

ProtSi: Prototypical Siamese Network with Data Augmentation for Few-Shot Subjective Answer Evaluation

Yining Lu and Jingxi Qiu and Gaurav Gupta

AI-LAB

Wenzhou Kean University

{yiningl, qiujing, ggupta}@kean.edu

Abstract

Subjective answer evaluation is a time-consuming and tedious task, and the quality of the evaluation is heavily influenced by a variety of subjective personal characteristics. Instead, machine evaluation can effectively assist educators in saving time while also ensuring that evaluations are fair and realistic. However, most existing methods using regular machine learning and natural language processing techniques are generally hampered by a lack of annotated answers and poor model interpretability, making them unsuitable for real-world use. To solve these challenges, we propose ProtSi Network, a unique semi-supervised architecture that for the first time uses few-shot learning to subjective answer evaluation. To evaluate students' answers by similarity prototypes, ProtSi Network simulates the natural process of evaluator scoring answers by combining Siamese Network which consists of BERT and encoder layers with Prototypical Network. We employed an unsupervised diverse paraphrasing model ProTAugment, in order to prevent overfitting for effective few-shot text classification. By integrating contrastive learning, the discriminative text issue can be mitigated. Experiments on the Kaggle Short Scoring Dataset demonstrate that the ProtSi Network outperforms the most recent baseline models in terms of accuracy and quadratic weighted kappa.

1 Introduction

The examination is a crucial method for evaluating the academic performance of students, whether objectively or subjectively. Using pre-defined correct answers, it is simple to evaluate objective multiple-choice questions. Subjective questions, on the other hand, allow students to provide descriptive responses, from which the evaluator can determine the students's level of knowledge and assign grades based on their opinion.

In general, subjective responses are usually length-

ier and include far more substance than objective answers. Therefore, it takes a lot of effort and time for evaluators to manually assess subjective responses. In addition, the marks assigned to the same subjective questions can be varied from evaluator to evaluator, depending on their ways of evaluating, moods at the evaluating moment, and even their relationship with students (Bashir et al., 2021). Therefore, automated evaluation for subjective answers with a consistent process becomes necessary. Not only can the human efforts in this repetitive task be saved and spent on other more meaningful education endeavors, but the evaluation results will also be fair and plausible for students. Due to the recent Covid-19 pandemic, most universities and colleges have shifted their examinations to the online mode. Therefore, developing a subjective answer evaluation model dovetails with the actual needs of the school and has vast application value. Subjective answer evaluation is a sub-field of text classification involving comparing student answers with the model answer and classifying the result into grade classes. Because questions are domain-specific, automatically evaluating answers is a challenging task because of the manually labeled data scarcity and heterogeneity of text (Iwata and Kumagai, 2020).

Most of the proposed evaluation methods (Patil et al., 2018; Sakhapara et al., 2019; Bhonsle et al., 2019; Johri et al., 2021; Girkar et al., 2021; Bhonsle et al., 2019) rely on traditional statistical and NLP techniques. Their evaluation methods consist of four main parts: preprocessing (tokenization, stopwords removal, chunking, POS tagging, lemmatizing), vectorizing (BoW, TF-IDF, Word2Vec, LSA), computing similarity score (cosine similarity), and classifying (KNN). However, these preprocessing steps are time-consuming and cumbersome to achieve significant results, especially for long texts which are heterogeneous and noisy. Since the invention of BERT (Devlin et al., 2018) and

transfer learning in NLP (Howard and Ruder, 2018; Radford et al.), recent deep learning models, which are robust and easier to train, outperform most approaches in NLP tasks with minimal feature engineering. Thus, there is still much room to improve existing automatic evaluation methods using advanced deep learning methods. Inspired by that, the efforts (Krishnamurthy et al., 2019) attempted to test the performance of CNN, LSTM, and BERT on shot answer scoring tasks, but without other model structure design and fine-tune. Additionally, the training dataset only has 1700 examples per question prompt, which is insufficient to train a deep neural network. Its highest accuracy comes from BERT and is only 71%. Recent research (Bashir et al., 2021) proposes a 2-step approach to evaluate descriptive questions using machine learning and natural language processing, and Sakhapara et al. (2019) applies information gained in computing the final score. Despite the effectiveness of these methods, Bashir et al. (2021) hired 30 annotators to prepare sufficient labeled data that requires high cost and is time-consuming in a real-world application, and Sakhapara et al. (2019) needs a large pre-graded training dataset to compute the score instead of using the model answer.

