In [1]: **import** pandas **as** pd import numpy as np import re import seaborn as sns import matplotlib.pyplot as plt from wordcloud import WordCloud from matplotlib import style style.use('ggplot') from textblob import TextBlob from nltk.tokenize import word_tokenize from nltk.stem import PorterStemmer import nltk from nltk.corpus import stopwords stop_words = set(stopwords.words('english')) from wordcloud import WordCloud from sklearn.feature_extraction.text import CountVectorizer from sklearn.model_selection import train_test_split from sklearn.linear_model import LogisticRegression from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, ConfusionMatrixDisplay df = pd.read_csv('vaccination_tweets.csv') In [2]: df.head() In [3]: id user_name user_location user_description user_created user_followers user_friends user_favourites user_verified date text hashtags sour Out[3]: Same folks Aggregator of Asian 2020-2009-04-08 La Crescentasaid daikon **0** 1340539111971516416 Rachel Roh 405 12-20 American news; 1692 3247 False ['PfizerBioNTech'] 17:52:46 paste could scanning di... 06:06:44 Andro treat a cyt... While the Marketing dude, tech 2020world has Twit 2009-09-21 **1** 1338158543359250433 Albert Fong geek, heavy metal & 834 666 178 False 12-13 been on the NaN W Francisco, CA 15:27:30 16:27:13 '80s ... wrong side of #coronavirus 2020-2020-06-25 #SputnikV ['coronavirus', 'SputnikV', **2** 1337858199140118533 heil, hydra 🖐 😊 10 88 155 Your Bed False 12-12 #AstraZeneca 'AstraZeneca', 'Pf... 23:30:28 Andro 20:33:45 #PfizerBio... Facts are Hosting 2020-Twit Vancouver, "CharlesAdlerTonight" Charles 2008-09-10 immutable, **3** 1337855739918835717 49165 3933 21853 True 12-12 NaN W BC - Canada Senator, even Adler 11:28:53 Global News Radi... 20:23:59 Α when you're... Explain to me Citizen 2020-Citizen News again why we Twit 2020-04-23 ['whereareallthesickpeople', NaN Channel bringing you 1473 **4** 1337854064604966912 152 580 need a False 12-12 News 17:58:42 'PfizerBioNTech'] Channel an alternati... 20:17:19 vaccine @Bor... df.info() In [4]: <class 'pandas.core.frame.DataFrame'> RangeIndex: 11020 entries, 0 to 11019 Data columns (total 16 columns): Non-Null Count # Column Dtype ------ - ------0 11020 non-null int64 id 11020 non-null user_name object user_location 8750 non-null object user_description 10340 non-null object user_created 11020 non-null object 5 user_followers 11020 non-null int64 6 11020 non-null int64 user_friends user_favourites 11020 non-null int64 user_verified 11020 non-null bool 8 9 11020 non-null object date 11020 non-null object 10 text 8438 non-null object 11 hashtags 11019 non-null object 12 source 13 retweets 11020 non-null int64 14 favorites 11020 non-null int64 11020 non-null bool 15 is_retweet dtypes: bool(2), int64(6), object(8) memory usage: 1.2+ MB In [5]: df.columns Out[5]: 'date', 'text', 'hashtags', 'source', 'retweets', 'favorites', 'is_retweet'], dtype='object') text_df = df.drop(['id', 'user_name', 'user_location', 'user_description', 'user_created', 'user_followers', 'user_friends', 'user_favourites', 'user_verified', 'date', 'hashtags', 'source', 'retweets', 'favorites', 'is_retweet'], axis=1) text_df.head() Out[6]: text 0 Same folks said daikon paste could treat a cyt... 1 While the world has been on the wrong side of ... 2 #coronavirus #SputnikV #AstraZeneca #PfizerBio... Facts are immutable, Senator, even when you're... 4 Explain to me again why we need a vaccine @Bor... print(text_df['text'].iloc[0],"\n") print(text_df['text'].iloc[1], "\n") print(text_df['text'].iloc[2], "\n") print(text_df['text'].iloc[3],"\n") print(text_df['text'].iloc[4],"\n") Same folks said daikon paste could treat a cytokine storm #PfizerBioNTech https://t.co/xeHhIMg1kF While the world has been on the wrong side of history this year, hopefully, the biggest vaccination effort we've ev... https://t.co/dlCHrZjkhm #coronavirus #SputnikV #AstraZeneca #PfizerBioNTech #Moderna #Covid_19 Russian vaccine is created to last 2-4 years... https://t.co/ieYlCKBr8P Facts are immutable, Senator, even when you're not ethically sturdy enough to acknowledge them. (1) You were born i... https://t.co/jqgV18kch4 Explain to me again why we need a vaccine @BorisJohnson @MattHancock #whereareallthesickpeople #PfizerBioNTech... https://t.co/KxbSRoBEHq text_df.info() In [8]: <class 'pandas.core.frame.DataFrame'> RangeIndex: 11020 entries, 0 to 11019 Data columns (total 1 columns): # Column Non-Null Count Dtype --- ----- ------ -----0 text 11020 non-null object dtypes: object(1) memory usage: 86.2+ KB In [9]: def data_processing(text): text = text.lower() text = re.sub(r"https\S+|www\S+https\S+", '', text, flags=re.MULTILINE) text = re.sub(r'\@w+|\#','',text)
text = re.sub(r'[^\w\s]','',text) text_tokens = nltk.word_tokenize(text) filtered_text = [w for w in text_tokens if not w in stopwords.words('english')] return " ".join(filtered_text) text_df.text = text_df['text'].apply(data_processing) text_df = text_df.drop_duplicates('text') stemmer = PorterStemmer() def stemming(data): text = [stemmer.stem(word) for word in data.split()] return " ".join(text) text_df['text'] = text_df['text'].apply(lambda x: stemming(x)) text_df.head() In [10]: text Out[10]: folk said daikon past could treat cytokin stor... world wrong side histori year hope biggest vac... coronaviru sputnikv astrazeneca pfizerbiontech... fact immut senat even your ethic sturdi enough... 4 explain need vaccin borisjohnson matthancock w... print(text_df['text'].iloc[0],"\n") In [11]: print(text_df['text'].iloc[1], "\n") print(text_df['text'].iloc[2],"\n") print(text_df['text'].iloc[3],"\n") print(text_df['text'].iloc[4],"\n") folk said daikon past could treat cytokin storm pfizerbiontech world wrong side histori year hope biggest vaccin effort weve ev coronaviru sputnikv astrazeneca pfizerbiontech moderna covid_19 russian vaccin creat last 24 year fact immut senat even your ethic sturdi enough acknowledg 1 born explain need vaccin borisjohnson matthancock whereareallthesickpeopl pfizerbiontech In [12]: text_df.info() <class 'pandas.core.frame.DataFrame'> Index: 10543 entries, 0 to 11019 Data columns (total 1 columns): # Column Non-Null Count Dtype -----0 text 10543 non-null object dtypes: object(1) memory usage: 164.7+ KB def polarity(text): In [13]: return TextBlob(text).sentiment.polarity text_df['polarity'] = text_df['text'].apply(polarity) In [14]: text_df.head(10) In [15]: polarity Out[15]: text folk said daikon past could treat cytokin stor... -0.250000 world wrong side histori year hope biggest vac... -0.500000 coronaviru sputnikv astrazeneca pfizerbiontech... 0.000000 fact immut senat even your ethic sturdi enough... 0.000000 4 explain need vaccin borisjohnson matthancock w... 0.000000 anyon use adviceguid whether covid vaccin safe... 0.500000 bit sad claim fame success vaccin patriot comp... -0.100000 mani bright day 2020 best 1 bidenharri win ele... 0.833333 covid vaccin get covidvaccin covid19 pfizerbio... 0.000000 covidvaccin state start get covid19vaccin mond... 0.000000 def sentiment(label): if label <0:</pre> return "Negative" **elif** label ==0: return "Neutral" elif label>0: return "Positive" text_df['sentiment'] = text_df['polarity'].apply(sentiment) text_df.head() In [18]: Out[18]: text polarity sentiment 0 folk said daikon past could treat cytokin stor... Negative -0.25 world wrong side histori year hope biggest vac... -0.50 Negative coronaviru sputnikv astrazeneca pfizerbiontech... 