Cognitive Media Processing

Cognitive Media Processing #11

Nobuaki Minematsu







Spectrum envelope-based feature such as CEP: o

But *o* depends on all the three kinds of info. (ling, para-ling, extra-ling).

Generation How to suppress extra-linguistic variation in o?

 \bigcirc Feature normalization: transforming o to that of the standard speaker

- Model adaptation: modifying model parameters to fit to the input speaker
- Statistical independence: hiding these variation through sample collection
- Physical independence: pursuing features invariant to these variation

A claim found in classical linguistics

Theory of relational invariance [Jakobson+'79]
 Also known as theory of distinctive features
 Proposed by R. Jakobson

We have to put aside the accidental properties of individual sounds and substitute a general expression that is the common denominator of these variables.

Physiologically identical sounds may possess different values in conformity with the whole sound system, i.e. in their relations to the other sounds.





Roman Jakobson Linda R. Waugh

mouton de gruyter

THE SOU

LANGUA

Complete transform-invariance

Complete invariance between two spaces

- An assumption
- An event in a space should be represented as distribution.
 - Event p in space A is transformed into event P in space B
 - p and P are physically different (/a/ of speaker A and /a/ of speaker B)



Complete transform-invariance

Any general expression for invariance?[Qiao'10]
 BD is just one example of invariant contrasts.
 f-divergence is invariant with any kind of transformation.

- $\bigcirc f_{div}(p_1, p_2) = \int p_2(\boldsymbol{x}) g\left(\frac{p_1(\boldsymbol{x})}{p_2(\boldsymbol{x})}\right) d\boldsymbol{x}$

♀ Invariant features have to be f-divergence.

 \subseteq If $\oint M(p_1(\boldsymbol{x}), p_2(\boldsymbol{x})) d\boldsymbol{x}$ is invariant with any transformation,

• The following condition has to be satisfied. $M = p_2(\boldsymbol{x})g\left(\frac{p_1(\boldsymbol{x})}{p_2(\boldsymbol{x})}\right)$



Invariance in variability

Topological invariance [Minematsu'09]

♀ Topology focuses on invariant features wrt. any kind of deformation.











Invariant speech structure

Utterance to structure conversion using *f*-div. [Minematsu'06]



An event (distribution) has to be much smaller than a phoneme.

A claim found in classical linguistics

Theory of relational invariance [Jakobson+'79]
 Also known as theory of distinctive feature
 Proposed by R. Jakobson

We have to put aside the accidental properties of individual sounds and substitute a general expression that is the common denominator of these variables.

Physiologically identical sounds may possess different values in conformity with the whole sound system, i.e. in their relations to the other sounds.







THE SOUL

Roman Jakobson Linda R. Waugh

LANGUAG

mouton de gruyter

More classical claims in linguistics

Ferdinand de Saussure (1857-1913)

- Father of modern linguistics
- "Course in General Linguistics" (1916)
- What defines a linguistic element, conceptual or phonic, is the relation in which it stands to the other elements in the linguistic system.
- The important thing in the word is pot the sound alone but the phonic listinguish this word from the others.

Junguage is a system of only conceptual differences and phonic differences.

c_{2}	$\begin{bmatrix} d_{11} \\ d_{21} \\ d_{31} \\ \vdots \\ d_{N1} \end{bmatrix}$	$egin{array}{c} d_{12} \ d_{22} \ d_{22} \ d_{N2} \end{array}$	 $d_{1N} \\ d_{2N}$	C G F

Course in General Linguistics Ferdinand de Saussure

A new framework for "human-like" speech machines #3

Nobuaki Minematsu





Cognitive Media Processing

Title of each lecture

Theme-1

- Multimedia information and humans
- Multimedia information and interaction between humans and machines
- Multimedia information used in expressive and emotional processing
- A wonder of sensation synesthesia -
- Theme-2
 - Speech communication technology articulatory & acoustic phonetics -
 - Speech communication technology speech analysis -
 - Speech communication technology speech recognition -
 - Speech communication technology speech synthesis -
- Theme-3
 - A new framework for "human-like" speech machines #1
 - A new framework for "human-like" speech machines #2
 - A new framework for "human-like" speech machines #3
 - A new framework for "human-like" speech machines #4



abcde g

h jk mn

opqrstu

VWXYZ



Menu of the last four lectures

Robust processing of easily changeable stimuli Robust processing of general sensory stimuli Q Any difference in the processing between humans and animals? Human development of spoken language Infants' vocal imitation of their parents' utterances What acoustic aspect of the parents' voices do they imitate? Speaker-invariant holistic pattern in an utterance Completely transform-invariant features -- f-divergence --Implementation of word Gestalt as relative timbre perception Application of speech structure to robust speech processing **Radical but interesting discussion** A hypothesis on the origin and emergence of language What is your definition of "human-like" machines?

