theyre' 'cannot' lused' make'
	'instead','given','say','become','isnt','take','needs','thing','never',
	top , arways , ive , ways , way , use , becomes , come , many , rives , 'mods', 'alive', 'far', 'age', 'accos']
	#separate stopwords list for sentiment analysis - used later in the code
	<pre>sw_senti = nltk.corpus.stopwords.words('english') + ['cmw:', 'cmw', 'people','us','need','give', 'would','think','life',</pre>
	'someone','top','new','much','also','day','one','etc', 'instead','diven','say', 'become', 'become', 'come', 'many','lives','states'
	'mods','alive','far', ⁷ age','across','believe','go','programs','time','times','states']
	<pre>posts['title_withoutstop'] = posts['title_cleaned'].apply(lambda x: ' '.join([word for word in x.split() if word not in (sw)]))</pre>
	<pre>posts['body_withoutstop'] = posts['body_cleaned'].apply(lambda x: ' '.join([word for word in x.split() if word not in (sw)]))</pre>
	<pre>posts['tokenized_title'] = posts.apply(lambda column: nltk.word_tokenize(column['title_withoutstop']), axis=1)</pre>
	<pre>posts['tokenized_body'] = posts.apply(lambda column: nltk.word_tokenize(column('body_withoutstop']), axis=1)</pre>
	<pre>posts.drop(posts['flair'] == 'Removed - Submission Rule E'].index, inplace = True) posts.drop(posts['flair'] == 'Removed - Submission Rule A'] index inplace = True)</pre>
	<pre>posts.drop(posts[posts['flair'] == 'Removed - Submission Rule A].index, inplace = True) posts.drop(posts['flair'] == 'Removed - Submission Rule B'].index, inplace = True)</pre>
	<pre>posts.drop(posts['flair'] == 'META'].index, inplace = True)</pre>
	#Checking values for flair

final = pd.DataFrame(posts,columns=['flair','score','num_comments','title_withoutstop','body_withoutstop','tokenized_title','tokenized_body']) final['flair'].unique()
mapping = {'Delta(s) from OP': 1, 'None': 0,'Delta(s) from OP - Fresh Topic Friday': 1,'Fresh Topic Friday':0} final["flair"] = [mapping[item] for item in final.flair]

no_delta = final[final.flair == 0]
with_delta = final[final.flair == 1]
final.head()

Out[5]:	fl	flair sco		num_comments	title_withoutstop	body_withoutstop	tokenized_title	tokenized_body
	0	0	168	123	student loans guaranteed government	surprised see topic debated sub recently light	[student, loans, guaranteed, government]	[surprised, see, topic, debated, sub, recently
	1	1	440	201	believe isolation lockdown mentally unhealthy	edit note already awarded deltas recommended f	[believe, isolation, lockdown, mentally, unhea	[edit, note, already, awarded, deltas, recomme
	2	1	1892	638	colleges waste students time useless mandatory	went competitive college astounded number abso	[colleges, waste, students, time, useless, man	[went, competitive, college, astounded, number
	3	0	15	85	wrong abortion whatsoever	treat something sacred countries certain crimi	[wrong, abortion, whatsoever]	[treat, something, sacred, countries, certain,
	4	1	1270	140	western countries invest subsidize semiconduct	admittedly know complexities matters mind coul	[western, countries, invest, subsidize, semico	[admittedly, know, complexities, matters, mind

Descriptive Stats

In [43]:	<pre>#number of delta vs final.sum(axis=None</pre>	. non-delta posts in dataset) #344 out of 754
0+[43]+	flair	344
00011000	score	95312
	num_comments	83415
	title withoutstop	student loans guaranteed governmentbelieve iso
	body withoutstop	surprised see topic debated sub recently light
	tokenized title	[student, loans, quaranteed, government, belie
	tokenized body	[surprised, see, topic, debated, sub, recently
	dtype: object	
In [49]:	#average reddit sco	re, number of comments for delta and non-delta posts
	print(with delta["s	core"].mean())
	print(with delta["n	um comments"].mean())

print(no_delta["score"].mean())
print(no_delta["num_comments"].mean())