To handle the problem of limited labeled data, we first consider subjective answer evaluation as a meta-learning problem: the model is trained to classify students' answer scores from a consecutive set of small tasks called episodes. Each episode contains support and query two sets with a limited number of N classes and a limited number of K labeled data for each class. That is usually known as N -way K -shots setup in meta-learning. Recently, meta-learning has been widely applied in NLP tasks, including emotion recognition in conversation (Guibon et al., 2021), named entity recognition (Das et al., 2021) or hypernym detection (Yu et al., 2020). The score evaluating process must involve verifying students' answers with model answers (answers from teachers). This step perfectly matches the function of the Siamese network commonly used in the Computer Vision field for image matching (Melekhov et al., 2016). We adopt the Siamese network to compare students' and model answers, then output a similarity vector for further classification. There are several metric-based meta-learning algorithms: Matching network (Vinyals et al., 2016) finds the weighted nearest neighbors; Prototypical Networks (Snell et al., 2017) com-

putes the class representations (prototypes) and classifies query based on Euclidean distance; Relation network (Sung et al., 2018) replaces the Euclidean distance by the deep neural network. In this work, we design a Prototypical network-based model, because it is the most effective when training data is insufficient in reality (Al-Shedivat et al., 2021) and it is the state-of-the-art for some NLP tasks, such as intent detection (Dopierre et al., 2021a), when equipped with a BERT text encoder. One challenge from meta-learning is that the models can easily overfit the biased distribution introduced by a few training examples (Ippolito et al., 2019). Inspired by the recent work PROTAUGMENT (Dopierre et al., 2021b), we utilize a pre-trained paraphrasing model that generates various paraphrases for the original randomly sampled text as data augmentation. Through data augmentation, we can incorporate more knowledge into the model by computing an unsupervised loss, thus making the model more robust. Another common challenge in subjective answer evaluation is the discriminative text problem which indicates two answers might have similar semantic representations but belong to different classes. Through data augmentation, we can tackle this issue by contrastive learning. The followings are the main contributions of this paper:

1. To the best of our knowledge, our work is the first of its kind to attempt to apply meta-learning to the evaluation of subjective answers.
2. We propose a novel semi-supervised approach, ProtSi Network, for evaluating students' answers from model answers.
3. We utilize the data augmentation to deal with overfit and heterogeneous text problems.
4. Through extensive experimental analysis, we demonstrate that the proposed approaches outperform the state-of-the-art on questions from the Kaggle dataset.

2 Problem Formulation

The meta-learning of few-shot text classification aims to learn from a set of small tasks or episodes and then transfer the learned knowledge to the testing tasks which are unseen during training.