Invariant speech structure

Utterance to structure conversion using *f*-div. [Minematsu'06]



An event (distribution) has to be much smaller than a phoneme.

A simple framework for isolated word recognition



🗳 Two big problems

Too strong invariance (two different words can be the same.)stream 1
 Multi-Stream Structuralization to constrain the invariance [Asakawa'0&ppstrur
 Too high dimension (N events leads to an NC2 dimensional vector.)
 2-stage LDA to reduce the dimension effectively [Asakawa'08]
 The invariance only wrt. speaker differences
 A mathematical model for VTL differences [Pitz,05]
 The invariance only wrt. any kind of band matrix (c' = Ac)



BD calc.

Structur

VTLC

Vocal tract length dif

Solution Θ Can be approximated as multiplication of matrix A in cep. domain. A is represented with warping parameter α .



BD calc.



Second se

 \bigcirc Word = V1V2V3V4V5 such as /eoaui/, PP = 120 (CL=0.8%)

- Solution Word-based HMMs (20 states) vs. word-based structures (20 events)
 - \bigcirc Training = 4M+4F adults, testing = other 4M+4F with various VTLs
- \bigcirc 4,130-speaker triphone HMMs are also tested with 0.30.
 - Some the speaker-independent HMMs widely used as baseline model in Japan



Second se

 \bigcirc Word = V1V2V3V4V5 such as /eoaui/, PP = 120 (CL=0.8%)

- Solution Word-based HMMs (20 states) vs. word-based structures (20 events)
 - \bigcirc Training = 4M+4F adults, testing = other 4M+4F with various VTLs
- - Some the speaker-independent HMMs widely used as baseline model in Japan



Second se

 \bigcirc Word = V1V2V3V4V5 such as /eoaui/, PP = 120 (CL=0.8%)

- Solution Word-based HMMs (20 states) vs. word-based structures (20 events)
 - \bigcirc Training = 4M+4F adults, testing = other 4M+4F with various VTLs
- - Some the speaker-independent HMMs widely used as baseline model in Japan



Second Se

- \bigcirc Mora-based length of words = 3 to 7
- Solution Word-based HMMs (25 states) vs. word-based structures (25 events)
 - \bigcirc Training = 15M+15F adults, testing = other 15M+15F with various VTLs



Application to more realistic ASR tasks [Suzuki+'15]
 Digits recognition and LVCSR (dictation)
 Use of structural features in discriminative reranking
 Str. scores and ASR scores are combined with average perceptron.



Continuous digits recognition

- Language = Japanese
- Baseline = GMM-HMM ASR
- Reranking = averaged perceptron
- ♀ Error reduction rate = 30%



Large vocabulary continuous speech recognition

- Language = Japanese
- Baseline = DNN-HMM ASR
- Reranking = averaged perceptron
- \bigcirc Error reduction rate = 5%

Many errors are due to a large number of homonyms in Japanese.

Table 6: CERs of the LVCSR experiment.BaselineProposedRelative improvement2.67%2.53%5.24%

Transformer model



Self-attention mechanism

- Relatedness (similarity) of the current input token
 - to the other tokens in the input sequence and
 - to the tokens in the output sequence generated so far.
- A token is converted to its three components.
 - Value vector, key vector, and query vector.



 $egin{aligned} {m E} \ (ec{e_1}, ec{e_2}, ..., ec{e_V}) \left(egin{aligned} 0 \ 1 \ 0 \ ec{e_1} \ 0 \end{array}
ight) = ec{e_2} \end{aligned}$

https://towardsdatascience.com/illustrated-self-attention-2d627e33b20a

Similarity \propto 1/Distance Distance \propto 1/Similarity

Utterance to

div. [Minematsu'06]

re



An event (distribution) has to be much smaller than a phoneme.

More classical claims in linguistics

Ferdinand de Saussure (1857-1913)

- Father of modern linguistics
- "Course in General Linguistics" (1916)
- What defines a linguistic element, conceptual or phonic, is the relation in which it stands to the other elements in the linguistic system.
- The important thing in the word is pot the sound alone but the phonic listinguish this word from the others.

