

Post-Workshop Research Proposal

Off-Topic Detection: Metaconversation and Small Talk

Robin Stewart (Williams College)

Supervisor: Yang Liu (ICSI & UT-Dallas)

Facilitator: Andrea Danyluk (Williams College)

Example

(Topic: Personal Habits)

...

R: Uh, I'm in college so, like, my drinking is pretty cheap.
Maybe like five bucks a week.

L: Oh, that's not bad.

R: [LAUGH] Yeah, it's pretty cheap.

L: Mhm.

Wait, what college do you go to by the way?

R: University of Illinois.

L: Really, in Champagne?

R: Yeah. In Champagne.

L: Oh, wow.

R: And you live in New York?

L: Yeah.

R: Interesting.

L: Yeah.

But - um - So anyways I guess we're off topic again [LAUGH].

R: [LAUGH] Yeah

L: Um- [LAUGH] um, what were the other things on the list?

Oh yeah, overeating.

See, you know what I heard about, um, overeating is that - or - or just in
general, like, you know, obesity and everything is that - um -

Right now smoking is the number one cause of death in the country.

But then pretty soon it's going at - um - switch over to obesity.

R: Yeah. I've - I've heard about that too.

Example

(Topic: Personal Habits)

- ...
- R: Uh, I'm in college so, like, my drinking is pretty cheap.
■ Maybe like five bucks a week.
- L: Oh, that's not bad.
- R: [LAUGH] Yeah, it's pretty cheap.
- L: Mhm.
Wait, what college do you go to by the way?
- R: University of Illinois.
- L: Really, in Champagne?
- R: Yeah. In Champagne.
- L: Oh, wow.
- R: And you live in New York?
- L: Yeah.
- R: Interesting.
- L: Yeah.
But - um - So anyways I guess we're off topic again [LAUGH].
- R: [LAUGH] Yeah
- L: Um- [LAUGH] um, what were the other things on the list?
Oh yeah, overeating.
See, you know what I heard about, um, overeating is that - or - or just in
general, like, you know, obesity and everything is that - um -
Right now smoking is the number one cause of death in the country.
But then pretty soon it's going at - um - switch over to obesity.
- R: Yeah. I've - I've heard about that too.

On topic

Example

(Topic: Personal Habits)

- ...
- R: Uh, I'm in college so, like, my drinking is pretty cheap.
■ Maybe like five bucks a week.
- L: Oh, that's not bad.
- R: [LAUGH] Yeah, it's pretty cheap.
- L: Mhm.
- Wait, what college do you go to by the way?
- R: University of Illinois.
- L: Really, in Champagne?
- R: Yeah. In Champagne.
- L: Oh, wow.
- R: And you live in New York?
- L: Yeah.
- R: Interesting.
- L: Yeah.
- But - um - So anyways I guess we're off topic again [LAUGH].
- R: [LAUGH] Yeah
- L: Um- [LAUGH] um, what were the other things on the list?
Oh yeah, overeating.
See, you know what I heard about, um, overeating is that - or - or just in
general, like, you know, obesity and everything is that - um -
Right now smoking is the number one cause of death in the country.
But then pretty soon it's going at - um - switch over to obesity.
- R: Yeah. I've - I've heard about that too.

On topic

Small talk

Example

(Topic: Personal Habits)

- ...
- R: Uh, I'm in college so, like, my drinking is pretty cheap.
■ Maybe like five bucks a week.
- L: Oh, that's not bad.
- R: [LAUGH] Yeah, it's pretty cheap.
- L: Mhm.
- Wait, what college do you go to by the way?
- R: University of Illinois.
- L: Really, in Champagne?
- R: Yeah. In Champagne.
- L: Oh, wow.
- R: And you live in New York?
- L: Yeah.
- R: Interesting.
- L: Yeah.
- But - um - So anyways I guess we're off topic again [LAUGH].
- R: [LAUGH] Yeah
- L: Um- [LAUGH] um, what were the other things on the list?
- Oh yeah, overeating.
- See, you know what I heard about, um, overeating is that - or - or just in general, like, you know, obesity and everything is that - um - Right now smoking is the number one cause of death in the country. But then pretty soon it's going at - um - switch over to obesity.
- R: Yeah. I've - I've heard about that too.

On topic

Small talk

Meta-
conversation

Example

(Topic: Personal Habits)

- ...
- R: Uh, I'm in college so, like, my drinking is pretty cheap.
■ Maybe like five bucks a week.
- L: Oh, that's not bad. **On topic**
- R: [LAUGH] Yeah, it's pretty cheap.
- L: Mhm.
- Wait, what college do you go to by the way?
- R: University of Illinois.
- L: Really, in Champagne?
- R: Yeah. In Champagne.
- L: Oh, wow. **Small talk**
- R: And you live in New York?
- L: Yeah.
- R: Interesting.
- L: Yeah.
- But - um - So anyways I guess we're off topic again [LAUGH].
- R: [LAUGH] Yeah **Meta-conversation**
- L: Um- [LAUGH] um, what were the other things on the list?
- Oh yeah, overeating.
- See, you know what I heard about, um, overeating is that - or - or just in
■ general, like, you know, obesity and everything is that - um -
■ Right now smoking is the number one cause of death in the country.
■ But then pretty soon it's going at - um - switch over to obesity.
- R: Yeah. I've - I've heard about that too.

Definitions

- **Small Talk:** Conversation that is not related to or not contributing to the assigned topic.
- **Metaconversation:** Conversation about the assigned topic, the task, and the phone call.
- **On-Topic:** Everything else.

