

Hankook University of Foreign Studies ELLT Graduate School Colloquium May. 12, 2022

> Modeling Emotional States from Speech Expressions

Tae-Jin Yoon Department of English Language and Literature







CONTENTS



Corpus-based phonology



Modeling Emotional States from Speech Expressions





Challenges of modern phonological research

- Without the ability to collect and analyze speech sounds in various environments and contexts, it can be a significant obstacle to trying to do the latest phonetics/phonology research.
- Learning a research methodology that can collect and process a variety of phonological and metainformation and speech data would be necessary to participate in various academic discussions, which is not an easy task to learn in a short time.



-카이스트가 왜 인문학 위기를 걱정하나.

"우리 사회엔 분명 인문학이 필요하다. 앞으론 문·이과 융합 인재 수요가 폭발적으로 늘어날 거다. 스스로 새 사상을 창조하는 사람이 리더가 된다. 대표적인 분야가 'AI와 협력'이다. 지난 2017년 바둑기사가 AI와 한 조를 이뤄 대결하는 복식경기에서 중국 롄샤오 8단이 구리 9단을 상대로 예상 밖 불계승을 거뒀다. 일대일로 붙으면 8단 기 사가 지겠지만, AI와 협력을 더 잘한 롄샤오가 구리를 이겼다. 바둑뿐 아니라 앞으로 모든 영역에서 AI와 협력을 잘하는 방법을 연구하는 인문학이 필요할 것이다. 인문학 도 4차 산업혁명 시대엔 연구 툴(방식)을 바꿔야 한다. 지금의 인문계 교육으론 '문송 (문과라서 죄송)'은 사라지지 않는다."





In The Handbook of Corpus Phonology (2012)



Opportunities

- Several benefits to doing phonetic research on a large database
 - □ the ability to do robust testing of linguistic hypotheses
 - □ the use of data from naturalistic interactions
 - the possibility of emergence of macroscopic structure in large enough databases.

(Goldstein, 2011)



Challenges

- How can we (semi)automatically extract the kinds of representation and features from large corpora?
- What kinds of tests can we apply to the extracted features to test our phonologically motivated questions?





Vowel Length merger in Seoul Korean



Kang, Yoonjung, Tae-Jin Yoon, and Sungwoo Han (2015). Frequency effects on the vowel length merger in Seoul Korean. *Laboratory Phonology*, 6(3-4): 469–503.



VOT and frequency over time

Low word freq

High word freq



Higher frequency words: smaller VOT difference (p = 0.025)

Hye-Young Bang, Morgan Sonderegger, Yoonjung Kang, Megan Clayards, and Tae-Jin Yoon (2018) The emergence, progress, and impact of sound change in progress in Seoul Korean: implications for mechanisms of tonogenesis. *Journal of Phonetics*. 66. 120-144.



Modeling of Nucleus F0 on Korean Accentual Phrase



Tae-Jin Yoon (2017) Growth Curve Modeling of Nucleus F0 on Korean Accentual Phrase. Phonetics and Speech Sciences 9(3), 17-23.



Corpus







음성/자연어	헬스케어	안전		
^{데이터 48종}	데이터 35종	데이터 21종		
비전	<mark>자율주행</mark>	농축수산	국토환경	
데이터 36종	^{데미터 23중}	데이터15종	데이터 128	
			교육 ^{ផのគ138}	

015 015	옷 알 옷 옷 배를 준비하는 소중한 우리말 자 알 휴 신청	원 (1) (1) (1) (1) (1) (1) (1) (1) (1) (1)	}	
	총 32 건	<u> 27</u> 지식하요가 ~ -		=
신규	신문 말뭉치 2021	(버전 1.0) 중합지, 전문지, 인터넷 가만 신문 매체의 기사(2020년)로 구성된 말랑치입니다.	1	신청하기 🕣
신규	국회 회의록 말뭉치 2021	(버진 1.0) 국회 소위원회 회의록(2003-2020년)으로 구성된 말봉치입니다.	(?)	신청하기 (+)
신규	추론_확신성 분석 말뭉치 2021	(버진 1.0) 내포문에서 주출한 가설에 화자의 확신성을 추론한 정보를 부탁한 말당치입니다.	?	신청하기 (+)
신규	맞춤법 교정 말뭉치 2021		3	신청하기 (+)
신규	속성 기반 감성 분석 말뭉치 2021	(버전 1.0) 국립국어원 감성 분석 말뭉치 2020(1.0)과 동일한 문서에 속성 기반 감성 정보를 부착한 말뭉치입니다.	1	신청하기 (+)
신규	개체명 분석 말뭉치 2021	(버전 1.0) 문장에 나티난 개체영의 경계를 표시하고 분석 표지를 부착한 말 당치입니다.	?	신청하기 🕀
신규	개체명 분석 말뭉치 개체 연결 2021	(버전 1.0) 개체명 분석 알용치에 위키피디아 정보를 부착한 자료입니다.	?	신청하기 (+)
신규	온라인 대화 말뭉치 2021	(버전 1.0) 두 명 이상의 대화 참여자가 온라인 공간에서 주고받은 대화 자료로 구성된 말중치입니다.	?	신청하기 (+)
	추론_확신성 분석 말뭉치 2020	(버전 1.0) 내포문에서 주출한 가설에 화자의 확신성을 추론한 정보를 부작한 말등차입니다.	?	신청하기 (+)
	의미역 분석 말뭉치	(버전 1.0) 문장의 술어가 가지는 논향을 분석하고 의미 역할을 부착한 말문처입니다.	?	신청하기 (+)

