	Vicente De Leon Deep Learning Principles: Homework 3 UID: 2001014594
	LSTM Recurrent Neural Networks for Classification
	Dependencies TensorFlow
	if gpus:
	<pre>print("GPU available.") else: print("GPU not available.") GPU available. TensorBoard</pre>
In []: In []:	!pip install -q -U tensorboard==2.12.0 #!pip install tensorboard==2.12.0 import datetime
In []:	from tensorflow.keras.callbacks import TensorBoard Pandas, Numpy, Keras
In []:	<pre>import pandas as pd import numpy as np from sklearn.model_selection import train_test_split from sklearn.preprocessing import LabelEncoder from tensorflow.keras.preprocessing.text import Tokenizer from tensorflow keras preprocessing sequence import pad sequences</pre>
	from tensorflow.keras.preprocessing.sequence import pad_sequences from tensorflow.keras.models import Sequential from keras.utils.vis_utils import plot_model from keras.layers import LSTM from keras.layers import Embedding from keras.layers import Dropout from tensorflow.keras.layers import Dense
	Exploring the IMDB Data Pete some from Kerrale IIMDB Detect of 50k Marie Parioural for I STM Binary Classification purposes
In []:	Data came from Kaggle "IMDB Dataset of 50k Movie Reviews" for LSTM Binary Classification purposes. Kaggle: https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews?resource=download df = pd.read_csv('IMDB Dataset.csv') df.head()
Out[]:	review sentiment O One of the other reviewers has mentioned that positive A wonderful little production. b /> > The positive I thought this was a wonderful way to spend ti positive
	3 Basically there's a family where a little boy negative 4 Petter Mattei's "Love in the Time of Money" is positive Summary statistics
<pre>In []: Out[]:</pre>	count 50000 50000
	top Loved today's show!!! It was a variety and not positive freq 5 25000
In []: Out[]:	Checking for missing values df.isnull().sum() review 0 sentiment 0 dtype: int64
	Preprocessing Label Encoding is used to handle categorical variables (column "sentiment" -> positive vs negative). Using this technique will allow me to convert (encode) the positive and negative sentiment into class 1 and class 0. This is for binary classification purposes. I am also splitting the data using SKlearn train_test_split() by setting the test size to 0.2 and random state to 42. These are the numbers I used for almost all the train_test_split steps last semester. I' m also creating a validation set, so its loss can be displayed in TensorBoard
In []:	along the training loss. X = df['review'] y = df['sentiment'] label_encoder = LabelEncoder()
In []:	<pre> Tabel_encoder = LabelEncoder() y = label_encoder.fit_transform(y) # This is how I used to split data for my APML CNN Model test size 20% and random state 42. # Also, a validation set was implemented. X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train, test_size=0.2, random_state=42) X_train, X_valid, y_train, y_valid = train_test_split(X_train, y_train, test_size=0.2, random_state=42)</pre>
Out[]:	<pre>X_train.shape (32000,) X_test.shape</pre>
	I am going to be using TensorFlow Tokenizer() to pre-process the text corpus by tokenizing text into words. This class allows to vectorize a text corpus, by turning each text into either a sequence of integers (each integer being the index of a token in a dictionary) or into a vector where the coefficient for each token could be binary, based on word count, based on tf-idf
	By default, all punctuation is removed, turning the texts into space-separated sequences of words (words may include the 'character). These sequences are then split into lists of tokens. They will then be indexed or vectorized. O is a reserved index that won't be assigned to any word. Reference: https://dev.to/balapriya/tokenization-and-sequencing-in-tensorflow-3p0n
	Back when I was taking "Intro to NLP" for Data Science with professor Olga Scrivner, I was using the below codes to normalize and preprocess data (course material) to create a Streamlit Sentiment Analyzer using Spotify data. However, it seems TensorFlow Tokenizer() seems to be more powerful, plus it goes hand to hand with the tutorial I am currently using.
