

Does Informativeness Matter? Active Learning for Educational Dialogue Act Classification

Wei Tan¹, Jionghao Lin^{1,2(✉)}, David Lang³, Guanliang Chen¹, Dragan Gasevic¹, Lan Du¹, and Wray Buntine^{4,1}

¹ Monash University, Clayton, Australia
 {wei.tan2, jionghao.lin1, lan.du, guanliang.chen, dragan.gasevic}@monash.edu

² Carnegie Mellon University, Pittsburgh, USA

³ Stanford University, Stanford, USA
 dnlang86@stanford.edu

⁴ VinUniversity, Hanoi, Vietnam
 wray.b@vinuni.edu.vn

Abstract. Dialogue Acts (DAs) can be used to explain what expert tutors do and what students know during the tutoring process. Most empirical studies adopt the random sampling method to obtain sentence samples for manual annotation of DAs, which are then used to train DA classifiers. However, these studies have paid little attention to sample informativeness, which can reflect the information quantity of the selected samples and inform the extent to which a classifier can learn patterns. Notably, the informativeness level may vary among the samples and the classifier might only need a small amount of low informative samples to learn the patterns. Random sampling may overlook sample informativeness, which consumes human labelling costs and contributes less to training the classifiers. As an alternative, researchers suggest employing statistical sampling methods of Active Learning (AL) to identify the informative samples for training the classifiers. However, the use of AL methods in educational DA classification tasks is under-explored. In this paper, we examine the informativeness of annotated sentence samples. Then, the study investigates how the AL methods can select informative samples to support DA classifiers in the AL sampling process. The results reveal that most annotated sentences present low informativeness in the training dataset and the patterns of these sentences can be easily captured by the DA classifier. We also demonstrate how AL methods can reduce the cost of manual annotation in the AL sampling process.

Keywords: Informativeness · Active Learning · Dialogue Act Classification · Large Language Models · Intelligent Tutoring Systems.

1 Introduction

Traditional one-on-one tutoring involves human participants (*e.g.*, a course tutor and a student), which has been widely acknowledged as an effective form of instruction [16, 24]. Understanding what expert tutors do and students know

during the tutoring process is a significant research topic in the field of artificial intelligence in education [5], which can contribute to the design of dialogue-based intelligent tutoring systems (*e.g.*, AutoTutor [14]) and the practice of human tutoring [12, 24]. A typical method used to understand the tutoring process is to map the sentences from tutors and students onto dialogue acts (DAs) which can manifest the intention behind sentences [24, 16]. For example, a tutor’s sentence (*e.g.*, “*No, it is incorrect!*”) can be coded as the **Negative Feedback** DA. The mapping process often needs a pre-defined DA scheme developed by educational experts, and the DA scheme is used to annotate the sentences with the DAs [17, 16]. However, manually annotating the DAs for millions of tutoring dialogues is impractical for conducting research related to educational dialogue since the annotating process is often time-consuming and costly [15]. Therefore, researchers typically annotate a small amount of sentences from tutoring dialogues and further use the annotated sentences to train machine learning classifiers to automate the annotation process for DAs [16, 12].

Though previous works demonstrated the potential of automating the DA classification [15, 17, 12, 3, 6], little attention has been given to investigate the extent to which the DA classifiers can learn the patterns from the annotated sentences, which is important to build a robust classifier to facilitate tutoring dialogue analysis. The ability of the DA classifiers that learn the relation between the model inputs (*e.g.*, sentences) and outputs (*e.g.*, DAs) is defined as learnability [21, 10]. To improve the learnability of a DA classifier, it is necessary to train the DA classifier on the instances that can help the classifier capture the underlying patterns and further improve the generalization of the DA classifier on the full dataset [10]. These instances are considered highly informative samples in the dataset [10]. It should be noted that the informativeness level might be different among the samples, and the classifier can learn the patterns of the low informative samples with a small amount of training data [21, 10]. Therefore, annotating more high informative samples and less low informative samples could further save human annotation costs. However, few works investigated sample informativeness on the task of educational DA classification, which motivated us to focus on the **first research goal**, *i.e.*, *investigate the sample informativeness levels of annotated DAs*. Additionally, the previous works annotated the sentences by randomly sampling from the archived tutoring dialogue corpora [17, 15, 16, 12], which might not provide sufficient high informative samples for training the DA classifier. We argue that blindly annotating the dialogue sentences to train the DA classifier might consume excessive human annotation costs and not enhance the classifier’s learnability on the full dataset [19]. Instead, a promising solution is to use the statistical active learning (AL) methods, which aim to select highly informative samples to train the classifier [19]. To provide details about the efficacy of AL methods, the **second research goal** of the current study was to *investigate the extent to which statistical AL methods support the DA classifier training on the DAs with different informativeness levels*.

