from keras.models import Sequential from keras.layers import Input, LSTM, Dense, Bidirectional, Activation from tensorflow.keras.layers import Embedding import keras.backend as K import joblib In []: ##Scrape data from twitter #attrs = [] #limit = 0#for i, tweet in enumerate(sntwitter.TwitterSearchScraper('Valheim since:2022-11-01 until:2022-12-16').get_items()): if limit>4999: # if(str(guess(tweet.content)) == "en"): attrs.append([tweet.user.username, tweet.date, tweet.likeCount, tweet.sourceLabel, tweet.content]) limit += 1## Creating a dataframe to load the list #df = pd.DataFrame(attrs, columns=["User", "Date Created", "Number of Likes", "Source of Tweet", "Tweet"]) #df['Date Created'] = pd.to_datetime(df['Date Created']).dt.date #pd.to_csv("scrapped_tweet.csv") In [2]: #Load twitter data and show first lines df = pd.read_csv('scrapped_tweet.csv') df = df.iloc[:,1:] df.head() User Date Created Number of Likes Source of Tweet Out[2]: Tweet 0 ElPibeCombo 2022-12-15 Twitter Web App @BS_artsss Oh, this is so good, i love it! can... MrsTadertaut 2022-12-15 Twitter Web App Hello lovelies!\n\nTonight we are jumping back... Meganoip 2022-12-15 0 Twitter Web App Storyline: Noi found Block fighting with giant... 3 TaesTeahouse 2022-12-15 Twitter Web App Ayo I'm totally not biased or anything but yal... 4 John90192410 2022-12-15 1 Twitter for Android The Queen Wants it all... \nDid Mistlands solo... In [3]: df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 5000 entries, 0 to 4999 Data columns (total 5 columns): Column Non-Null Count Dtype -----5000 non-null object 0 User 5000 non-null object Date Created 1 Number of Likes 5000 non-null int64 Source of Tweet 5000 non-null object Tweet 5000 non-null object dtypes: int64(1), object(4) memory usage: 195.4+ KB In [4]: #Create new list with tweet as its values tweet = (df["Tweet"].values.tolist()) In [5]: #Initialize lemmatizer wordnet_lemmatizer = WordNetLemmatizer() In [6]: #Create new variable to be filled with English stop-words stop_words = nltk.corpus.stopwords.words('english') In [8]: **#Fix** contractions contractions.add('cant', 'can not') In [9]: #Cleaning the tweets #Create a new list to store the cleaned tweets clean_tweet = [] for i in range(0,len(tweet)): **#Undercase all letters** temp = p.clean(tweet[i]).casefold() #Remove all punctuations temp = temp.translate(str.maketrans('', '', string.punctuation)) #Create new list to store filtered tokenized words filtered = [] #Create new variable for tokenized sentence words_temp = word_tokenize(temp) #Create new variable to store string that has its contractions removed str_temp_fixed = "" #Clean the words from contractions for j in words_temp: j = contractions.fix(j) str_temp_fixed += " " + j #Create new variable for the cleaned tokenized sentence words_temp_fixed = word_tokenize(str_temp_fixed) for c in words_temp_fixed: #Check if a word is a stop word if c not in stop_words: #Lemmatize the word, then stores it into the "filtered" list c = wordnet_lemmatizer.lemmatize(c) filtered.append(c) #Create a new variable to store the filtered words that are re-combined back into a sentence cleaned = " ".join(filtered) #Insert the "cleaned" variable to "clean_tweet" list clean_tweet.append(cleaned) #Create a new dataframe filled with the cleaned tweets clean_tweet = pd.DataFrame(clean_tweet, columns=["Tweet"]) clean_tweet.head() Out[9]: oh good love wait get home play valheim 0 1 hello loveliestonight jumping back made copper... storyline noi found block fighting giant valhe... ayo totally biased anything go watch lil si pl... 4 queen want mistlands solo cheating died time b... In [10]: #Manual stop words removal for confirmation def word count(str): counts = dict() words = str.split() for word in words: if word in counts: counts[word] += 1 else: counts[word] = 1return counts #Take 20 top words a = ' '.join(map(str, clean_tweet["Tweet"])) counted = {k: v for k, v in sorted(word_count(a).items(), key=lambda item: item[1], reverse=True)} counted = dict(Counter(counted).most_common(15)) plt.barh(range(len(counted)), list(counted.values()), align='center') plt.yticks(range(len(counted)), list(counted.keys())) plt.gca().invert_yaxis() plt.show() valheim mistlands new game live update playing time play get stream going like server 1000 1500 2000 2500 3000 The graph shows some words that are stopwords but haven't removed already. The next lines will remove the words from the tweet. In [11]: manual_stop = ['valheim', 'mistlands', 'new', 'game', 'live', 'update', 'playing', 'time', 'play', 'get', 'stream',