	<pre>stop_words = nltk.corpus.stopwords.words('english') def normalize_document(doc): doc = re.sub(r'[^a-zA-20-9\s]', '', doc, re.I re.A) doc = doc.lower() doc = doc.strip() tokens = nltk.word_tokenize(doc) filtered_tokens = [token for token in tokens if token not in stop_words] doc = '.join(filtered_tokens) return doc</pre>
	<pre>normalize_corpus = np.vectorize(normalize_document) df = pd.read_csv(csv_file) df[df.Review.str.strip() == ''].shape[0] norm_corpus = normalize_corpus(df['Review']) df['Clean Review'] = norm corpus</pre>
	<pre>df = df[['Review', 'Clean Review']] df.replace(r'^(\s?)+\s', np.nan, regex=True) df.dropna().reset_index(drop=True) df['Clean Review'] = norm_corpus df = df[['Review', 'Clean Review']] st.write(df.head())</pre>
	OOV -> Out of Vocab will replace any unkown words with this special token. limiting vocabulary size to 10000 in the dataset. Only the top num_words will be considered.
In []:	tokenizer = Tokenizer(num_words=10000,oov_token='<00V>') tokenizer.fit_on_texts(X_train) word_index -> Used to see how tokens have been created and the indices assigned to words. word_index = tokenizer.word_index print(word_index)
	The Tokenizer is then fit on the training data only (X_train) to learn the vocabulary. Later, the texts_to_sequences method is applied separately to the training and testing sets to convert the text into sequences. Finally, the sequences are padded to a consistent length using pad_sequences. tokenizer.texts_to_sequences() will be used to convert the text into sequences. Padding will also be implemented to the sequence to same length using pad_sequences.
In []:	This is a crucial step of preprocessing because as we know, machines can't understand text data as we do. In order to make the machine understand what we want is to convert the data into integers (sequence of numbers). Also, padding = 'post' will help us preserve the original order of the words. X_train_sequences = tokenizer.texts_to_sequences(X_train)
	LSTM Model This model was taken from tutorial Google Colab tutorial "IMDB Movie review.ipynb". Extra ideas were taken from Applied Machine Learning course, Introduction to NLP, and internet sources which helped me complete the model regarding regularization techniques, TensorBoard, model structure etc. Parameters like max_length, epochs, dropout numbers were randomly chosen after many experiments. Why regularization techniques using L2 and dropout? My model was suffering from cases of very low accuracy scores, high training
	losses and low validation losses. My model was mostly overfitting and sometimes underfitting. After playing with the max_length, epochs, and integers within the dropout I managed to get the highest accuracy score and best validation/training loss behavior. Validation and Training behavior can be seeing using the TensorBoard dashboard below. Other options like histograms, graphs etc. can be selected within the dashboard. Very useful tool to see real time experimentation. I got the idea from Applied Machine Learning course. Colab tutotial: https://colab.research.google.com/github/dipanjannC/Deep-Learning/blob/master/IMDB_Movie_review.ipynb#scrollTo=L0y6pKlrHfHS
	 The tutorial is following Keras IMDB movie dataset and I am following the Kaggle IMDB Movie 50k Reviews. LSTM Layer is known as the Long Short Term Memory layer (type of RNN model). LSTM is a type of recurrent neural network but is better than traditional recurrent neural networks in terms of memory. Having a good hold over memorizing certain patterns LSTMs perform fairly better. This type of model is extremely useful in memorizing important information. The model will be able to find out the actual meaning in input string and will give the most accurate output class.
