

Prosody-Based Unsupervised Speech Summarization with Two-Layer Mutually Reinforced Random Walk

Sujay Kumar Jauhar Yun-Nung (Vivian) Chen Florian Metze

{sjauhar, yvchen, fmetze}@cs.cmu.edu

The 6th International Joint Conference on Natural Language Processing – Oct. 14-18, 2013



Language Technologies Institute School of Computer Science Carnegie Mellon University







Extractive Summarization



Extractive Summarization

- Speech Summarization
 - Spoken documents are more difficult to browse than texts
 - \rightarrow easy to browse, save time, easily get the key points
- Prosodic Features
 - Speakers may use prosody to implicitly convey the importance of the speech



Æ Extractive Summarization

Extractive Summarization (1/2)

- *Extractive Speech Summarization*
 - Select the indicative utterances in a spoken document
 - Cascade the utterances to form a summary



Extractive Summarization (2/2)

- Selection of Indicative Utterances
 - Each utterance U in a spoken document d is given an importance score I(U, d)
 - Select the indicative utterances based on I(U,d)
 - The number of utterances selected as summary is decided by a predefined ratio





9

- Prosodic Feature Extraction
- Ø Graph Construction
- Two-Layer Mutually Reinforced Random Walk



- Prosodic Feature Extraction
 - Oraph Construction
 - Two-Layer Mutually Reinforced Random Walk

Prosodic Feature Extraction

- For each pre-segmented audio file, we extract
 - number of syllables
 - number of pauses
 - duration time: speaking time including pauses
 - o phonation time: speaking time excluding pauses
 - speaking rate: #syllable / duration time
 - o articulation rate: #syllable / phonation time
 - fundamental frequency measured in Hz: avg, max, min
 - energy measured in Pa²/sec
 - intensity measured in dB



- Prosodic Feature Extraction
- Graph Construction
 - Two-Layer Mutually Reinforced Random Walk

Graph Construction (1/3)

- OUtterance-Layer
 - Each node is the utterance in the meeting document



13

Graph Construction (2/3)

- OUtterance-Layer
 - Each node is the utterance in the meeting document
- Prosody-Layer
 - Each node is a prosodic feature



14

Graph Construction (3/3)

- OUtterance-Layer
 - Each node is the utterance in the meeting document
- Prosody-Layer
 - Each node is a prosodic feature
- Between-Layer
 Relation



The weight of the edge is the normalized value of the prosodic feature extracted from the utterance



Prosodic Feature Extraction

Ø Graph Construction

Two-Layer Mutually Reinforced Random Walk

Mathematical Formulation

utterance scores at (t+1)-th iteration

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$



the second second second second second

Mathematical Formulation

$$\begin{cases} F_{U}^{(t+1)} = (1-\alpha)F_{U}^{(0)} + \alpha \cdot L_{UP}F_{P}^{(t)} \\ F_{P}^{(t+1)} = (1-\alpha)F_{P}^{(0)} + \alpha \cdot L_{PU}F_{U}^{(t)} \end{cases}$$

- Original importance
 - Outterance: equal weight



18

Mathematical Formulation

scores propagated from prosody nodes weighted by prosodic values

19

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$

- Original importance
 - Outterance: equal weight



Mathematical Formulation

$$\begin{cases} F_{U}^{(t+1)} = (1-\alpha)F_{U}^{(0)} + \alpha \cdot L_{UP}F_{P}^{(t)} \\ F_{P}^{(t+1)} = (1-\alpha)F_{P}^{(0)} + \alpha \cdot L_{PU}F_{U}^{(t)} \end{cases}$$

prosody scores at (t+1)-th iteration

- Original importance
 - Outterance: equal weight



Mathematical Formulation

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$

original importance of prosodic features

- Original importance
 - O Utterance: equal weight
 - Prosody: equal weight



Mathematical Formulation

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$

- Original importance
 - Outterance: equal weight
 - Prosody: equal weight

scores propagated from utterances weighted by prosodic values



Mathematical Formulation

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$

Utterance node U can get higher score when

 More important prosodic features with higher weights corresponding to utterance U

Mathematical Formulation

$$\begin{cases} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \cdot L_{UP}F_P^{(t)} \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \cdot L_{PU}F_U^{(t)} \end{cases}$$

Utterance node U can get higher score when

- More important prosodic features with higher weights corresponding to utterance U
- Prosody node P can get higher score when
- More important utterances have higher weights corresponding to the prosodic feature P

→ Unsupervised learn important utterances/prosodic features



- *O* Experimental Setup
- Evaluation Metrics
- Results
- Analysis



- - Evaluation Metrics
 - Results
 - Analysis

Experimental Setup

- CMU Speech Meeting Corpus
 - ✓ 10 meetings from 2006/04 2006/06
 - #Speaker: 6 (total), 2-4 (each meeting)
 - 0 WER = 44%
- Reference Summaries
 - Manually labeled by two annotators as three "noteworthiness" level (1-3)
 - REDUNDANCY Ø Extract utterances with level 3 as reference summar
- $\begin{array}{c} \checkmark \quad \text{Parameter Setting} \\ \land \quad \alpha = 0 \\ \end{array} \left\{ \begin{array}{c} F_U^{(t+1)} = (1-\alpha)F_U^{(0)} + \alpha \\ F_P^{(t+1)} = (1-\alpha)F_P^{(0)} + \alpha \end{array} \right\} L_{UP}F_P^{(t)} \\ L_{PU}F_U^{(t)} \end{array}$
 - Extractive summary ratio = 10%, 20%, 30%



- Experimental Setup
- O Evaluation Metrics
 - Results
 - Analysis

Evaluation Metrics

- 🖉 ROUGE
 - ROUGE-1
 - F-measure of <u>matched unigram</u> between extracted summary and reference summary
 - ROUGE-L (Longest Common Subsequence)
 - F-measure of <u>matched LCS</u> between extracted summary and reference summary
- Average Relevance Score
 - Average noteworthiness scores for the extracted utterances



- *i* Experimental Setup
- Evaluation Metrics

Results

Analysis

Baseline

- //>
 Longest
 - the longest utterances based on #tokens
- 🟉 Begin
 - the utterances that appear in the beginning
- / Latent Topic Entropy (LTE)
 - Stimate the "focus" of an utterance
 - Lower topic entropy represents more topically informative
- TFIDF
 - Average TFIDF scores of all words in the utterances

10%

Results



32

For 10% summaries, Begin performs best and proposed performs comparable results

10% & 20% **Results**



33

For 20% summaries, proposed approach outperforms all of the baselines

10% & 20% & 30% **Results**

2.30

2.25

2.20

Longest



34

For 30% summaries, proposed approach outperforms all of the baselines

LTE

TFIDF

Proposed

Begin



- *O* Experimental Setup
- Evaluation Metrics
- Results

Analysis

Analysis

Ø Based on converged scores for prosodic features

36

- Predictive features
 - number of pauses
 - 🧑 min pitch
 - 🧷 avg pitch
 - 🧑 intensity
- Least predictive features
 - the duration time
 - the number of syllables
 - 🙋 the energy



- Two-layer mutually reinforced random walk integrates prosodic knowledge into an unsupervised model for speech summarization
- We show the first attempt at performing unsupervised speech summarization without using lexical information
- Compared to some lexically derived baselines, the proposed approach outperforms all of them but one scenario



Thanks for your attention! © Q&A