Let the student answer dataset D be split into dis-

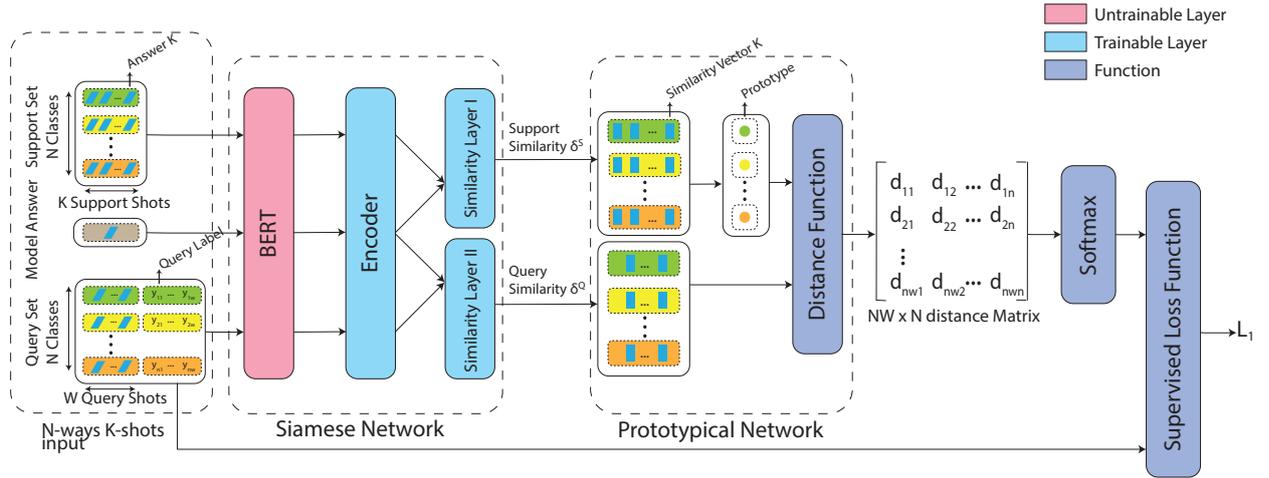


Figure 1: Framework of Supervised ProtSi Network

joint D_{train} , D_{val} , D_{test} . At each episode, a task consists of support set S and query set Q is drawn from either D_{train} , D_{val} , and D_{test} for training, validation, or testing. In one episode, N classes are randomly sampled from score class set (called rubric range). For each of the N classes, K students' answers are randomly sampled to compose the support set; and W different student's answers are sampled to compose the query set. Let a pair (x_i^S, y_i^S) denote the i^{th} item of total $N \times K$ items in support set. (x_i^Q, y_i^Q) with similar meaning denotes query set. The problem becomes training a meta-learner through consecutive small tasks that attempt to classify the students' answers in the query set Q based on the given model answer m and a few scored answers in the support set S .

3 Methodology

In this section, we present our semi-supervised model ProtSi Network. In order to assess student answers based on similarity prototypes, ProtSi Network combines the Siamese network, which is constituted of a BERT and encoder layer, with the Prototypical network. It also incorporates unsupervised paraphrasing loss and unsupervised contrastive loss to address the overfitting problem and discriminative text problem.

3.1 Supervised ProtSi Network

3.1.1 Siamese network

For a few-shot subjective answer evaluation, before assigning marks to students' answers, a comparison between students' answers and a model answer M is needed to map the textual and contextual difference into a similarity vector so that the further

metrics for classification can be computed. Unlike regular Siamese network consisting of twin networks that accept two distinct inputs (Koch et al., 2015), our Siamese network accepts three inputs x_{nk}^S, x_{nw}^Q, M and performs dual comparisons for M with x_{nk}^S and x_{nw}^Q separately, where x_{nk}^S and x_{nw}^Q are answers from support set and query set.

Specifically, in Siamese network, a text embedding layer is needed to map the textual data onto a high dimensional vector space where the similarity function between texts can be applied. Recently, more and more NLP tasks utilize pre-trained language models, such as BERT, to obtain text representations and achieve promising results. Following the work (Chen et al., 2022; Dopierre et al., 2021b; Hui et al., 2020; Mehri et al., 2020), we also use the pre-trained $BERT_{base}$ (Peters et al., 2018) as the embedding layer. The BERT model encodes the word information more effectively and flexibly because the context is considered. Thus, the numeric representations of a word can differ among sentences. For later use, we denote the $BERT_{base}$ embedding as function $e(x)$.

After obtaining word representation, a context encoder layer is needed to extract the whole sentence information and outputs a sentence feature vector. A recurrent neural network (RNN) is widely used in the NLP field because it can memorize the information from prior inputs and influence the current inputs and outputs. We express the embedding and encoding operations, including both the BERT and encoder layers, as the following equation:

$$x' = f(e(x)|\theta) \quad (1)$$

where θ denotes the trainable parameters in the encoding function $f(\cdot) : \mathbb{R}^E \rightarrow \mathbb{R}^M$, which convert the E-dimensional input x to the M-dimensional encoded output x' . To be noticed, support and query set share the same BERT and Encoder layers. The encoded support and query set are x'_i^S and x'_i^Q and encoded model answer is m' . Then we use two different similarity layers $s_1(\cdot, \cdot | \theta_1), s_2(\cdot, \cdot | \theta_2) : \mathbb{R}^M \rightarrow \mathbb{R}^M$ to draw the similarity representations between model answer with support set and query set separately. The similarity representations of answers in support set and query set can be computed as:

$$\begin{aligned}\delta_i^S &= s_1(x'_i^S, m' | \theta_1) \\ \delta_i^Q &= s_2(x'_i^Q, m' | \theta_2)\end{aligned}$$

The cardinality of set $\{\delta_i^S\}$ is NK and of set $\{\delta_i^Q\}$ is NW .

