# Tensorflow libraries # Tensorflow libraries import tensorflow as tf from tensorflow import keras from tensorflow.keras.models import Sequential from tensorflow.keras.preprocessing import text, sequence # from tensorflow.keras.preprocessing.sequence import pad\_sequences from tensorflow.keras import regularizers import tensorflow\_hub as hub # sklearn libraries from sklearn.model\_selection import train\_test\_split from sklearn.utils import class\_weight from sklearn.metrics import roc\_auc\_score from gensim.models import Word2Vec # Word2Vec module from gensim.parsing.preprocessing import preprocess\_string, strip\_tags, strip\_punctuation, remove\_stopwords, strip\_numeric, stem\_text import os for dirname, \_, filenames in os.walk('/kaggle/input'): for filename in filenames: print(os.path.join(dirname, filename)) # You can write up to 5GB to the current directory (/kaggle/working/) that gets preserved as output when you create a version using "Save & Run All" # You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session /kaggle/input/fake-news-content-detection/train.csv /kaggle/input/fake-news-content-detection/test.csv /kaggle/input/fake-news-content-detection/sample submission.csv **Load Dataset** In [203... train\_data = pd.read\_csv("/kaggle/input/fake-news-content-detection/train.csv") test\_data = pd.read\_csv("/kaggle/input/fake-news-content-detection/test.csv") submission\_data = pd.read\_csv("/kaggle/input/fake-news-content-detection/sample submission.csv") In [204... # Sample data from training data train\_data.sample(3) Out [204... Labels Text Text\_Tag 4 Rick Santorum says Rick Perry requested 1,200 ... congress,federal-budget 8026 0 Studies suggest the 2017 College Football Play... jobs,sports 1 The debt comes up all the time in town meeting... debt,federal-budget In [205... # Dataset information train\_data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 10240 entries, 0 to 10239 Data columns (total 3 columns): # Column Non-Null Count Dtype 0 Labels 10240 non-null int64 10240 non-null object 1 Text 2 Text\_Tag 10238 non-null object dtypes: int64(1), object(2) memory usage: 240.1+ KB In [206... train\_data[train\_data.duplicated(['Text'])] Labels Text Text\_Tag 1014 2 On abortion abortion, candidates-biography On support for gay marriage. 1814 civil-rights,families,gays-and-lesbians,marriage 1846 Obama says Iran is a 'tiny' country, 'doesn't ... foreign-policy debates, elections, states 2697 On repealing the 17th Amendment Four balanced budgets in a row, with no new ta... job-accomplishments,jobs,state-budget,state-fi... 2846 3256 On a cap-and-trade plan. cap-and-trade,climate-change,environment 4386 On the Trans-Pacific Partnership. trade 4839 2 During Sherrod Browns past decade as a D.C. po... economy,job-accomplishments,jobs 4940 On changing the rules for filibusters on presi... congressional-rules 6759 2 On torture. human-rights,terrorism 6784 On support for the Export-Import Bank trade 7248 On the status of illegal immigrants immigration 7647 5 Six justices on the U.S. Supreme Court have be... congress,legal-issues,supreme-court 8906 Says Mitt Romney flip-flopped on abortion. abortion,message-machine-2012 9400 Twenty million Americans are out of work. jobs 9642 On changing the rules for filibusters on presi... congressional-rules 9750 Some 20,000 Delphi salaried retirees lost up t... corporations, economy train\_data = train\_data.drop\_duplicates(['Text']) In [208... train\_data.sample(3) Labels Text Text\_Tag Nobody is leaving Memphis. Thats a myth. 1767 4 census,population,taxes 3299 2 Says black women are fastest-growing demograph... guns Ken Buck wants to outlaw abortion, even in cas... abortion, message-machine 4179 ## word cloud for text\_tag cloud = WordCloud(width=1300, height=950).generate(" ".join(train\_q.astype(str))) plt.figure(figsize=(15, 10)) plt.imshow(cloud) plt.axis('off') from wordcloud import WordCloud cloud = WordCloud(width=1300, height=950).generate(" ".join(train\_qs.astype(str))) plt.figure(figsize=(15, 10)) plt.imshow(cloud) plt.axis('off') United State Barack Obama Cost In [252... X = train\_data['Processed'] y = train\_data['Labels'] y\_category = keras.utils.to\_categorical(y, 6) # Split data into Train and Holdout as 80:20 ratio X\_train, X\_valid, y\_train, y\_valid = train\_test\_split(X, y\_category, shuffle=True, test\_size=0.33, random\_state=111) print("Train shape : {}, Holdout shape: {}".format(X\_train.shape, X\_valid.shape)) Train shape : (6849,), Holdout shape: (3374,) In [253... def word\_embedding(train, test, max\_features, max\_len=200): # Keras Tokenizer class object tokenizer = text.Tokenizer(num\_words=max\_features) tokenizer.fit\_on\_texts(train) train\_data = tokenizer.texts\_to\_sequences(train) test\_data = tokenizer.texts\_to\_sequences(test) # Get the max\_len vocab\_size = len(tokenizer.word\_index) + 1 # Padd the sequence based on the max-length x\_train = sequence.pad\_sequences(train\_data, maxlen=max\_len, padding='post') x\_test = sequence.pad\_sequences(test\_data, maxlen=max\_len, padding='post') # Return train, test and vocab size return tokenizer, x\_train, x\_test, vocab\_size except ValueError as ve: raise(ValueError("Error in word embedding {}".format(ve))) In [308... max\_features = 5000  $max_len = 128$ output\_dim = len(np.unique(y)) # Test data X\_test = test\_data['Processed'] tokenizer, x\_pad\_train, x\_pad\_valid, vocab\_size = word\_embedding(X\_train, X\_valid, max\_features) In [309... # Test data X\_test = test\_data['Processed'] tokenizer.fit\_on\_sequences(X\_test) X\_test\_seq = tokenizer.texts\_to\_sequences(X\_test) x\_pad\_test = sequence.pad\_sequences(X\_test\_seq, maxlen=max\_len, padding='post') In [310... **def** compute\_classweights(target): Computes the weights of the target values based on the samples :param target: Y-target variable :return: dictionary object # compute class weights class\_weights = class\_weight.compute\_class\_weight('balanced', np.unique(target), target) # make the class weight list into dictionary weights = {} # enumerate the list for index, weight in enumerate(class\_weights): weights[index] = weight return weights # Get the class weights for the target variable weights = compute\_classweights(y) In [311... weights Out[311... {0: 1.0307521677757612, 1: 0.8574903539674551, 2: 0.8078868342026236, 3: 0.868859425463199, 4: 2.0307906237584428, 5: 1.017821585025886} In [312... X\_train.sample(3) Out[312... 3960 taxes new funding state sent floridians billio... 6085 tens texas voters rick perrys wendy defeat law... 6655 voted baldwin candidate extreme increase budge... Name: Processed, dtype: object In [313... rnn\_model = Sequential([ keras.layers.Embedding(vocab\_size,128, input\_length=max\_len), keras.layers.BatchNormalization(), keras.layers.Dense(128, activation='relu', kernel\_regularizer=tf.keras.regularizers.L2(0.002)), keras.layers.GlobalMaxPool1D(), # Remove flatten layer keras.layers.Dense(64, activation='relu', kernel\_regularizer=tf.keras.regularizers.L2(0.002)), keras.layers.Dropout(0.3), keras.layers.Dense(32, activation='relu', kernel\_regularizer=tf.keras.regularizers.L2(0.002)), keras.layers.Dropout(0.3), keras.layers.Dense(output\_dim, activation='softmax') return rnn\_model In [314... rnn\_model.summary() Model: "sequential\_25" Layer (type) Output Shape Param # embedding\_26 (Embedding) (None, 128, 128) 1515520 batch\_normalization\_26 (Batc (None, 128, 128) 512 16512 dense\_99 (Dense) (None, 128, 128) global\_max\_pooling1d\_24 (Glo (None, 128) 8256 dense\_100 (Dense) (None, 64) dropout\_51 (Dropout) (None, 64) 2080 dense\_101 (Dense) (None, 32) dropout\_52 (Dropout) (None, 32) dense\_102 (Dense) (None, 6) 198 \_\_\_\_\_ Total params: 1,543,078 Trainable params: 1,542,822 Non-trainable params: 256 In [315... # Compile the model rnn\_model.compile(optimizer=tf.keras.optimizers.Adam(1e-3), loss=keras.losses.CategoricalCrossentropy(from\_logits=True), metrics=[tf.metrics.AUC()]) In [316... history = rnn\_model.fit(x\_pad\_train, y\_train, batch\_size=512, epochs=20, verbose=1, validation\_data=(x\_pad\_valid, y\_valid), class\_weight=weights) Epoch 1/20 Epoch 3/20 Epoch 4/20 Epoch 5/20 Epoch 6/20 Epoch 7/20 Epoch 9/20 Epoch 10/20 Epoch 11/20 Epoch 13/20 Epoch 14/20 Epoch 15/20 Epoch 17/20 Epoch 19/20 Epoch 20/20 In [317... results = rnn\_model.evaluate(x\_pad\_valid, y\_valid) In [318... y\_preds = rnn\_model.predict\_proba(x\_pad\_test, batch\_size=256) In [319... y\_preds[:,0] Out[319... array([0.10799249, 0.10470885, 0.20419715, ..., 0.11288042, 0.07156285, 0.0953638 ], dtype=float32) In [320... final\_df = pd.DataFrame({'0': y\_preds[:,0], '1': y\_preds[:,1], '2': y\_preds[:,2], '3': y\_preds[:,3], '4': y\_preds[:,4], '5': y\_preds[:,5]}, index=test\_data.index) In [321... final\_df