223.42732558139534 179.53488372093022 80.2304347826087 94.15217391304348

Topic Analysis

In [6]:
 #creating a vectorizer instance for topic analysis and the models
 count_vectorizer = CountVectorizer(stop_words = sw)

Dataset with Delta Awarded

Top ten words and topic analysis

count_datal = count_vectorizer.fit_transform(with_delta['title_withoutstop']) def plot_10_most_common_words1(count_data1, count_vectorizer): import mast_common_worksi(count_data; count import mast_pollib.pyplot as plt words = count_vectorizer.get_feature_names() total_counts = np.zeros(len(words)) for t in count_datal: total_counts+=t.toarray()[0] x_pos = np.arange(len(words)) plt.bar(x_pos, counts, align = 'center', color = '#ff4500')
plt.xticks(x_pos, words, rotation = 90, fontname = 'helvetica')
plt.xlabel('Words', fontname = 'helvetica')
plt.ylabel('Counts', fontname = 'helvetica') plt.title('10 Most Common Words for Delta awarded', fontname = 'helvetica')
plt.show()

plot_10_most_common_words1(count_data1, count_vectorizer)



In [8]:	<pre>def print_topics(model, count_vectorizer, n_top_words): words = count_vectorizer.get_feature_names() for topic_idx, topic in enumerate(model.components_): print(^`\nTopic #dd' % topic_idx) print(" ".join([words[i]</pre>
	<pre>for i in topic.argsort()[:-n_top_words - 1:-1]]))</pre>

lda = LDA(n_components = number_topics, random_state = 42) lda.fit(count datal) print("Topics found via LDA:")
print_topics(lda, count_vectorizer, number_words)

Topics found via LDA:

Topic #0:

pay exist women states companies right ukraine elon Topic #1:

gender reason useless media time wing death trans

Topic #2: men american important believe society sports champion body

Dataset with No Delta Awarded

In [10]:

count data2 = count vectorizer.fit transform(no delta['title withoutstop']) def plot_10_most_common_words2(count_data2, count_vectorizer): import matplotlib.pyplot as plt words = count_vectorizer.get_feature_names() total_counts = np.zeros(len(words)) for t in count_data2: total_counts+=t.toarray()[0] count_dict = (zip(words, total_counts)) count_dict = sorted(count_dict, key = lambda x:x[1], reverse = True)[0:10] words = [w[0] for w in count_dict] counts = [w[1] for w in count_dict] x_pos = np.arange(len(words)) plt.bar(x_pos, counts, align = 'center', color = '#ff4500')
plt.xticks(x_pos, words, rotation = 90, fontname = 'helvetica')
plt.xlabel('Words',fontname = 'helvetica')
plt.ylabel('counts',fontname = 'helvetica')
plt.title('10 Most Common Words for No Delta awarded', fontname = 'helvetica') plt.show() plot 10 most common words2 (count data2, count vectorizer)



number_topic = 3
number_words = 8

lda = LDA(n_components = number_topics, random_state = 42)
lda.fit(count_data2)

print("Topics found via LDA:")
print_topics(lda, count_vectorizer, number_words)

Topics found via LDA:

Topic #0: companies student government reason russia society loan human

Topic #1:

world state america sexual time student western wrong

Topic #2: men gender death change women body ukraine society

Complete dataset (with and without delta)

count_data3 = count_vectorizer.fit_transform(final['title_withoutstop'])

```
def plot_10_most_common_words3(count_data3, count_vectorizer):
    import matplotlib.pyplot as plt
    words = count_vectorizer.get_feature_names()
    total_counts = np.zeros(len(words))
    for t in count_data3:
                    total_counts+=t.toarray()[0]
          count_dict = (zip(words, total_counts))
count_dict = sorted(count_dict, key=lambda x:x[1], reverse=True)[0:10]
words = [w[0] for w in count_dict]
counts = [w[1] for w in count_dict]
x_pos = np.arange(len(words))
          plt.bar(x_pos, counts, align = 'center', color = '#ff4500')
plt.xticks(x_pos, words, rotation = 90, fontname = 'helvetica')
plt.xlabel('Words', fontname = 'helvetica')
plt.ylabel('Counts', fontname = 'helvetica')
plt.title('10 Most Common Words for Complete Dataset', fontname = 'helvetica')
           plt.show()
plot_10_most_common_words3(count_data3, count_vectorizer)
```



number_topic = 3
number_words = 8

lda = LDA(n_components = number_topics, random_state = 42) lda.fit(count data3)

print("Topics found via LDA:")
print_topics(lda, count_vectorizer, number_words)