3.1.2 Prototypical network

The Prototypical network (Snell et al., 2017) is a few-shot classification model that computes prototypes for each class in the support set and assigns samples in the query set to classes based on their distances with prototypes. The prototypes are computed by averaging all the instance representations belonging to the same class in the support set. For our ProtSi Network with N-ways K-shots input, the prototypes are computed as:

$$\mathbf{c}_n = \frac{1}{K} \sum_{k=1}^K \delta_{nk}^S \quad (2)$$

where $\delta_{nk}^S \in \mathbb{R}^M$ is the similarity vector of answers k in the support set class n with model answer and similarity prototype $\mathbf{c}_n \in \mathbb{R}^M$ is the M-dimensional representation for class n . Given a distance function $d(\cdot, \cdot) : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}^1$, we can compute the distance $d(\mathbf{c}_n, \delta_i^Q)$ between prototypes \mathbf{c}_n and the i^{th} similarity vectors in the query set. Then the probability over classes for a query answer x_i^Q based on a softmax over distances to the prototypes can be computed as the following equation:

$$p(y = n | x_i^Q) = \frac{\exp(-d(\mathbf{c}_n, \delta_i^Q))}{\sum_{j=1}^N \exp(-d(\mathbf{c}_j, \delta_i^Q))} \quad (3)$$

Given the computed probability distribution p_i over score classes for the i^{th} answer in the query set and

the one-hot truth label y_i^Q , the supervised cross entropy loss can be computed as:

$$\mathcal{L}_1 = -\frac{1}{NW} \sum_{i=1}^{NW} y_i^Q \cdot \log(p_i) \quad (4)$$

where the inner product is used. The supervised ProtSi Network is shown in Figure 1.

3.2 Unsupervised Regularization

To alleviate the overfitting problem caused by a few training examples in one episode and the discriminative text problem, we propose to train the supervised ProtSi Network under the regularization of unsupervised paraphrasing loss and unsupervised contrastive loss. The overall unsupervised part of ProtSi Network is shown in Figure 2.

3.2.1 Data Augmentation

Along with the labeled answers randomly sampled at each episode, we also randomly choose S unlabeled answers from the whole dataset, including $D_{train}, D_{val}, D_{test}$. Then in the data augmentation step, we obtain Z paraphrases for each unlabeled answer. We use x_s to denote the s^{th} sampled unlabeled answer in the unlabeled sample set $U = \{x_1, x_2, \dots, x_S\}$ and \tilde{x}_s^z to denote the z^{th} paraphrase of x_s . The recent work (Dopierre et al., 2021b) proposes a diverse paraphrasing model PROTAUGMENT for data augmentation and it has shown to be effective in few-shot text classification. We apply PROTAUGMENT to do data augmentation in our unsupervised learning.

3.2.2 Paraphrasing Loss

After generating paraphrases for each unlabeled answer, we put the original answer, paraphrases, and model answer through the same ProtSi Network as shown in Figure 2 to create unlabeled similarity prototypes. Specifically, let δ_s denote the similarity vector of unlabeled answer x_s with model answer and $\tilde{\delta}_s^z$ denote the similarity vector of the z^{th} paraphrase of x_s . In Prototypical network, we derive unlabeled similarity prototypes of x_s by averaging all similarity vectors of its paraphrases:

$$\mathbf{c}_{x_s} = \frac{1}{Z} \sum_{z=1}^Z \tilde{\delta}_s^z \quad (5)$$

After obtaining the unlabeled similarity prototypes, we can compute the distance between δ_s and \mathbf{c}_{x_s} using the function $d(\cdot, \cdot)$. Given the computed $S \times$

