

Deep Neural Ranking Models for Argument Retrieval

Master's Thesis by Saeed Entezari

Referees: Prof. Stein, PD. Dr. Jakoby

Supervisor: Michael Völske

Faculty of Media
Bauhaus Universität Weimar

September 16, 2020

Outline

Introduction

Arguments

Ranking Task

Dataset and Models

Experiments and Results

Conclusion

Why Argument Retrieval

- Different types of opinions toward controversial topics
- Getting an overview of every opinion is an exhaustive and time consuming task
- Automated decision making
- Opinion Summarization

Argument components

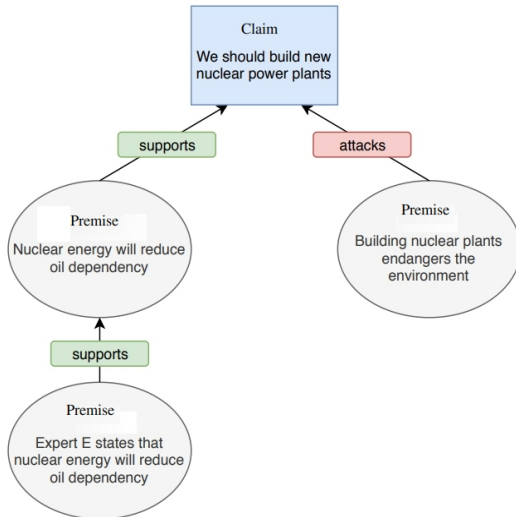


Figure: The relation between the argument units ([Dumani(2019)])

Outline

Introduction

Arguments

Ranking Task

Dataset and Models

Experiments and Results

Conclusion

Outline

Introduction

Dataset and Models

Touché Shared Task Dataset

Preprocessing and Visualisation

Query Relevance Information

Training and Validation sets

Deep Neural Ranking Models

Experiments and Results

Conclusion

Args.me Corpus

387740 annotated arguments in total from crawling 4 debate portals (json format):

- Debatewise (14000 arguments)
- IDebate.org (13000 arguments)
- Debatepedia (21000 arguments)
- Debate.org (338000 arguments)

Information for each argument:

- unique ID
- claim
- premise
- source of crawling
- time of crawling
- stance of premise regard to claim

Outline

Introduction

Dataset and Models

Touché Shared Task Dataset

Preprocessing and Visualisation

Query Relevance Information

Training and Validation sets

Deep Neural Ranking Models

Experiments and Results

Conclusion

Preprocessing and Visualisation: Claims

- Forming normalized claims
 - punctuation removal and case sensitivity
 - stop words removal
- Visualization and Statics
 - 66473 unique claims
 - 29970 unique tokens

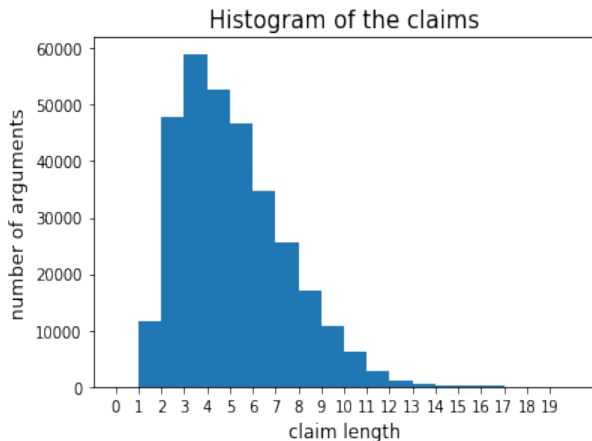


Figure: Histogram of the unique claims based on the number of tokens.

