# CS410 : US Midterm Election – Twitter Sentiment & Topic Analysis

Unnati Hasija (uhasija2@illinois.edu)

Prateek Dhiman (pdhiman2@illinois.edu)

# Contents



1) **Project Background**
2) **Implementation Details**
   a. Requirements
   b. Setup & Usage Instructions with Results and Visualizations
   c. Errors and Solution
3) **Conclusion and Future Work**

# Project Background

**Background**

As discussed in our initial project proposal we took **Twitter** and performed the following for **US Midterm Election** as a key area.

- **Sentiment Analysis** and prediction
    - With **Naïve Bayes** Classification
    - With **K-Nearest Neighbor** Classification
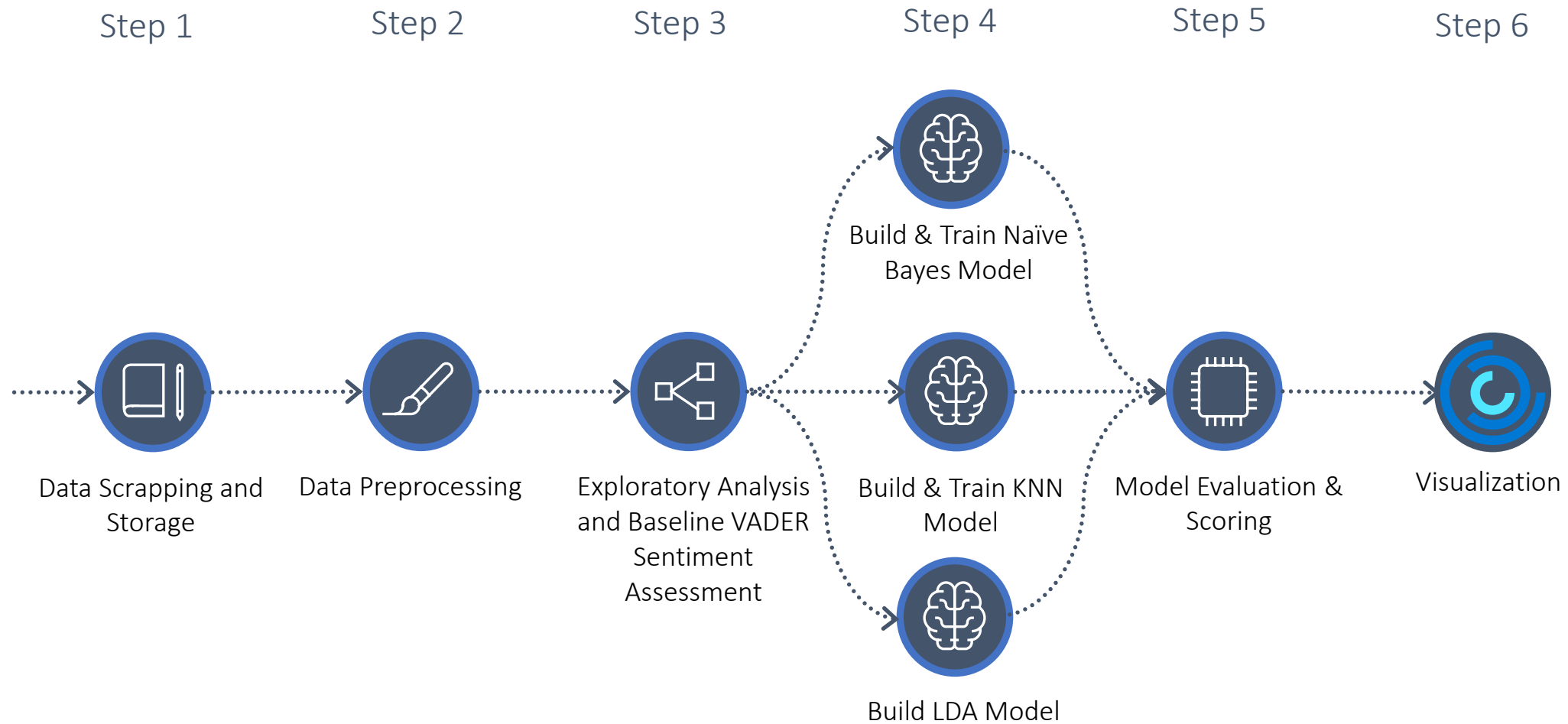- **Topic Analysis** with **Latent Dirichlet Allocation**