Topics found via LDA:

Topic #0: time states companies american pay men exist women

Topic #1: death media ukraine state wrong person world speech

Topic #2: gender student important society trans body government believe

Sentiment Analysis

In [14]:

#creating a new datasent for sentiment analysis that is cleaned using different stopword #creating a new datasent for sentiment analysis that is cleaned using different stopword
senti = posts
senti = senti.drop(columns =["tokenized_title", "tokenized_body", "title_withoutstop", "body_withoutstop","id","url"], axis = 1)
senti['title_withoutstopsenti'] = senti['title_cleaned'].apply(lambda x: '.join((word for word in x.split() if word not in (sw_senti)]))
senti['tokenized_title_senti'] = senti['body_cleaned'].apply(lambda x: '.join((word for word in x.split() if word not in (sw_senti)]))
senti['tokenized_title_senti'] = senti.apply(lambda column: nltk.word_tokenize(column['title_withoutstopsenti']), axis=1)
senti['tokenized_body_senti'] = senti.apply(lambda column: nltk.word_tokenize(column['body_withoutstopsenti']), axis=1)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule E"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule E"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace = True)
senti.drop(senti[senti['flair"] == "Removed - Submission Rule B"].index, inplace B"].index].senti.head()
senti['flair"] = [mapping[item] f

senti

tokenized_body_senti	tokenized_title_senti	body_withoutstopsenti	title_withoutstopsenti	body_cleaned	title_cleaned	title body flair score num_comments		t[14]:			
[surprised, see, topic, debated, sub, recently	[student, loans, guaranteed, government]	surprised see topic debated sub recently light	student loans guaranteed government	i was surprised to see that this topic has not	cmv student loans should not be guaranteed by	123	168	0	I was surprised to see that this topic has not	CMV: Student Loans Should Not be Guaranteed by	0
[edit, note, ive, already, awarded, deltas, re	[isolation, lockdown, mentally, unhealthy, mus	edit note ive already awarded deltas recommend	isolation lockdown mentally unhealthy must log	edit note while ive already awarded deltas its	cmv if you believe isolation from lockdown was	201	440	1	edit note: \nWhile I've already awarded delta	CMV: *IF* you believe isolation from lockdown	1
[went, competitive, college, astounded, number	[colleges, waste, students, useless, mandatory	went competitive college astounded number abso	colleges waste students useless mandatory classes	i went to a very competitive college in the us	cmv us colleges should not waste students time	638	1892	1	I went to a very competitive college in the US	CMV: US Colleges should not waste student's ti	2
[treat, something, sacred, countries, certain,	[nothing, wrong, abortion, whatsoever]	treat something sacred countries certain crimi	nothing wrong abortion whatsoever	we do not treat life as something that is sacr	cmv there is nothing wrong with abortion whats	85	15	0	We do not treat life as something that is sacr	CMV: There is Nothing Wrong with Abortion What	3
[admittedly, know, complexities, matters, mind	[western, countries, invest, subsidize, semico	admittedly know complexities matters mind coul	western countries invest subsidize semiconduct	admittedly i do not know much about the comple	cmv western countries should invest in and sub	140	1270	1	Admittedly, I do not know much about the compl	CMV: Western countries should invest in and su	4
[big, aspect, lot, seem, skip, discussing, fin	[financially, things, get, better, saved, money]	big aspect lot seem skip discussing financial	financially things get better saved money	one big aspect a lot of people seem to skip wh	cmv financially things would get better if mor	31	0	1	One big aspect a lot of people seem to skip wh	CMV: Financially, things would get better if m	576
[let, put, bias, world, view, table, upfront,	[united, level, outlaw, abortion, gay, marriag	let put bias world view table upfront im strai	united level outlaw abortion gay marriage inte	let me put my bias and world view on the table	cmv the united states will at some level outla	46	0	1	Let me put my bias and world view on the table	CMV: The United States will, at some level, ou	577
[view, wish, changed, ukraine, can, not, win, 	[matter, russia, wins, ukraine]	view wish changed ukraine cannot win war curre	matter russia wins ukraine	the view i wish to have changed is that ukrain	cmv it is only a matter of time until russia w	862	1511	1	The view I wish to have changed is that Ukrain	CMV: It is only a matter of time until Russia	578
[philosophy, incredibly, important, subject, c	[philosophy, important, academic, subject, man	philosophy incredibly important subject core c	philosophy important academic subject mandator	i think that philosophy is an incredibly impor	cmv philosophy is the most important academic	99	21	1	I think that philosophy is an incredibly impor	CMV: Philosophy is the most important academic	579
[arent, aware, weeks, ago, reddit, revamped, b	[abuse, reddits, revamped, block, feature, out	arent aware weeks ago reddit revamped block fe	abuse reddits revamped block feature outweigh	for those who arent aware a few weeks ago redd	cmv the abuse of reddits revamped block featur	214	80	0	For those who aren't aware, a few weeks ago Re	CMV: The abuse of Reddit's revamped block feat	580