Our results showed that most annotated instances from the training samples presented low informativeness, which indicated that the patterns of these samples

could be sufficiently captured by our DA classifier. Additionally, the DA classifier needs to train on more high informativeness samples to improve the classification performance on the unseen dataset, which further requires support by statistical AL methods. We compared the state-of-the-art AL method (*i.e.*, CoreMSE [22]) with other commonly-used AL methods (*i.e.*, Least Confidence and Maximum Entropy) and random baseline. We found that the CoreMSE AL method could select more high informative instances for the DA classifier at the early stage of the classifier training process and gradually reduced the selection of the low informative instances, which could alleviate the cost of manual annotation.

2 Related Work

2.1 Educational Dialogue Act Classification

To automate classifying educational DAs, previous studies in the educational domain have employed machine learning models to train on annotated sentences [15, 1, 18, 17, 12]. For example, Boyer *et al.* [1] employed a Logistic Regression model to train on linguistic features (*e.g.*, N-grams) from the annotated sentences in 48 programming tutoring sessions. They [1] achieved 63% accuracy on classifying thirteen DAs. Later on, Samei *et al.* [18] used Decision Trees and Naive Bayes models to train on linguistic features and contextual features (*e.g.*, the speaker’s role of previous sentences) from 210 annotated tutoring sentences on the science-related topic. Their work [18] achieved the accuracy of 56% in classifying seven educational DAs. These previous works trained DA classifiers by using the sentences which were randomly sampled from the larger tutoring dialogue corpora. Though these works demonstrated the feasibility of automating DA classification [15, 1, 18, 17, 12], it still remains largely unknown whether the DA classifier learned the patterns sufficiently from the randomly selected sentences, which is important for building robust and reliable DA classifiers.

2.2 Sample Informativeness

Training a DA classifier on a sufficient amount of annotated sentences can help the classifier achieve satisfactory classification performance. However, manual annotation is time-consuming and expensive [12, 15, 17]. To mitigate the human annotation cost and also help the classifier achieve satisfied classification performance, Du *et al.* [4] suggested annotating the most informative samples, which can reduce the generalization error and uncertainty of the DA classifier on the unseen data. A recent work [21] proposed the *Data Maps* framework which used **Confidence**, **Variability**, and **Correctness** to measure the informativeness. **Confidence** denotes the averaged probability of the classifier’s prediction on the correct labels across all training epochs; **Variability** denotes the spread of probability across the training epochs, which capture the uncertainty of the classifier; **Correctness** denotes the fraction of the classifier correctly predicted labels over all training epochs [21]. Building upon the work introduced in [21], Karamcheti *et al.* [10] further used the **Correctness** to categorize the samples into four groups: *Easy*, *Medium*, *Hard*, and *Impossible* [21] and these groups

indicated the extent to which the classifier can learn the patterns from the annotated instances. Inspired by [4, 21, 10], we argue that when scrutinizing the sampling process for training the DA classifier, if samples are randomly selected, there will be redundant for overly sampling the *Easy* samples and insufficient for sampling high informative samples, which might consume the human annotation budget and lead to poor generalizability of the DA classifier. As a remedy, it is important to select the most suitable samples for training the classifier. AL can offer promising methods to select the most suitable samples.

2.3 Statistical Active Learning

Recent studies on educational DA classification [12, 15] have agreed that the high demand for the annotated dataset was still an issue for the DA classification task. To alleviate this issue, a promising solution is to use AL methods which can select the high informative samples from the unlabeled pool and send them to Oracle (*e.g.*, human annotator) for annotation [19]. Traditionally, there are three typical scenarios of AL methods: 1) *membership query synthesis*, which focuses on generating artificial data-point for annotation rather than sampling the data-point from the real-world data distribution, 2) *stream-based sampling*, which focuses on scanning through a sequential stream of non-annotated data-points and make sampling query decision individually, and 3) *pool-based sampling* which focuses on selecting the most informative samples from the non-annotated data pool and send them to the oracle for annotation [19].

