

Automatic Explanation Quality Assessment in Online Learning Environments

Sameer Bhatnagar¹, Amal Zouaq¹, Michel C. Desmarais¹, and Elizabeth Charles²

¹ Ecole Polytechnique Montreal

{sameer.bhatnagar,amal.zouaq,michel.desmarais}@polymtl.ca

² Dawson College echarles@dawsoncollege.qc.ca

Abstract. *Asynchronous Peer Instruction* is an increasingly popular item type in online learning environments. Students first respond to a question item, but they must also provide an explanation for their reasoning. They are then presented alternative explanations as written by their peers, and given the opportunity to change their initial answer choice, based on those they find most convincing. The peer-explanations that students find most convincing represent valuable data, for teachers to better grasp their students’ understanding, and for the learning environment itself, as higher quality explanations can be shown to students as examples to compare their work to. This study reports on the application of argument mining methods in the context of asynchronous peer instruction, with the objective of automatically identifying high quality student explanations. Our results offer the potential to inform the design of “learnersourcing” systems, where the content is both *generated* and *evaluated* by novices. These design choices are critical as these systems scale, especially with respect to providing pedagogically insightful reports to teachers, and presenting engaging alternative explanations to students to promote higher order thinking.

Keywords: Argument mining · Learnersourcing · Peer Instruction

1 Introduction

2 Related Work

2.1 Learnersourcing & Comparative Peer Assessment

Ripple[7], AXIS[16] Juxtapeer[1]

2.2 Argument Quality & Convincingness

Conventional argument-mining pipelines include several successive components, starting with the automatic detection of argumentative units, classification of these units into types (e.g. major claim, minor claim, premise), and identification

of argumentative relations (which evidence units support which claim). Such pipelines are essential in question-answering systems [8] and are at the heart of the IBM Project Debater initiative.

Work in the area of automatic evaluation of argument quality finds its roots in detecting evidence in legal texts[10], but has accelerated in recent years as more datasets become available in everyday contexts, and focus shifts to modelling more qualitative measures, such as *convincingness*.

Some of the earlier efforts included work on automatically scoring of persuasive essays [12] and modelling persuasiveness in online debate forums [14]. However, evaluating argument *convincingness* with an absolute score can be challenging, which has led to significant work in adopting a pairwise approach, where data consists of pairwise observations of two arguments, labelled with which of the two is most convincing.

rationale	chosen_rationale
The graph shows constant positive acceleration and then constant negative acceleration. This means that the velocity-time graph should have a positive slope and then a negative slope, and graph C is the only option that satisfies those requirements.	1st, acceleration is +, indicating an increase in velocity.\r\nthen, acceleration is suddenly and without warning negative, and velocity is reduced.

Table 1: Example of instance in pairwise comparison task, where two students explanations are compared, and one is chosen as more convincing

In [5], the authors propose a feature-rich support vector machine, as well as an end-to-end neural approach based on pre-trained Glove vectors and a bidirectional Long-Short-Term Memory network for the pairwise classification task. This is extended in [4], where the authors build a Siamese network architecture, where each leg is a BiLSTM, taking as input the pair of explanations as Glove embeddings [11], in order to detect which of argument in a pair has the most convincing evidence. Finally, based on the success of transformer models such as BERT[3], the authors of [15] release a dataset of argument pairs and show that these models accurately predict the most convincing argument in a pair.

- [15] Assessment of argument quality, with a dataset that has both individual scores and pairwise-ranked data

3 Methods

3.1 Data

The dataset is comprised of pairs of student explanations for a particular answer choice to a given question. The first explanation is always the one written by

Table 2: Descriptive statistic for each dataset of argument pairs, with last rows showing *dalite* split by discipline

dataset	N_{pairs}	N_{topics}	N_{args}	\overline{wc} (SD)	$\overline{\Delta wc}$ (SD)
IBM	9125	11	3474	23 (7)	3 (2)
UKP	11650	16	1052	49 (28)	30 (23)
dalite	8551	102	8942	17 (15)	12 (7)
dalite:Biology	3919	49	4116	15 (14)	10 (6)
dalite:Chemistry	1666	24	1758	20 (14)	12 (7)
dalite:Physics	2966	29	3068	19 (15)	15 (7)

the learner-annotator, while the second is an alternative which they either chose as more convincing, or not. The data is filtered so as to only keep observations where the explanations are within half a standard deviation in length of each other. To ensure internal reliability, we only keep explanations that were chosen at least 5 times. To ensure that the explanations in each pair are of comparable length, we keep only those with word counts that are within 25 words of each other. This leaves us a dataset with 8551 observations, spanning 2216 learner annotators having completed, on average, 4.0 items each, from a total of 109 items across three disciplines, with at least 50 explanation-pairs per item.