574 rows × 11 columns

sentiments = senti.drop(columns = ["tokenized_title_senti", "tokenized_body_senti"], axis = 1)

vader = SentimentIntensityAnalyzer()

scores = [vader.polarity_scores(body_withoutstopsenti) for body_withoutstopsenti in sentiments.body_withoutstopsenti]

Convert the list of dicts into a DataFrame body feelings = pd.DataFrame(scores)

pd.set_option('display.max_rows', 100)

#Appending the dataframe with just the sentiments to the data.frame (final2). sentiments = sentiments.join(body_feelings)

#Creating a column that displays a binary valence column
sentiments['valence'] = sentiments['compound'].apply(lambda c: 'pos' if c >=0 else 'neg')

sentiments = sentiments.dropna()
sentiments = sentiments.drop(columns = ['body_withoutstopsenti'])

sentiments.head()

Out[15]:		title	body	flair	score	num_comments	title_cleaned	body_cleaned	title_withoutstopsenti	neg	neu	pos	compound	valence
	0	CMV: Student Loans Should Not be Guaranteed by	I was surprised to see that this topic has not	0	168	123	cmv student loans should not be guaranteed by	i was surprised to see that this topic has not	student loans guaranteed government	0.125	0.589	0.286	0.9684	pos
	1	CMV: *IF* you believe isolation from lockdown	edit note: \nWhile I've already awarded delta	1	440	201	cmv if you believe isolation from lockdown was	edit note while ive already awarded deltas its	isolation lockdown mentally unhealthy must log	0.352	0.422	0.225	-0.9902	neg
	2	CMV: US Colleges should not waste student's ti	I went to a very competitive college in the US	1	1892	638	cmv us colleges should not waste students time	i went to a very competitive college in the us	colleges waste students useless mandatory classes	0.114	0.678	0.208	0.9934	pos
	3	CMV: There is Nothing Wrong with Abortion What	We do not treat life as something that is sacr	0	15	85	cmv there is nothing wrong with abortion whats	we do not treat life as something that is sacr	nothing wrong abortion whatsoever	0.384	0.409	0.207	-0.9834	neg
	4	CMV: Western countries should invest in and su	Admittedly, I do not know much about the compl	1	1270	140	cmv western countries should invest in and sub	admittedly i do not know much about the comple	western countries invest subsidize semiconduct	0.051	0.725	0.224	0.9844	pos

