# Text Analysis
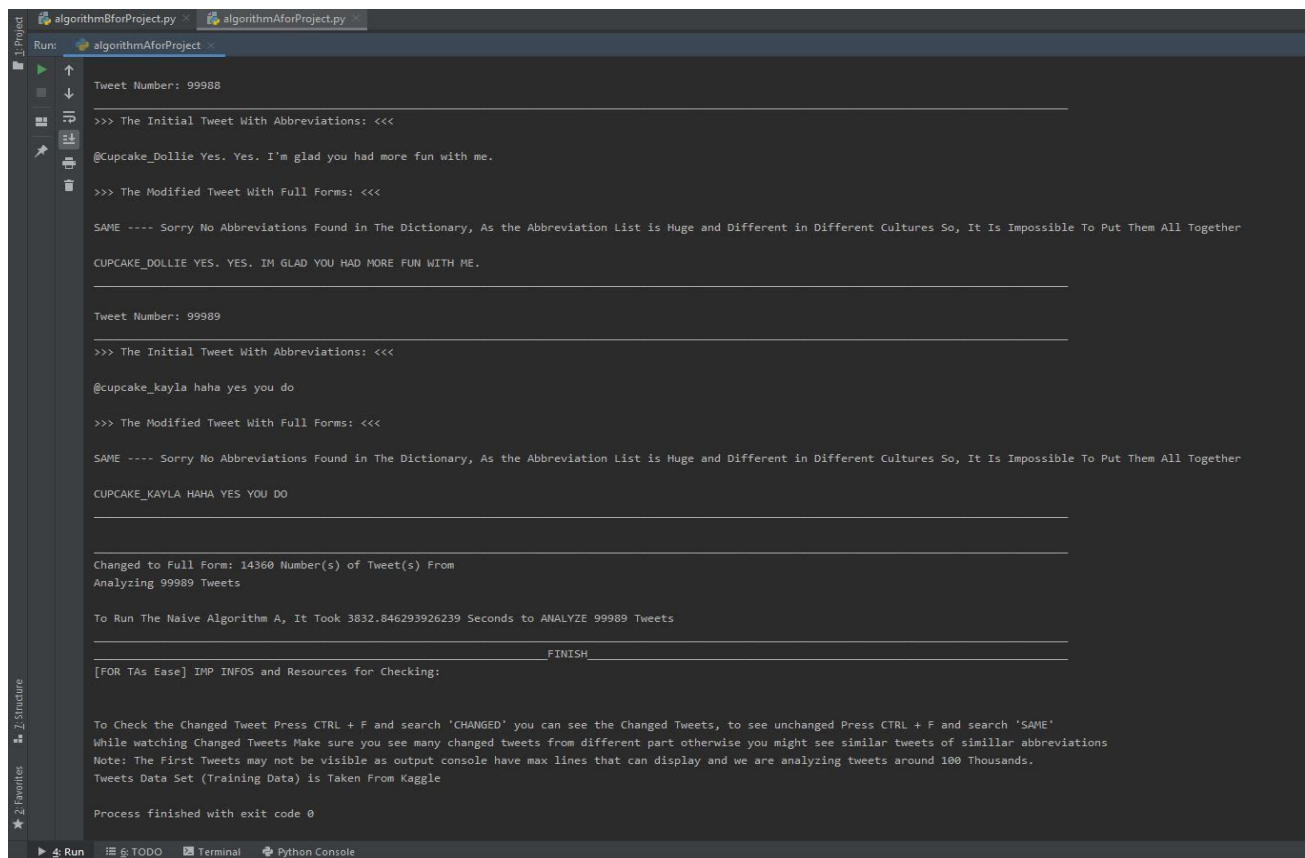# Reading Twitter Data Set and Replacing the Abbreviation (Slang words) with Full Form.
ReadMe + Instruction Guide

## File Descriptions:

1. algorithmAforProject: is the python file for the naive algorithm. It took around 1 hours to run the program, as it is the most inefficient solution (mentioned on project topic pdf) and analyzing 100 thousand tweets against 235 abbreviations list.

Summary of Output in Screenshot:



\* Screenshot is provided in the repository also

(Note: Running algorithmAforProject, the naive one took around 1 hour, so I would suggest run the Program when you will be taking meals or shower or watching movies)

2. <u>algorithmBforProject:</u> is the python file for the efficient algorithm. It took around 30 seconds to run the program, as it is the most efficient solution and it is analyzing 100 thousand tweets against 235 abbreviations hash: dictionary.

Summary of Output in Screenshot:



* Screenshot is provided in the repository also

3. <u>train.csv:</u> is the dataset for the tweets collected from kaggle. It contains around 100,000 (100 thousands) Tweets, among 100,000 tweets 14,360 Tweet's abbreviation is updated to full form, as the abbreviations and slang list is huge and it depends on many different cultures it is impossible to collect the whole list of abbreviations/slang, managed to collect 235 nos of abbreviation/slang

4. <u>finalSlang.csv:</u> combined collection of all the dataset of the abbreviations and slang list used in the python program.

5. <u>Book1.csv:</u> initial collection of the abbreviations and slang list, not used in the program, then it is updated to more collections and changed to finalSlang.csv which is used in the program.

6. <u>Venv + .idea:</u> python build files along with the libraries.

# Efficiency of The Programs:

algorithmAforProject Takes approx. 1 hours to run completely where algorithmBforProject takes 30 seconds to run completely. Both are analyzing same no. of tweets using the same no. of abbreviations. So the efficient solution algorithmBforProject is 120 times faster than the naive solution algorithmAforProject.

# Instructions to Run The Program:

1.  Make sure you have python installed with proper IDE.

2.  Clone the repository.

3.  Open the repository/python file (algorithmAforProject, algorithmBforProject) with IDE and run it

4.  It will take around 1 hour to run algorithmAforProject as I mentioned above, so suggestion would be to run algorithmBforProject before.

5.  See the last lines of the console after running the program to see some summary of the programs.