



TOPIC MODELING ON SUMMARIZED NEWS ARTICLES

GIOVANNI NOÈ - 881765

FRANCESCO VOLPI GHIRARDINI - 933735

TEXT MINING AND SEARCH

INTRODUCTION



AIM OF THE PROJECT

This project investigates the interaction between topic modeling & extractive text summarization

DATASET

CNN/Daily Mail, news articles (subset of 25,000 documents from each)



DATA CLEANING AND EXPLORATION

CLEANING AND EXPLORATION




CLEANING

- Headers and footers are removed
- Web links are removed
- Escape characters are removed
- Missing full stops are added
- Data is shuffled
- Documents and summaries are separated

EXPLORATION

Both the documents and summaries have a skewed distribution for the sentence count, the extreme values are removed. This distribution was the same in the subset and in the original data

The background features a light gray pattern of horizontal lines and rectangular blocks, resembling a newspaper layout. A large, hand-drawn red oval is centered on the page, framing the title text. The title is written in a bold, black, serif font.

TOPIC MODELING ON FULL ARTICLES

TEXT PROCESSING



TEXT PROCESSING

- Strip non-alphabetic characters, lowercasing, tokenization, lemmatization
- Remove stop words (both general and domain-specific)
- Build dictionary (types in 0.1–80% of docs, max 20,000 types)

TEXT REPRESENTATION

- TF (Term Frequency): used by LDA, normalized for pLSA
- TF-IDF: used by LSA
- Word embeddings (BERT): used by BERTopic