As the annotated DAs were available in our study and the dataset was not collected in a sequential stream, we consider our study fits well with the *pool-based sampling* scenario. The pool-based AL methods can both reduce the computational cost of model training and maintain the performance of the model trained on the annotated dataset [19]. Many studies employed the pool-based AL methods (*e.g.*, Least Confidence [20, 8] and Maximum Entropy [11]) on the educational tasks (*e.g.*, student essay classification [8], and educational forum post classification [20]) and their results have demonstrated the promise of AL methods on alleviating the demand for annotated datasets. However, it still remains largely unknown about the extent to which AL methods can select the informative samples to support the automatic classification of educational DA.

3 Methods

3.1 Dataset

The current study obtained ethics approval from the Monash University Human Research Ethics Committee under application number 26156. The dataset used in our study was provided by an educational technology company that operated online tutoring services and collected the data from tutors and students along with informed consent allowing the use of the de-identified dataset for research. The dataset contained detailed records of the tutoring process where tutors and students collaboratively solve various problems (the subjects including mathematics, physics, and chemistry) via textual message. In total, our dataset contained 3,626 utterances (2,156 tutor utterances and 1,470 student utterances)

from 50 tutorial dialogue sessions. The average number of utterances per tutorial dialogue session was 72.52 ($min = 11$, $max = 325$); tutors averaged 43.12 utterances per session ($min = 5$, $max = 183$), and students 29.40 utterances per session ($min = 4$, $max = 142$). We provided a sample dialogue in the digital appendix via <https://github.com/jionghaolin/INFO>.

3.2 Educational Dialogue Act Scheme and Annotation

Identifying the DAs in tutorial dialogues often relies on a pre-defined educational DA coding scheme [9]. By examining the existing literature, we employed the DA scheme introduced in [24] whose effectiveness in analyzing online one-on-one tutoring has been documented in many previous studies (*e.g.*, [7, 12, 23]). The DA scheme developed in [24] characterizes the DAs into a two-level structure. To discover more fine-grained information from tutor-student dialogue, in our study, we decided to annotate the tutoring dialogues by using the second-level DA scheme. Notably, some utterances in dialogues contained multiple sentences, and different sentences can indicate different DAs. To address this concern, Vail and Boyer [24] suggested partitioning the utterances into multiple sentences and annotating each sentence with a DA. After the utterance partition, we then removed the sentences which only presented meaningless symbols or emoji. Lastly, we recruited two human coders to annotate DAs, and we obtained Cohen’s κ score of 0.77 for the annotation. The annotations achieved a substantial agreement between the two coders, and we recruited a third educational expert to resolve the inconsistent cases. The full DA scheme can be found in a digital appendix via <https://github.com/jionghaolin/INFO>, which contains 31 DAs. Due to the space limit, we only presented students’ DAs in Table 1.

Table 1. The DA scheme for annotating student DAs. The DAs were sorted based on their frequency (*i.e.*, the column of **Freq.**) in the annotated dataset.

Dialogue Acts (DAs)	Sample Sentences	Freq.
Confirmation Question	“So that’d be 5?”	4.93%
Request Feedback by Image	[Image]	4.34%
Understanding	“Oh, I get it”	1.46%
Direction Question	“Okay what do we do next?”	1.20%
Information Question	“Isn’t there a formula to find the n th term?”	1.06%
Not Understanding	“I don’t know.”	0.24%
Ready Answer	“Yep, ready to go.”	0.07%

3.3 Identifying Sample Informativeness via Data Maps

To identify the informativeness of annotated instances, we need to train the annotated dataset on a classifier. Building upon our previous work [13], we used ELECTRA [2] as the backbone model for classifying 31 DAs, which is effective in capturing nuanced relationships between sentences and the DAs. The dataset (50 dialogue sessions) was randomly split to *training* (40 sessions) and *testing set* (10 sessions) in the ratio of 80%:20% for training the classifier. The classifier achieved accuracy of 0.77 and F1 score of 0.76 on the testing set. Then, we applied

