



Affective computing from speech: Towards robust recognition of emotions in ecologically valid situations

Fabien Ringeval Laboratoire d'Informatique de Grenoble Université Grenoble-Alpes

Motivation

• Interactive technology is getting more and more ubiquitous







Motivation

- Engineering research focused on human-centered technology
 - Automatic human behaviour understanding
 - General interaction dynamics
 - Developmental, psychiatric disorders
 - Depression, suicide risk, marital therapy
 - Solitude, gerontechnology
 - Overlap with many research domains
 - Social & Human Sciences
 - Psychology, Psychiatry
 - Signal processing
 - Machine learning





Human Behaviour

• Push-Pull effect (Scherer) — WYSIWYG





Grenoble Alpes

K. Scherer, *The evolutionary origin of multimodal synchronization in emotional expression*, Journal of Anthropological Sciences, Vol. 91 (2013), pp. 185-200.

Human Behaviour Understanding

UNIVERSITÉ

Grenoble Alpes





Human Behaviour Understanding from Speech

• Paralinguistic – reading between the lines: « How » vs « What »



Human Behaviour Understanding from Speech

• The voice tells a lot about ourselves

• Age, gender, emotion, culture, personality, role, mental state, physical state, ...



Human Behaviour Understanding from Speech

• Paralinguistic: from short-term states to long-term traits



Collecting Affective Data

Elicitation methods

- Acted, certainty +, naturalness (EMODB, GEMEP)
- Induced, certainty +--, naturalness +-- (eNTERFACE, CPESD)
- Spontaneous: certainty -, naturalness + (SEMAINE, RECOLA)





Alpes



GEneva Multimodal Emotion Portrayal



Child Pathological & Emotional Speech Database



REmote **COL**laborative and **A**ffective interaction



Quantifying Emotions

- Discrete framework
 - Semantic descriptions
 - "Hard-coded" in the brain
 - Common categories:
 - Anger
 - Fear
 - Happiness
 - Sadness



R. Plutchik, *Emotions: A general psychoevolutionary theory*, in K. R. Scherer and P. Ekman [Eds], *Approaches to Emotion*, Erlbaum, Hillsdale (NJ), 1984.





Fabien Ringeval

Quantifying Emotions

- Continuous framework
 - Description of properties
 - Common dimensions
 - Arousal: calm vs. excited
 - Valence: negative vs. positive



J. A. Russell. *A circumplex model of affect*. Journal of Personality and Social Psychology, 39(6):1161–1178, 1980.





Fabien Ringeval

Quantifying Emotions

- Continuous framework
 - Description of properties
 - Common dimensions
 - Arousal: calm vs. excited
 - Valence: negative vs. positive
 - Many categories can be mapped onto the VA space





G. Paltoglou and M. Thelwall. Seeing stars of valence and arousal in blog posts. IEEE TAFC, 4(1):116–123, 2013.

Data Collection: RECOLA database

• Collect multimodal data in ecologically valid situation (collaboration)



Data Collection: RECOLA database

• First dyad (P13 — P14)





F. Ringeval, Introducing the RECOLA Multimodal Corpus of Remote Collaborative and Affective Interactions, in Proc. of EmoSPACE, IEEE International Conference on Face & Gestures 2013.

14 L I G

Fabien Ringeval

• Web-based annotation interface (open source); 6 raters (3F/3M)

valence - SEQUENCE 2.mp4 social dimension annotation arousal - SEQUENCE 3.mp4 valence - SEQUENCE 3.mp4 social dimension annotation

arousal - SEQUENCE 2.mp4

arousal - SEQUENCE 4.mp4 valence - SEQUENCE 4.mp4 social dimension annotation

arousal - SEQUENCE_5.mp4 valence - SEQUENCE_5.mp4 social dimension annotation

arousal - SEQUENCE 6.mp4 valence - SEQUENCE 6.mp4 social dimension annotation

UNIVERSITÉ

Alpes

renoble







• Raw data



G

 Interpolated data (piecewise Hermite interpolation; d²/ dt not continuous)



• Interpolated data for all raters



G

• Can we improve by post-processing? (E.g., zero-centring)



- Can we improve by post-processing? (E.g., zero-centring)
 - Looks great and improves Inter-Rater Agreement (averaged ρ_c)
 - BUT, mean perceived emotion is forced to neutral





20



Fabien Ringeval

- Can we improve by post-processing?
 - Proposed method: centre to mean weighted by the inter-rater agreement
 - 1. Compute inter-rater agreement $\bar{\rho}_d(i)$ for a dimension d and an evaluator $i \in [1, N_e]$
 - 2. Average ratings $y_d^{e_i}$ of all evaluators e_i and weight by $\bar{\rho}_d(i)$ to obtain the centring value \bar{y}_d
 - 3. Centre ratings and average to obtain a **gold-standard** trace $g_d(t)$

$$\bar{\rho}_{d}(i) = \frac{1}{N_{e} - 1} \sum_{i=1, j \neq i}^{N_{e}} \rho_{d}(i, j)$$
(1)

$$\bar{y}_{d} = \frac{1}{\sum_{i=1}^{N_{e}} \bar{\rho}_{d}(i)} \sum_{i=1}^{N_{e}} \frac{1}{T} \sum_{t} y_{d}^{e_{i}}(t) \bar{\rho}_{d}(i) \quad (2)$$







