

CSE713

Pattern Recognition

Course Instructor:

Annajiat Alim Rasel

Research Assistant:

Humaion Kabir Mehedi

Student Tutor:

Abid Hossain

Team-19:

23366009 Ripa Sarkar



Paper Title

“Using Natural Language Processing to Classify Serious Illness
Communication with Oncology Patients”

by Anahita Davoudi, Hegler Tissot, Abigail Doucette, Peter E. Gabriel, Ravi Parikh
Danielle L. Mowery & Stephen P. Miranda.



Table of Contents:

- Introduction
- Objectives
- What is Serious Illness Communication (SIC)?
- Challenges in SIC Documentation
- Natural Language Processing (NLP) for SIC
- Methods
- Dataset and Schema
- Serious Illness Communication Classifier Development and Evaluation
- Serious Illness Communication Subdomain Characterization
- Algorithms
- Results
- Limitations
- Future Work
- Discussion
- Conclusion
- Bibliography



Introduction

- Importance of serious illness communication (SIC) in patient-centered decision-making.
- Challenges in measuring and evaluating the quality of SIC.
- Potential of Natural Language Processing (NLP) for identifying and evaluating SIC documentation.



Objectives

- Develop NLP algorithms to identify and characterize SIC with oncology patients.
- Demonstrate the potential for using NLP-based metrics in oncology and other serious illness care settings.



What is Serious Illness Communication (SIC)?

- Importance of SIC for patient-centered decision-making.
- Impact on quality of life and goal-concordant care.
- SIC aims to address end-of-life care, prognosis, treatment options, and patient preferences.
- Inadequate SIC associated with psychosocial distress and incongruent end-of-life care.
- Need for evaluating SIC documentation as a core quality measure.



Challenges in SIC Documentation

- Underutilization and inconsistency of traditional forms of SIC documentation.
- Difficulty in tracking SIC across inpatient and outpatient settings.
- Limited ability of keyword-based approaches to capture nuanced documentation about patient priorities and prognostic communication.



Natural Language Processing (NLP) for SIC

- Overview of NLP as an efficient and accurate alternative for identifying SIC in electronic health records (EHR).
- Current approaches relying on keyword-based algorithms.
- Introduction of machine learning approaches for more accurate and automatic identification.
- Importance of expanding beyond keywords to capture critical SIC domains.
- Need for more sophisticated NLP approaches to capture nuanced SIC documentation.



Methods

- Collection of a weakly annotated dataset of free-text entries containing SIC documentation.
- Training of machine learning algorithms (Logistic Regression, XGBoost, BERT, Bio+Clinical BERT) to classify SIC documentation by domain and subdomain.
- Classification of SIC documentation by domain and subdomain.
- Characterization of features associated with each SIC subdomain.



Dataset and Schema

- Description of the dataset from the University of Pennsylvania Abramson Cancer Center.
- Structure of the "Serious Illness Conversation" note template for documenting SIC.
- Random split of the dataset for training and testing.

Prognosis Domains			
Subdomain	Prompt	Responses	Comment
Prognostic Understanding (PU)	What is your understanding now of where you are with your illness?	<i>Overestimates prognosis;</i> <i>Accurate understanding of prognosis;</i> <i>Underestimates prognosis;</i> <i>No understanding of prognosis;</i>	"He knows he only has weeks to live."
Information Preferences (IP)	How much information about what is likely to be ahead with your illness would you like from me?	<i>Patient wants to be fully informed;</i> <i>Patient wants to be informed of big picture, but not details;</i> <i>Patient wants some information, but no "bad news";</i> <i>Patient prefers information to be shared with ***</i>	"She prefers weekly prognosis updates."
Prognostic Communication (PC)	Information shared with patient about prognosis	<i>Uncertain prognosis;</i> <i>Possibility of getting sick quickly;</i> <i>Limited time, may be as short as</i> <i>May never get stronger or regain function</i>	"He had questions about prognosis."

Goal Domains			
Subdomain	Prompt	Responses	Comment
Main Goals (MG)	If your health situation worsens, what are your most important goals?	<i>Live as long as possible;</i> <i>Pursue every available treatment;</i> <i>Avoid hospitalizations/maximize time at home;</i> <i>Not be a burden/maintain independence;</i> <i>Be physically comfortable;</i> <i>Be mentally aware;</i> <i>Spent time with family</i>	"The patient wants to live to see his daughter's wedding."
Fears/Worries (FW)	What are your biggest fears and worries about the future with your health?	<i>Pain or other symptoms;</i> <i>Loss of control or dignity;</i> <i>Burdening others;</i> <i>Family concerns;</i> <i>Financial concerns</i>	"He worries about becoming dependent."
Strengths (ST)	What gives you strength as you think about the future with your illness?	<i>Friends/family;</i> <i>Faith/spirituality;</i> <i>Prior experience with adversity</i>	"Support of family and friends."



Serious Illness Communication Classifier Development and Evaluation

- Performance of machine learning algorithms on the test set (F1-score, precision, recall).
- Comparison of algorithms (Logistic Regression, XGBoost, BERT, Bio+Clinical BERT).
- Importance of selecting the appropriate algorithm for different SIC domains.

SIC classifier performance by SIC domain on the test set.