the *Data Maps*⁵ to the DA classifier to analyze the behaviour of the classifier on learning individual instance during the training process. Following the notation of *Data Maps* in [21], the training dataset denotes $Dataset = \{(x, y^*)_i\}_{i=1}^N$ across E epochs where the N denotes the size of the dataset, i th instance is composed of the pair of the observation of x_i and true label y_i^* in the dataset. **Confidence** ($\hat{\mu}_i$) was calculated by $\hat{\mu}_i = \frac{1}{E} \sum_{e=1}^E p_{\theta^{(e)}}(y_i^* | x_i)$ where $p_{\theta^{(e)}}$ is the model’s probability with the classifier’s parameters $\theta^{(e)}$ at the end of the e^{th} epoch. **Variability** ($\hat{\sigma}_i$) was calculated by the variance of $p_{\theta^{(e)}}(y_i^* | x_i)$ across epochs: $\hat{\sigma}_i = \sqrt{\frac{\sum_{e=1}^E (p_{\theta^{(e)}}(y_i^* | x_i) - \hat{\mu}_i)^2}{E}}$. **Correctness** (Cor) was the fraction of the classifier correctly annotated instance x_i over all epochs. A recent study [10] further categorized **Correctness** into *Easy* ($Cor \geq 0.75$), *Medium* ($0.75 > Cor \geq 0.5$), *Hard* ($0.5 > Cor \geq 0.25$), and *Impossible* ($0.25 > Cor \geq 0$). In line with the studies [21, 10], we mapped the training instances along two axes: the y – *axis* indicates the confidence and the x – *axis* indicates the variability of the samples. The colour of instances indicates **Correctness**.

3.4 Active Learning Selection Strategies

We aimed to adopt a set of AL methods to examine the extent to which AL methods support the DA classifier training on the DAs with different informativeness levels. We employed the state-of-the-art AL method (CoreMSE [22]) to compare with two commonly used AL methods (*i.e.*, Least Confidence [19] and Maximum Entropy [25]) and the random baseline. Following the algorithms in the papers, we re-implemented them in our study. **Random Baseline** is a sampling strategy that samples the number of instances uniformly at random from the unlabeled pool [19]. **Maximum Entropy** is an uncertainty-based method that chooses the instances with the highest entropy scores of the predictive distribution from the unlabeled pool [25]. **Least Confidence** is another uncertainty-based method that chooses the instances with the least confidence scores of the predictive distribution from the unlabeled pool [19]. **CoreMSE** is a pool-based AL method proposed by [22]. The method involves both diversity and uncertainty measures via the sampling strategy. It selects the diverse samples with the highest uncertainty scores from the unlabeled pool, and the uncertainty scores are estimated by the reduction in classification error. These diverse samples can provide more information to train the model to achieve high performance.

3.5 Study Setup

In the AL training process, the DA classifier was initially fed with 50 training samples randomly selected from the annotated dataset. For each AL method, the training process was repeated six times using different random seeds. A batch size of 50 was specified. We set the maximum sequence length to 128 and fine-tuned it for 30 epochs. The AdamW optimizer was used with the learning rate of 2e-5 to optimize the training of the classifier. All experiments were implemented on RTX 3090 and Intel Core i9 CPU processors.

⁵ <https://github.com/allenai/cartography>

4 Results

4.1 Estimation of Sample Informativeness

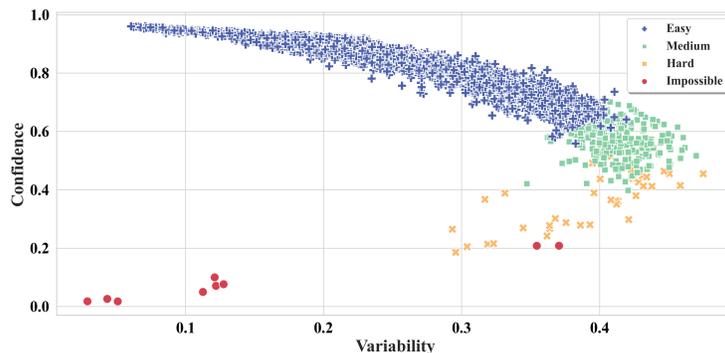


Fig. 1. The sample informativeness on the annotated sentences.