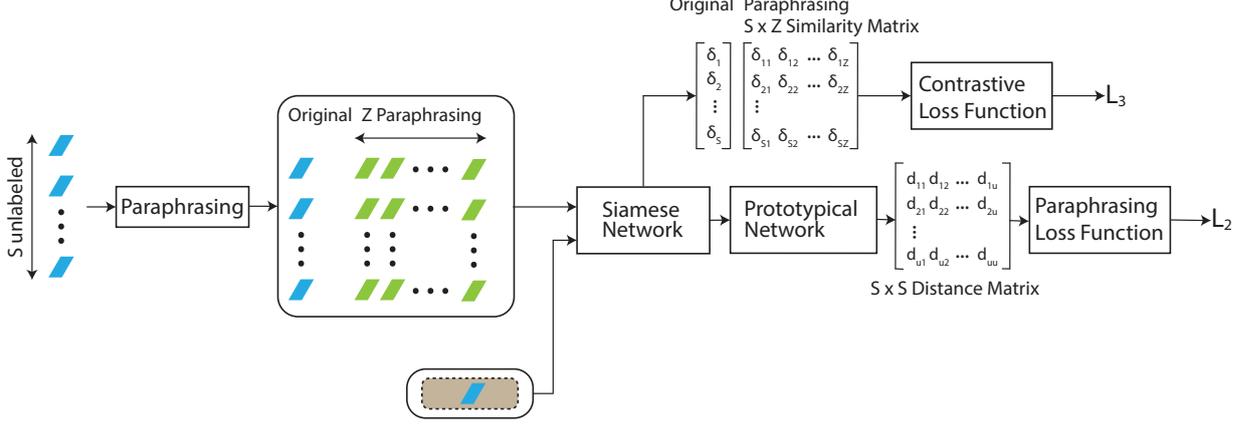


Figure 2: Framework of Unsupervised ProtSi Network

S distance matrix, we model the probability of each unlabeled similarity vector being assigned to each unlabeled similarity prototype as the Equation 3 but unsupervised:

$$p(s = t | x_s) = \frac{\exp(-d(\mathbf{c}_{x_t}, \delta_s))}{\sum_{x_s \in U} \exp(-d(\mathbf{c}_{x_s}, \delta_s))} \quad (6)$$

Given the probability distribution p_{x_s} of the unlabeled sample be assigned to unlabeled prototypes, we can compute an unsupervised entropy loss \mathcal{L}_2 using the following equation as the entropy regularization (Mnih et al., 2016):

$$\mathcal{L}_2 = -\frac{1}{S} \sum_{x_s \in U} p_{x_s} \cdot \log(p_{x_s} + \beta) \quad (7)$$

where β is a hyperparameter to prevent $\log(p_{x_s})$ being negative infinity when p_{x_s} is close to zero. \mathcal{L}_2 encourages the answers' similarity representation closer to its paraphrases' similarity prototypes and further away from the prototypes of other unlabeled similarity representations. According to Mnih et al. (2016) and Dopierre et al. (2021b), adding this unsupervised paraphrasing loss improves the model's exploration ability. It brings more knowledge and noise from unlabeled answers to the model, thus making it more robust to alleviate the overfitting problem.

3.2.3 Contrastive Loss

Classifying discriminative text is a big challenge in text classification tasks. A Prototypical network is unsatisfactory in learning discriminative texts with similar semantics and representations belonging to different classes. It simply learns the representations but ignores the relationship among texts.

To tackle the discriminative text problem in our ProtSi Network, we propose to use unsupervised contrastive loss that encourages learning discriminative text representations via contrastive learning motivated by its successful application in the few-shot text classification (Chen et al., 2022).