Prognosis	Recall	Precision	F1-score
Logistic Regression (baseline)	0.81	0.85	0.83
XGBoost	0.85	0.86	0.86
BERT	0.80	0.64	0.71
Bio+Clinical BERT	0.86	0.80	0.83
Goals	Recall	Precision	F1-score
Logistic Regression (baseline)	0.91	0.89	0.90
XGBoost	0.92	0.91	0.91
BERT	0.80	0.90	0.84
Bio+Clinical BERT	0.88	0.92	0.90

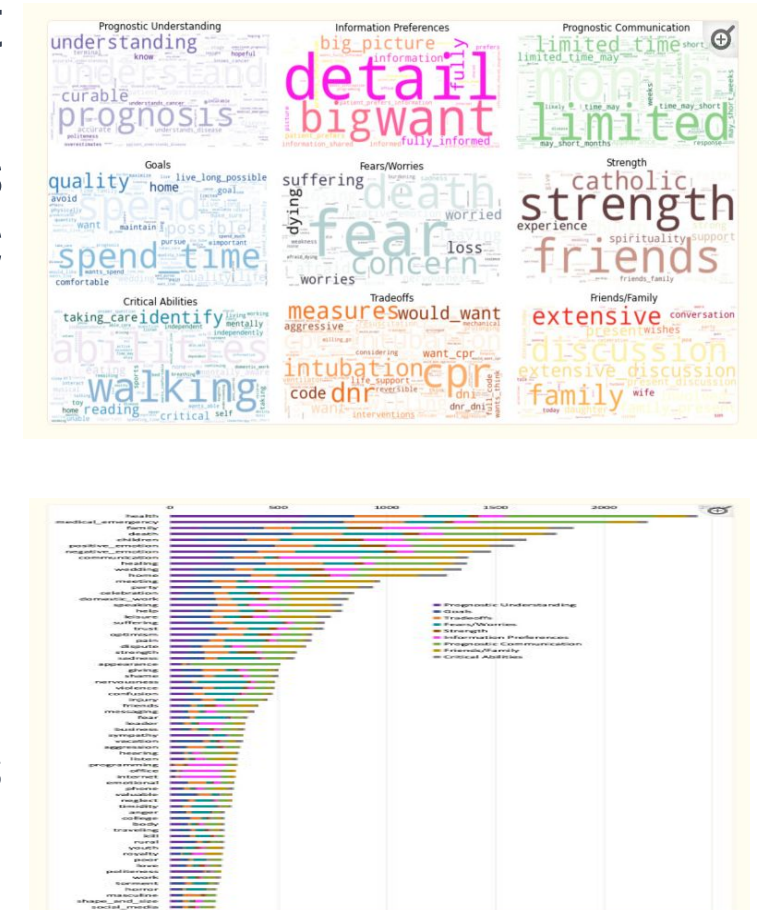
Logistic Regression SIC classifier performance by SIC subdomain on the test set.

Prognosis	Recall	Precision	F1-score
Prognostic Understanding	0.58	0.64	0.61
Information Preferences	0.44	0.42	0.43
Prognostic Communication	0.57	0.63	0.60
Goals	Recall	Precision	F1-score
Main Goals	0.52	0.68	0.59
Fears/Worries	0.62	0.40	0.49
Strengths	0.75	0.58	0.65
Critical Abilities	0.70	0.71	0.71
Tradeoffs	0.60	0.65	0.63
Friends/Family	0.47	0.27	0.35



Serious Illness Communication Subdomain Characterization

- Identification of the most informative n-grams and Empath categories associated with each SIC subdomain.
- Visualization of associated features using WordCloud.
- Distribution of Empath categories across subdomains.



Algorithms

- Logistic Regression: Learns a logit regression model to explain the relationship between features and classes.
- XGBoost: Gradient descent algorithm that predicts residual errors and minimizes loss for class prediction.
- BERT: Pretrained deep bidirectional representations fine-tuned using a masked language model.
- Bio+Clinical BERT: BERT model initialized from BioBERT and fine-tuned using clinical notes.



Results

- Overview of the study's findings on the identification and classification of SIC documentation.
- Distribution of comments by subdomain.
- Predictive performance of machine learning algorithms on test set.
- Comparison of F1-scores, precision, and recall for prognosis and goals domains.
- Performance of the SIC classifier and identification of informative features.



Limitations

- Semi-structured Epic EHR modules.
- Replicability in other settings.
- Free-text clinical notes.
- Discrimination between relevant and irrelevant text.
- Performance in population-level datasets.
- Limited number of clinicians at one institution.
- Generalizability.
- Patient preferences evolution.
- SIC across gender, race, ethnicity, and culture.



Future Work

- Enhance discrimination within each SIC domain to improve classifier performance in free-text clinical notes.
- Test and validate the algorithm's ability to identify and classify SIC in undifferentiated clinical notes.
- Further explore the performance of algorithms on longer free-text entries.



Discussion

- Importance of accurate and scalable identification of SIC in the EHR.
- Evaluation of the developed NLP algorithm's performance.
- Potential for further improvement and exploration in classifying SIC in free-text clinical notes.



Conclusion

- Successful development of an NLP algorithm for classifying SIC documentation.
- Potential for utilizing NLP-based metrics in measuring and improving the quality of oncology care.
- Future work to enhance algorithm performance, expand applications, and integrate into clinical practice.
- Future directions and implications for research in serious illness communication.



Bibliography

A. Davoudi, H. Tissot, A. Doucette, P. E. Gabriel, R. Parikh, D. L. Mowery & S. P. Miranda. (2022). Using Natural Language Processing to Classify Serious Illness Communication with Oncology Patients. AMIA Jt Summits Transl Sci Proc. 2022; 2022: 168–177, PMCID: PMC9285137, PMID: 35854756.