**Goal**

Key Goals for the project that we set and achieved were

- **Sentiment Analysis** : Identify overall positive , negative and neutral sentiments towards elections and respective parties in US
- **Prediction** : Predict tweets for sentiments using Naïve Bayes and KNN. Also perform a comparison which Algorithm performs better and impact of data cleaning techniques on precision and accuracy of predictions.
- **Topic Analysis** : Identify top topics in tweets and predict the dominant topic and its percentage for each tweet.
- **Visualize :** Lastly compare our findings and present them as part of our final project.

# Implementation Details

Step 1      Step 2      Step 3      Step 4      Step 5      Step 6



Build & Train Naïve Bayes Model

Data Scrapping and Storage

Data Preprocessing

Exploratory Analysis and Baseline VADER Sentiment Assessment

Build & Train KNN Model

Model Evaluation & Scoring

Visualization

Build LDA Model

# Implementation Details – Requirements

- Requirements
  - Local Jupyter Server / VSCODE for running Jupyter notebooks
  - Python 3.9+
  - git
  - Github Account
  - Snscrape
  - Numpy
  - Pandas
  - NLTK
  - Matplotlib
  - Wordcloud
  - Emot
  - Gensim
  - pyLDAVis

- Setup Instructions
  - Latest Python version and VScode can be installed via Anaconda distribution https://www.anaconda.com/
  - Details of PIP modules are available in in requirement.txt in the project repository. They will be installed as part of setup
  - If you are unable use the requirement.txt file , please use the PIP commands as mentioned in README.MD
  - Follow Screenshots from next slides to setup the project locally.

# Setup and Usage Instructions

1. Install Anaconda with Python 3.9 from
   https://www.anaconda.com

2. Ensure you also install VSCODE via Anaconda Navigator

3. In VSCODE , please ensure you install the python extension

# Setup and Usage Instructions

3. Install Git in your local machine if it does not exist from [Git (git-scm.com)](https://git-scm.com)

4. Perform the initial git setup if not done in past

*$ git config --global user.name "<your github username> "*

*$ git config --global user.email <your github email>*

# Setup and Usage Instructions

5. Create a folder in your machine where you would clone the project files.

6. Open git bash and run the following

*git clone https://github.com/pd-illinois/CS410Googolplex.git*

# Setup and Usage Instructions

7. Open the folder in VSCODE and go to source folder.

8. Open Terminal in VSCODE and run

pip install -r requirements.txt

Wait for the installation to complete. If installation fails due to genism please see section – Expected Errors and Solution

# Setup and Usage Instructions

9. Ensure you have the python 3.9 selected as your main python interpreter .

10. Now Open **Twitter_Sentiment_Naive.ipynb** and Run the Jupyter Notebook

11. Validate the successful run and results.

# Setup and Usage Instructions

12. Key Naïve Bayes Results

Naïve Bayes



```
          naïvebayesCV()

[5]

...   Accuracy Score for ngrams = 1 is: 75.0%
      Recall for ngrams = 1 is: 75.0%
      Precision for ngrams = 1 is: 75.0%
      F1-score for ngrams = 1 is: 75.0%
      Accuracy Score for ngrams = 2 is: 51.0%
      Recall for ngrams = 2 is: 50.0%
      Precision for ngrams = 2 is: 55.00000000000001%
      F1-score for ngrams = 2 is: 50.0%
      Accuracy Score for ngrams = 3 is: 34.0%
      Recall for ngrams = 3 is: 32.0%
      Precision for ngrams = 3 is: 46.0%
      F1-score for ngrams = 3 is: 28.000000000000004%
      Accuracy Score for ngrams = 5 is: 25.0%
      Recall for ngrams = 5 is: 23.0%
      Precision for ngrams = 5 is: 48.0%
      F1-score for ngrams = 5 is: 21.0%
```

