Comparing emotion feature extraction approaches for predicting depression and anxiety

Hannah A. Burkhardt¹, Michael D. Pullmann¹, Thomas D. Hull², Patricia A. Areán¹, Trevor Cohen¹

¹University of Washington, Seattle ²Talkspace, New York, NY

Background

- Counseling is informed by patient symptoms
- Emotional state changes are symptoms; differ by emotion (e.g. pride), not just sentiment (pos./neg.)
- Patient-generated text from digital mental health services can be used to develop automatic assessments via measuring emotion
- There are well-established tools, but they have limitations

GoEmotions captures clinically relevant nuance

Mixed-effects linear regressions showed significant associations and high R2s. Findings reflect current understanding:

- lower reactivity in depression \rightarrow less **excitement**
- **grief** often co-occurs with depression
- lower self-image in depression \rightarrow less **pride**

UW Medicine

UW SCHOOL OF MEDICINE

BIOMEDICAL INFORMATICS AND MEDICAL EDUCATION

Feature importance differs between depression and anxiety

- Random Forest models trained w/ 80-20 data split (by patient), binary targets (PHQ-9≥10: major depressive disorder, MDD; GAD-7≥10: general anxiety disorder, GAD), using all features
- SHAP feature importance: fear ranked higher for predicting GAD than MDD. Sadness higher for predicting MDD than GAD.

Can NN-based emotion extraction methods (GoEmo) outperform established word counting methods (LIWC)?

Linguistic Inquiry and Word Count (LIWC)

- Established word-counting tool used for mental-health prediction tasks (Pennebaker et al., 2007)
- Counts words belonging to pre-defined categories
- Categories with known relationship to anxiety/depression: first-person singular pronouns ("I"), first-person plural pronouns ("we"), bio, health, sadness, anxiety, anger, pos. and neg. emotion. (Tausczik and Pennebaker, 2010)

GoEmotions

• BERT-based classifier pipeline trained on GoEmotions (annotated Reddit

• feeling of not fitting in, isolation in depression & anxiety \rightarrow perception of violations of social norms aka self-

disgust

• heightened self-criticism, interpersonal sensitivity in anxiety \rightarrow increased **disapproval** / decreased **approval**

Fi	ne-g	raine	d GoE	Emotic	ons (C	owen)	GAD	MDD
sadness			•						
admiration		=							
annoyance			=						
disappointment			•						
јоу		=							
pride									
excitement		-							
disapproval			-						
approval		=							
confusion			8						
relief		#							
neutral		(
anger			\$						
disgust			*						
optimism		=							
realization									
amusement		=							
fear			••						
nervousness			•						
caring		8							
gratitude		-	-						
embarrassment			\$						

		MDD	GAD	
	1 -	GE disappointment	GE negemo	PHQ rank
	2 -	GEE sadness	GEE negemo	
	3 –	GEnegemo	GEE posemo	
	4 –	GEE posemo	GEE joy	40
	5 –	GEE joy	GE posemo	
	6 –	LIWCi	GEE sadness	
	7 -	GEE negemo	GE sadness	30
	8 -	LIWC we	GE pride	
	9_		GE fear	
	10-	GF admiration	GE admiration	20
	11 -	GE excitement		
	12_	LIWC health	LIWC health	10
	12	GE posemo		10
	14	GE sadness	GE realization	
	14	CE prido	GE pervousposs	
	10-			
	10-			
	17-		GEE lear	
	18-	GE disapproval	GE JOY	
	19-	GE realization	GE disappointment	
	20-	GE amusement		
	21-	GE confusion	GEE anger	
	22-	GEjoy	LIWC i	
	23-	LIWC posemo	GE excitement	
×	24-	GE approval	GE disgust	
lan	25–	GE relief	LIWC anx	
	26-	GEE disgust	GEE disgust	
	27-	GE disgust	GE grief	
	28-	GE surprise	GE neutral	
	29-	GE desire	GE amusement	
	30-	GE optimism	GE relief	
	31-	GElove	GE surprise	
	32-	GE grief	GE confusion	
	33-	GE embarrassment	GE curiosity	
	34-	GE curiosity	GE embarrassment	
	35-	GE caring	GElove	
	36-	GE remorse	GE desire	
	37-	GE fear	GE gratitude	
	38-	GE anger	GE caring	
	39-	GE gratitu de	GEE neutral	
	40-	GE annoyance	GEE surprise	
	41-	GE nervousness	GE disapproval	
	42-	LIWC negemo	GE optimism	
	43-	GEE fear	GE approval	
	44-	GEE surprise	GE anger	
	45-	GE neutremo	GE annovance	
	46-	GEE neutral	GE neutremo	
	Δ7_	GE neutral	GE remorse	

posts) (Demszky et al., 2020)

• Granularity settings: **6 basic** emotions (Ekman, 1992) & 27 fine-grained emotions (Cowen and Keltner, 2017). Positive and negative emotion features were calculated.