Topic Analysis based on Valence

	<pre>#splitting dataset based on valence possentiment = sentiments[sentiments.valence == 'pos'] negsentiment = sentiments.valence == 'neg']</pre>
	Positive Sentiment
[n [17]:	<pre>count_data4 = count_vectorizer.fit_transform(possentiment['title_withoutstopsenti'])</pre>
	<pre>def plot_10_most_common_words4(count_data4, count_vectorizer): import matplotlib.pyplot as plt words = count_vectorizer.get_faature_names() total_counts = np.zeros(len(words)) for t in count_data4: total_counts+et.toarray()[0] count_dict = (zip(words, total_counts))</pre>
	<pre>count_dict = sorted(count_dict, key=lambda x:x[1], reverse=True)[0:10] words = [w[0] for w in count_dict] counts = [w[1] for w in count_dict] x_pos = np.arange(len(words)) </pre>
	<pre>pit.bar(x_pos, counts, align = center, color = #1f4500) plt.vicks(x_pos, words, rotation = 90, fontname = 'helvetica') plt.vlabel('Words', fontname = 'helvetica') plt.ylabel('Counts', fontname = 'helvetica') plt.tite('10 Most Common Words for Positive Valence', fontname = 'helvetica') plt.show()</pre>

10 Most Common Words for Positive Valence 14 12 10 Counts 8 2 men person trans media society reason gender vomen merica Words

In [18]:

number_topics = 3 number_words = 8

lda = LDA(n_components = number_topics, random_state = 42)
lda.fit(count_data4)

print("Topics found via LDA:")
print_topics(lda, count_vectorizer, number_words)

Topics found via LDA:

Topic #0: rights russia media gender care ukraine government student

Topic #1: society person exist problems useless state work important

Topic #2: trans reason argument elon gender valid america wrong

Negative Sentiment

In [19]: count_data5 = count_vectorizer.fit_transform(negsentiment['title_withoutstopsenti'])

def plot_10_most_common_words5(count_data5, count_vectorizer):
 import matplotlib.pyplot as plt count_dict = (zip(words, total_counts)) count_dict = sorted(count_dict, key=lambda x:x[1], reverse=True)[0:10] words = [w[0] for w in count_dict] counts = [w[1] for w in count_dict] x_pos = np.arange(len(words)) plt.bar(x_pos, counts, align = 'center', color = '#ff4500')
plt.xticks(x_pos, words, rotation = 90, fontname = 'helvetica')
plt.xlabel('Words', fontname = 'helvetica')
plt.ylabel('Counts', fontname = 'helvetica')

plt.title('10 Most Common Words for Negative Valence', fontname = 'helvetica') plt.show()

plot_10_most_common_words5(count_data5, count_vectorizer)



In [20]:

number topics = 3 number_words = 8

lda = LDA(n components = number topics, random state = 42) lda.fit(count_data5)

print("Topics found via LDA:")
print_topics(lda, count_vectorizer, number_words)

Topics found via LDA:

Topic #0:

men companies world anything school religion home change

Topic #1: society american wrong sports women ukraine wild problem

Topic #2:

death champion twitter gender penalty worst date white

Creating a Model Using Sentiment Scores as Features (Model 1)

sentiments.head() title body flair score num_comments title_cleaned body_cleaned title_withoutstopsenti neg neu pos compound valence 0 CMV: Student Loans Should I was surprised to see that 123 cmv student loans should i was surprised to see that student loans guaranteed 0 168 0.125 0.589 0.286 0.9684 pos Not be Guaranteed by ... this topic has not ... not be guaranteed by ... this topic has not ... government 201 cmv if you believe isolation edit note while ive already CMV: *IF* you believe edit note: \nWhile I've isolation lockdown mentally 1 440 0.352 0.422 0.225 -0.9902 1 neg isolation from lockdown ... already awarded delta ... from lockdown was... awarded deltas its ... unhealthy must log.. 638 cmv us colleges should not i went to a very competitive 2 CMV: US Colleges should not I went to a very competitive colleges waste students 1 1892 0.114 0.678 0.208 0.9934 pos waste student's ti... college in the US... waste students time ... college in the us... useless mandatory classes 3 CMV: There is Nothing Wrong nothing wrong abortion 0.384 0.409 0.207 We do not treat life as 85 cmv there is nothing wrong we do not treat life as 0 15 -0.9834 neg with Abortion What... something that is sacr... with abortion whats... something that is sacr... western countries invest 0.051 0.725 0.224 CMV: Western countries Admittedly, I do not know cmv western countries admittedly i do not know 1 1270 140 Should invest in and sub... 0.9844 4 pos should invest in and su... much about the compl... much about the comple... subsidize semiconduct...