Fabien Ringeval

• Post-processing by centring with IRA

- Same improvement as zero-centring for IRA
- Original (positive) skew better preserved

	Arousal	Valence
raw; ρ _c	0.28	0.37
raw;%pos.	59.0	70.5
zero-m. ; $ ho_c$	0.33	0.43
zero-m. ; % pos.	50.8	44.8
wgt-m. ; $ ho_c$	0.33	0.43
wgt-m.;% pos.	48.5	74.1





Emotion Recognition from Speech

Speech conveys many relevant cues

Alpes



Emotion Recognition from Speech

- Features extraction from speech openSMILE:)
 - Low-level descriptors and functionals
 - Sliding window shifted forward at a constant rate (25 Hz)



by audEERING[™]



Emotion Recognition from Speech

• Machine learning: sequential learning problem



25



- Modelling asynchronous ratings of emotion with (B)LSTM-RNN
 - LSTM-RNN can capture long-range contextual dependencies
 - Exploit this to model delay between raters (multi-task learning)









Fabien Ringeval, Prediction of Asynchronous Dimensional Emotion Ratings from Audiovisual and Physiological Data. Pattern Recognition Letters, 66:22–30, November 2015.

- Setup
 - RECOLA: 27 subjects (9/9/9) x 5 min. Features: ComParE 130 LLDs x 5 func.
 - Window size: from 0.5 s (micro-expressions) to 6 s (information loss)
 - (FF, (B)LSTM) **DNN**: $\sigma_i = 0.1$; 2 h.l. (128); 100 epochs; early stop. MSE
 - Effect of window size (best network topology, single/multi-task)



- Setup
 - RECOLA: 27 subjects (9/9/9) x 5 min. Features: ComParE 130 LLDs x 5 func.
 - Window size: from 0.5 s (micro-expressions) to 6 s (information loss)
 - (FF, (B)LSTM) **DNN**: $\sigma_i = 0.1$; 2 h.l. (128); 100 epochs; early stop. MSE
 - Impact of contextual modelling (averaged over all modalities)



- Setup
 - RECOLA: 27 subjects (9/9/9) x 5 min. Features: ComParE 130 LLDs x 5 func.
 - Window size: from 0.5 s (micro-expressions) to 6 s (information loss)
 - (FF, (B)LSTM) **DNN**: $\sigma_i = 0.1$; 2 h.l. (128); 100 epochs; early stop. MSE
 - Effect of multi-tasking (mean vs. individual ratings)

	Arousal	Valence
AUDIO - SINGLE	.732	.412
AUDIO - MULTI	.738	.343
VIDEO - SINGLE	.403	.339
VIDEO - MULTI	.427	.349



Experiments: CCC as objective function

- Neural networks typically trained using SSE-related criteria
- Goal: Directly optimize predictor for the evaluation metric instead
- ho as evaluation metric **not directly usable** for optimization
 - Scale (y_{sc}) and shift (y_{sh}) invariant (!!!)
 - Leads to **ill-posed** optimisation problem

$$\mathcal{O} = -\sum_{i\in\beta,f\in\mathcal{F}}\rho_c$$



Felix Weninger, Fabien Ringeval, Erik Marchi, and Björn Schuller. Discriminatively trained recurrent neural networks for continuous dimensional emotion recognition from audio. In Proc. of IJCAI 2016, pp. 2196–2202, New York City (NY), USA, July 2016

Experiments: CCC as objective function

- ρ_c : CCC over the whole evaluation set
 - Expect good CCC when predicting the correct mean for each sequence
- $\Sigma
 ho_c$: Sum the CCC in each sequence
 - Need to get the variation within each sequence right, but sensitive to noise in the gold-standard

•
$$\rho_c$$
 = 0.37, $\Sigma \rho_c$ = 0



Experiments: CCC as objective function

- Setup
 - **RECOLA**: 46 subjects (16/15/15) x 5 min. **Features**: eGeMAPS (AV+EC 2015)
 - (FF, (B)LSTM) **DNN**: $\sigma_i = 0.1$; 2 h.l. (128); 100 epochs; early stop. MSE
 - Significant improvements with ${f \Sigma}
 ho_c$ objective; trade-off between RMSE and ho_c

	Arousal			Valence		
О	RMSE	ρ _c	$E\{ ho_c \}$	RMSE	ρ _c	$E\{ ho_c \}$
SSE	.128	.097	.161	.108	.131	.052
ρ_c	.193	.254	.212	.130	.155	.080
Σρ _c	.200	.350	.268	.192	.199	.139





Experiments: Robustness to noise

- How non-stationary noise and reverberation degrade performance?
- Noises: CHiME'15, Smartphone (Nexus One), Hall, Train
- Various SNR (0dB, 3dB, 6dB, 9dB, 12dB)
- Method: signal/feature enhancement with auto-encoders



UNIVERSITÉZixing Zhang, Fabien Ringeval, Jin Han, Jun Deng, Erik Marchi, and Björn Schuller. Facing realism in spontaneous emotion recognition Grenoble from speech: Feature enhancement by autoencoder with LSTM neural networks. In *Proc. INTERSPEECH 2016*. 33

Experiments: Robustness to noise

Alpes

• Results on Arousal / CHiME'15 (MFCCs, SVR for prediction – AVEC'16)





Experiments: Adieu features?