Junguage is a system of only conceptual differences and phonic differences.

c_{2}	$\begin{bmatrix} d_{11} \\ d_{21} \\ d_{31} \\ \vdots \\ d_{N1} \end{bmatrix}$	$egin{array}{c} d_{12} \ d_{22} \ d_{22} \ d_{N2} \end{array}$	 $d_{1N} \\ d_{2N}$	C G F

Course in General Linguistics Ferdinand de Saussure

Transformer model

*y*₂

RNN

RNN

Decoder

Attention is all you need !!

- <u>https://arxiv.org/abs/1706.03762</u>
- Explicit modeling of the relations (similarities) of the current input token to other ones in the input sequence and to the tokens in the output sequence generated so far.

Self-attention mechanism

Encoder

RNN

X2

RNN

 x_1

 h_3

 x_3

RNN



Structure-based transformer

Comparison between the two transformers [Wang+'22]

CTC weight	Transform	er (6)		Proposed model (c)			
	CER-dev	CER-test	training time	CER-dev	CER-test	training time	
0	28.2	29.5	163134	28.4	29.9	107156	
0.3	5.5	6.2	132001	5.5	6.2	113071	
0.6	6.1	6.6	143398	5.9	6.6	112970	
1.0	5.4	5.9	120568	5.4	5.9	91890	







Language acquisition through vocal imitation

Utterance \rightarrow symbol sequence \rightarrow production of each sym.







Phonemic awareness is too poor to decompose an utterance.

Several answers from developmental psychology

- General Holistic/related sound patterns embedded in utterances
 - Holistic wordform [Kato'03]
 - Word Gestalt [Hayakawa'06]
 - Related spectrum pattern [Lieberman'80]

No mathematical formulation

The patterns have to include no speaker information in themselves. 9

- If they do it, children have to try to impersonate their fathers.
- What is the speaker-invariant and holistic pattern in an utterance?

Structure-to-speech conversion

Speech representation with extra-ling. features removed
 Speaker-specific vocal tract features are removed.
 With them, we can identify speakers by hearing voices.

Structure-to-speech conversion







How to implement the vocal imitation?

Acoustic instances are searched for in the voice space.
 Initial conditions : a few acoustic instances given from an infant
 Constrained conditions : speech Gestalt (distance matrix)



How to implement the vocal imitation?

Geometrical interpretation of BD-based constraints

$$BD(p_1(x), p_2(x)) = \frac{1}{8}(\mu_1 - \mu_2)^T \Sigma_{12}^{-1}(\mu_1 - \mu_2) + \frac{1}{2} \ln \frac{|\Sigma_{12}|}{|\Sigma_1||\Sigma_2|}$$

Search for a new target using BD(1, new), BD(2, new), BD(3, new)...

 $\subseteq \Sigma_{new}$ is given. Only μ_{new} is searched for in the current paper.



An experiment with real vocal imitation

Demonstration with my wife and daughter

Constraint conditions are given by my wife.Initial conditions are given by my daughter.













An experiment with real vocal imitation

Demonstration with my wife and daughter

Constraint conditions are given by my wife.Initial conditions are given by my daughter.



A big problem in CALL development

A very important and requisite function for CALL systems

- The system has to be able to ignore speaker differences.
 - Age and gender (the size and length of the vocal tube)
 - But no current system can ignore speaker differences well enough.
- Requirement of "acoustic matchedness" bet. HMMs and learners
 - Collection of children's speech or speaker adaptation of adult HMMs
 - Q : Learning to pronounce is learning to impersonate?



Mismatch problem

Speech model for another separation
 Separation between source and filter
 Separation between ling. and extra-ling.



A big solution for CALL development

For which does Minematsu's normal English sound closer ?

speaker	USA/F12	Minematsu O Minematsu
gender	female	male o male
age	?	37 0 37
mic	Sennheiser	cheap mic O cheap mic
room	recording room	living room on living room
AD	SONY DAT	PowerBook PowerBook
proficiency	perfect	good X Japanized



A big solution for CALL development **Proficiency estimation based on P(M | o) = GOP** $P(M|o) = P(p_1, ..., p_N|o)$ $= \frac{P(o|p_1, ..., p_N)P(p_1, ..., p_N)}{\sum_{p_i} P(o|p_1, ..., p_N)P(p_1, ..., p_N)}$ USA matsu nized) $\approx \frac{P(o|p_1,...,p_N)}{\sum_{p_i} P(o|p_1,...,p_N)}$ $P(o|p_1, \dots, p_N)$ \approx $\overline{\max_{p_i} P(o|p_1, \dots, p_N)}$ P(o|M)matsu USA $\overline{\max_M P(o|M)}$ nized) GOP (Goodness Of Pronunciation)

A big solution for CALL development

Proficiency estimation based on structural distance











Evaluation is done not based on whether each vowel sound has adequate acoustic property independently of others but based on whether a good vowel system underlies a learner's pronunciation.