Naïve Bayes with TD-IDF
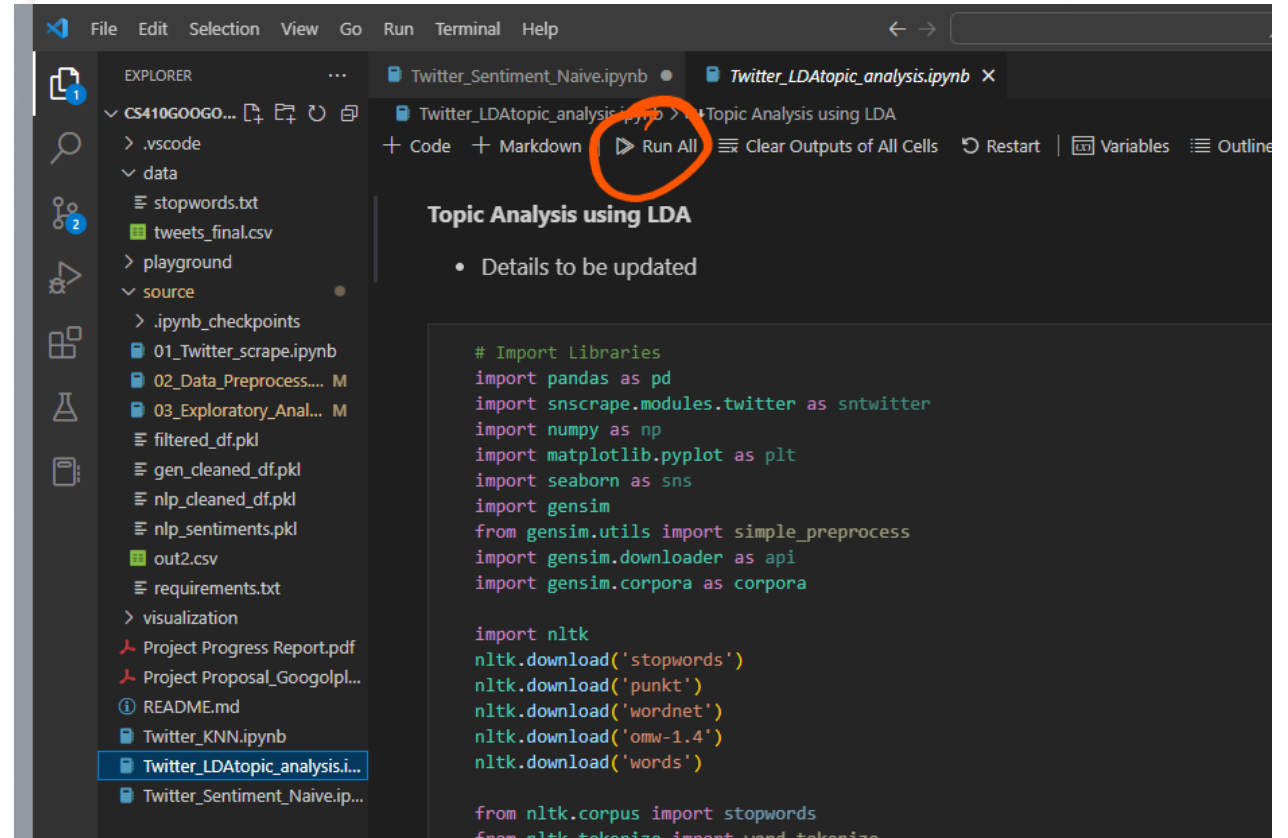
# Setup and Usage Instructions