Preprocessing and Visualisation: Premises

- Tokenizing punctuation
 - for static embedding: god exists. ⇒ god exists <PERIOD>
 - for contextualized embedding is not required!
- Removing consecutive repetitive tokens
 - !!!!!!! ⇒ <EXCLAMATIONMARK>
 - yes yes yes ⇒ yes
- Mapping digits to words
 - 95 ⇒ ninety-five
- Removing the URLs
 - <http://example.net/achiever.html?boy=armyauthority=beginner>

Preprocessing and Visualisation: Premises

- Statistics of the premises:
 - vocabulary size: 586796
 - 85% of the premises have the length of less than 200 words
- Arguments with the premise length of less than 15 tokens are removed

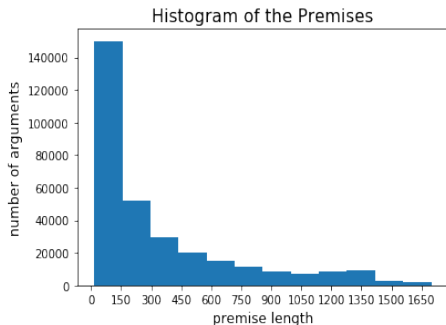


Figure: Histogram of the premises based on their length (number of tokens separated by white space)

Outline

Introduction

Dataset and Models

Touché Shared Task Dataset

Preprocessing and Visualisation

Query Relevance Information

Training and Validation sets

Deep Neural Ranking Models

Experiments and Results

Conclusion

Learning to Rank

- Learning goal: related documents over the unrelated ones
- Pairwise hinge cost function
- Relevant and irrelevant Query-Document pairs are required and are **missing in the corpus**
- A model to produce the similarity scores (We use Deep ranking models)

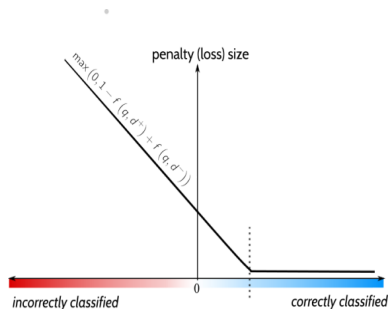


Figure: Hinge as a pairwise cost function

Binary Query Relevance Generation

RQ.1: Useful dataset for ad-hoc task

- Distant Supervision Approach
 - Claims \Rightarrow Queries
 - Premises \Rightarrow Related Documents
- Unrelated premise for each query
 - qrel files contain also unrelated query-document pairs
 - similarity measure: fuzzy similarity
 - premise of an unrelated claims could be an unrelated document to our claims
- A binary query relevance is formed \Rightarrow Exploitation of deep ranking models in the context of argument retrieval is possible now!

Dataset Ready for Ad-hoc Task

Data collection ready for the ad-hoc task (for static and contextualized embedding) with the following columns:

Important Note: Different arguments may have same claims and different premises

id	claim	norm-claim	premise	unrelated id	unrelated premise
arg1
arg2

Outline

Introduction

Dataset and Models

Touché Shared Task Dataset

Preprocessing and Visualisation

Query Relevance Information

Training and Validation sets

Deep Neural Ranking Models

Experiments and Results

Conclusion

Training and Validation Sets

- Training set: 312248 arguments with one unrelated documents each
- Validation set: 4885 arguments: 20 unrelated documents each