In [22]: #Using score, num_comments, and compound as predictors of Delta

#dropping columns from sentiments to form dataset for Model 1 numericmodel = sentiments.drop(columns = ['title withoutstopsenti', "title cleaned", "body cleaned"]) numericmodel.head()

ut[22]:		title	body	flair	score	num_comments	neg	neu	pos	compound	valence
	0	CMV: Student Loans Should Not be Guaranteed by	I was surprised to see that this topic has not	0	168	123	0.125	0.589	0.286	0.9684	pos
	1	CMV: *IF* you believe isolation from lockdown	edit note: \nWhile I've already awarded delta	1	440	201	0.352	0.422	0.225	-0.9902	neg
	2	CMV: US Colleges should not waste student's ti	I went to a very competitive college in the US	1	1892	638	0.114	0.678	0.208	0.9934	pos
	3	CMV: There is Nothing Wrong with Abortion What	We do not treat life as something that is sacr	0	15	85	0.384	0.409	0.207	-0.9834	neg
	4	CMV: Western countries should invest in and su	Admittedly, I do not know much about the compl	1	1270	140	0.051	0.725	0.224	0.9844	pos

scaler = MinMaxScaler() features = data.loc[:,'score':'compound']
features = scaler.fit_transform(features) target = data['flair'] x_train, x_test, y_train, y_test = train_test_split(features, target, test_size = test_size, random_state = 808) init model = None if model == 'knn': init_model = knn(n_neighbors = 1)
elif model == 'randomforest':
 init_model = RandomForestClassifier(random_state = 808) elif model == 'boost': init_model = GradientBoostingClassifier(random_state = 808)
elif model == 'sym':
 init_model = SVC(random_state = 808)
elif model == 'decisiontree': init_model = DecisionTreeClassifier(random_state = 808) elif model == 'bagging': init model = BaggingClassifier(random_state = 808)
elif model == 'logistic':
 init_model = LogisticReression(random_state = 808)
elif model == 'passive_aggressive':

init_model = PassiveAggressiveClassifier(random_state = 808)
elif model == 'radius_neighbors':
 init_model = RadiusNeighborsClassifier(radius = 100, outlier_label = "most_frequent") fitted_model = init_model.fit(x_train ,y_train)
test_predictions = fitted_model.predict(x_test)
accuracy_score = fitted_model.score(x_test,y_test)

cm = confusion_matrix(y_test, test_predictions, labels = fitted_model.classes_)
disp = ConfusionMatrixDisplay(confusion_matrix = cm,display_labels = fitted_model.classes_)
disp.plot(cmap = 'Oranges')
plt.title('Confusion Matrix of Observation Counts for {}'.format(model), fontname = 'helvetica')
plt.show()

print("Accuracy Score for {}: {:.2%}".format(model, accuracy_score))

evaluate_model(numericmodel, model = 'decisiontree')
evaluate_model(numericmodel, model = 'logistic')
evaluate_model(numericmodel, model = 'boost')
evaluate_model(numericmodel, model = 'boost')
evaluate_model(numericmodel, model = 'randomforest')
evaluate_model(numericmodel, model = 'sem')
evaluate_model(numericmodel, model = 'radius_neighbors')







Accuracy Score for logistic: 63.16%



Accuracy Score for bagging: 56.14% Confusion Matrix of Observation Counts for boos



Accuracy Score for boost: 59.65% Confusion Matrix of Observation Counts for randomforest



Accuracy Score for randomforest: 64.91%



Creating a Model Using The Vectorized Text to Predict (Model 2)

