

Sentiment Analysis Using Python

Submitted By - Akshat Charan (209301286)
Pranshu Raghuwanshi (209301252)
Ansh Harjai (209301140)

Supervised And Unsupervised Learning

Supervised Learning	Unsupervised Learning
Supervised learning algorithms are trained using labeled data.	Unsupervised learning algorithms are trained using unlabeled data.
Supervised learning model takes direct feedback to check if it is predicting correct output or not.	Unsupervised learning model does not take any feedback.
Supervised learning model predicts the output.	Unsupervised learning model finds the hidden patterns in data.
In supervised learning, input data is provided to the model along with the output.	In unsupervised learning, only input data is provided to the model.
The goal of supervised learning is to train the model so that it can predict the output when it is given new data.	The goal of unsupervised learning is to find the hidden patterns and useful insights from the unknown dataset.
Supervised learning needs supervision to train the model.	Unsupervised learning does not need any supervision to train the model.
Supervised learning can be categorized in Classification and Regression problems.	Unsupervised Learning can be classified in Clustering and Associations problems.
Supervised learning can be used for those cases where we know the input as well as corresponding outputs.	Unsupervised learning can be used for those cases where we have only input data and no corresponding output data.
Supervised learning model produces an accurate result.	Unsupervised learning model may give less accurate result as compared to supervised learning.
Supervised learning is not close to true Artificial intelligence as in this, we first train the model for each data, and then only it can predict the correct output.	Unsupervised learning is more close to the true Artificial Intelligence as it learns similarly as a child learns daily routine things by his experiences.
It includes various algorithms such as Linear Regression, Logistic Regression, Support Vector Machine, Multi-class Classification, Decision tree, Bayesian Logic, etc.	It includes various algorithms such as Clustering, KNN, and Apriori algorithm.

Natural Language Processing

Natural Language Processing or NLP is a branch of Artificial Intelligence which deal with bridging the machines understanding humans in their Natural Language. Natural Language can be in form of text or sound, which are used for humans to communicate each other. NLP can enable humans to communicate to machines in a natural way.

Text Classification

Text Classification is a process involved in Sentiment Analysis. It is classification of peoples opinion or expressions into different sentiments. Sentiments include Positive, Neutral, and Negative, Review Ratings and Happy, Sad. Sentiment Analysis can be done on different consumer centered industries to analyse people's opinion on a particular product or subject.

Data Preprocessing

Tweet texts often consists of other user mentions, hyperlink texts, emoticons and punctuations. In order to use them for learning using a Language Model. We cannot permit those texts for training a model. So we have to clean the text data using various preprocessing and cleansing methods.

Stemming/ Lemmatization

For grammatical reasons, documents are going to use different forms of a word, such as write, writing and writes. Additionally, there are families of derivationally related words with similar meanings. The goal of both stemming and lemmatization is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form.

Stemming usually refers to a process that chops off the ends of words in the hope of achieving goal correctly most of the time and often includes the removal of derivational affixes.

Lemmatization usually refers to doing things properly with the use of a vocabulary and morphological analysis of words, normally aiming to remove inflectional endings only and to return the base and dictionary form of a word

Hyperlinks and Mentions

Twitter is a social media platform where people can tag and mentions other people's ID and share videos and blogs from internet. So the tweets often contain lots of Hyperlinks and twitter mentions.

Twitter User Mentions - Eg. @arunrk7, @andrewng

Hyperlinks - Eg. <https://keras.io>, <https://tensorflow.org>

Stopwords

Stopwords are commonly used words in English which have no contextual meaning in an sentence. So therefore we remove them before classification.