'u', 'going', 'server', 'am', 'pm', 'est', 'tonight', 'today', 'tomorrow', 'yesterday'] #Create a new list to store the cleaned tweets manual_clean_tweet = [] for i in range(0,len(tweet)): temp = clean_tweet.iloc[i,0] #Create new list to store filtered tokenized words #Create new variable for tokenized sentence words_temp = word_tokenize(temp) for c in words_temp: #Check if a word is a stop word if c not in manual_stop: #Stores it into the "filtered" list filtered.append(c) #Create a new variable to store the filtered words that are re-combined back into a sentence cleaned = " ".join(filtered) #Insert the "cleaned" variable to "manual_clean_tweet" list manual_clean_tweet.append(cleaned) #Create a new dataframe filled with the cleaned tweets clean_tweet["Tweet"] = manual_clean_tweet clean_tweet.head() Out[11]: **Tweet** oh good love wait home 1 hello loveliestonight jumping back made copper... storyline noi found block fighting giant drago... ayo totally biased anything go watch lil si 4 queen want solo cheating died biome felt impos... In [12]: #Add new columns to the dataframe for the result of the analyzer sentiments = SentimentIntensityAnalyzer() clean_tweet["Positive"] = [sentiments.polarity_scores(i)["pos"] for i in clean_tweet["Tweet"]] clean_tweet["Neutral"] = [sentiments.polarity_scores(i)["neu"] for i in clean_tweet["Tweet"]] clean_tweet['Compound'] = [sentiments.polarity_scores(i)["compound"] for i in clean_tweet["Tweet"]] clean_tweet.head() Out[12]: Tweet Positive Neutral Compound 0 oh good love wait home 0.703 0.297 0.7964 1 hello loveliestonight jumping back made copper... 1.000 0.0000 storyline noi found block fighting giant drago... 0.000 0.772 -0.7717 ayo totally biased anything go watch lil si 0.000 0.745 -0.3384 4 queen want solo cheating died biome felt impos... 0.188 0.348 -0.8360 Originally, the SentimentIntensityAnalyzer has 4 columns Positive, Neutral, Negative, Compound. For the sake of simplicity, the Neutral result are dismissed from the column. Compound column tells which "way" does the sentiment probably headed to. A positive compound means the sentence is more likely be a positive one, and vice versa for the negative compound. True 0 compound means the sentence is completely neutral, but in this research it's recognized as a negative sentence. In [13]: #Create a new column in the dataframe to store the sentiment of each tweets using Compound value score = clean_tweet["Compound"].values sentiment = [] for i in score: **if** i > 0 : sentiment.append(1) sentiment.append(0) clean_tweet["Sentiment"] = sentiment clean_tweet.head() Out[13]: Tweet Positive Neutral Compound Sentiment oh good love wait home 0.703 0.297 0.7964 1 1 hello loveliestonight jumping back made copper... 0 0.000 1.000 0.0000 storyline noi found block fighting giant drago... 0.000 0.772 -0.7717 0 ayo totally biased anything go watch lil si 0.000 0.745 -0.3384 0 4 queen want solo cheating died biome felt impos... 0.188 0.348 -0.8360 0 In [14]: #Create a tokenizer, and fit the Tweet to the tokenizer. t = Tokenizer() t.fit_on_texts(clean_tweet["Tweet"]) In [15]: #Encode the tweets then pad it encoded_tweets = t.texts_to_sequences(tweet) padded_tweets = pad_sequences(encoded_tweets, padding='post', maxlen=100) Manual Analysis In [16]: ax = sns.countplot(x=clean_tweet["Sentiment"]) ax.bar_label(ax.containers[0]) [Text(0, 0, '2676'), Text(0, 0, '2324')] Out[16]: 2676 2500 2324 2000 1500 1000 500 1 Sentiment More than half of the sentiments against Valheim are negative (2676), and the rest are positive (2324). In [17]: #Initialize and show the WordCloud for positive sentiments WC_pos = WordCloud(width = 1000, height = 1000, background_color ='white', min_font_size = 10).generate(' '.join(map(str, (clean_tweet.loc[clean_tweet['Sentiment'] == 1, 'Tweet'])))) plt.figure(figsize = (8, 8), facecolor = None) plt.imshow(WC_pos) plt.axis("off") plt.tight_layout(pad = 0) plt.show() amp later cool survival A hard viking Look ldi save enough next startingrelease night even awesome bui please though let guy .ot pretty people make help bit ready wait stuff house le excited ast home come join first finally getting oh adventure amazing playthrough Work definitely chat content know video In [18]: #Initialize and show the WordCloud for negative sentiments