To estimate the sample informativeness level of each training instance, we employed the *Data Maps*. In Fig. 1, the y-axis represents the **Confidence** measure, the x-axis represents **Variability**, and the colours of instances represent **Correctness**, which are detailed in Sec. 3.3. The samples in the *Easy* group (*i.e.*, blue scatters) always presented a high confidence level, which indicates that most sample points could be easily classified by our DA classifier. Then, the samples in the *Medium* and *Hard* groups presented generally lower confidence and higher variability compared with the *Easy* group, which indicates that the samples in *Medium* and *Hard* were more informative than those in *Easy*. Lastly, the samples in the *Impossible* group always presented low confidence and variability. According to [21], the samples in *Impossible* could be the reasons for mis-annotation or insufficient training samples (*e.g.*, insufficient training samples for a DA can impede the classifier’s ability to capture its classification pattern) for the DA classifier. Most training samples were considered *Easy* to be classified, and only a small number of samples are *Impossible*. The distribution of *Easy* and *Impossible* points in Fig. 1 indicates that our dataset presented high-quality annotation, which is generally considered learnable for the DA classifier.

We further investigated the distribution of informativeness level for the annotated DAs⁶. As described in Sec. 3.2, we mainly present the distribution of the informativeness level for each DA in Fig. 2. We sorted the DAs based on their frequency in Table 1. We observed that most samples in the **Confirmation Question** and **Request Feedback by Image** DAs were in the *Easy* group, which indicates that both DAs are easy for the DA classifier to learn. Then, in the middle of Fig.2, the **Understanding**, **Direction Question**, and **Information Question** DAs had roughly 1% frequency in Table 1. Compared with the first two

⁶ Due to the space limit, we only present a part of the analysis results, and the full results can be accessed at <https://github.com/jionghaolin/INFO>.

DAs in Fig. 2, the middle three DAs had a higher percentage of *Medium*, which indicates that the DA classifier might be less confident in classifying these DAs. Lastly, the DAs *Not Understanding* and *Ready Answer* had the frequency lowered than 1% in Table 1. The samples in *Not Understanding* and *Ready Answer* were considered *Hard* and *Impossible* to be classified, respectively. The reasons might be that the annotated samples of *Not Understanding* and *Ready Answer* in our dataset were insufficient for the DA classifier to learn the patterns.

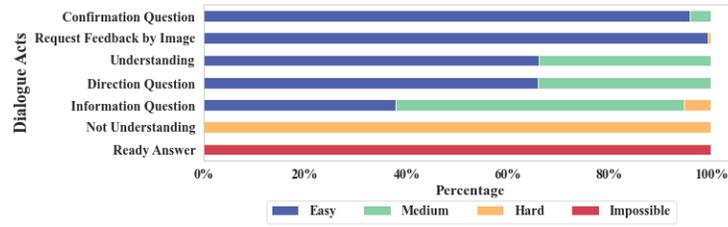


Fig. 2. Distribution of informativeness level for each DA (Student only)

4.2 Efficacy of Statistical Active Learning Methods

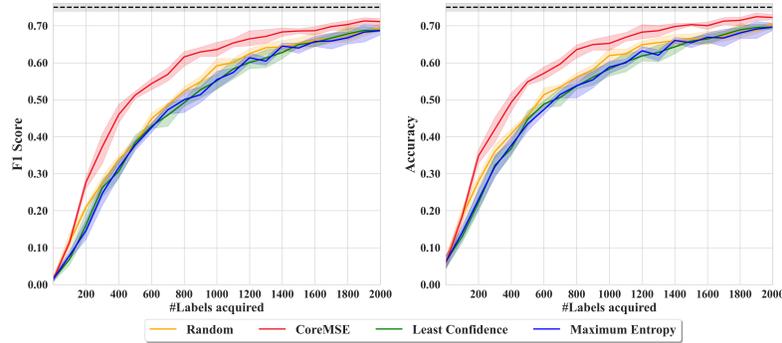


Fig. 3. Learning Curve of AL methods. The sampling batch for each AL method is 50. The classification performance was measured by F1 score and accuracy.