In []:	Training and TensorBoard from tensorflow.keras.regularizers import 12 Extra Information:
	• The embedding layer has weights that are learned. This layer is responsable for converting the sequences (input) into embedding vectors and learns the word embedding during training. • Sigmoid Activation Function Will always returns a value between 0 and 1. $Sigmoid(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$
	• Binary Cross Entropy Computes the cross-entropy loss between true labels and predicted labels. $\mathcal{L}_{\mathrm{BCE}}(y,\hat{y}) = -\frac{1}{N}\sum_{i=1}^N y_i \log(\hat{y_i}) + (1-y_i)\log(1-\hat{y_i})$
In []:	<pre>vocab_size = len(tokenizer.word_index) + 1 embedding_dim = 32 input_length = max_length batch = 32 # Keras batch size default "APML" epochs = 15 # 10</pre>
	<pre>with tf.device('/GPU:0'): # using GPU for faster run model = Sequential() model.add(Embedding(vocab_size, embedding_dim, input_length = max_length)) model.add(Dropout(0.5)) #5, 5, 5, 5, 5, 5, 3 model.add(LSTM(100, kernel_regularizer=12(0.01))) model.add(Dropout(0.1)) #1, 15, 45, 5, 4, 4</pre>
	<pre>model.add(Dense(1, activation = 'sigmoid')) model.compile(loss='binary_crossentropy', # Loss also used within my CNN model from last semester</pre>
	<pre>with tf.device('/GPU:0'): history = model.fit(X_train_padded, y_train,</pre>
	Epoch 1/15 1000/1000 [==================================
	1000/1000 [============] - 16s 16ms/step - loss: 0.6264 - accuracy: 0.6681 - val_loss: 0.5331 - val_accuracy: 0.7830 Epoch 4/15 1000/1000 [==============] - 13s 13ms/step - loss: 0.6560 - accuracy: 0.6180 - val_loss: 0.6617 - val_accuracy: 0.6201 Epoch 5/15 1000/1000 [===============] - 14s 14ms/step - loss: 0.6255 - accuracy: 0.6800 - val_loss: 0.6663 - val_accuracy: 0.5761 Epoch 6/15 1000/1000 [================] - 15s 15ms/step - loss: 0.5366 - accuracy: 0.7643 - val_loss: 0.6165 - val_accuracy: 0.7159 Epoch 7/15
	1000/1000 [==================================
	Epoch 11/15 1000/1000 [==================================
In []:	Epoch 14/15 1000/1000 [============] - 12s 12ms/step - loss: 0.3457 - accuracy: 0.8637 - val_loss: 0.3578 - val_accuracy: 0.8529 Epoch 15/15 1000/1000 [==================================
	Layer (type) Output Shape Param #
	lstm_1 (LSTM) (None, 100) 53200 dropout_3 (Dropout) (None, 100) 0 dense_1 (Dense) (None, 1) 101
In []: Out[]:	Total params: 3,294,965 Trainable params: 3,294,965 Non-trainable params: 0 plot_model(model, to_file='model_plot.png', show_shapes=True, show_layer_names=True)
	InputLayer output: [(None, 250)]
	embedding_1 input: (None, 250) Embedding output: (None, 250, 32) dropout_2 input: (None, 250, 32)
	dropout_2 input: (None, 250, 32) Dropout output: (None, 250, 32) lstm_1 input: (None, 250, 32)
	LSTM output: (None, 100) dropout_3 input: (None, 100)
	Dropout output: (None, 100) dense_1 input: (None, 100)
In []:	Dense output: (None, 1) *tensorboardlogdir logs/fit
	TensorBoard TIME SERIES SCALARS GRAPHS DISTRIBUTIONS HISTOGRAMS □ Show data download links □ Ignore outliers in chart scaling □ Ignore outliers in chart s
	Tooltip sorting method: default Smoothing 0.6 © One of the prochaccuracy tags epoch_accuracy tags epoch
	Horizontal Axis STEP RELATIVE WALL Runs Write a regex to filter runs O.5
	Write a regex to filter runs □ 20230530-042207/train □ 20230530-042207/validation □ 20230530-042754/validation □ 30230530-042754/validation □ 40500000000000000000000000000000000000
	© 20230530-043137/rain C O 20230530-043137/validation TOGGLE ALL RUNS logs/fit 0.6 0.55 0.5 0.5 0.5 0.2 4 6 8 10 12 14
In []:	
	<pre>min_loss = min(history.history['loss']) # APML Homework min_val_loss = min(history.history['val_loss']) max_tr_acc = max(history.history['accuracy']) max_val_acc = max(history.history['val_accuracy'])</pre>
	min_val_loss = min(history.history['val_loss']) max_tr_acc = max(history.history['accuracy']) max_val_acc = max(history.history['val_accuracy']) print(f'The lowest training loss after 15 epochs is: {min_loss}') print(f'The lowest val loss after 15 epochs is: {min_val_loss}') print(f'The hightest training accuracy after 15 epochs is: {max_tr_acc}') print(f'The hightest validation accuracy after 15 epochs is: {max_val_acc}') The lowest training loss after 15 epochs is: 0.30322331190109253