13. Open **Twitter_KNN.ipynb** and Run the Jupyter Notebook

11. Validate the successful run and results.

# Setup and Usage Instructions

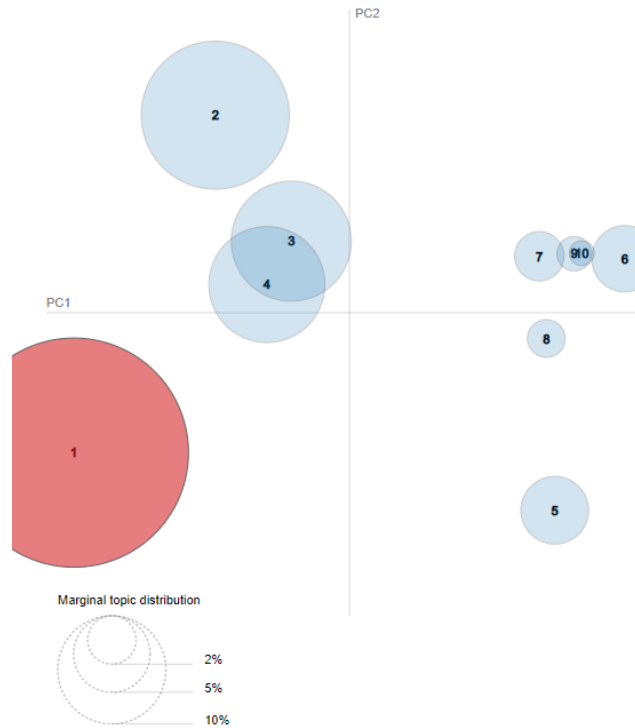## 12. Key KNN Results

```
Confusion Matrix for k = 1 is:

[[ 174  697   81]
 [  12 1238   28]
 [  74  764  669]]

Classification Report for k = 1 is:

              precision    recall  f1-score   support

    negative       0.67      0.18      0.29       952
     neutral       0.46      0.97      0.62      1278
    positive       0.86      0.44      0.59      1507

    accuracy                           0.56      3737
   macro avg       0.66      0.53      0.50      3737
weighted avg       0.67      0.56      0.52      3737

Accuracy Score for k = 1 is: 56.00000000000001%
Macroaveraged Recall for k = 1 is: 53.0%
Macroaveraged Precision for k = 1 is: 66.0%
Macroaveraged F1-score for k = 1 is: 50.0%
...
```

KNN

# Setup and Usage Instructions

13. Open **Twitter_LDAtopic_Analysis.ipynb** and Run the Jupyter Notebook

11. Validate the successful run and results.

# Setup and Usage Instructions

## 12. Key LDA Results



```
from gensim.models import CoherenceModel
# Compute Coherence Score
coherence_model_lda = CoherenceModel(model=lda_model, texts=data_words, dictionary=id2word, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('\nCoherence Score: ', coherence_lda)


Coherence Score:  0.41753944190016695
```

# Expected Errors and Solutions

Error : Gensim fails due to missing C++ Build

- Solution Install Microsoft Visual C++ Build Tools.

# Expected Errors and Solutions

https://visualstudio.microsoft.com/visual-cpp-build-tools/

# Conclusion and future work

❑ - Accuracy improves if the data has been cleansed properly. We removed the words that are not a part of nltk.words. This helped in improving the accuracy for Naïve Bayes from 65% to 75%.

❑ - For KNN, difficult to build a feature vector of huge number of tweets.

❑ - For sentiment analysis on tweets data, we found Naïve Bayes to show better accuracy over K-NN.

❑ - It was more difficult to train K-NN as compared with Naïve Bayes. Hence, we conclude to use Naïve Bayes for sentiment analysis on Tweets data.

❑ - For future work, the same model could be evaluated for different subjects other than US midterm elections by simply scraping the data on another topic from Twitter using the code provided and plugging that data for Topic Modeling and Text Categorization and could be evaluated for accuracy.

❑ - Future work could also include modeling other text categorization techniques like SVM, Deep Learning (LSTM), etc. with word embedding techniques such as Word2vec could also be applied and compared for accuracy.

Thank you