

# Group 13

## Authorship attribution based on Stylometry

### *Minutes of Meeting (MoM)*

<b>Week 1</b> February 8 February 11	<ol style="list-style-type: none"><li>1. Discussed various topics which can be pitched.</li><li>2. Ideas discussed:<ol style="list-style-type: none"><li>a. Sarcasm detection (Vivek)</li><li>b. Authorship Attribution based on Stylometry (Sarvani)</li><li>c. Coronavirus sentiment detection (Lyn)</li></ol></li></ol>
<b>Week 2</b> February 17	<ol style="list-style-type: none"><li>1. More ideas discussed<ol style="list-style-type: none"><li>a. An idea around creating a Chatbot (Piyush)</li><li>b. Analysis of music lyrics to make new recommendations to users (Nikhil)</li><li>c. Google QnA labelling on Kaggle (Lyn)</li><li>d. On input, find similar questions like on QnA websites (StackOverflow)</li></ol></li></ol>
<b>Week 3</b> February 25	<ol style="list-style-type: none"><li>1. Finalized the topic for the essay after discussing with the professor (Feb 21)</li><li>2. Formed research questions</li><li>3. Discussed various possible datasets:<ul style="list-style-type: none"><li>• Kaggle</li><li>• Pan Competitions</li><li>• Project Gutenberg</li></ul></li><li>4. Found Parts of Speech features for the defined text corpus.</li></ol>
<b>Week 3</b> March 4 March 8	<ol style="list-style-type: none"><li>1. Worked on literature review and writing various parts of the essay</li><li>2. Kaggle dataset is finalized.</li><li>3. Preliminary data analysis and feature extraction.</li><li>4. Discussion of possible features and extraction tools.</li><li>5. Intermediate results are presented.</li></ol>
<b>Week 4</b> March 11	<ol style="list-style-type: none"><li>1. Discussion of Peer Reviews received.</li><li>2. Minor changes to the essay to address smaller feedbacks.</li><li>3. Piyush: "<i>We are gonna get the reviews back. How about we go authorship attribution on that</i>". Good idea. :-)</li></ol>
<b>Week 5</b> March 18	<ol style="list-style-type: none"><li>1. Created a list of all the changes to be made in the final group document based on the reviews received</li><li>2. Looked into several approaches for extracting Text-based and meta-features using different libraries and tools.</li><li>3. Discussed valuable inputs received from discussion with other teams.</li></ol>

<b>Week 6</b> March 25	<ol style="list-style-type: none"> <li>1. Meta Feature-based approach implementation</li> <li>2. Decided to use multiple classification models: Naive Bayes, Light GBM and XGBoost.</li> <li>3. List of important features is selected among all.</li> </ol>
<b>Week 7</b> April 5	<ol style="list-style-type: none"> <li>1. Content-based features extraction (finally used TF-IDF)</li> <li>2. Getting feature importance from different classifiers.</li> <li>3. Metadata and content-based approaches are implemented parallely to get accuracy.</li> <li>4. Decided to use at least two of the classification models from Meta-Features based approach</li> <li>5. Agreed on Incorporating Ridge Classifier, Decision Tree and Logistic Regression also to Meta-feature based approach.</li> </ol>
<b>Week 8</b> April 13	<ol style="list-style-type: none"> <li>1. Wrote the final essay</li> <li>2. Formatted the essay in Latex</li> <li>3. Generated visualizations for easy understanding</li> <li>4. Got reviews on our final work from two other groups (Group 12 and Group 2), made changes to address the problems.</li> </ol>