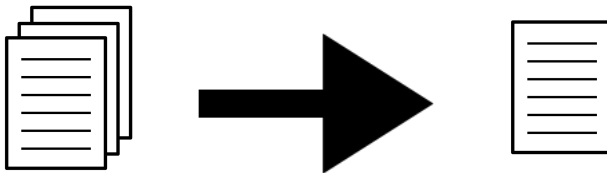


Automatic Detection of Linguistic Quality Violations

Jonathan Oberländer

Bachelor Thesis Defense
Universität des Saarlandes
21.08.2014

Automatic Summarization



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)

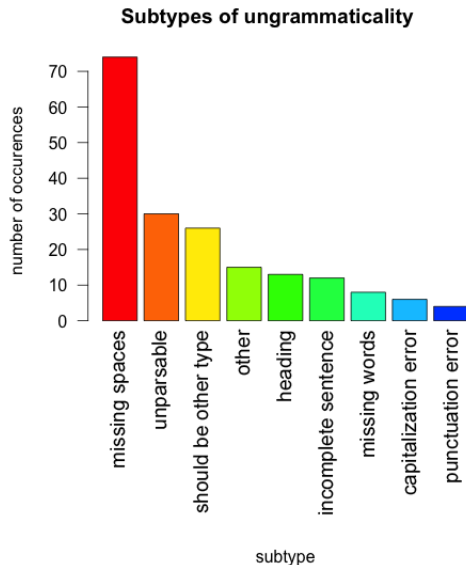
Annotated results of TAC 2011 Guided Summarization task
(Owczarzak and Dang, 2011)

- ▶ 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)

Ungrammaticality (OTHRUNGR+INCOMPLSN)



Detecting missing spaces

*“A strong earthquake measuring 7.8 magnitude struck **Wenchuancounty** of Sichuan Province on Monday, leaving at least **12,000people** died and thousands more injured.”*

*“Virginia Tech reported a campus shooting Monday and told **studentsto** stay inside their residences and away from windows.”*

*“A gunman opened fire on classrooms at Virginia Tech University **onMonday** morning, killing at least 30 people before turning his **gunon** himself in the bloodiest school shooting in US history.”*

Unknown Tokens

Idea:

Sentence contains violation iff any word \notin known tokens

Idea:

Sentence contains violation iff any word \notin known tokens

known tokens?

- ▶ Source documents available? \rightarrow all tokens in source documents = **UnknownTokens**_{source}
- ▶ Otherwise \rightarrow **UnknownTokens**_{general}

- ▶ Tokens from (parts of) Gigaword = **UnknownTokens_{gw}**

- ▶ Tokens from (parts of) Gigaword = **UnknownTokens_{gw}**
- ▶ + Heuristics (Capitalized words) = **UnknownTokens_{gw+heur}**

- ▶ Tokens from (parts of) Gigaword = **UnknownTokens_{gw}**
- ▶ + Heuristics (Capitalized words) = **UnknownTokens_{gw+heur}**
- ▶ + NER (Finkel et al., 2005) = **UnknownTokens_{gw+heur+ner}**

- ▶ Tokens from (parts of) Gigaword = **UnknownTokens_{gw}**
- ▶ + Heuristics (Capitalized words) = **UnknownTokens_{gw+heur}**
- ▶ + NER (Finkel et al., 2005) = **UnknownTokens_{gw+heur+ner}**
- ▶ + Wikipedia = **UnknownTokens_{gw+heur+ner+wiki}**

UnknownTokens: Evaluation

	Missing spaces			No missing spaces		
	P	R	F	P	R	F
Baseline	0.0	0.0	0.0	94.8	100	97.3
UT_{gw}	15.0	98.7	26.0	99.9	69.1	81.7
UT_{gw+heur}	30.5	97.3	46.5	99.8	87.8	93.4
UT_{gw+heur+ner}	35.5	97.3	52.0	99.8	90.3	94.8
UT_{gw+heur+ner+wiki}	70.3	96.0	81.2	99.8	97.8	98.8
UT_{source}	95.9	94.6	95.2	99.7	99.7	99.7

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

foo bar

- Finkel, J. R., Grenager, T., and Manning, C. (2005). Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 363–370. Association for Computational Linguistics.
- 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.