Positi	ve	Negat	Ambiguous	
admiration 👋	joy 😃	anger 😡	grief 😢	confusion 😕
amusement 😂	love 🤎	annoyance 😒	nervousness 😬	curiosity 🤔
approval 📥	optimism 🤞	disappointment	remorse 😔	realization 💡
caring 🤗	pride 😌	disapproval 👎	sadness 😞	surprise 😲
desire 😍	relief 😅	disgust 🤮		
excitement 🤩		embarrassment 😳		
gratitude 🙏		fear 😨		

GoEmotions taxonomy: Includes 28 emotion categories, including "neutral" https://ai.googleblog.com/2021/10/goemotions-dataset-for-fine-grained.html

Data

- 13,000 documents labeled with PHQ-9 and **GAD-7 scores** created from **>337,000 messages** from message-based therapy sessions from **>6,500 unique** patients collected via Talkspace (Hull et al., 2020)
- Patients and clinicians gave consent; IRB approved; data handled securely.





GoEmotions Ekman



GE = GoEmotions (fine-grained). GEE = GoEmotions Ekman

GoEmotions features are collectively more predictive

- Trained Random Forest models with combinations of feature sets as input.
- Fine-grained GoEmotions features combined with LIWC's syntactic and topic (non-emotion) features were most predictive.

AUROCs, F1 score (positive class), precision, and recall of models trained with combinations of non-emotion and emotion features for predicting MDD/GAD

	MDD				GAD				
	ROC	F1	Pr	Rc	ROC	F1	Pr	Rc	
LIWC non-emo	0.577	0.413	0.525	0.341	0.549	0.290	0.478	0.209	
LIWC emo	0.621	0.471	0.561	0.405	0.613	0.405	0.541	0.324	
GoEmo Ekman	0.643	0.493	0.583	0.427	0.643	0.443	0.550	0.371	
GoEmo Cowen	0.662	0.522	0.613	0.455	0.652	0.444	0.565	0.366	
LIWC non-emo+									
LIWC emo	0.640	0.484	0.569	0.420	0.617	0.401	0.529	0.324	
GoEmo Ekman	0.655	0.498	0.585	0.434	0.637	0.441	0.548	0.369	
GoEmo Cowen	0.671	0.514	0.615	0.441	0.654	0.451	0.568	0.374	



95% confidence intervals shown. Variables that were not significant ($p \ge 0.05$) are shown in gray.

slope

0.5

1.0

0.68 0.72 0.70 0.75

R2

R2

0.0

-0.5

we

bio

- Obtainable from all 3 sources: anger, sadness, positive & negative emotion
- Significant associations between these features and the PHQ-9/GAD-7 scores, and comparable predictive power (measured in R2).
- LIWC emotion features performed well, indicating that these features remain a good choice, e.g. if computational constraints preclude NN models.



Acknowledgements & References



This work was supported by the National Library of Medicine (grant number 67-3780) and by Innovation Grant "Informatics-Supported Authorship for Caring Contacts (ISACC)" from the Garvey Institute for Brain Health Solutions. Conflicts of interest: TDH is an employee of the platform that provided the data.

Alan S. Cowen and Dacher Keltner. 2017. Self-report captures 27 distinct categories of emotion bridged by continuous gradients. Proceedings of the *National Academy of Sciences of the United States of America*, 114(38):E7900–E7909, September.

Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A Dataset of Fine-<u>Grained Emotions., May.</u>

Paul Ekman. 1992. Are There Basic Emotions? Psychological Review, 99(3):550-553.

Thomas D. Hull, Matteo Malgaroli, Philippa S. Connolly, Seth Feuerstein, and Naomi M. Simon. 2020. Two-way messaging therapy for depression and anxiety: longitudinal response trajectories. *BMC Psychiatry*, 20(1):297, December.

James W. Pennebaker, R. J. Booth, and M. E. Francis. 2007. Linguistic Inquiry and Word Count: LIWC.

Yla R. Tausczik and James W. Pennebaker. 2010. The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods. Journal of Language and Social Psychology, 29(1):24–54, March.

All three 0.520 0.612 0.453 **0.657 0.456** 0.567 **0.382** 0.671

Summary

- LIWC's emotion features are as predictive as GoEmotions features \rightarrow still a good choice.
- 2. GoEmotions features capture emotional state comprehensively, yielding additional clinically relevant nuance and benefitting predictive performance.
- Limitations: Non-diverse patient sample lacksquare(79% ≤ 35 y.o., 79% female, 75% BS or higher)
- Future work: Clinical decision support ullettools. **Interpretability is key**: models based on interpretable emotion features are preferred over black-box models
- Ethics: Monitoring may be considered lacksquareinvasive - informed consent is paramount. Further research & applications must take ethical considerations into account.