To investigate how AL methods select different informative samples for training the DA classifier, we first evaluated the overall performance of the DA classifier with the support of AL methods. In Fig. 3, the x-axis represents the training sample size after selecting samples and the y-axis represents the classification performance measured by F1 score and classification accuracy. We compared the CoreMSE method to the baseline methods, *i.e.*, Maximum Entropy (ME), Least Confidence (LC), and Random. The results in Fig. 3 demonstrate the learning

curves of the models where the CoreMSE method could help the DA classifier achieve better performance with fewer training samples compared with the baseline methods. For example, when acquiring 600 samples from the annotated data pool, the CoreMSE method could achieve roughly the F1 score of 0.55, which was equivalent to the classifier performance training on 900 samples with the use of baseline methods. It indicates that the CoreMSE method could save 30% human annotation costs compared to the baseline methods. Whereas the efficacy of LC and ME methods was similar to that of the random baseline in both F1 score and accuracy value; this indicates that the traditional uncertainty-based AL methods (*i.e.*, LC and ME) might not be effective on our classification task.

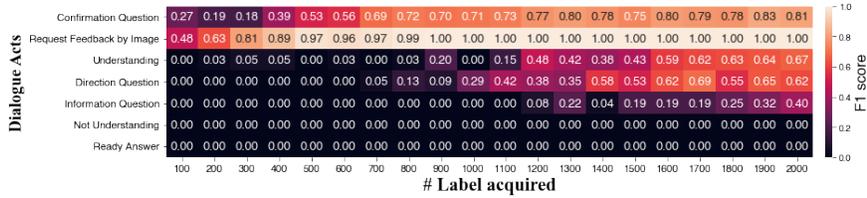


Fig. 4. F1 score for the dialogue acts only specific to students (Random)

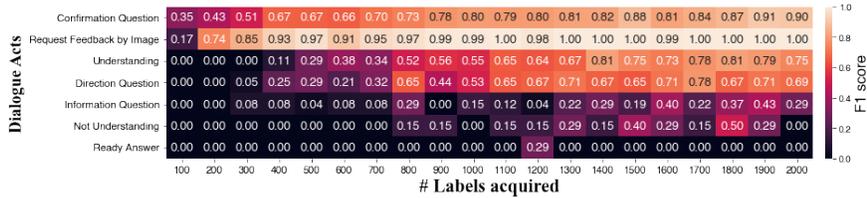


Fig. 5. F1 score for the dialogue acts only specific to students (CoreMSE)

Based on the results in Fig. 3, we observed that the random baseline performed slightly better than the ME and LC methods. In the further analysis of each DA, we decided to compare the efficacy between CoreMSE and the random baseline⁷. As shown in Fig. 4 and Fig. 5, we sorted DAs based on their frequency in the dataset. Fig. 4 shows the changes of F1 scores for each DA with the use of the random baseline; this indicates that the DA classifier performed quite well for classifying the **Confirmation Question** and **Request Feedback by Image** DAs as the training sample size increased. Regarding the **Understanding** and **Direction Question**, the DA classifier needed more than 1,600 annotated samples to achieve a decent performance. For the DAs **Information Question**, **Not Understanding**, and **Ready Answer**, the DA classifier almost failed to learn the patterns from the samples selected by the Random method. In comparison, Fig. 5 shows the classification performance for each DA with the support of the CoreMSE method. The results indicate that CoreMSE could support the DA

⁷ Due to the space limit, we documented our full results in a digital appendix, which is accessible via <https://github.com/jionghaolin/INFO>.

classifier to make accurate predictions for the DAs (e.g, **Understanding**) when acquiring fewer annotation samples than the random baseline. Then, though the classification performance for the DAs **Information Question** and **Not Understanding** was not sufficient, CoreMSE demonstrated the potential to improve the accuracy for both DAs as more labels were acquired. Lastly, both Random and CoreMSE methods failed to support the DA classifier to make an accurate prediction on the **Ready Answer** DA. The reason might be the fact that the **Ready Answer** DA was rare in our annotated dataset.

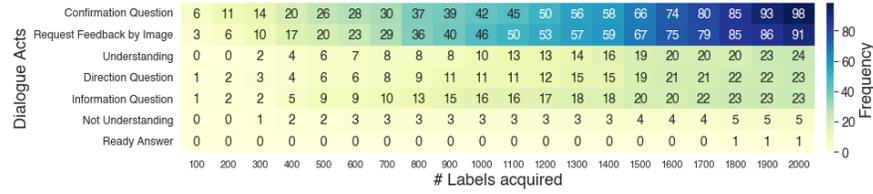


Fig. 6. The distribution of sampling frequency for each dialogue acts (**Random**)