Some stopwords are...

```
> stopwords("english")  
[1] "i"          "me"          "my"          "myself"      "we"  
[6] "our"        "ours"        "ourselves"   "you"         "your"  
[11] "yours"      "yourself"    "yourselves"  "he"          "him"  
[16] "his"        "himself"     "she"         "her"         "hers"  
[21] "herself"    "it"          "its"         "itself"      "they"  
[26] "them"       "their"       "theirs"      "themselves"  "what"  
[31] "which"      "who"         "whom"        "this"        "that"  
[36] "these"      "those"       "am"          "is"          "are"  
[41] "was"        "were"        "be"          "been"        "being"  
[46] "have"       "has"         "had"         "having"      "do"  
[51] "does"       "did"         "doing"       "would"       "should"  
[56] "could"      "ought"       "i'm"         "you're"      "he's"  
[61] "she's"      "it's"        "we're"       "they're"     "i've"  
[66] "you've"     "we've"       "they've"     "i'd"         "you'd"  
[71] "he'd"       "she'd"       "we'd"        "they'd"      "i'll"  
[76] "you'll"     "he'll"       "she'll"      "we'll"       "they'll"  
[81] "isn't"      "aren't"      "wasn't"      "weren't"     "hasn't"  
[86] "haven't"    "hadn't"      "doesn't"     "don't"       "didn't"  
[91] "won't"      "wouldn't"    "shan't"      "shouldn't"   "can't"  
[96] "cannot"     "couldn't"    "mustn't"     "let's"       "that's"  
[101] "who's"      "what's"      "here's"      "there's"     "when's"  
[106] "where's"    "why's"       "how's"       "a"           "an"
```

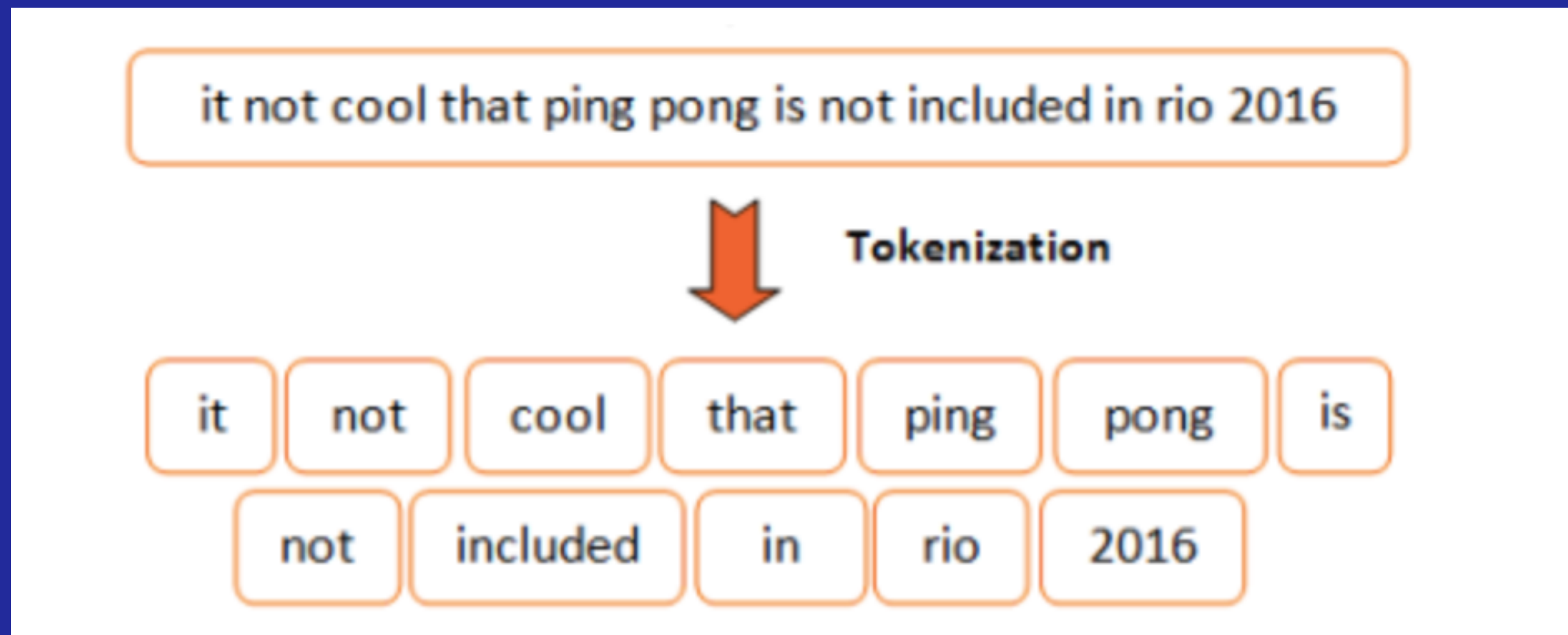

Train and Test Split

The train-test split procedure is used to estimate the performance of machine learning algorithms when they are used to make predictions on data not used to train the model.

