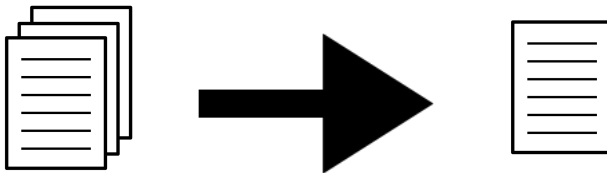


Automatic Detection of Linguistic Quality Violations

Jonathan Oberländer

Bachelor Thesis Defense
Universität des Saarlandes
21.08.2014

Automatic Summarization



- ▶ **Single-Document:** One document
- ▶ **Multi-Document:** Multiple documents on the same topic

Automatic Summarization

- ▶ **Single-Document:** One document
- ▶ **Multi-Document:** Multiple documents on the same topic
- ▶ **Abstractive:** Internal semantic representation + generation
- ▶ **Extractive:** New summary from source sentences

Automatic Summarization

- ▶ **Single-Document:** One document
- ▶ **Multi-Document:** Multiple documents on the same topic
- ▶ **Abstractive:** Internal semantic representation + generation
- ▶ **Extractive:** New summary from source sentences

	Single-document	Multi-document
Abstractive		
Extractive		

Summarization systems should produce coherent and grammatical output.

Summarization systems **don't produce coherent and grammatical output.**

Why?

- ▶ It's hard.

Summarization systems **don't produce coherent and grammatical output.**

Why?

- ▶ It's hard.
- ▶ Evaluation: content, information density

Summarization systems **don't produce coherent and grammatical output.**

Why?

- ▶ It's hard.
- ▶ Evaluation: content, information density

⇒ LQVCorpus (Friedrich et al., 2014)

Annotated results of TAC 2011 Guided Summarization task
(Owczarzak and Dang, 2011)

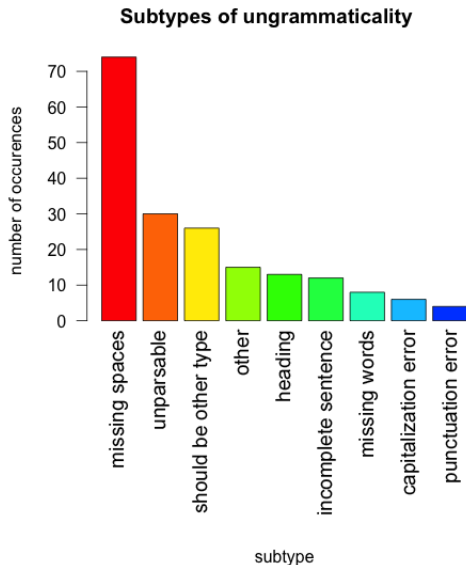
Annotated results of TAC 2011 Guided Summarization task
(Owczarzak and Dang, 2011)

- ▶ Entity level:
 - ▶ FM-EXPL, SM+EXPL
 - ▶ DNP-REF, INP+REF
 - ▶ PRN-ANT, PRN+MISLA
 - ▶ ACR-EXPL

Annotated results of TAC 2011 Guided Summarization task
(Owczarzak and Dang, 2011)

- ▶ Entity level:
 - ▶ FM-EXPL, SM+EXPL
 - ▶ DNP-REF, INP+REF
 - ▶ PRN-ANT, PRN+MISLA
 - ▶ ACR-EXPL
- ▶ Clause level:
 - ▶ **incomplete sentence (INCOMPLSN)**
 - ▶ **inclusion of datelines (INCLDATE)**
 - ▶ **other ungrammatical form (OTHRUNGR)**
 - ▶ no semantic relatedness (NOSEMREL)
 - ▶ **redundant information (REDUNINF)**
 - ▶ no discourse relation (NODISREL)

Detecting Ungrammaticality (OTHRUNGR)



detecting ungrammaticality:subtypes

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

- Friedrich, A., Valeeva, M., and Palmer, A. (2014). Lqvsumm: A corpus of linguistic quality violations in multi-document summarization.
- Owczarzak, K. and Dang, H. T. (2011). Overview of the tac 2011 summarization track: Guided task and aesop task. In *Proceedings of the Text Analysis Conference (TAC 2011)*, Gaithersburg, Maryland, USA, November.