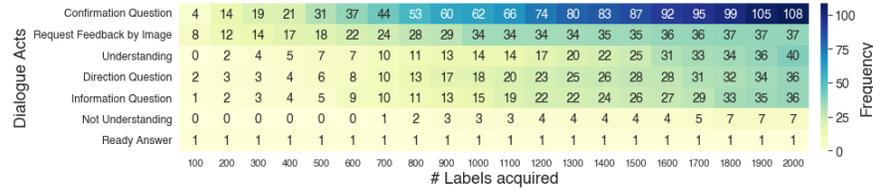


Fig. 7. The distribution of sampling frequency for each dialogue acts (**CoreMSE**)

Next, we investigated the sampling process between the CoreMSE and random methods to learn more details about how the CoreMSE method saved the annotation budget. For each sampling batch, we counted the cumulative frequency for each DA in Figs.6 and 7. For example, in the first 200 acquired annotations (Fig.6), the random method acquired 11 **Confirmation Question**, among which 6 of them were from the first 100 samples and 5 from the subsequent 100 acquired labels. We observed that compared with the random baseline (Fig.6), the CoreMSE method (Fig.7) gradually reduced sampling **Request Feedback by Image** instances from 700 labels acquired. This result indicates that the random baseline retained sampling the DA **Request Feedback by Image** instances even the F1 score achieved satisfactory performance, which might consume the budget for the human manual annotation, whereas, the CoreMSE method could alleviate the manual annotation cost when the DA classifier sufficiently learned the patterns. Then, compared with the random baseline (Fig. 6), the CoreMSE method (Fig. 7) generally selected more instances of **Understanding**, **Direction Question**, and **Information Question** DAs for training the DA classifier across the sampling process, which could explain the reason why the CoreMSE method supported the classifier achieving better performance than the random baseline.

5 Conclusion

Our study demonstrated the potential value of using a well-established framework *Data Maps* [21] to evaluate the informativeness of instances (*i.e.*, the tutorial sentences annotated with dialogue acts) in automatic classification of educational DAs. We found that most instances presented low informativeness in the training dataset, which was easy-to-learn for the dialogue act (DA) classifier. To improve the generalizability of the DA classifier on the unseen instances, we proposed that the classifier should be trained on the samples with high informativeness. Since the annotation of educational DA is extremely time-consuming and cost-demanding [15, 12, 17], we suggest avoiding the use of random sampling for annotation. Our study provided evidence of how the state-of-the-art statistical AL methods (*e.g.*, CoreMSE) could select informative instances for training the DA classifier and gradually reduce selecting the easy-to-learn instances, which can alleviate the cost of manual annotation. We acknowledged that some instances were quite rare in our original annotation, so we plan to employ the AL methods to select more of these rare instances for annotation in future research. Lastly, the AL methods might also be useful for costly evaluation tasks in the education field (*e.g.*, automated classroom observation), which requires the educational experts to annotate the characteristics of behaviours among the students and teachers. Thus, a possible extension of the current work would be to develop an annotation dashboard for human experts to be used in broader educational research.

Bibliography

- [1] Boyer, K., Ha, E.Y., Phillips, R., Wallis, M., Vouk, M., Lester, J.: Dialogue act modeling in a complex task-oriented domain. In: Proceedings of the SIGDIAL 2010 Conference. pp. 297–305 (2010)
- [2] Clark, K., Luong, M.T., Le, Q.V., Manning, C.D.: Electra: Pre-training text encoders as discriminators rather than generators. In: ICLR (2019)
- [3] D’Mello, S., Olney, A., Person, N.: Mining collaborative patterns in tutorial dialogues. *Journal of Educational Data Mining* **2**(1), 1–37 (2010)
- [4] Du, B., Wang, Z., Zhang, L., Zhang, L., Liu, W., Shen, J., Tao, D.: Exploring representativeness and informativeness for active learning. *IEEE transactions on cybernetics* **47**(1), 14–26 (2015)
- [5] Du Boulay, B., Luckin, R.: Modelling human teaching tactics and strategies for tutoring systems: 14 years on. *IJAIED* **26**(1), 393–404 (2016)
- [6] Ezen-Can, A., Boyer, K.E.: Understanding student language: An unsupervised dialogue act classification approach. *JEDM* **7**(1), 51–78 (2015)
- [7] Ezen-Can, A., Grafsgaard, J.F., Lester, J.C., Boyer, K.E.: Classifying student dialogue acts with multimodal learning analytics. In: Proceedings of the Fifth LAK. pp. 280–289 (2015)
- [8] Hastings, P., Hughes, S., Britt, M.A.: Active learning for improving machine learning of student explanatory essays. In: International Conference on Artificial Intelligence in Education. pp. 140–153. Springer (2018)