Table 3: Observations of students choosing a peer explanation as more convincing than their own, or not, aggregated by discipline and whether they started and finished with the correct answer

	rr	rw	wr	ww
Biology	2459	124	733	603
Chemistry	1151	51	228	236
Physics	2288	66	278	334

Table 3 highlights one key difference between the modelling task of this study, and related work in argument mining, where annotators are presented pairs of arguments that are always for the same stance, in order to limit bias due to their opinion on the motion when evaluating which argument is more convincing. In a *Peer Instruction* learning environment, other pairings are possible and pedagogically relevant. In this dataset, the majority of students keep the same answer choice between the two steps of the prompt, and so they are comparing two explanations that are either both correct (“rr”) or incorrect (“wr”). However, there is 17 % of the observations in this dataset are for students who not only choose an explanation more convincing than their own, but also switch answer choice, either from the incorrect to correct, or the reverse. These pairs add a different level of complexity to the model, but are very pertinent in the

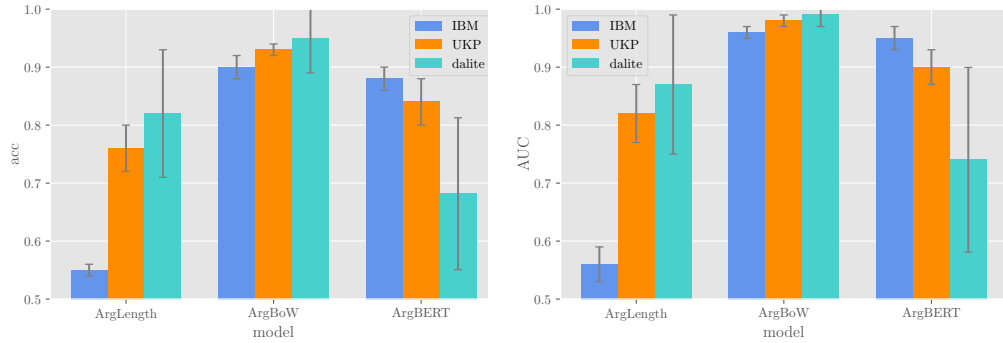
pedagogical context: what are the argumentative features which can help students remediate an initial wrong answer choice (“*wr*”)? What are the features that might be responsible for getting students to actually move away from the correct answer choice (“*rw*”)?

3.2 Models

In pairwise classification tasks so 50% is the baseline performance. In Table ?? we begin by comparing the *ArgLongest* model baseline, where we predict that students simply choose the longer explanation of the pair. We also include two baseline models on *Bag of Words* models: *ArgBoWGen* where the term-document-matrix is built from an open-source textbook from the corresponding discipline³, and *ArgBoWItemSpec*, where the term-document matrix is built from the words students have used for the item (pertinent when no reference text is available for a discipline). As this is a new context for these argument mining methods, we include the same baselines on the carefully curated *IBMArgPair* dataset, which is of the same format.

For our experiments, we begin by following the line of work proposed by [5], and experiment with a feature-rich linear SVM classifier for the pairwise classification task. We use a similar feature set, which we categorize as **lexical**, **syntactic**, and **semantic**, as described in Table ?. we begin by computing the feature vector for each explanation, and compute the difference for each pairwise ranking instance as per the well established SVM-Rank algorithms [6], training the model to learn which of the pair is the more convincing argument.

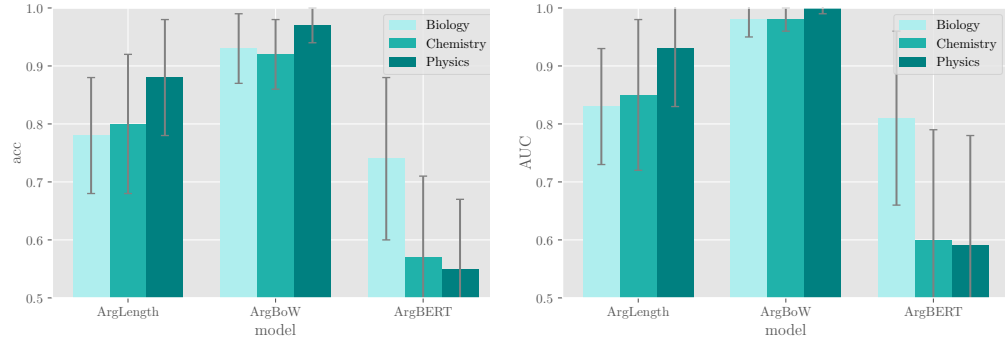
4 Results



(a) Pairwise ranking accuracy for different models across datasets

(b) Pairwise ranking classification ROC-AUC for different models across datasets

³ <https://openstax.org/>



(a) Pairwise ranking accuracy for different models in myDalite dataset, across disciplines

(b) Pairwise ranking classification ROC-AUC for different models in myDalite dataset, across disciplines

5 Discussion

In [9], for task of pairwise ranking of newspaper articles based on “quality”, the authors achieve a similar result: when comparing the performance of SVM-rank models using different input feature sets (e.g. *use of visual language*, *use of named entities*, *affective content*), their top performing models achieve “same-topic” pairwise ranking accuracy of 0.84 using a combination of content and writing features, but also a 0.82 accuracy with the content words as features alone.