THE MODELS



LSA

- Uses SVD on TF-IDF matrix
- Optimal topics*: 10

PLSA

- Probabilistic model on TF matrix
- Optimal topics*: 30

LDA

- Probabilistic model on normalized TF matrix
- Optimal topics*: 20

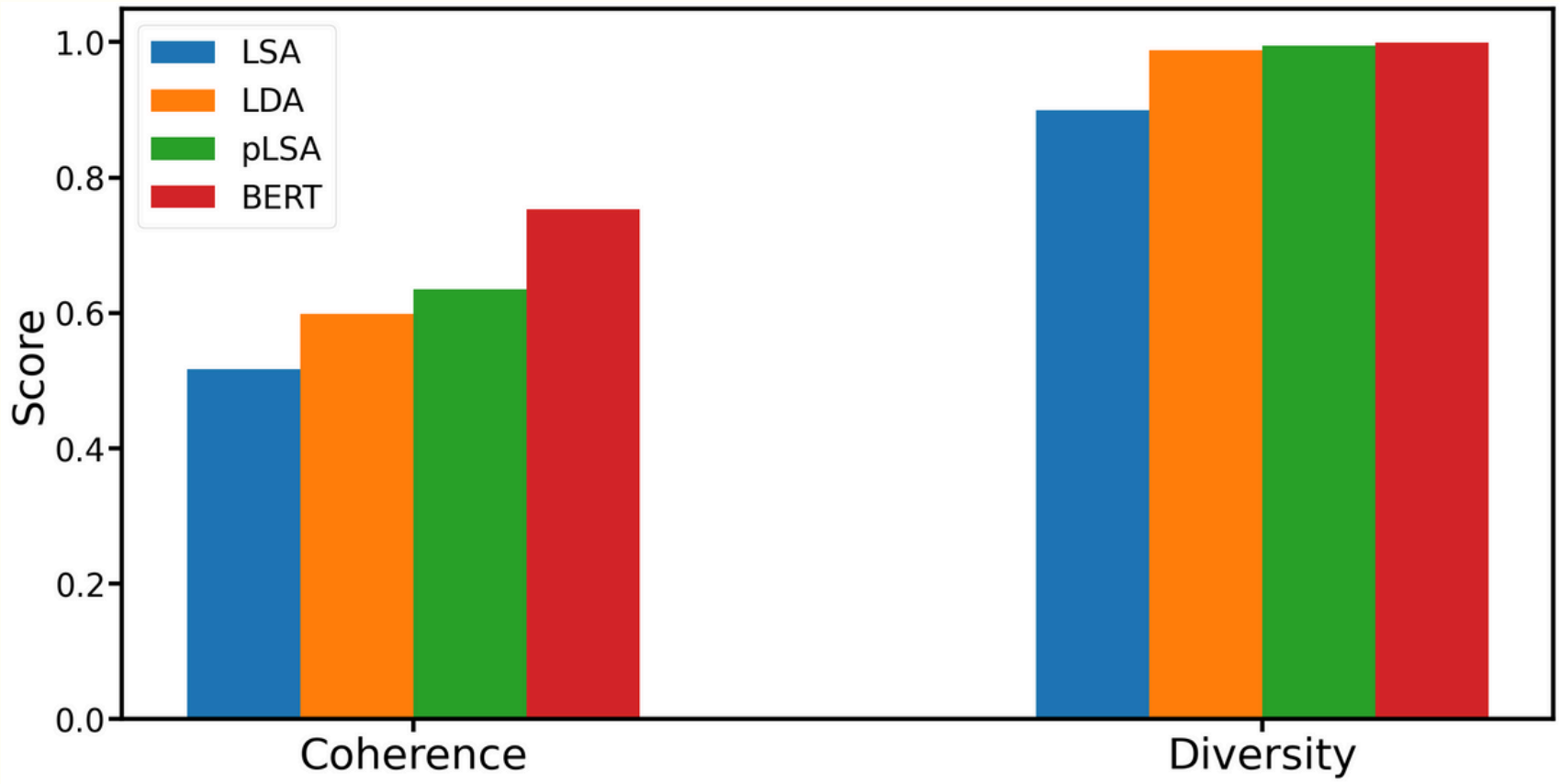
BERTOPIC

- Uses BERT embeddings
- No fine-tuning required (automatically determines topics)

*Estimated through tuning tests

EVALUATION

METRICS



VISUAL INSPECTION

Visual inspection shows that LSA performs the worst, while pLSA, LDA, and BERTopic produce more meaningful topics

TIME COMPARISON

Model	Tokenization	Tuning	Fitting	Total
LSA	3m	37m	3s	40m
LDA	3m	77m	7m	88m
pLSA	3m	122m	25m	150m
BERTopic	---	---	41m	41m



TEXT SUMMARIZATION

THE METHOD



KEY POINTS

- Extractive summarization
- 5 sentences per document kept
- 85–15% train–test split

UPPER BOUND ESTIMATE

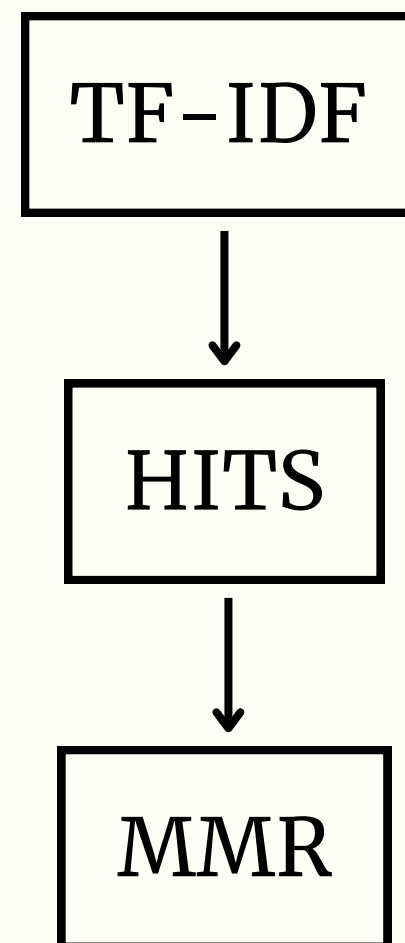
MMR with ROUGE scores
and cosine similarity on
BERT embeddings