- CNN-DNN can learn appropriate representation from raw input
- How well can it perform for **emotion recognition**?



- Sergpæyæirtstælkællukiowi (FKH2460tsplæsseltingesæltplædtees) @ 5ms)
- Maxing calicrog sections e (pharonins ize (2) to rside v20) sampling @ 8kHz

Grenoble Alpes

George Trigeorgis, Fabien Ringeval, Raymond Brueckner, Erik Marchi, Mihalis Nicolaou, Björn Schuller, and Stefanos Zafeiriou. Adieu features? End-to-end speech emotion recognition using a deep convolutional recurrent network. In *Proc ICASSP 2016*.



Experiments: Adieu features?

- CNN-DNN can learn appropriate representation from raw input
- How well can it perform for **emotion recognition**?







Experiments: Adieu features?

• Results

Predictor	Features	Arousal (ρ_c)	Valence (ρ_c)	
a. Mean squared error objective				
SVR	eGeMAPS	.318 (.489)	.169 (.210)	
SVR	ComParE	.366 (.491)	.180 (.178)	
BLSTM	eGe <mark>M</mark> APS	.300 (.404)	.192 (.187)	
BLSTM	ComParE	.132 (.221)	.117 (.152)	
Proposed	raw signal	.684 (.728)	.249 (.312)	

b. Concordance correlation coefficient objective

BLSTM	eGeMAPS	.316 (.445)	.195 (.190)
BLSTM	ComParE	.382 (.478)	.187 (.246)
Proposed	raw signal	.686 (.741)	.261 (.325)







Cooperative Regression Models (CRMs)

G

Arianna Mencattini, Eugenio Martinelli, Fabien Ringeval, Björn Schuller, and Corrado Di Natale. Continuous estimation of emotions in speech by dynamic cooperative speaker models. *IEEE Transactions on Affective Computing*, 2016.





• Frequency of Single Speaker Regression Model inclusion



noble

Alpes







Experiments: Bag of audio words

• Learn dictionary of acoustic words instead of using functionals (MFCC)



openWORD

$$TF'(w) = \lg(TF(w) + 1)$$





Maximilian Schmitt, Fabien Ringeval, and Björn Schuller. At the border of acoustics and linguistics: Bag-of-Audio-Words for the recognition of emotions in speech. In *Proc. INTERSPEECH 2016*, San Fransisco (CA), USA, September 2016.

Experiments: Bag of audio words

• Results: learning time-delay







Experiments: Bag of audio words

• Results: comparison with state-of-the-art

Model	Ref.	CCC			
		Arousal		Valence	
		Valid	Test	Valid	Test
BLSTM-RNN	[26]	.800		.398	
CNN (end-to-end)	[4]	.741	.686	.325	.261
Proposed (BoAW)	Table 2	.793	.753	.550	.430
Proposed (early fusion)	Table 4	.799	.738	.521	.465





Conclusion

- Affective computing is a "hot" topic
- Vast application fields, robotics, behaviomedics, media, ...
- Shifted from small and acted to large and spontaneous
- Context modelling with LSTM helps, CNN+LSTM performs well
- Simple but well defined systems can perform better than end-to-end
 => not (yet) the end of signal processing!



Perspectives

- Toward BIG DATA
 - ER community => 10 hours of data
 - ASR community => 10.000 hours of data



- Crawl the web vlogs, vimeo, youtube, ...
- Use novelty detection to identify non-neutral
- Use active learning to automatically label easiest





Perspectives

- Toward real-life data
 - Not "properly" recorded => issue
 - But that's real life!!!
 - Wearable sensors, smartphones, webcam
 - Face diversity of individuals (transfer learning)



Perspectives

- Toward multimodality
 - Datasets are not all with ECG, EDA, FAU, ...
 - Use pseudo-multimodality to cross fertilise the missing modalities
 - Model context with available sensors (e.g., GPS, movements, logs, calls)





Commercial break – AVEC'16 16th October



The 6th Audio/Visual Emotion Challenge and Workshop



@ ACM-Multimedia 2016, Amsterdam, The Netherlands



Commercial break – ACM TOMM 31st October

Special Section on Multimedia Computing and Applications of Socio-Affective Behaviors in the Wild

Fabien Ringeval, Björn Schüller, Michel Valstar, Jonathan Gratch, Roddy Cowie, Maja Pantic



ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)

Formerly known as TOMCCAP, TOMM focuses on multimedia computing, multimedia communications, and multimedia applications.





Thank you!

Imperial College London





UNIVERSITÉ DE FRIBOURG UNIVERSITÄT FREIBURG







JNIVERSITÉ Grenoble

Alpes



Technische Universität München