<u>ering</u> (33333333333333 Preparation of data -- 96 simulated learners --**12** Japanese students who are returnees from US (A to L) See English words of /b-V-t/ and Japanese words of /b-V-to/ AE vowels : 1 word utterance per vowel J vowels : 5 word utterances per vowel Vowel segments are extracted automatically to estimate a vowel system. **Replacement of some AE vowels with J vowels** ♀ 12 speakers [A-L] x 8 pronunciations [1-8] = 96 learners α æ Ð ð Ι υ u 3 Э Λ ĺ S1E E E E E S2J a, æ, ʌ, ə, ð a J Ε E E E E E $\mathbf{S3}$ J J I, 1 1 E E Ε Е Е $\mathbf{S4}$ Ε .] $\mathbf{S5}$ E E E E E J υ, u u . J J J Ε **S6** E Ε E E J .) .) .) e 3 $\mathbf{S7}$ E E E E E E J J E **S**8 E E E E E E E E E Ε 0 0

Structure-to-structure distance measure

Euclidian distance between two distance matrices





Can approximate the structural distance after shift and rotation



Minimum of the total distances between corresponding points

§ 96 x 96 large distance matrix (12 spk. x 8 pron.)

Speakers: A to LProns: 1 to 8

Pronunciation clustering





Speaker clustering

Search Another distance measure between two structures

- Contrast-based comparison
- Substance-based comparison

$$\sqrt{\frac{1}{M} \sum_{i < j} (S_{ij} - T_{ij})^2}$$

Contrast-based comparison

Search Another distance measure between two structures

- Contrast-based comparison
- Substance-based comparison

$$\sqrt{\frac{1}{M} \sum_{i < j} (S_{ij} - T_{ij})^2}$$

Contrast-based comparison

1111113133333333336666666666666666665552555552528222222585584444413444447777777777772888288888788

GFFCEDHKKGK

Which vowels to correct at first?

Global difference between Student and Teacher

Euclidian distance between two distance matrices

Can be decomposed into local differences

Contribution of individual vowels to the global difference

$$d(v) = \sqrt{\frac{1}{M} \sum_{j=1}^{M} (S_{vj} - T_{vj})^2}$$

 $\$ Vowels of larger d(v) are should be corrected at first!!!

Which vowels to correct at first?

Estimation of the order of vowel correction

Only with given two matrices without the replacement table

	α	æ	Λ	Э	9r	Ι	i	ប	u	3	Э
S 1	J	J	J	J	J	J	J	J	J	J	J
S2	E	Е	Е	Е	Е	J	J	J	J	J	J
S3	J	J	J	J	J	Е	Е	Е	Е	Е	Е
S4	E	Е	J	J	J	Е	Е	J	J	Ε	Е
$\mathbf{S5}$	J	J	Е	Е	Е	J	J	Е	Е	J	J
S6	E	J	Е	J	Е	J	J	J	J	Е	Е
S7	J	Е	J	Е	J	Е	Е	Е	Ε	J	J
S 8	E	Е	Е	Е	Е	Е	Е	Е	Ε	Е	Е
							_				

S1 - S7 : Japanese English
S8 : American English
: replaced by J vowels

: no replacement

Very motivating interface for CALL

Select your favorite teacher!!

Shortest cut to your model!

Menu of the last four lectures

Robust processing of easily changeable stimuli Robust processing of general sensory stimuli Q Any difference in the processing between humans and animals? Human development of spoken language Infants' vocal imitation of their parents' utterances What acoustic aspect of the parents' voices do they imitate? Speaker-invariant holistic pattern in an utterance Completely transform-invariant features -- f-divergence --Implementation of word Gestalt as relative timbre perception Application of speech structure to robust speech processing **Radical but interesting discussion** A hypothesis on the origin and emergence of language What is your definition of "human-like" machines?

The final lecture of CMP

ssignme

will be given on next Tuesday (Jan 16).
The final assignment will also be given on that day.
If you cannot attend it, you should view the video.