It is a fast and easy procedure to perform, the results of which allow you to compare the performance of machine learning algorithms for your predictive modeling problem. Although simple to use and interpret, there are times when the procedure should not be used, such as when you have a small dataset and situations where additional configuration is required, such as when it is used for classification and the dataset is not balanced.

Tokenization

Given a character sequence and a defined document unit, tokenization is the task of chopping it up into pieces, called tokens, perhaps at the same time throwing away certain characters, such as punctuation. The process is called Tokenization.



Label Encoding

We are building the model to predict class in encoded form (0 or 1 as this is a binary classification). We should encode our training labels to encodings.

Word Embedding

In Language Model, words are represented in a way to intend more meaning and for learning the patterns and contextual meaning behind it.

Word Embedding is one of the popular representation of document vocabulary. It is capable of capturing context of a word in a document, semantic and syntactic similarity, relation with other words, etc.

Basically, it's a feature vector representation of words which are used for other natural language processing applications.

We could train the embedding ourselves but that would take a while to train and it wouldn't be effective. So going in the path of Computer Vision, here we use Transfer Learning. We download the pre-trained embedding and use it in our model.

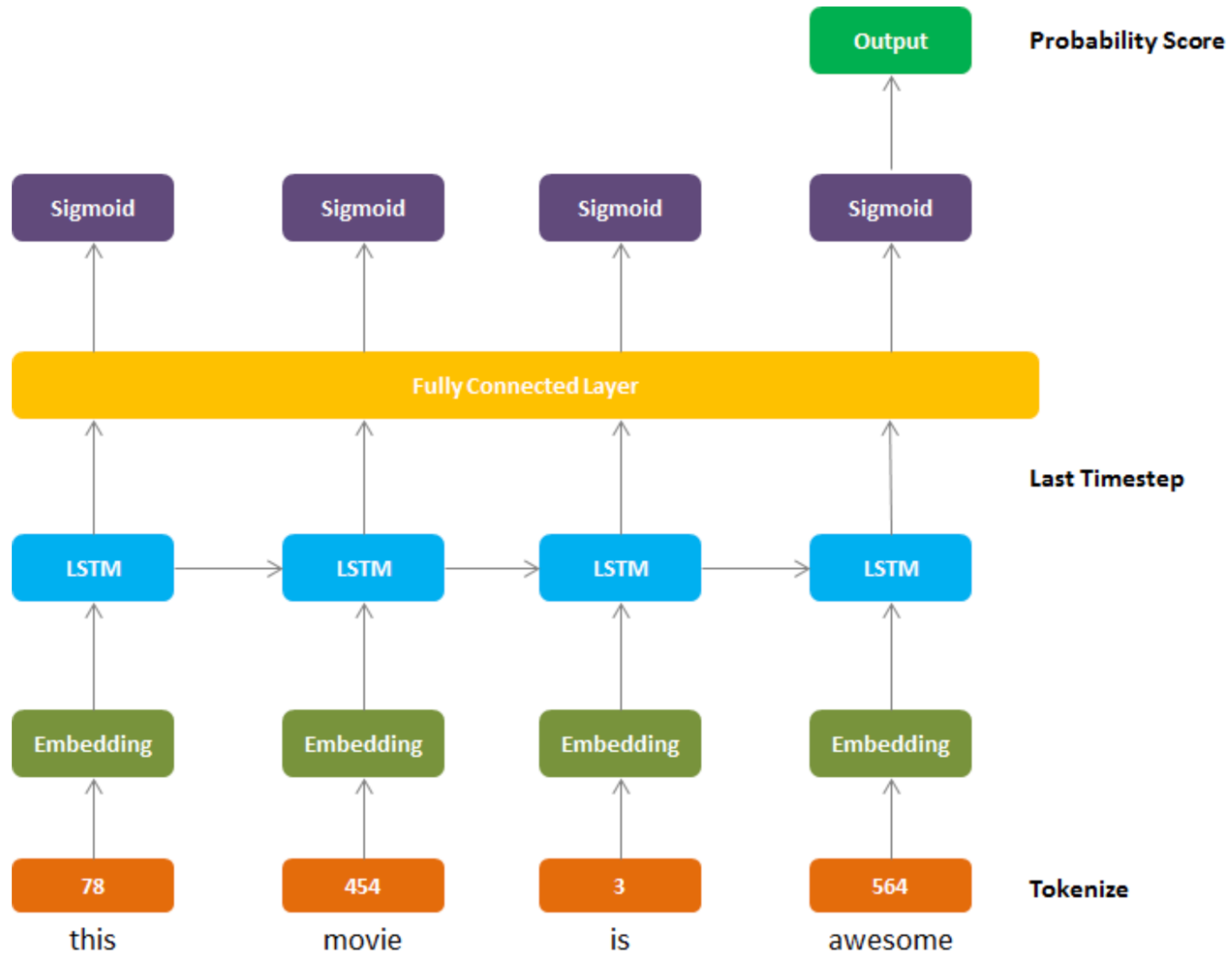
The pretrained Word Embedding like GloVe & Word2Vec gives more insights for a word which can be used for classification. If you want to learn more about the Word Embedding, please refer some links that I left at the end of this notebook.