0.00 Neutral fact immut senat even your ethic sturdi enough... Neutral 4 explain need vaccin borisjohnson matthancock w... 0.00 Neutral In [20]: fig = plt.figure(figsize=(5,5)) sns.countplot(x='sentiment', data = text_df) <Axes: xlabel='sentiment', ylabel='count'> Out[20]: 6000 5000 4000 count 3000 2000 1000 0 Positive Negative Neutral sentiment fig = plt.figure(figsize=(7,7)) colors = ("yellowgreen", "gold", "red") wp = {'linewidth':2, 'edgecolor':"black"} tags = text_df['sentiment'].value_counts() explode = (0.1, 0.1, 0.1)tags.plot(kind='pie', autopct='%1.1f%%', shadow=True, colors = colors, startangle=90, wedgeprops = wp, explode = explode, label='') plt.title('Distribution of sentiments') Text(0.5, 1.0, 'Distribution of sentiments') Distribution of sentiments Negative 8.5% Positive 30.3% 61.2% Neutral In [22]: pos_tweets = text_df[text_df.sentiment == 'Positive'] pos_tweets = pos_tweets.sort_values(['polarity'], ascending= False) pos_tweets.head() Out[22]: text polarity sentiment 2004 fulli vaccin covid 19 best gift 2021 thank hel... 1.0 Positive 5889 prguy17 scottythequeuejump resignaustralian de... 1.0 Positive 9216 despit israel pfizerbiontech vaccin consid bes... Positive 1.0 9317 best way get merrygoround pfizer pfizerbiontec... 1.0 Positive 5636 thank europ develop best vaccin avail let us s... 1.0 Positive text = ' '.join([word for word in pos_tweets['text']]) In [23]: plt.figure(figsize=(20,15), facecolor='None') wordcloud = WordCloud(max_words=500, width=1600, height=800).generate(text) plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.title('Most frequent words in positive tweets', fontsize=19) Most frequent words in positive tweets felt major report covid suppl: neg_tweets = text_df[text_df.sentiment == 'Negative'] neg_tweets = neg_tweets.sort_values(['polarity'], ascending= False) neg_tweets.head() Out[24]: polarity sentiment 7715 got first dose less wait time airport vaccin c... -0.005556 Negative 7157 nas_k27 second dose due end next month well fa... -0.006250 Negative 575 u think bill gate evil becoim rich nice person... -0.008333 Negative 3738 due experi norway suggest covid19 vaccin may r... -0.012500 Negative 8138 long american british isra feed final got shot... -0.012500 Negative text = ' '.join([word for word in neg_tweets['text']]) In [25]: plt.figure(figsize=(20,15), facecolor='None') wordcloud = WordCloud(max_words=500, width=1600, height=800).generate(text) plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.title('Most frequent words in negative tweets', fontsize=19) plt.show() Most frequent words in negative tweets Э go 0 scienc oos medic expect know dose neutral_tweets = text_df[text_df.sentiment == 'Neutral'] In [26]: neutral_tweets = neutral_tweets.sort_values(['polarity'], ascending= False) neutral_tweets.head() text polarity sentiment Out[26]: 2 coronaviru sputnikv astrazeneca pfizerbiontech... 0.0 Neutral 7364 test pfizerbiontech vaccin hospit worker cambr... 0.0 Neutral 7377 1 1 go pfizerbiontech pr 0.0 Neutral 7375 break moderna pfizerbiontech vaccin highli eff... 0.0 Neutral 7374 eu mam vaccinat covid19 covidwar pfizerbiontech 0.0 Neutral text = ' '.join([word for word in neutral_tweets['text']]) In [27]: plt.figure(figsize=(20,15), facecolor='None') wordcloud = WordCloud(max_words=500, width=1600, height=800).generate(text) plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.title('Most frequent words in neutral tweets', fontsize=19) plt.show() Most frequent words in neutral tweets deathbooster side effect obsei moderna data∞ feel make studi In [29]: vect = CountVectorizer(ngram_range=(1,2)).fit(text_df['text']) X = text_df['text'] Y = text_df['sentiment'] X = vect.transform(X) In [30]: x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=42) In [31]: print("Size of x_train:", (x_train.shape)) print("Size of y_train:", (y_train.shape)) print("Size of x_test:", (x_test.shape)) print("Size of y_test:", (y_test.shape)) Size of x_train: (8434, 71682) Size of y_train: (8434,) Size of x_test: (2109, 71682) Size of y_test: (2109,) In [32]: logreg = LogisticRegression() logreg.fit(x_train, y_train) logreg_pred = logreg.predict(x_test) logreg_acc = accuracy_score(logreg_pred, y_test) print("Test accuracy: {:.2f}%".format(logreg_acc*100)) Test accuracy: 90.14%