# Using Natural Language Processing to Classify Suicide Notes

**John P. Pestian\*, Pawel Matykiewicz, Jacqueline Grupp-Phelan,**
Sarah Arszman Lavanier, Jennifer Combs, and Robert Kowatch
Cincinnati Children's Hospital Medical Center
Cincinnati, OH 45220, USA
john.pestian@cchmc.org

## Abstract

We hypothesize that machine-learning algorithms (MLA) can classify completer and simulated suicide notes as well as mental health professionals (MHP). Five MHPs classified 66 simulated or completer notes; MLAs were used for the same task. Results: MHPs were accurate 71% of the time; using the sequential minimization optimization algorithm (SMO) MLAs were accurate 78% of the time. There was no significant difference between the MLA and MPH classifiers. This is an important first step in developing an evidence based suicide predictor for emergency department use.

## 1 Problem

Suicide is the third leading cause of death in adolescents and a leading cause of death in the United States[1]. Those who attempt suicide usually arrive at the Emergency Department seeking help. These individuals are at risk for a repeated attempt, that may lead to a completed suicide[2]. We know of no evidence-based risk assessment tool for predicting repeated suicide attempts. Thus, Emergency Medicine clinicians are often left to manage suicidal patients by clinical judgment alone. This research focuses on the initial stage for constructing such an evidence based tool, the Psychache[3] Index. Our efforts herein posit that suicide notes are an artifact of a victim's thoughts and that the thoughts between completers and attempters are different. Using natural language processing we attempt to distinguish between completer notes and notes that have been simulated by individuals who match the profile of the completer. Understanding how to optimize classification methods between these types of notes prepares us for future work that can include clinical and biological factors.

## 2 Methods

Suicidal patients are classified into three categories: ideators —those who think about committing suicide, attempters —those who attempt suicide, and completers —those who complete suicide. This research focuses on the completers and a group of individuals called *simulators*. These simulators were matched to completers by age, gender and socioeconomic status and asked to write a suicide note[4]. Suicide notes from 33 completers and 33 simulators were annotated with linguistic characteristics using a perl-program with the EN:Lingua:Tagger module. Emotional characteristics were annotated by assigning terms in the note to a suicide-emotion ontology that was developed from a meta analysis of 2,166 suicide related manuscripts and validated with expert opinion. This ontology includes such classes as: affection, anger, depression, and worthlessness. Each class had multiple concepts, i.e, affection→ love, concern for others, and gratitude. Three MHPs read each note and tagged emotion-words found in the notes with the appropriate classes and concepts. Analysis of variance between structures was conducted to insure that there actually was a difference that could be detected. Emotional annotations were used for machine-learning.

We then tested the hypothesis that MLAs could distinguish between completer and simulated notes as well as MHPs. Copies of the notes were given to five MHPs who classified them as either written by a completer or an simulator. MLA feature space was defined by matrix of selected characteristics from four sources: words, parts of speech, concepts, and readability indexes. Collinearity was eliminated by removing highly correlated features. The final feature space included: specific words (such as "love", "life", "no"), specific parts of speech (such as, personal pronouns, verbs) Kincaid readability index and emotional concepts (such as anger, and hopelessness). We then tested the following algorithms' ability to distinguish between completer and simulator notes: *decision trees* - J48, C4.5, LMT, DecisionStump, M5P; *classification rules* - JRip, M5, OneR, PART; *function models* - SMO, logistic builds, multinomial logistic regression, linear regression; *lazy learners* and *meta learners*[5].

## 3 Results

A significant difference was found between the linguistic and emotional characteristics of the notes. Linguistic differences (completer/simulated): word count 120/66 p=0.007, verbs 25/13 p=0.012, nouns 28/12 p=0.0001, and prepositions 20/10 p=0.005. This difference justified testing the classification hypothesis. Emotionally, completers gave away their possessions 20% of the time, simulators, never did. Mental health experts accurately classified the notes 71% of the time. The MLAs were accurate 60-79% of the time with SMO giving the highest results when the word count, part-of-speech, and readability vectors were included. Performance weakened when the emotional vector was included, yet the emotional vector was the primary source of data for the MHPs.

## 4 Conclusion

Machine learning methods for classifying suicide and non-suicide notes are promising. Future efforts to represent the thoughts of the suicidal patient will require larger sample sizes, inclusion of attempters response to open-ended questions, biological and clinical characteristics.

## 5 Acknowledgements

**References:**
[1] Jeffrey A Bridge, Tina R Goldstein, and David A Brent. Adolescent suicide and suicidal behavior. *J Child Psychol Psychiatry*, 47(3-4):372–394, 2006.
[2] P M Lewinsohn, P Rohde, and J R Seeley. Psychosocial risk factors for future adolescent suicide attempts. *J Consult Clin Psychol*, 62(2):297–305, 1994.
[3] E S Shneidman. Suicide as psychache. *J Nerv Ment Dis*, 181(3):145–147, 1993.
[4] ES Shneidman and NL Farberow. *Clues to Suicide*. McGraw Hill Paperbacks, 1957.
[5] I.H. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools ad Techniques*. Morgan Kaufman, 2nd edition, 2005.