BASELINE

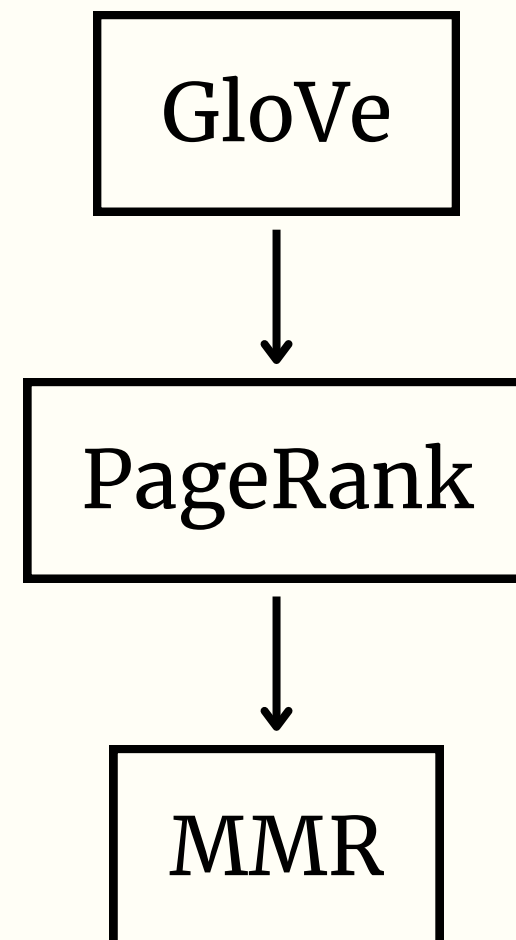
Random summarizer

THE MODELS

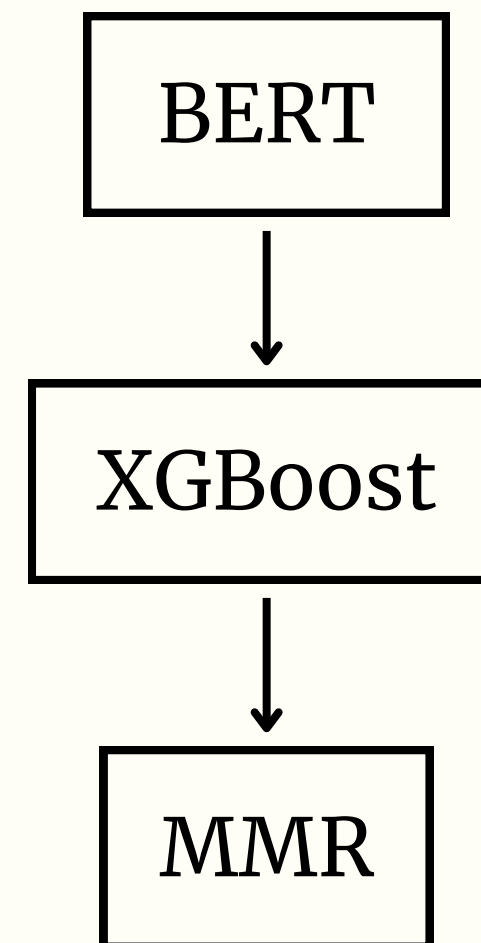
FIRST ALGORITHM



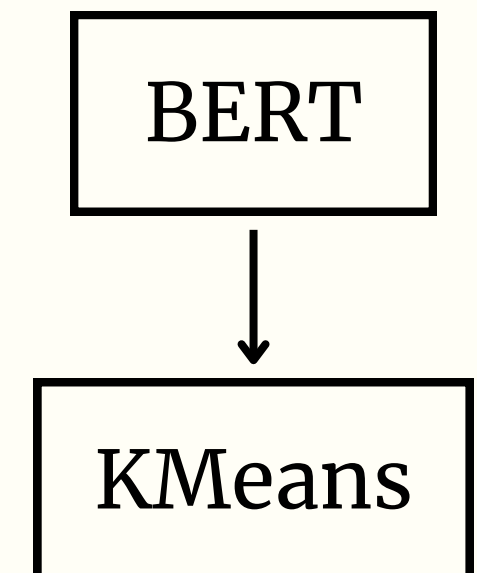
SECOND ALGORITHM



THIRD ALGORITHM