In [202... import re

Out [321...

1 2 3 4

**0** 0.107992 0.101322 0.150519 0.378211 0.134148 0.127807

1 0.104709 0.158395 0.332498 0.163373 0.189925 0.051100

**2** 0.204197 0.133102 0.249838 0.167857 0.143589 0.101418

**3** 0.167783 0.325967 0.100499 0.088033 0.130143 0.187575

**4** 0.174848 0.282983 0.110808 0.102282 0.210029 0.119051

**1262** 0.125772 0.097854 0.083867 0.338114 0.165196 0.189197

**1263** 0.142004 0.104889 0.095647 0.265502 0.085052 0.306905

**1264** 0.112880 0.168333 0.251629 0.280145 0.070717 0.116296

**1265** 0.071563 0.260155 0.113663 0.235579 0.170863 0.148177

**1266** 0.095364 0.361185 0.258039 0.103341 0.052368 0.129703

import warnings

import numpy as np # linear algebra

warnings.filterwarnings('ignore')

from nltk.stem import WordNetLemmatizer

import pandas as pd # data processing, CSV file I/O (e.g. pd.read\_csv)

# Input data files are available in the read-only "../input/" directory

# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

1267 rows × 6 columns

1207 TOWS ^ O COIDITIES

In [322... final\_df.to\_csv("fake\_news\_ann\_08.csv", index=False)