After getting the similarity matrix $\mathcal{D} \in \mathbb{R}^{S \times Z}$ of paraphrases from the Siamese network, we randomly sample L columns from \mathcal{D} so that each original similarity vector δ_s is paired with l new similarity vectors from paraphrasing. We combine all δ_s and sampled columns as a training batch $\{\delta\}$ of total $S(L+1)$ similarity vectors. δ_i denotes the i^{th} similarity vector and δ'_i denotes the paired similarity vector of δ_i in the batch $\{\delta\}$. The contrastive loss is:

$$\mathcal{L}_3 = -\frac{1}{S(L+1)} \sum_{\delta_i \in \{\delta\}} \log \frac{\sum_{\delta_i = \delta'_i} \exp(\frac{\delta_i \cdot \delta'_i}{\tau})}{\sum_{\delta_i = \delta'_i} \exp(\frac{\delta_i \cdot \delta'_i}{\tau}) + \sum_{\delta_i \neq \delta'_i} \exp(\frac{\delta_i \cdot \delta'_i}{\tau})} \quad (8)$$

Each similarity representation is paired and compared to all other representations. Thus the contrastive loss \mathcal{L}_3 encourages different similarity representations to be distant from each other and pull closer similarity representations belonging to the same original answers, which prevents the different answers with similar representations from being assigned the same score and alleviate the discriminative text problem.

3.2.4 Final Objective

The supervised Prototypical network in Figure 1 aligns the query similarity vectors to the similarity prototypes and generates an $NW \times N$ distance matrix by $d(\cdot, \cdot)$. Then the score of query answer

$x_i^Q \in Q$ is determined as the class of prototypes c_n who has the minimum distance with x_i^Q :

$$\hat{y}_i^Q = \arg \min_n d(c_n, \delta_i^Q) \quad (9)$$

During training, we combine supervised loss \mathcal{L}_1 and unsupervised loss \mathcal{L}_2 and \mathcal{L}_3 into the final loss \mathcal{L} using the following equation:

$$\mathcal{L} = t^\alpha \times \mathcal{L}_1 + (1 - t^\alpha) \times \mathcal{L}_2 + \gamma \times \mathcal{L}_3 \quad (10)$$

where α and γ are hyperparameters and t is the monotonically decreasing function with respect to the epoch. We use a loss annealing scheduler which will gradually incorporate the noise from the unsupervised loss \mathcal{L}_2 as training progresses.

To summarize, the final loss mainly uses the supervised loss so that the model learns to classify. Then incorporating more and more noise and knowledge from paraphrasing data using paraphrasing loss makes the model more robust to unseen answers. Finally, the contrastive loss is introduced to learn more distinct similarity representations of different answers to deal with discriminative text problems.

3.3 Overall Algorithm

The overall algorithm is shown in Algorithm 1.

4 Experiment

4.1 Datasets

We used the Short Answer Scoring dataset by the Hewlett Foundation from Kaggle. The dataset consists of 10 question prompts, which are from different domains. Questions 1, 2, 10 are from the Science subject, Questions 5, 6 are from Biology subject, Questions 7, 8, 9 are from English subject, and Questions 3, 4 are from English Language Art subject. There are around 1700 labeled answers per question with an average length of 50 words. Different questions have different model answers and rubric ranges.

4.2 Baselines

We compare our method with the following competitive baseline methods released in the recent two years¹:

¹Because the source code of **BERT**, **IG_WN**, **ASSESS**, and **SACS** are not released, we re-implemented all these four methods on the same Kaggle dataset by closely following the methods described in (Krishnamurthy et al., 2019), (Sakhapara et al., 2019), (Johri et al., 2021), and (Girkar et al., 2021) respectively.

BERT (Krishnamurthy et al., 2019) uses BERT base model with a dense layer for score classification using the labeled answers.

IG_WN (Sakhapara et al., 2019) uses information gain with WordNet to vectorize textual answers and apply KNN on the cosine similarity between labeled answer dataset and input answer.

ASSESS (Johri et al., 2021) proposes an Assessment algorithm consisting of semantic learning to compute the similarity score between model answer and input answer.

SACS (Girkar et al., 2021) proposes a Subjective Answer Checker System using NLP and Machine Learning to compute the final score from model answer and input answer.