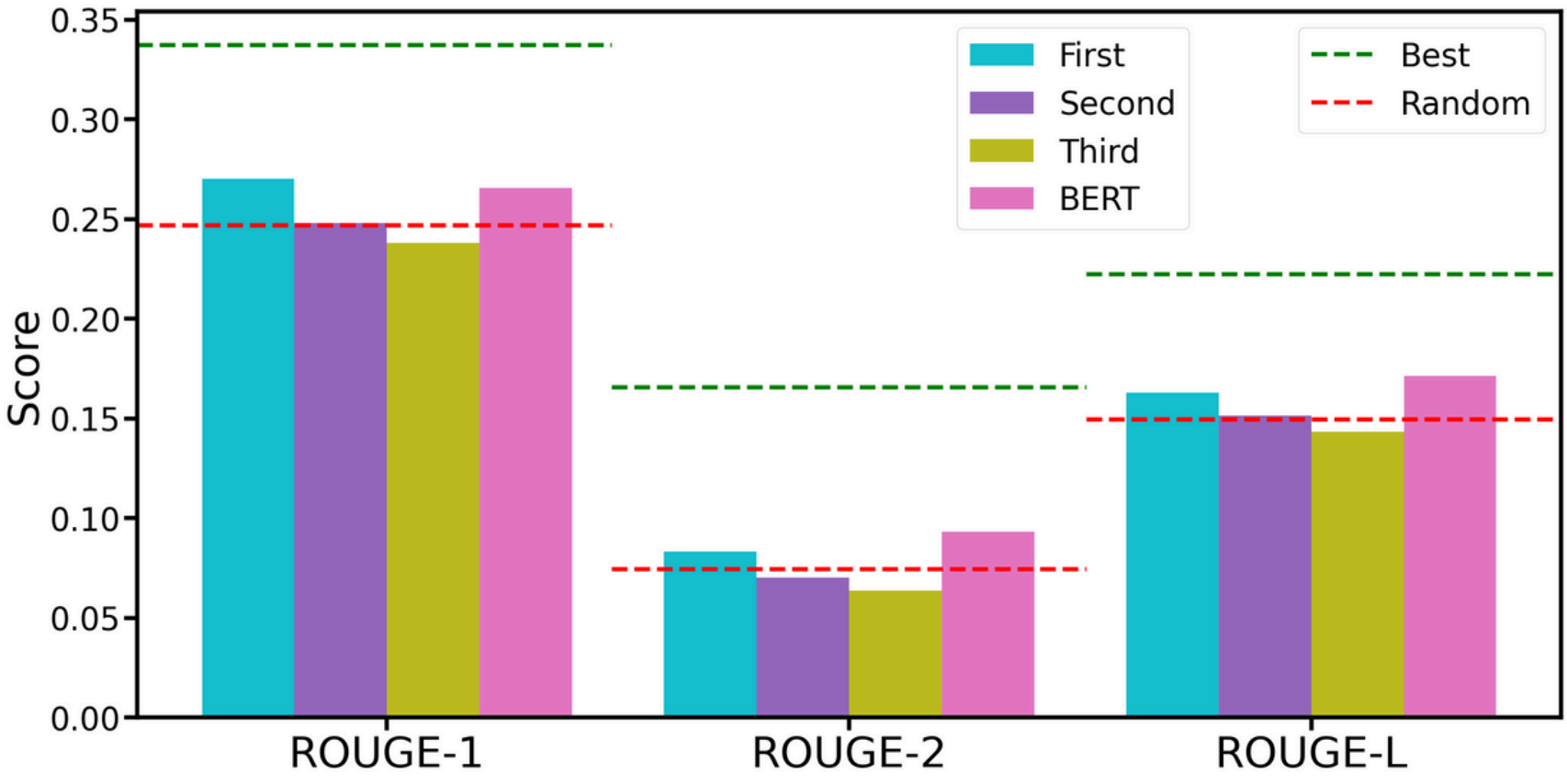


BERT SUMMARIZER



EVALUATION

METRICS



TIME COMPARISON

Model	Embedding	Training	Test	Total
First	---	---	1m	1m
Second	---	---	2m	2m
Third	16m	31m	12s	47m
BERT	---	---	22m	22m



