



Connecting Suicidal and Help-Seeking Behaviors on Social Media

Kushal Gevaria kgg5247@rit.edu

Advisor: Prof. Christopher Homan cmh@cs.rit.edu



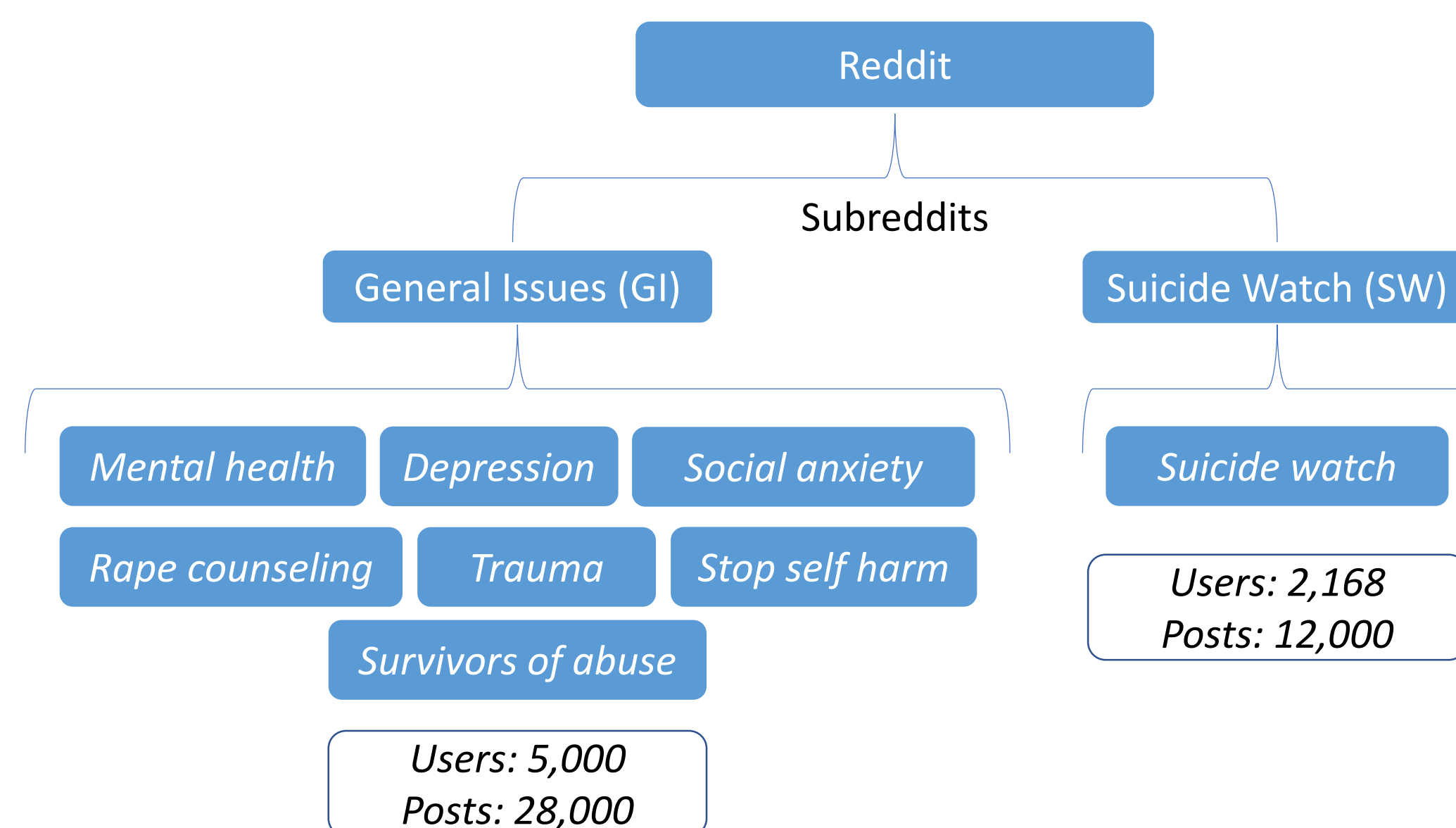
INTRODUCTION

- Social media contributes an enormous amount of user-generated content on the topics of mental illness and suicidal
- Analyzing these data helps to identify the individuals showing mental concerns that lead to the thought of committing suicide
- **Goal:** To differentiate between individuals who have suicidal thoughts due to depression, anxiety or mental health issues and individuals who just exhibit general issues but don't have suicidal thoughts on social media