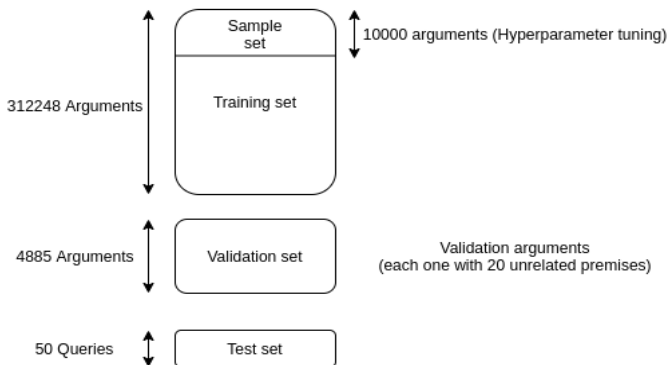


Figure: Different datasets and their number of arguments

Validation Arguments

RQ.1: Forming an appropriate training and validation dataset

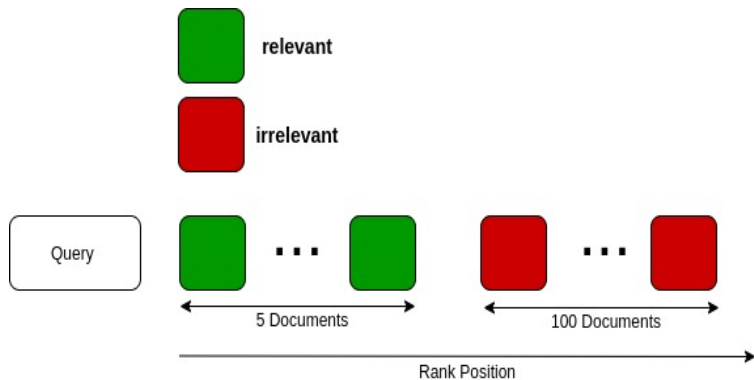


Figure: An ideal ranking for a validation query

Outline

Introduction

Dataset and Models

Touché Shared Task Dataset
Preprocessing and Visualisation
Query Relevance Information
Training and Validation sets
Deep Neural Ranking Models

Experiments and Results

Conclusion

Siamese Network

Model	type	embedding	re-rank
GRU	rep	static	yes
DRMM	int	static	yes
KNRM	int	static	yes
CKNRM	int	static	yes
Vanilla BERT	int	contx	yes
DRMM BERT	int	contx	yes
KNRM BERT	int	contx	yes
SNRM	rep	static	no

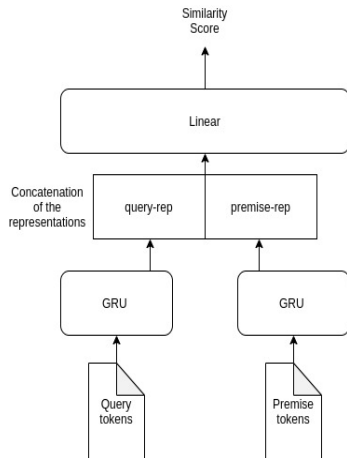


Figure: Similarity scores using recurrent neural network

KNRM: Kernel-based Neural Ranking Model

Model	type	embedding	re-rank
GRU	rep	static	yes
DRMM	int	static	yes
KNRM	int	static	yes
CKNRM	int	static	yes
Vanilla BERT	int	contx	yes
DRMM BERT	int	contx	yes
KNRM BERT	int	contx	yes
SNRM	rep	static	no

- Another strategy for encoding the input pair interaction
- Forming translation matrix: elements are the cos similarity of the term embedding
- Applying the RBF as the kernels and forming the input features for fully connected network
- A linear layer learns the score similarity of the input pairs

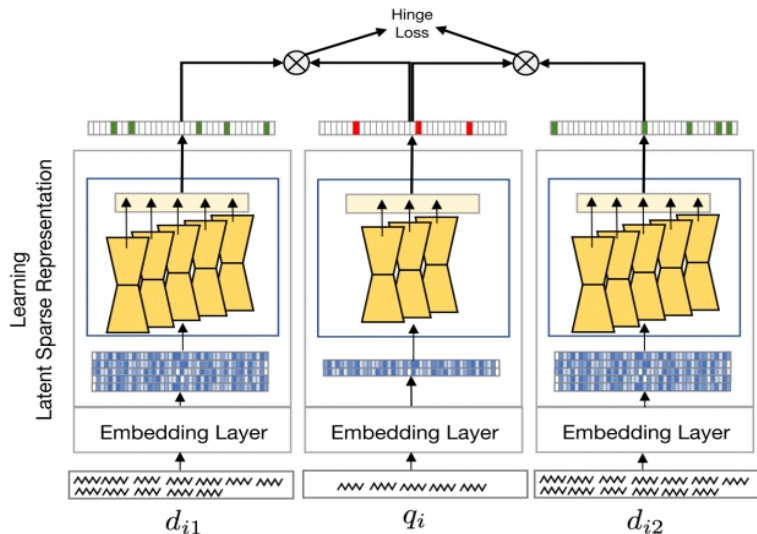


Figure: Training process of SNRM
 ([Zamani et al.(2018)Zamani, Dehghani, Croft, Learned-Miller, and Kamps])

Outline

Introduction

Dataset and Models

Experiments and Results

Training and Validation Phase

Test Phase

Model Output Analysis

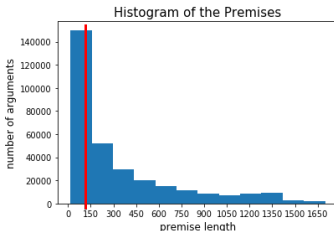
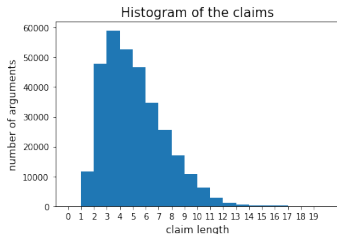
Aggregation

Test Results

Conclusion

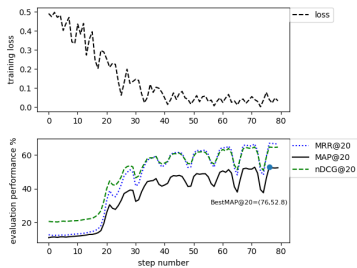
Train and Validation Phase

- 10000 sample data for hyper-parameter tuning and debug the codes so that the models run correctly
- Query length: 20 and Document length: 100

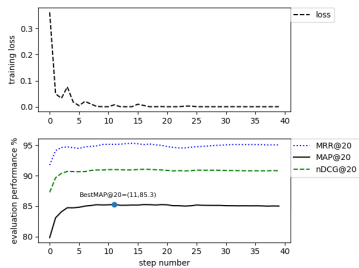


- Each batch: 32 argument
- Train the models
 - static embedding: 10 epochs
 - contextualized embedding: 5 epochs
- Validation run for 8 times within a training epoch
 - Top 20 hits among the 105 validation documents for each query
 - Validation metrics: MRR@20, MAP@20, and nDCG@20
 - For binary qrel: MAP@20 more stable validation scores

Sample Training and Validation Curves



(a) DRMM



(b) Vanilla BERT

Outline

Introduction

Dataset and Models

Experiments and Results

Training and Validation Phase

Test Phase

Model Output Analysis

Aggregation

Test Results

Conclusion

Re-ranking Candidate Arguments

- 50 test queries provided in the Touché task
- 100 first hits by each model for each test query is saved

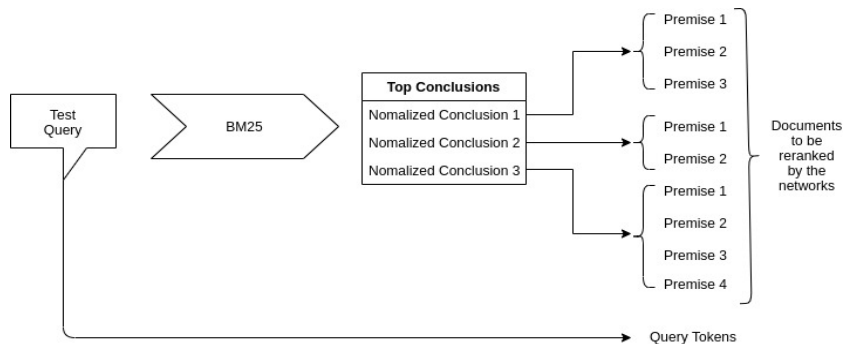


Figure: Candidate documents to be re-ranked in the test phase

Inference in SNRM

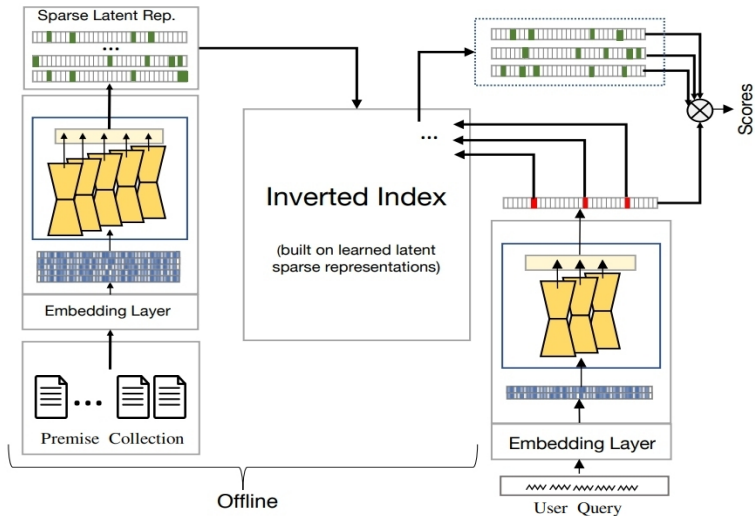


Figure: Document retrieval process
([Zamani et al.(2018)Zamani, Dehghani, Croft, Learned-Miller, and Kamps])

RQ3. Aggregation Strategy

- Why to aggregate?
 - Performance improvement
 - Aggregation of the different model principles
- How to aggregate?
 - Using regression between the *normalized* model scores
- What do we need to know before the regression?
 - How diverse the model results are.
 - Models with outlier results. Assumption: Outlier results belong to weak models!

Outline

Introduction

Dataset and Models

Experiments and Results

Training and Validation Phase

Test Phase

Model Output Analysis

Aggregation

Test Results

Conclusion

Model Output Analysis

- The model results are vectors: retrieved documents as dimensions and scores are the values in each dimension
 - retrieved documents are not the same for the models
- Jaccard and Spearman Coefficients for measuring the similarity of the ranking results
 - Jaccard: portion of the documents in common
 - Spearman: correlation of the ranking scores of the common documents
- The average of the coefficients over 50 test queries are calculated

Jaccard Coefficient as Similarity Measure

Jaccard: portion of the documents in common $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$

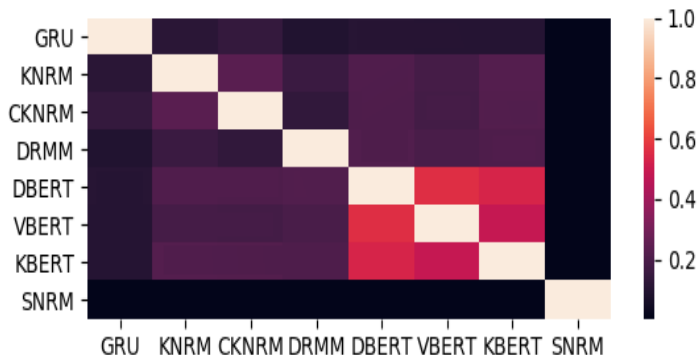


Figure: The heat map of the Jaccard coefficient for the 50 test queries

Outline

Introduction

Dataset and Models

Experiments and Results

Training and Validation Phase

Test Phase

Model Output Analysis

Aggregation

Test Results

Conclusion

Linear Regression as an Aggression Strategy

- We assume SNRM results as outlier data (Based on the similarity results)
- Regression model is trained on validation set (1 related and 1 unrelated document)
 - $2 * 4885$ data points for training the regression with the dimension of 7
- union of the retrieved documents by models are scored by the regression model
 - If a model did not retrieve a document, 0 is assigned to the corresponding dimension

Outline

Introduction

Dataset and Models

Experiments and Results

Training and Validation Phase

Test Phase

Model Output Analysis

Aggregation

Test Results

Conclusion

Argument Quality Dimensions

- **Logical**: acceptable and relevant premises to the arguments
- **Rhetorical**: the ability of convince the audiences
- **Dialectical** (utility): the ones by which a stance can be built
- Our concern in this study: Focusing on the **Logical** aspect

Test Results

- nDCG@5 score is calculated over the retrieved arguments
- Manually annotation is done by human annotators based on the different quality dimensions of the arguments

Model	type	embedding	re-rank	MAP@20	nDCG@5
GRU	rep	static	yes	0.241	x
DRMM	int	static	yes	0.528	x
KNRM	int	static	yes	0.727	0.684
CKNRM	int	static	yes	0.733	x
Vanilla BERT	int	contx	yes	0.88	0.404
DRMM BERT	int	contx	yes	0.881	0.371
KNRM BERT	int	contx	yes	0.902	0.319
SNRM	rep	static	no	0.701	x
Aggregation	x	x	x	x	0.372

Test Results

- KNRM (our best performing model) ranked 4th in the competition
- Most of the competitors got less score than the baseline (Dirichlet LM)
 - Argument retrieval meeting the quality dimensions is not an easy task
- Validation results and test results were not correlated
 - related arguments \neq good arguments (meeting the argument quality dimensions)
 - Relevance is a required but not enough condition for a good argument
- Interaction-focused network outperformed representation-focused networks
 - Representation focused networks' results are not shown in the table
- Aggregation model has been trained on the validation set and its MAP@20 score on the validation set is useless.

Outline

Introduction

Dataset and Models

Experiments and Results

Conclusion

Summary

Future Works

Outline

Introduction

Dataset and Models

Experiments and Results

Conclusion

Summary

Future Works

What's next...

- Providing a concrete mathematical definition of the argument quality dimensions to be included in the cost function of the networks
- Working on strategies to map the interaction of the input pairs
- Devising more intuitive structures to create sparse representation for end-to-end models

Evaluation Metrics: Mean Reciprocal Rank (MRR)

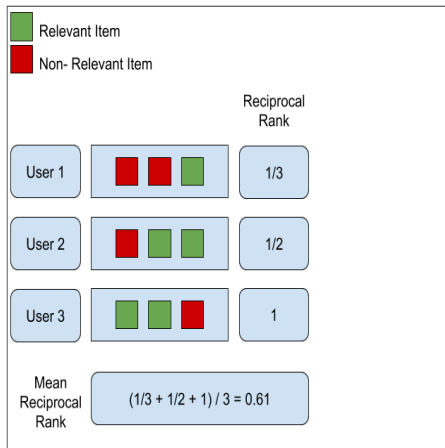


Figure: An example of MRR calculation

Evaluation Metrics: Mean Average Precision (MAP)

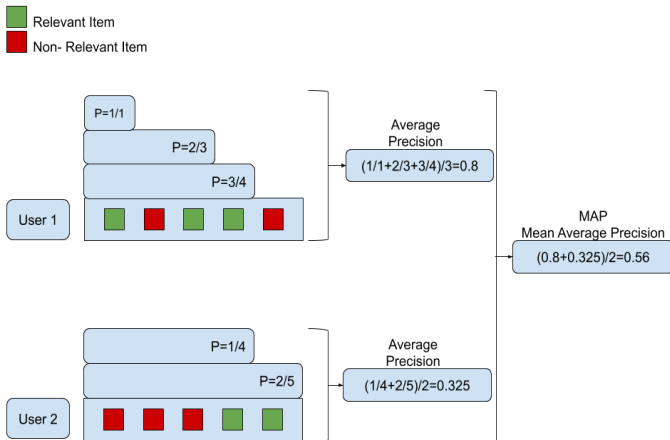


Figure: An example of MAP calculation

Evaluation Metrics: Normalized Discounted Cumulative Gain (nDCG)

$$DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)} \quad (2)$$

$$nDCG_p = \frac{DCG_p}{IDCG_p}. \quad (3)$$



Lorik Dumani.

Good premises retrieval via a two-stage argument retrieval model.

In *Grundlagen von Datenbanken*, pages 3–8, 2019.



Richard D Rieke, Malcolm Osgood Sillars, and Tarla Rai Peterson.

Argumentation and critical decision making.

Longman New York, 1997.



Hamed Zamani, Mostafa Dehghani, W Bruce Croft, Erik Learned-Miller, and Jaap Kamps.

From neural re-ranking to neural ranking: Learning a sparse representation for inverted indexing.

In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 497–506, 2018.