TOPIC MODELING ON SUMMARIZED TEXT

THE METHOD



KEY POINTS

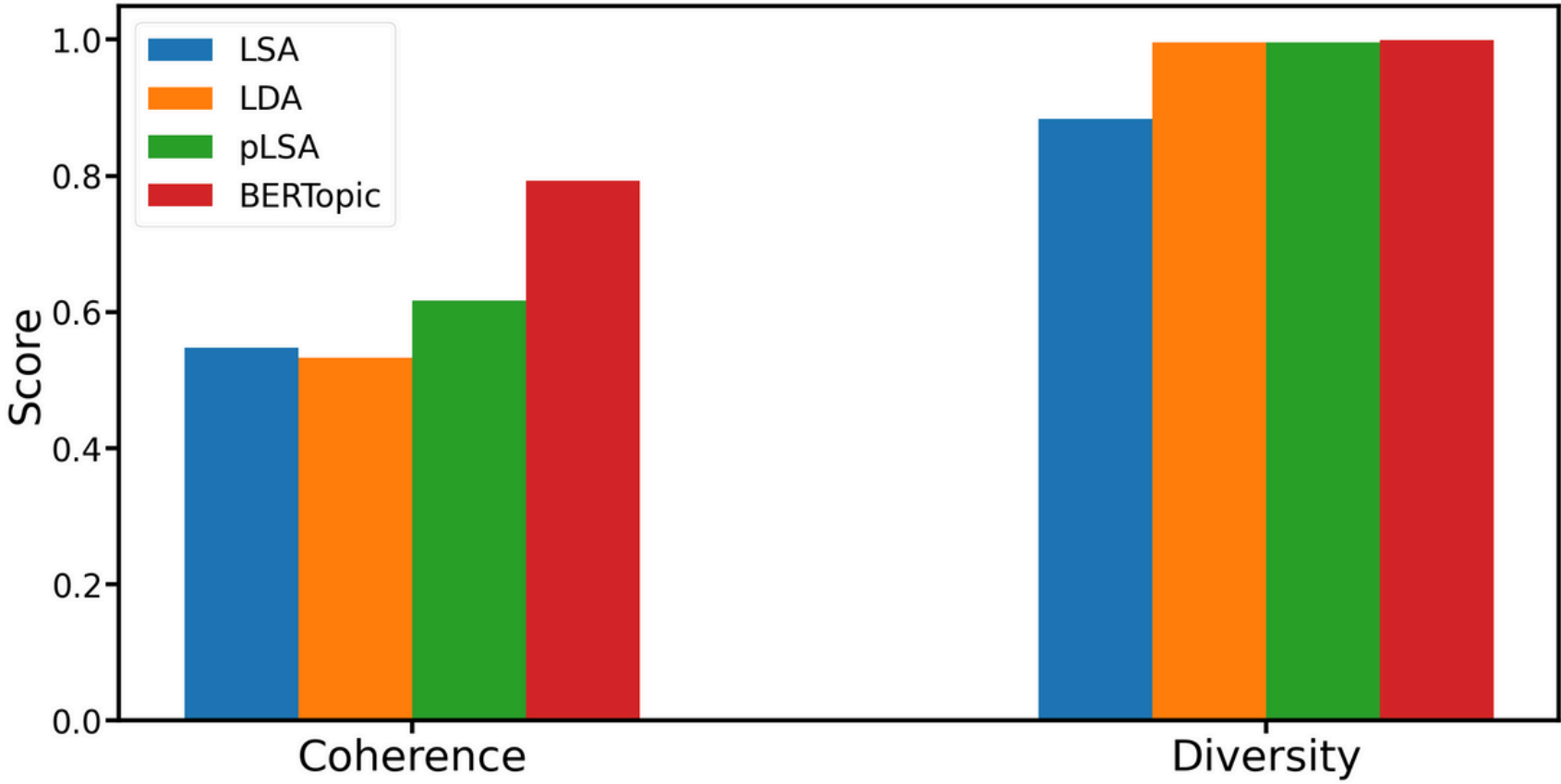
- Apply topic modeling to the summarized corpus
- Use the same pipeline as for original articles for a fair comparison

PROS AND CONS

- Potential benefits:
 - Less noise and redundancy
 - Lower computational costs
- Possible drawback:
 - Loss of context

EVALUATION

METRICS



VISUAL INSPECTION

By visual inspection BERTopic remains the best-performing model, producing the most coherent and meaningful topics, while LDA performs much worse, pLSA slightly worse, and LSA only slightly better

TIME COMPARISON

Model	Tokenization	Tuning	Fitting	Total
LSA	44s	47s	1s	2m
LDA	44s	8m	3m	12m
pLSA	44s	7m	3m	11m
BERTopic	---	---	39m	39m



CONCLUSIONS

CONCLUSIONS



GENERAL CONSIDERATIONS

Both the topic modeling and text summarization experiments worked, with models employing BERT embeddings performing better than traditional algorithms on both tasks

ON THE AIM OF THE PROJECT

The topic modeling on the summarized text gave mixed results, with some models achieving better topic quality while others greatly decreased their processing time

**THANKS FOR
YOUR ATTENTION**