In [40]: text_predictor = final

#Create a Target
flair = text_predictor.flair

X_train, X_test, y_train, y_test = train_test_split(text_predictor['title_withoutstop'], flair, test_size=0.2, random_state= 42)

#Intialize the Count Vectorizer count_vectorizer = CountVectorizer(stop_words=sw) count_train = count_vectorizer.fit_transform(X_train) count_test = count_vectorizer.transform(X_test)

Initialize tfidf vectorizer
tfidf_vectorizer = TfidfVectorizer(stop_words=sw)

Create tfidf train and test variables
tfidf_train = tfidf_vectorizer.fit_transform(X_train)
tfidf_test = tfidf_vectorizer.transform(X_test)

#Creating Naive Bayes Model Using TFIDF
tfidf_nb = NultinomialN8()
tfidf_nb_redict(tfidf_train, y_train)
tfidf_nb_redict(tfidf_test)
tfidf_nb_redict(tfidf_test)
tfidf_nb_recore = metrics.accuracy_score(y_test, tfidf_nb_pred)
em_tfidf = Confusion_matrix(y_test, tfidf_nb_pred, labels = tfidf_fitted_model.classes_)
disp_tfidf = Confusion_matrix(s Context for Naive Bayes TFIDF', fontname = 'helvetica')
plt.show()
print('NaiveBayes Model Using Count Vectorizer
count_nb_redict(count_train, y_train)
count_nb_red = confusion_matrix(y_test, count_nb_pred)
em_conf_nb_score = metrics.accuracy_score(y_test, count_nb_pred)
count_nb_red = cont_sion_matrix(y_test, count_nb_pred)
fcreating Maive Bayes Model Using Count Vectorizer
count_nb_red = count_nb.redict(count_train, y_train)
count_nb_red = confusion_matrix(y_test, count_nb_pred)
count_nb_red = confusion_matrix(y_test, count_nb_pred)
em_count = Confusion_matrix(y_test, count_nb_red, labels = count_fitted_model.classes_)
disp_count.plot(map = 'Oranges')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('confusion_matrix(y_test, count_nb_red, labels = count_fitted_model.classes_)
disp_count.plot(map = 'Oranges')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('confusion Natrix of Observation Counts for Naive Bayes CountVectorizer', fontname = 'helvetica')
plt.tile('NaiveBayes Count Score: ', count_nb_score)



Predicted label

NaiveBayes Tfidf Score: 0.5739130434782609 Confusion Matrix of Observation Counts for Naive Bayes CountVectorizer



NaiveBayes Count Score: 0.5304347826086957

```
In [41]:
```

tfidf_svc = LinearSVC()
tfidf_svc.fit(tfidf_train, y_train)
tfidf_svc_pred = tfidf_svc.predict(tfidf_test)
tfidf_svc_score = metrics.accuracy_score(y_test, tfidf_svc_pred)

print("LinearSVC Score: %0.3f" % tfidf_svc_score)

LinearSVC Score: 0.539

In [42]: # delta_post = "I do not feel very strongly about veganism, so please try to convince me" # nondelta_post = "I will never change my view about veganism, but come at me"

delta_post = "I do not feel very strongly about Alex winning Eurovision, so please try to convince me" # nondelta_post = "I will never change my view about Alex winning Eurovision, but come at me"

delta_post = "I do not think Alex will win Eurovision, so convince me otherwise"
nondelta_post = "I will not change my view about Alex winning Eurovision, but come at me"

delta_post_vectorized = tfidf_vectorizer.transform([delta_post])
nondelta post vectorized = tfidf vectorizer.transform([nondelta post])

delta_post_pred = tfidf_svc.predict(delta_post_vectorized)
nondelta_post_pred = tfidf_svc.predict(nondelta_post_vectorized)

print('Predicted Post that Led to Changed View', delta_post_pred)
print('Predicted Post that Did Not Lead to View Change', nondelta_post_pred)

Predicted Post that Led to Changed View [1] Predicted Post that Did Not Lead to View Change [0]

In []