6 Future Work

In this study we do not ever infer which are, overall, the most convincing student explanations for any given item. Inferring a gold standard of global rankings, starting from these pairwise preference data can be accomplished using research from the information retrieval community[2]. Work on deriving point wise scores for argument pairs is proposed as a Gaussian Process Preference Learning task by [13]. Seeing the lack of pointwise labels for overall convincingness, [15] released a dataset where they collect this data as well. A comparable source of data inside the myDALITE platform are the feedback scores teachers can optionally provide to students on their explanations.

References

1. Cambre, J., Klemmer, S., Kulkarni, C.: Juxtapeer: Comparative Peer Review Yields Higher Quality Feedback and Promotes Deeper Reflection. In: Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18. pp. 1–13. ACM Press, Montreal QC, Canada (2018). <https://doi.org/10.1145/3173574.3173868>

2. Chen, X., Bennett, P.N., Collins-Thompson, K., Horvitz, E.: Pairwise ranking aggregation in a crowdsourced setting. In: Proceedings of the sixth ACM international conference on Web search and data mining. pp. 193–202 (2013)
3. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
4. Gleize, M., Shnarch, E., Choshen, L., Dankin, L., Moshkovich, G., Aharonov, R., Slonim, N.: Are you convinced? choosing the more convincing evidence with a Siamese network. arXiv preprint arXiv:1907.08971 (2019), <https://www.aclweb.org/anthology/P19-1093/>
5. Habernal, I., Gurevych, I.: Which argument is more convincing? Analyzing and predicting convincingness of Web arguments using bidirectional LSTM. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1589–1599. Association for Computational Linguistics, Berlin, Germany (Aug 2016). <https://doi.org/10.18653/v1/P16-1150>, <https://www.aclweb.org/anthology/P16-1150>
6. Joachims, T.: Optimizing search engines using clickthrough data. In: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 133–142 (2002)
7. Khosravi, H., Kitto, K., Williams, J.J.: Ripple: A crowdsourced adaptive platform for recommendation of learning activities. arXiv preprint arXiv:1910.05522 (2019)
8. Lippi, M., Torroni, P.: Argumentation mining: State of the art and emerging trends. ACM Transactions on Internet Technology (TOIT) **16**(2), 10 (2016)
9. Louis, A., Nenkova, A.: What Makes Writing Great? First Experiments on Article Quality Prediction in the Science Journalism Domain. Transactions of the Association for Computational Linguistics **1**, 341–352 (2013)
10. Moens, M.F., Boiy, E., Palau, R.M., Reed, C.: Automatic detection of arguments in legal texts. In: Proceedings of the 11th international conference on Artificial intelligence and law. pp. 225–230 (2007)
11. Pennington, J., Socher, R., Manning, C.D.: Glove: Global Vectors for Word Representation. In: EMNLP. vol. 14, pp. 1532–1543 (2014)
12. Persing, I., Ng, V.: End-to-End Argumentation Mining in Student Essays. In: HLT-NAACL. pp. 1384–1394 (2016)
13. Simpson, E., Gurevych, I.: Finding Convincing Arguments Using Scalable Bayesian Preference Learning. Transactions of the Association for Computational Linguistics **6**, 357–371 (2018), <https://www.aclweb.org/anthology/Q18-1026>
14. Tan, C., Niculae, V., Danescu-Niculescu-Mizil, C., Lee, L.: Winning Arguments: Interaction Dynamics and Persuasion Strategies in Goodfaith Online Discussions. arXiv:1602.01103 [physics] pp. 613–624 (2016). <https://doi.org/10.1145/2872427.2883081>, <http://arxiv.org/abs/1602.01103>, arXiv: 1602.01103
15. Toledo, A., Gretz, S., Cohen-Karlik, E., Friedman, R., Venezian, E., Lahav, D., Jacovi, M., Aharonov, R., Slonim, N.: Automatic Argument Quality Assessment-New Datasets and Methods. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 5629–5639 (2019), <https://www.aclweb.org/anthology/D19-1564.pdf>
16. Williams, J.J., Kim, J., Rafferty, A., Maldonado, S., Gajos, K.Z., Lasecki, W.S., Heffernan, N.: AXIS: Generating Explanations at Scale with Learnersourcing and

Machine Learning. In: Proceedings of the Third (2016) ACM Conference on Learning @ Scale - L@S '16. pp. 379–388. ACM Press, Edinburgh, Scotland, UK (2016). <https://doi.org/10.1145/2876034.2876042>