Model Training - LSTM

We are clear to build our Deep Learning model. While developing a DL model, we should keep in mind of key things like Model Architecture, Hyperparameter Tuning and Performance of the model.

As you can see in the word cloud, the some words are predominantly feature in both Positive and Negative tweets. This could be a problem if we are using a Machine Learning model like Naive Bayes, SVD, etc.. That's why we use Sequence Models.

Sequence Model



Recurrent Neural Networks can handle a sequence of data and learn a pattern of input sequence to give either sequence or scalar value as output. In our case, the Neural Network outputs a scalar value prediction.

For model architecture, we use

- 1) Embedding Layer - Generates Embedding Vector for each input sequence.
- 2) Conv1D Layer - Its using to convolve data into smaller feature vectors.
- 3) LSTM - Long Short Term Memory, its a variant of RNN which has memory state cell to learn the context of words which are at further along the text to carry contextual meaning rather than just neighbouring words as in case of RNN.
- 4) Dense - Fully Connected Layers for classification

Optimization Algorithm

Our code uses Adam, optimization algorithm for Gradient Descent.

Callbacks

Callbacks are special functions which are called at the end of an epoch. We can use any functions to perform specific operation after each epoch. I used two callbacks here,

LRScheduler - It changes a Learning Rate at specific epoch to achieve more improved result. In this notebook, the learning rate exponentially decreases after remaining same for first 10 Epoch.

ModelCheckPoint - It saves best model while training based on some metrics. Here, it saves the model with minimum Validity Loss.

Model Evaluation

Now that we have trained the model, we can evaluate its performance. We will use some evaluation metrics and techniques to test the model.

Let's start with the Learning Curve of loss and accuracy of the model on each epoch.

The model will output a prediction score between 0 and 1. We can classify two classes by defining a threshold value for it. In our case, I have set 0.5 as THRESHOLD value, if the score above it. Then it will be classified as POSITIVE sentiment.

Confusion Matrix

Confusion Matrix provide a nice overlook at the model's performance in classification task. In the field of machine learning and specifically the problem of statistical classification, a confusion matrix, also known as an error matrix, is a specific table layout that allows visualization of the performance of an algorithm, typically a supervised learning one.

	precision	recall	f1-score	support
Negative	0.79	0.78	0.78	160542
Positive	0.78	0.79	0.78	159458
accuracy			0.78	320000
macro avg	0.78	0.78	0.78	320000
weighted avg	0.78	0.78	0.78	320000

Thank
You