	<pre>min_val_loss = min(history.history['val_loss']) max_tr_acc = max(history.history['accuracy']) max_val_acc = max(history.history['val_accuracy']) print(f'The lowest training loss after 15 epochs is: {min_loss}') print(f'The lowest val loss after 15 epochs is: {min_val_loss}') print(f'The lowest val loss after 15 epochs is: {min_val_loss}') print(f'The hightest training accuracy after 15 epochs is: {max_tr_acc}') print(f'The hightest validation accuracy after 15 epochs is: {max_val_acc}')</pre>
	min_val_loss = min(history.history('val_loss')) max_ral_scc = max(history.history('val_accuracy')) max_val_scc = max(history.history('val_accuracy')) max_val_scc = max(history.history('val_accuracy')) max_val_scc = max(history.history('val_accuracy')) mrin(f'the lowest training loss after 15 epochs is: {min_loss}') print(f'the hightest training accuracy after 15 epochs is: {max_val_acc}') print(f'the hightest validation accuracy after 15 epochs is: {max_val_acc}') print(f'the hightest validation accuracy after 15 epochs is: {max_val_acc}') The lowest training loss after 15 epochs is: 0.3222331190109253 The lowest val loss after 15 epochs is: 0.3223331190109253 The hightest training accuracy after 15 epochs is: 0.8712649804490765 As mentioned before, our highest train accuracy is of 88% after experimenting with different cases, integers, and techniques. Since this is my first time working with LSTM Model, I wanted it to try and keep it simple regarding accuracy score and epochs. Many tutorials found on internet has pretty good accuracy scores, however many models were extremely overfitting. The following image, shows one of my first experiments where my model was sufflering from underfitting:
	min yal_loss = min(history/history[val_loss']) max trace = max(history.history[accuracy']) max_val_ace = max(history.history[accuracy']) max_val_ace = max(history.history[val_accuracy']) max_val_ace = max(history.history[val_accuracy']) print(f'The lowest training loss after 15 epochs is: {min_val_loss}') print(f'The highest training accuracy after 15 epochs is: {max_val_acc}') print(f'The highest validation accuracy after 15 epochs is: {max_val_acc}') print(f'The highest training loss after 15 epochs is: 0.30322331190109253 The lowest val loss after 15 epochs is: 0.303229319916167 The highest training accuracy after 15 epochs is: 0.8812187314033508 The highest validation accuracy after 15 epochs is: 0.8812187314033508 As mentioned before, our highest train accuracy is of 88% after experimenting with different cases, integers, and techniques. Since this is my first time working with LSTM Model, I wanted it to try and keep it simple regarding accuracy score and epochs. Many tutorials found on internet has pretty good accuracy scores, however many models were extremely overfitting. The following image, shows one of my first experiments where my model was sufffering from underfitting:
	main_vai_tous = main_history_history vectory tourner() main_vai_tous = main_history_history vectory vectory main_vai_tous = main_history_history vectory vectory main_vai_tous = main_history_history vectory vectory main_vai_tous = main_history_history vectory main_vai_tous = main_history vectory vectory main_vai_tous = main_history vectory vectory main_vai_tous = main_history vectory vectory vectory vec
	max for soc = max[history_history] (secure(y)) max vi acc = max[history,history] (secure(y)) max vi acc = max[history
In []:	The algorithms recovered by the control of the cont
In []:	mis_tal_dos= _mantateop_aistery(_val_sets_y) mantateop_aistery(_val_sets_y) mantateop_aistery(_
	Evaluation and Testing Evaluation and Testing Proc. Pro
	The field company and the street of the stre
	The second secon
	State of the state processing the case () State of the state processing the case () State of the state of t
	Figure 1. The contract of the
	For the set of the property states of the set of the se
	Residuacy Andreasty Market Service (1) in a control (1) i
	Secretary and the state of the
	Section 1. The section of the section 1. The sectio
	## The Property of the Propert
	The state of the s
	The state of the s
	### Property of the Company of the C
	Section of the state of the sta