Definitions

- **Small Talk:** Conversation that is not related to or not contributing to the assigned topic.
- **Metaconversation:** Conversation about the assigned topic, the task, and the phone call.
- **On-Topic:** Everything else.

**Goal: Automatically classify sentences
in recorded telephone speech**

Motivations

- Just as “edit” regions can be removed to improve parsing, “small talk” regions could be removed to improve **information extraction**.
(someone searching for weather information shouldn’t get audio clips of “so, how’s the weather?”)
- Both metaconversation and small talk regions may help to identify changes in topic for **new topic detection**.
 - Meta: “Now we’re supposed to talk about US public schools...”
 - Small talk: fills the gap between more-significant topics

Motivations

- Can also be applied to:
 - Meeting corpora
 - (“You should have seen the traffic today..”)
 - (“Let’s talk about the quarterly revenue report.”)
 - Broadcast news
 - (“I’m glad I’m safe inside the studio!”)
 - (“We now go live to Jim for an update.”)
 - Surreptitiously recorded telephone conversations
 - (“We had mac and cheese again tonight”)
 - (“So I was calling you because...”)
 - Lectures, etc.

Related Work

- “Off-talk” detection for human-machine interaction
(University of Munich)
 - “Oh, I have to click on that with the mouse”
- Social dialogue with conversational agents
(Northwestern, MIT Media Lab)
 - Generating and responding to small talk with human users
- **NIST** Topic Detection and Tracking benchmark tasks
(1998-2004)
 - Supervised and unsupervised classification techniques
 - Evaluation metrics

Proposal

- Weakly supervised classification of sentence units
- Local classification techniques:
 - Naive Bayes (“bag of words”) classifier
 - Maximum-entropy (MaxEnt) classifier
 - Support Vector Machine (SVM) classifier
- Sequence decoding:
 - Hidden Markov Model (HMM)
 - Conditional Random Field (CRF)
- Train the classifier on a small set, use it to automatically “annotate” a much larger corpus, then iteratively re-train on the larger corpus

Feature Extraction

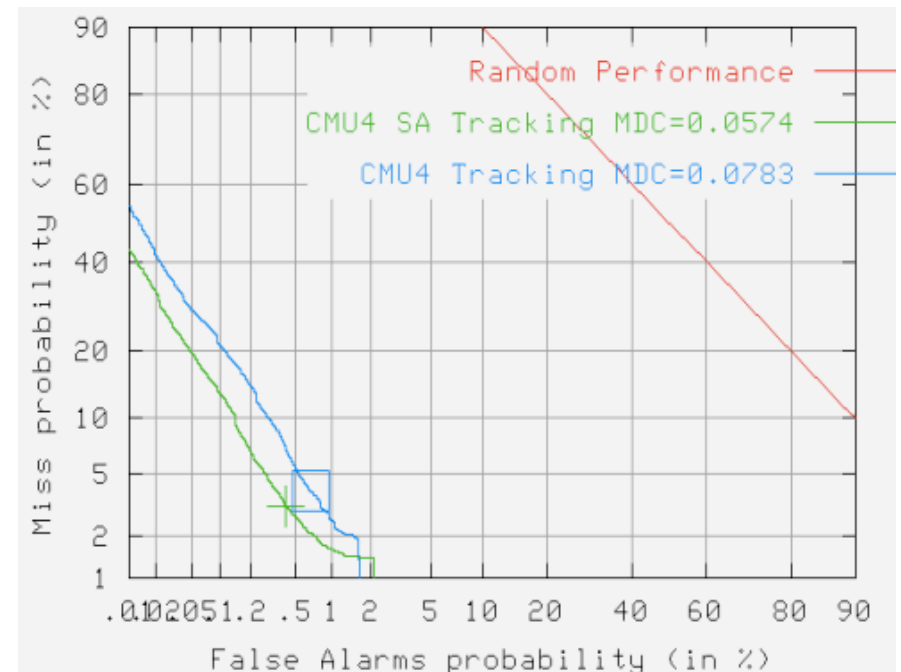
- Similar to our metadata reranking system
- Features which might prove useful:
 - Bigram or trigram language model
 - Key words such as filled pauses and discourse markers
 - Speaker changes and overlap
 - Duration of pauses
 - Frequency of awkward laughs
 - Etc.
- Easily extracted from our corpus

Annotation

- I've fully annotated 5 conversations, and looked over many others.
 - The time it takes to annotate is at *most* twice the length of the conversation.
 - We expect high annotator agreement.
- Weakly supervised learning techniques minimize the amount of annotation needed.
 - Need ~ 3 hours of training data (30 conversations) and another 3 hours for evaluation
 - 2 annotators for each conversation, plus a “tiebreaker”
 - ~ 30 hours of work = feasible
- Create annotation spec

Evaluation

- Accuracy - % of sentences correctly identified
- NIST metrics for Detection Evaluation
- Detection Error Tradeoff curves
 - uses probability estimates to graph the tradeoff between misses and false alarms



We will find out:

- How well can off-topic regions be detected using standard machine learning techniques?
- How much training data is needed?
- Which machine learning algorithms work well?
- What features are effective?
- What is the effect of ASR and MDE errors?
- How well do ASR and MDE systems perform in on-topic vs. off-topic regions?

Conclusion

- **Useful**
 - Improve Information Extraction and New Topic Detection
- **Generalizable**
 - Meetings, Broadcast News, Phone Calls, ...
- **Feasible**
 - Builds on NIST TDT benchmark tasks
 - Small amount of annotation