

# BIG DATA PROJECT REPORT

## Machine Learning with Spark Streaming

### Sentimental Analysis Dataset BD\_036\_186\_425\_462

Aishwarya B S	PES1UG19CS036
Harshitha C	PES1UG19CS186
Samhitha Harish	PES1UG19CS425
Shreya Srinivas	PES1UG19CS462

Our dataset is the sentimental analysis. Which basically involves preprocessing the streamed data in batches and classifying using models we have used in our context. We have used decision trees and random forest classifiers as our 2 classification algorithms.

Starting from the streaming part of it we created a streamterminal.py file that takes in the argument which takes in the given batchsize and streams it according to the number given. Argparse() module helps in taking the values and parsing through. Time.sleep(2) indicates that it takes a timestamp of 2 seconds to process every batch

Preprocessing: The initial part of the data analysis starts with preprocessing. This part focuses to create basis for building our model. We have used 2 transformations here by importing tokenize and IF-TDF modules from the pyspark MLLIB features. We have localized the names "tweetcol" and "sentimentcol" for easier preprocessing. We tokenize the words then use idf for inverse column matrix which makes it easier to calculate the transpose of the words for it to be generated according to our criteria to convert it into numerical values as models work well with numerical values.

#### CLASSIFIER ALGORITHMS:

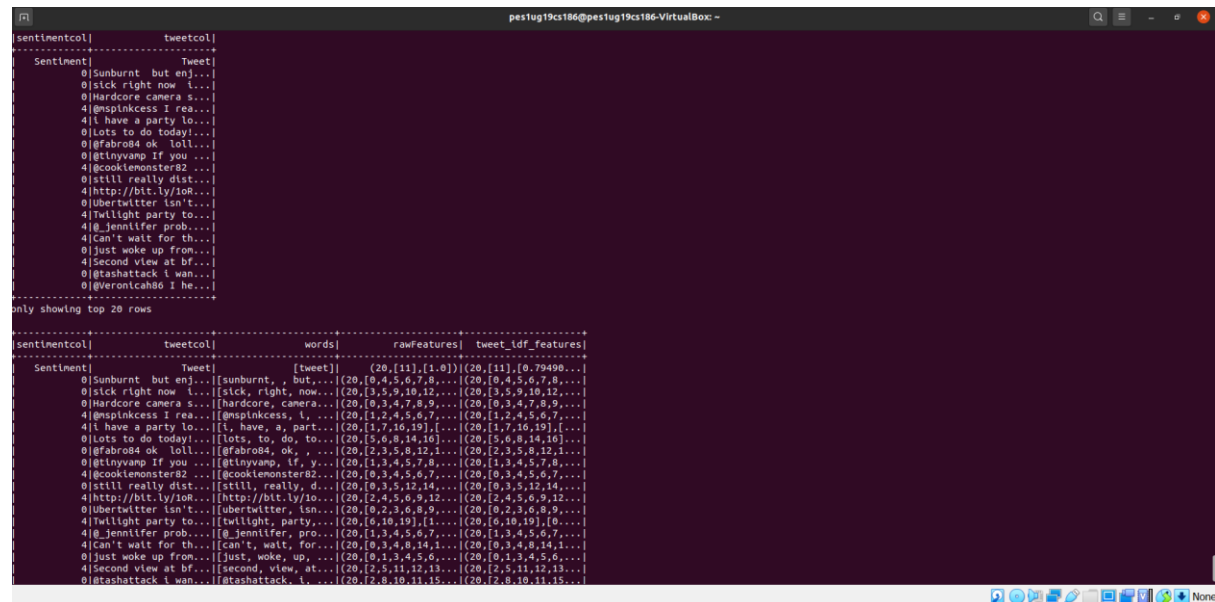
##### 1) Random Forest:

The Random forest or Random Decision Forest is a supervised Machine learning algorithm used for classification, regression, and other tasks using decision trees. The Random forest classifier creates a set of decision trees from a randomly selected subset of the training set. It is basically a set of decision trees (DT) from a randomly selected subset of the training set and then collects the votes from different decision trees to decide the final prediction. The libraries used are sklearn.metrics , accuracy score sklearn.model selection ,(train\_test\_split) sklearn.preprocessing (OneHotEncoder).