WC_neg = WordCloud(width = 1000, height = 1000, background_color ='white', min_font_size = 10).generate(' '.join(map(str, (clean_tweet.loc[clean_tweet['Sentiment'] == 0, 'Tweet'])))) plt.figure(figsize = (8, 8), facecolor = None) plt.imshow(WC_neg) plt.axis("off") plt.tight_layout(pad = 0) plt.show() christmas minute streaming twitch think chill death right streaming little play trying release later exploring Vorkwaitingswamp know al twitch check community everything run animated biomeplan got thing next kill night take feel house let might building maybe starting build want beat coming gon na yet something sick hang week tree start working a content hey ill even die video let go yagluth chat come make boy shit way wai explore amp probably around bonemass -- Positive Wordcloud Analysis The most obvious positive words that shows up are friend, fun, good, and love. As an online and mostly cooperative game, the game experience would most likely be better when there are friends that are playing together. --Negative Wordcloud Analysis The most obvious negative words that shows up are back, day, one. This could be correlated to the new major update that contains a new update called "Mistland biome". Its difficulties might be the one making players frustated by having their items lost in the new biome, and they are unable to retrieve it back. Modelling In [19]: #Creating feature and target variables X = padded_tweets y = clean_tweet["Sentiment"] In [20]: #Splitting 70:30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=26) In [21]: #Creating the BiLSTM model model = Sequential() model.add(Embedding((len(t.word_index) + 1), 10, input_length = 100)) model.add(Bidirectional(LSTM(5, return_sequences=True))) model.add(Bidirectional(LSTM(5))) model.add(Activation('relu')) model.add(Dense(1,activation='sigmoid')) model.compile(loss=['binary_crossentropy'] , optimizer='Adam', metrics=['accuracy']) model.summary() Model: "sequential" Layer (type) Output Shape Param # embedding (Embedding) 80760 (None, 100, 10) bidirectional (Bidirectiona (None, 100, 10) 640 bidirectional_1 (Bidirectio (None, 10) 640 activation (Activation) (None, 10) dense (Dense) (None, 1) 11 ______ Total params: 82,051 Trainable params: 82,051 Non-trainable params: 0 In [22]: #Compiling the model with 10 epochs model.compile(optimizer='Adam',loss='binary_crossentropy',metrics=['accuracy']) history = model.fit(X_train, y_train, validation_data=(X_test, y_test), batch_size=35, epochs=10) Epoch 1/10 Epoch 2/10 Epoch 3/10 Epoch 4/10 Epoch 5/10 Epoch 6/10 Epoch 7/10 Epoch 8/10 Epoch 9/10 100/100 [===== Epoch 10/10 100/100 [===== =========] - 5s 49ms/step - loss: 0.1371 - accuracy: 0.9820 - val_loss: 0.5058 - val_accuracy: 0.8020 In [23]: #Creating the model's train loss against validation loss plt.plot(history.history['loss']) plt.plot(history.history['val_loss']) plt.title('model train vs validation loss') plt.vlabel('loss') plt.xlabel('epoch') plt.legend(['train', 'validation'], loc='upper right') plt.show() model train vs validation loss 0.7 validation 0.6 0.5 S 0.4 0.3 0.2 The model starts to overfit after the 4 epoch, the validation loss rises starting from the 5th epoch. The next lines will create the same model and trains it only until the 4 epoch. In [27]: #Creating the BiLSTM model model = Sequential() model.add(Embedding((len(t.word_index) + 1), 10, input_length = 100)) model.add(Bidirectional(LSTM(5, return_sequences=True))) model.add(Bidirectional(LSTM(5))) model.add(Activation('relu')) model.add(Dense(1,activation='sigmoid')) model.compile(loss=['binary_crossentropy'] , optimizer='Adam', metrics=['accuracy']) model.summary() Model: "sequential_2" Output Shape Layer (type) Param # embedding_2 (Embedding) (None, 100, 10) 80760 bidirectional 4 (Bidirectio (None, 100, 10) 640 nal) bidirectional_5 (Bidirectio (None, 10) 640 nal) activation_2 (Activation) (None, 10) 0 dense_2 (Dense) (None, 1) 11 _____ Total params: 82,051 Trainable params: 82,051 Non-trainable params: 0 In [28]: #Compiling the model with 4 epochs model.compile(optimizer='Adam',loss='binary_crossentropy',metrics=['accuracy']) history = model.fit(X_train,y_train,validation_data=(X_test,y_test),batch_size=35,epochs=4) Epoch 1/4 Epoch 2/4 Epoch 3/4 Epoch 4/4 In [29]: #Creating the model's train loss against validation loss plt.plot(history.history['loss']) plt.plot(history.history['val_loss']) plt.title('model train vs validation loss')