- [9] Hennessy, S., Rojas-Drummond, S., Higham, R., Márquez, A.M., Maine, F., Ríos, R.M., García-Carrión, R., Torreblanca, O., Barrera, M.J.: Developing a coding scheme for analysing classroom dialogue across educational contexts. *Learning, Culture and Social Interaction* **9**, 16 – 44 (2016)
- [10] Karamcheti, S., Krishna, R., Fei-Fei, L., Manning, C.D.: Mind your outliers! Investigating the negative impact of outliers on active learning for visual question answering. In: *Proceedings of the 59th ACL*. pp. 7265–7281 (2021)
- [11] Karumbaiah, S., Lan, A., Nagpal, S., Baker, R.S., Botelho, A., Heffernan, N.: Using past data to warm start active machine learning: Does context matter? In: *LAK21: 11th LAK*. pp. 151–160 (2021)
- [12] Lin, J., Singh, S., Sha, L., Tan, W., Lang, D., Gašević, D., Chen, G.: Is it a good move? mining effective tutoring strategies from human–human tutorial dialogues. *Future Generation Computer Systems* **127**, 194–207 (2022)
- [13] Lin, J., Tan, W., Du, L., Buntine, W., Lang, D., Gašević, D., Chen, G.: Enhancing educational dialogue act classification with discourse context and sample informativeness. *IEEE TLT* (in press)
- [14] Nye, B.D., Graesser, A.C., Hu, X.: Autotutor and family: A review of 17 years of natural language tutoring. *IJAIED* **24**(4), 427–469 (2014)
- [15] Nye, B.D., Morrison, D.M., Samei, B.: Automated session-quality assessment for human tutoring based on expert ratings of tutoring success. *International Educational Data Mining Society* (2015)
- [16] Rus, V., Maharjan, N., Banjade, R.: Dialogue act classification in human-to-human tutorial dialogues. In: *Innovations in smart learning*, pp. 185–188. Springer (2017)
- [17] Rus, V., Maharjan, N., Tamang, L.J., Yudelson, M., Berman, S., Fancsali, S.E., Ritter, S.: An analysis of human tutors’ actions in tutorial dialogues. In: *The Thirtieth International Flairs Conference* (2017)
- [18] Samei, B., Li, H., Keshtkar, F., Rus, V., Graesser, A.C.: Context-based speech act classification in intelligent tutoring systems. In: *International conference on intelligent tutoring systems*. pp. 236–241. Springer (2014)
- [19] Settles, B.: *Active learning*. Synthesis digital library of engineering and computer science, Morgan & Claypool, San Rafael, Calif. (2012)
- [20] Sha, L., Li, Y., Gasevic, D., Chen, G.: Bigger data or fairer data? Augmenting BERT via active sampling for educational text classification. In: *Proceedings of the 29th COLING*. pp. 1275–1285 (2022)
- [21] Swayamdipta, S., Schwartz, R., Lourie, N., Wang, Y., Hajishirzi, H., Smith, N.A., Choi, Y.: Dataset cartography: Mapping and diagnosing datasets with training dynamics. In: *Proceedings of EMNLP. ACL, Online* (2020)
- [22] Tan, W., Du, L., Buntine, W.: Diversity enhanced active learning with strictly proper scoring rules. *NeurIPS* **34** (2021)
- [23] Vail, A.K., Grafsgaard, J.F., Boyer, K.E., Wiebe, E.N., Lester, J.C.: Predicting learning from student affective response to tutor questions. In: *ITS*. pp. 154–164. Springer (2016)
- [24] Vail, A.K., Boyer, K.E.: Identifying effective moves in tutoring: On the refinement of dialogue act annotation schemes. In: *ITS*. Springer (2014)
- [25] Yang, Y., Loog, M.: Active learning using uncertainty information. In: *Proc. 23rd Int. Conf. Pattern Recognit.* pp. 2646–2651 (2016)