## 2) Decision Tree:

Decision tree is a type of supervised learning algorithm that can be used for both regression and classification problems. The algorithm uses training data to create rules that can be represented by a tree structure. Like any other tree representation, it has a root node, internal nodes, and leaf nodes. The internal node represents condition on attributes, the branches represent the results of the condition and the leaf node represents the class label. The libraries used are numpy, itertools combinations sklearn.metrics accuracy\_score.

### DATA BEFORE AND AFTER PREPROCESSING FOR 1 STREAMED BATCH



The screenshot shows a Jupyter Notebook interface with two code cells. The first cell displays a DataFrame with two columns: 'sentimentcol' and 'tweetcol'. The second cell displays a DataFrame with five columns: 'sentimentcol', 'tweetcol', 'words', 'rawfeatures', and 'tweet\_idf\_features'. The data is shown for the top 20 rows.

sentimentcol	tweetcol
0	Sunburnt but enj...
0	sick right now i...
0	Hardcore camera s...
4	@mspinkess I rea...
4	i have a party to...
0	Lots to do today!
0	@fabro84 ok lol!
0	@tinyvamp if you ...
4	@cookenmonster2 ...
0	still really dist...
4	http://bit.ly/1oR...
0	UberTwitter isn't...
4	Twilight party to...
4	@jennifer prob...
4	Can't wait for th...
0	Just woke up from...
4	Second view at bf...
0	@tashattack i wan...
0	@Veronica86 I he...

sentimentcol	tweetcol	words	rawfeatures	tweet_idf_features
0	Sunburnt but enj...	[Sunburnt, but, ...]	(20, [11], [1, 0])	(20, [11], [0.79490...])
0	sick right now i...	[sick, right, now, ...]	(20, [3, 5, 9, 10, 12], ...)	(20, [3, 5, 9, 10, 12], ...)
0	Hardcore camera s...	[hardcore, camera, ...]	(20, [0, 3, 4, 7, 8, 9], ...)	(20, [0, 3, 4, 7, 8, 9], ...)
4	@mspinkess I rea...	[@mspinkess, i, ...]	(20, [1, 2, 4, 5, 6, 7], ...)	(20, [1, 2, 4, 5, 6, 7], ...)
4	i have a party to...	[i, have, a, part...	(20, [1, 7, 10, 19], ...)	(20, [1, 7, 10, 19], ...)
0	Lots to do today!	[lots, to, do, to...	(20, [5, 6, 8, 14, 16], ...)	(20, [5, 6, 8, 14, 16], ...)
0	@fabro84 ok lol!	[@fabro84, ok, , ...]	(20, [2, 3, 5, 8, 12, 1...])	(20, [2, 3, 5, 8, 12, 1...])
0	@tinyvamp if you ...	[@tinyvamp, if, y...	(20, [1, 3, 4, 5, 7, 8], ...)	(20, [1, 3, 4, 5, 7, 8], ...)
4	@cookenmonster2 ...	[@cookenmonster2...	(20, [0, 3, 4, 5, 6, 7], ...)	(20, [0, 3, 4, 5, 6, 7], ...)
0	still really dist...	[still, really, d...	(20, [0, 3, 5, 12, 14], ...)	(20, [0, 3, 5, 12, 14], ...)
4	http://bit.ly/1oR...	[http://bit.ly/1o...	(20, [2, 4, 5, 6, 9, 12], ...)	(20, [2, 4, 5, 6, 9, 12], ...)
0	UberTwitter isn't...	[UberTwitter, isn...	(20, [0, 2, 3, 6, 8, 9], ...)	(20, [0, 2, 3, 6, 8, 9], ...)
4	Twilight party to...	[twilight, party...	(20, [6, 10, 19], [1...])	(20, [6, 10, 19], [0...])
4	@jennifer prob...	[@jennifer, pro...	(20, [1, 3, 4, 5, 6, 7], ...)	(20, [1, 3, 4, 5, 6, 7], ...)
4	Can't wait for th...	[can't, wait, for...	(20, [0, 3, 4, 8, 14, 1...])	(20, [0, 3, 4, 8, 14, 1...])
0	Just woke up from...	[just, woke, up, ...]	(20, [0, 1, 3, 4, 5, 6], ...)	(20, [0, 1, 3, 4, 5, 6], ...)
4	Second view at bf...	[second, view, at...	(20, [2, 5, 11, 12, 13], ...)	(20, [2, 5, 11, 12, 13], ...)
0	@tashattack i wan...	[@tashattack, i...	(20, [2, 8, 10, 11, 15], ...)	(20, [2, 8, 10, 11, 15], ...)