https://www.researchgate.net/figure/BM-big-data-characteristics-3V-Adopted-from-4_fig6_303562879



Children's developmental pattern of AP









(

Modeling Emotional States from Speech Expressions





Data-driven approach



https://christophm.github.io/interpretable-ml-book/terminology.html



Why Emotion in phonetics?

- Emotions: psychological states variously associated with thoughts, feelings, behavioral responses, and a degree of pleasure or displeasureurce: https://en.wikipedia.org/wiki/Emotion
- Modulation of pitch, loudness, duration, and voice quality across syllables in an utterance
 - ightarrow conveys both linguistic and non-linguistic information

prominence, prosodic phrasing, ...

age, gender, speaker's emotional status, etc...



Two theories on emotion

□ The discrete emotion theory

- Basic discrete emotions
 - (1) surprise, (2) interest, (3) joy, (4) rage, (5) fear, (6) disgust, (7) shame, (8) anguish
- Individual emotions have biological and neurological profiles

□ The dimensional theory

- Two emotional dimensional spaces distinguish emotions
 - (1) valence how positive or negative an emotion is
 - (2) arousal the intensity of an emotion



The discrete emotion approach

- Emotions are discrete, measurable, and physiologically distinct.
- □ Certain emotions appeared to be universally recognized.

→ Many studies have examined the vocal characteristics of speech in hope of defining a vocal signature for each basic emotion (Russell 2003)



The Dimensional approach

- So far, the strongest single association found for vocal acoustic have been with the sender's general **arousal** level.
- High-arousal emotions such as anger and joy have similar characteristics low arousal emotions such as sadness
 - □ greater loudness,
 - higher pitch
 - □ faster speech
- Few works have concentrated on distinguishing emotions between positive- and negative- valence emotions such as anger and joy.



Eerola, T., & Vuoskoski, J. K. (2010). *A comparison of the discrete and dimensional models of emotion in music. Psychology of Music, 39(1), 18–49.* doi:10.1177/0305735610362821



Research topic

- □ F0 contours (not a summary F0 such as meanF0 or F0 range) contains discriminatory information about emotions.
- Very few can be found in the literature that made the efforts to describe the shape of f0 contours directly in classifying emotions

The Ryerson Audio-Visual Database of Emotional Speech 🧭 실선여자대학교 and Song (RAVDESS)

- The RAVDESS dataset is a multimodal validated English dataset that contains speech, song, and video files that represent 8 emotions.
- ❑ The portion of the dataset that I use in this study is the speech audio files that are represented by 1440 wave file.
- □ Twenty-four professional actors (12 female and 12 male) with 60 trials for each actor produced the 1440 wave files (24×60 =1440).

	Fname	Emotions	Intensity	Repetition	Actor	Gender	Statement
0	03-01-05-01-02-01-16	angry	normal	1	16	female	Dogs are sitting by the door
1	03-01-06-01-02-02-16	fear	normal	2	16	female	Dogs are sitting by the door
2	03-01-06-02-01-02-16	fear	strong	2	16	female	Kids are talking by the door
3	03 - 01 - 05 - 02 - 01 - 01 - 16	angry	strong	1	16	female	Kids are talking by the door
4	03 - 01 - 07 - 01 - 01 - 01 - 16	disgust	normal	1	16	female	Kids are talking by the door

Livingstone, S.R. and F.A. Russo. (2018). The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS ONE* 13(5), e0196391.

The Ryerson Audio-Visual Database of Emotional Speech 🛞 생선여자대학교 and Song (RAVDESS) (2)

□ The actors vocalized two sentences in a neutral North American accent.

- □ "Kids are talking by the door"
- Dogs are sitting by the door"
- The emotions included in this dataset are
 - □ neutral, calm, happy, sad, angry, fearful, surprise, and disgust
- Each expression is produced at two levels of emotional intensity (normal and strong) except for the neutral emotion that is recorded in a normal intensity only.

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) (3)

성신여자대학교



neutral calm happy sad angry fearful disgust surprised

https://www.kaggle.com/datasets/uwrfkaggler/ravdess-emotional-speech-audio



RAVDESS Examplars

RAVDESS Exemplars

Emotional Speech

Livingstone & Russo, 2018



Validity of the RAVDESS dataset

- 247 untrained research participants from North America



Livingstone, S.R. and F.A. Russo. (2018). The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PLOS ONE* 13(5), e0196391.