plt.ylabel('loss')
plt.xlabel('epoch')

plt.show()

0.70

0.65

0.60

0.55 <u>8</u> 0.50

0.45

0.40

0.35

In [33]:

0.0

#Saves the model

...layers\activation

...layers\bidirectional

....vars

....vars

.....vars
.....1

.....vars
.....1

....vars

....vars

.....vars
......0
.....1

.....vars
.....1
.....2

.....vars
...layers\dense
....vars
.....0
.....1

.....vars
.....0
...metrics\mean
....vars
.....0

. 0 1 ...optimizervars 0 110 13 14 1516 17 18 19 2 2021 22 23 24 2526 27 2829 3

......5
.......6
......8
.....9

File Name

Out[33]:

config.json

metadata.json

variables.h5

Keras model archive saving:

['valheim_bilstm_mistlands_twitter.sav']

Modified

2022-12-19 13:58:05

2022-12-19 13:58:05

2022-12-19 13:58:05

Size

3455

1043592

64

...layers\embedding

0.5

1.0

...layers\bidirectional\backward_layer

...layers\bidirectional\forward_layer

...layers\bidirectional\layer

...layers\bidirectional_1

...layers\bidirectional\layer\cell

...layers\bidirectional_1\backward_layer

...layers\bidirectional_1\forward_layer

...layers\bidirectional_1\layer

...metrics\mean_metric_wrapper

...layers\bidirectional_1\layer\cell

...layers\bidirectional_1\backward_layer\cell

...layers\bidirectional_1\forward_layer\cell

...layers\bidirectional\backward_layer\cell

...layers\bidirectional\forward_layer\cell

plt.legend(['train', 'validation'], loc='upper right')

model train vs validation loss

1.5

joblib.dump(model, "valheim_bilstm_mistlands_twitter.sav")

Keras weights file (<HDF5 file "variables.h5" (mode r+)>) saving:

2.0

2.5

3.0

- train - validation

Sentiment Analysis of Valheim (Video Game) based on Twitter Dataset using Keras BiLSTM

Valheim is a video game with a nord-inspired theme. The game's main aspects are exploration and adventure. With the last announcement of the latest major update (Mistland update), the game has gained some attention among the gamer community in twitter. This project aims to analyze the sentiment of Valheim players on Twitter and uses it to train a BiLSTM model. The model would be evaluated

MISTLANDS

and its result should tell whether the model is usable for prediction or not.

Preprocessing

import pandas as pd
import numpy as np

import seaborn as sns
import contractions

import tensorflow as tf

import preprocessor as p

from keras import layers

import nltk

import string

import keras

#Importing essential libraries

import matplotlib.pyplot as plt

from collections import Counter
from wordcloud import WordCloud

from nltk.corpus import stopwords

from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer

import snscrape.modules.twitter as sntwitter

from guess_language import guess_language as guess

from tensorflow.keras.preprocessing.sequence import pad_sequences

from nltk.sentiment.vader import SentimentIntensityAnalyzer

from sklearn.model_selection import train_test_split

from keras.preprocessing.text import Tokenizer

In [1]: