

Modern Talking: Key-Point Analysis using Modern Natural Language Processing

Max Henze and Thi Kim Hanh Luu and Jan Heinrich Reimer

Martin Luther University Halle-Wittenberg, Germany

{max.henze, thi.luu, jan.reimer}@student.uni-halle.de

Abstract

In this paper, we present our two approaches for Track 1 Key Point Matching of the common Quantitative Summarization - Key Point Analysis task.

Our two approaches use two methods to determine whether a Key Point matches an argument under the given topic: Jaccard coefficients and binary classification using a fine-tuned language model RoBERTa. We also applied several classical preprocessing methods from the field of natural language processing to preprocess our input text data, such as stop words, lemmatization, and application of the WordNet lexical database.

With our approaches, we got the best score ... of all participants in the Shared Task Track 1 Key Point Matching.

1 Introduction

Test citation (?)

TODO

2 Related Work

2.1 Key Point Analysis

Given the first track of the shared task analyzing key points is crucial. In their work (?) Bar-Haim et al. propose an approach for summarising large argument collections to small sets of key points. Thus covering a sufficient amount of all arguments. They show that domain experts can very quickly create pro and con key points, which are able to "capture the gist" of the arguments on the given topic. All this without being exposed to the arguments themselves. Furthermore they develop the large-scale dataset ArgKP which is the foundation of this shared task. In a later work (?) Bar-Haim et al. construct an automatic method for key point extraction which can compete with key points created by human domain experts. The method consists of two aspects. Assuming that the key points

can be found among the given comments they first select short, high quality comments as key point candidates and then select the candidate with the highest data coverage. Using the HuggingFace transformer framework they fine-tune four different models from which ALBERT (?) has the best F1 score with 0.809 but RoBERTa (?) (F1 score of 0.773) is chosen for key point extraction since it has a 6 times faster inference time. In the paper (?), Egan et al. propose a summarising for informal arguments such as they occur in online political debates. By extracting verbs and their syntactic arguments they retrieve points which can make key content accessible. By grouping these points they propose to create discussion summaries.

2.2 Argument Clustering

Argument Clustering is a mighty tool that enables algorithms to assign multiple arguments, which address a similar key message to a given topic. In the paper (?), Reimers et al. make use of "contextualized word embeddings to classify and cluster topic-dependent arguments". Having performed argument classification they then compute similar and dissimilar pairs of arguments. Two approaches one with clustering and one without are being used. Clustering arguments is achieved by usage of agglomerative hierarchical clustering (?). Without clustering a fine-tuned BERT-base-uncased model reached a F1 mean score for similar and dissimilar arguments of 0.7401. Agglomerative hierarchical clustering being a strict partitioning algorithm, results for clustering perform worse by up to 7.64pp (Bert-large F1 mean score: 0.7135). Hence they conclude that "strict partitioning clustering methods introduce a new source of errors". Another approach proposed in a paper (?) by Ajjour et al. revolves around clustering arguments into so called frames which are "a set of arguments that focus on the same aspect". Thereby framing (?) only a specific information to present to the listeners and

convince them of your stance. They propose that an argument consists of two crucial parts. The topic and the frame. Hence their approach splits into three steps: First, all arguments are clustered into m topics. Second, topical features are extracted from all arguments and therefore from its cluster. Third, the arguments are reclustered into k non-overlapping frames. By utilizing k-means (?) for clustering and Term Frequency-Inverse Document Frequency (TF-IDF) for topic removal they achieved a F1 score of 0.28.

2.3 Stance Classification

By knowing the stance of an argument it becomes nearly effortless for a human to classify it as pro or con to a given topic.

3 Data

The dataset used in this shared task is the IBM Debater(R) - ArgKP v1 (?) which consists of over 24.000 argument and key point pairs labeled as matching/non-matching. They all belong to one of 28 controversial topics ranging from "Assisted suicide should be a criminal offence." to "We should abandon marriage." therefore drawing a clear line between supporting and non-supporting key points.

The section used for training has over 5500 arguments belonging to over 200 key points with 24 topics. This leaves the dev dataset with over 900 arguments and 36 key points for 4 topics.

3.1 Characteristics

TODO

4 Approach

In this section we present a rule-based baseline performing well, given it's simplicity. In addition to this we came up with two BERT and RoBERTa based approaches. Both improving our baseline.

4.1 Baseline

We start with a very simple baseline. Therefore choosing a Term Overlap baseline with preprocessed terms. Generally it can be assumed, that key-points are summarizing ideas of all associated arguments. We therefore came up with the idea that certain key-words, contained in a lot of arguments, are also very likely to be present in the associated key-point. This makes sense with our intuition, because rather than using completely new words for summarization of arguments, a human

would rather reuse certain important words, which have been already found in the arguments.

For example, the following argument exists in the ArgKP dataset:

People reach their limit when it comes to their quality of life and should be able to end their suffering. This can be done with little or no suffering by assistance and the person is able to say good bye.

In relation to this the following key-point can be found:

Assisted suicide reduces suffering.

It can already be seen, that an overlap with the words *suffering* exists. We can further increase this overlap by performing simple preprocessing steps. First of all, we utilize stop word removal for reducing the noise within all arguments. Initially this can be seen counter productive, because less words means less overlap and therefore worse performance. But at second glance this makes a lot of sense. A lot of arguments and key-points contain unnecessary words like *the, and, as etc..* Removing these gives us purer sentences and results in less confusion with the term-overlap algorithm. Furthermore the redundancy of language makes it possible to contain key-aspects in sentences, even with out these unnecessary stop words.

Secondly, Stemming reduces terms to their corresponding stems and thus achieves a better generalization, when comparing terms. For example, the word: *weakness* will be stemmed to *weak* using the Porter-Stemmer (?). Thus creating an overlap between those words and increasing the possibility that an argmunt containing *weakness* will be associated to a key-pont containing *weak*.

Thirdly, we increase the generalization of our term-overlap algorithm even further by creating lists of synonyms and antonyms and testing if checked words can be replaced with candidates from these to increase overlap.

For the actual similarity computation of given arguments and key-points we use the Jaccard similarity coefficient (?). Meaning a higher proportion of terms that appear in an argument as well as in a key-point will classify this key-point to be more likely to match.

5 Results

TODO

6 Conclusion and Future Work

TODO

6.1 Future Work

TODO

A Appendix

We release our source code online under the free MIT license.¹

¹<https://github.com/heinrichreimer/modern-talking>