I. Random Forest



Support = # of samples

성신여자대학교

- Precision: 양성으로 예측된 것 중 진짜 양성은 얼마나?
- Recall: 양성 샘플 중 양성 클래스로 분류된 샘플의 수는?
- F-score: precision과 recall의 조화 점수

Müller, A. C., & Guido, S. (2016). Introduction to machine learning with Python: a guide for data scientists. " O'Reilly Media, Inc.".



II. Deep Learning (CNN)



modified from https://github.com/marcogdepinto/emotion-classification-from-audio-files/blob/master/legacy_code/data_exploration/EmotionsRecognition.ipynb



model accuracy					precision		f1-score	support			
0.9	— tra	ain			-	alayan ta	neutral	0.26	0 40	0 31	25
0.8		51		at the second second			calm	0.65	0.59	0.62	61
0.7			ALL DESCRIPTION OF				happy	0.53	0.61	0.56	66
0.6		and the second second		All and the	L.S	-	sad	0.58	0.45	0.51	66
0.5		A Company of the second se	-	week and the second			angry	0.70	0.60	0.65	72
0.4		Land States	1				fearful	0.58	0.63	0.61	60
03	and the second second	e 1					disgust	0.50	0.55	0.52	60
0.5	1						surprised	0.60	0.56	0.58	66
0.2	1										
0.1	0	200	10.0	600	000	1000	accuracy			0.56	476
	0	200	400 en(600 ach	800	1000	macro avg	0.55	0.55	0.55	476
			cpt				weighted avg	0.58	0.56	0.57	476

de Pinto, M. G., Polignano, M., Lops, P., & Semeraro, G. (2020, May). Emotions understanding model from spoken language using deep neural networks and mel-frequency cepstral coefficients. In 2020 IEEE Conference on Evolving and Adaptive Intelligent Systems (EAIS) (pp. 1-5). IEEE.



III. Transfer Learning



Baevski, A., Zhou, Y., Mohamed, A., & Auli, M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, *33*, 12449-12460.



wav2vec2.0 for Children-Specific ASR





from unicode_jamo import *
join_jamos(processor.decode(pred_ids))

'거북인 이렇게 겁자난 엉그멍금'



<pre>***** train metrics *****</pre>			
epoch	=		30.0
total_flos	=		465552446523GF
train_loss	=		0.4268
train_runtime	=	14	days, 11:02:18.34
train_samples	=		161732
<pre>train_samples_per_second</pre>	=		3.884
<pre>train_steps_per_second</pre>	=		0.061



wav2vec2.0 based classification



	precision	recall	fl-score	support
angry	0.73	0.81	0.77	27
calm	0.84	0.78	0.81	27
disgust	0.72	0.67	0.69	27
fear	0.37	0.82	0.51	28
happy	0.50	0.07	0.13	27
neutral	0.00	0.00	0.00	14
sad	0.26	0.33	0.30	27
surprise	0.83	0.74	0.78	27
accuracy			0.56	204
macro avg	0.53	0.53	0.50	204
ghted avg	0.57	0.56	0.53	204

Support = # of samples



- **IV. Generalized Additive Mixed Modeling**
 - In Linear Model, the mean of data is modeled as a sum of linear terms

$$y_i = \beta_0 + \sum_j \beta x_{ji} + \varepsilon_i$$

 In Generalized Additive Mixed Model, the mean of data is modeled as a sum of *smooth* functions (= smooths)

$$y_i = \beta_0 + \sum_j s_j(x_{ji}) + \varepsilon_i$$







GAMM approach to the F0 contour modeling



성신여자대학교

Gamm Modeling



Emotionscalm-0.68600.6277-1.0930.274459Emotionsdisgust2.21080.62773.5220.000428***Emotionsfear7.33910.627711.692< 2e-16</td>***Emotionshappy5.75150.62779.163< 2e-16</td>***Emotionssad2.81230.62774.4807.46e-06***Emotionssurprise6.27530.62779.997< 2e-16</td>***Signif. codes:0'***'0.001'*'0.05'.'0.1'1

R-sq.(adj) = 0.649 Deviance explained = 64.9% fREML = 1.0445e+06 Scale est. = 20.228 n = 357120



Pair-wise comparison of contours



• F0 contour differentiates angry from happy.

• F0 contour hardly differentiates calm from neutral.



Conclusion

- I reviewed briefly what and how we can do with large-scaled datasets of spoken language.
- □ I attempted to compare approaches to emotional classification
- □ I put an emphasis on modelling emotions using F0 contours as an input to generalized additive model (GAM)
 - □ The present approach has predictive power.
 - □ The additive model provides visualized aids and makes us better understand validity of the data obtained from human labelers.



 \bigcap

Thank you