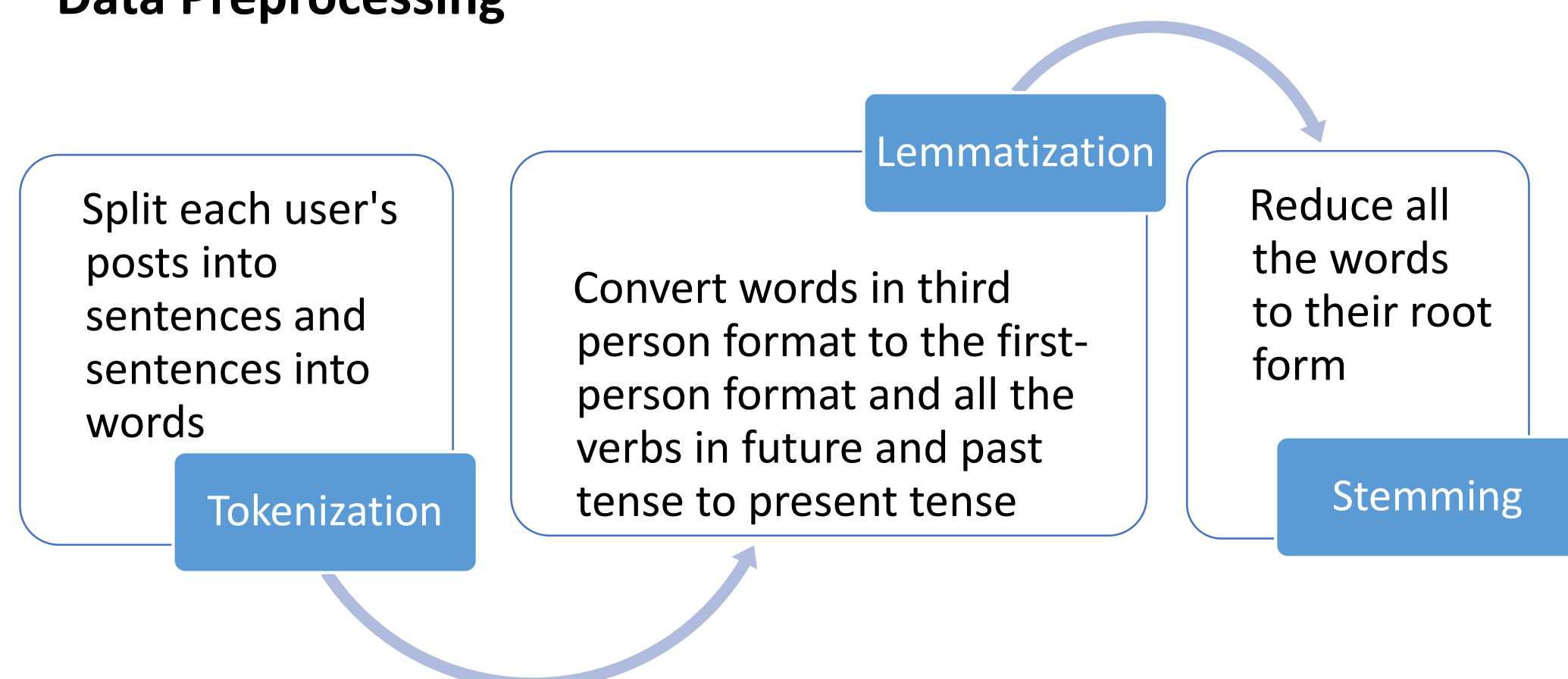
DATA

Data Selection & Extraction

- We used the *Pushshift* API to collect user information and their posts for year 2015



Data Preprocessing



MODEL

Latent Dirichlet Allocation

- A generative probabilistic model which is used to classify the collection of composites based on the parts used in it to the number of topics considered

	Topic 0 GI	Topic 1 SW
improve	0.271	0.004
kill	0.004	0.565
life	0.360	0.004
suicide	0.004	0.425

Table 1: Probability of choosing a word given the number of topics

	Topic 0 GI	Topic 1 SW
[stress, depress, ..., suicide]	0.065	0.935
[severe, pain, ...,improve]	0.924	0.076
[kill, suicide, ..., pain, quit]	0.246	0.754
[sick, never, ..., bed, sad]	0.345	0.355
[worry, improve, ..., life]	0.924	0.076

Table 2: Probability of choosing a sentence given the number of topics

Support Vector Machines

- A supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis
- Input feature: Matrix representation of documents v/s words
 - 20,000 posts from the training dataset
 - Each document in the matrix have a vector representation of the 4,211 words (features) with binary values
- Feature Reduction:
 - Extracted top 200 words as the feature list instead of all the 4211 words from the two topics
 - Re-trained the SVM model using reduced features

RESULTS

We compute the classification accuracy and results were evaluated by cross-validating whether GI users classified as SW subscribed to “**SuicideWatch**” subreddit.

	Actual GI 4200	Actual SW 1800
Predicted GI 4118	3122	996
Predicted SW 1182	1078	804

Table. 1: Confusion matrix for **LDA** with an accuracy of **65.04%**

	Actual GI 4200	Actual SW 1800
Predicted GI 3937	3345	592
Predicted SW 2063	855	1208

Table. 2: Confusion matrix for **support vector machine** algorithm with an accuracy of **75.88%**

	Actual GI 4200	Actual SW 1800
Predicted GI 4088	3563	525
Predicted SW 1912	637	1275

Table. 3: Confusion matrix for **support vector machine** algorithm with an accuracy of **80.63%** after **feature reduction**

CONCLUSION

- Successfully identified users exhibiting suicidal ideation
- It does not implicitly mean that the users classified to have suicidal thoughts are expressing the same
- Learned the use of topic modelling as a feature for the classification algorithm