4.3 Implementation Details

In the Encoder layer, we use three LSTM networks to extract the context information of sentences. Moreover, the encoded student’s answer and encoded model answer are concatenated to form the input of similarity layer 1 or similarity layer 2. Each similarity layer has one dense layer with a Relu activation function and is trained using batch normalization. We choose Euclidean distance as the distance function in the prototypical network. To demonstrate the efficiency of the proposed framework, we generated 8100 episodes of data for training and 2700 episodes for testing. For each episode or task, we set $N = 4$ for those questions with a rubric range from 0 to 3 and $N = 3$ for other questions ranging from 0 to 2. We set query shots $W = 1$ and support shots $K = 3$. Along with each episode, we randomly sample $S = 5$ unlabeled answers and generate $Z = 5$ paraphrases for each unlabeled answer with a 0.5 diversity penalty and 15 num beams for the Beam search algorithm. In the BERT layer, we set the maximum tokens per answer as 90 and the embedding size as 768. We set $L = 1$ for computing contrastive loss and $\alpha = 0.8$, $\beta = 0.01$, $\gamma = 0.02$ for computing the final loss. We evaluate our method and all baseline methods on the same dataset using the Accuracy (Acc) and the Quadratic Weighted Kappa (QWK) metrics.

4.4 Results

4.4.1 Main Results

Table 1 summarizes the testing results of four baseline methods and our ProtSi Network on the Short

Algorithm 1: ProtSi Network

Data: Model answer m , Student answer dataset D
Input: EPOCHS, EPISODES, N , K , W , S , Z , L
Output: Predicted score \hat{y}_i^Q

```
// Prepare Dataset
1 Split  $D$  into disjoint  $D_{train}$ ,  $D_{val}$ ,  $D_{test}$ ;
2 for  $episode \leftarrow 0$  to EPISODES do
3   Randomly sample  $N$ -ways  $K$ -shots support set  $S$  and  $N$ -ways  $W$ -shots query set  $Q$  from
   either  $D_{train}$ ,  $D_{val}$ , or  $D_{test}$  to create training dataset, validation dataset and testing dataset,
   respectively;
4   Randomly sample  $S$  unlabeled answers from the entire dataset  $D$ ;
5   Obtain  $Z$  paraphrases from the pre-trained paraphrasing model for each unlabeled answer;
6 end
// Training
7 for  $epoch \leftarrow 0$  to EPOCHS do
8   for  $x_i^S, x_i^Q, y_i^Q, x_s, \tilde{x}_s^z, m$  in training dataset, unlabel answer, paraphrases, model answer do
9      $\delta_i^S \leftarrow s_1(f(e(x_i^S)|\theta_0), f(e(m)|\theta_0)|\theta_1)$ ;
10     $\delta_i^Q \leftarrow s_2(f(e(x_i^Q)|\theta_0), f(e(m)|\theta_0)|\theta_2)$ ;
11     $\delta_s \leftarrow s_1(f(e(x_s)|\theta_0), f(e(m)|\theta_0)|\theta_1)$  // either  $s_1$  or  $s_2$ 
12     $\tilde{\delta}_s^z \leftarrow s_1(f(e(\tilde{x}_s^z)|\theta_0), f(e(m)|\theta_0)|\theta_1)$  // either  $s_1$  or  $s_2$ 
13     $\hat{y}_i^Q \leftarrow \delta_i^S, \delta_i^Q$  using equations 2, 9;
14     $\mathcal{L}_1 \leftarrow \delta_i^S, \delta_i^Q, y_i^Q$  using equations 2, 3, 4;
15     $\mathcal{L}_2 \leftarrow \delta_s, \tilde{\delta}_s^z$  using equations 5, 6, 7;
16    Randomly sample  $L$  columns from the matrix  $\mathcal{D}$  of  $\tilde{\delta}_s^z$  and combine with  $\delta_s$  as  $\{\delta\}$ ;
17     $\mathcal{L}_3 \leftarrow \{\delta\}$  using equation 8;
18     $\mathcal{L} \leftarrow \mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3$  using equation 10;
19    Perform backpropagation
20   end
21 end
// Testing
22 Perform the same steps from line 9-18 and compute accuracy and quadratic weighted kappa
```

Answer Scoring dataset with limited labeled data². Rather than taking advantage of model answers, IG_WN and BERT directly use labeled answers for supervised prediction, which is more like "learning for learning". Hence their Acc and QWK are generally higher than those of the other two baseline methods. However, ASSESS, SACS, and our ProtSi Network, which utilize the model answers for comparison and evaluation, are more dovetail with how real teachers grade answers. Furthermore, our ProtSi Network establishes new state-of-the-art performance on questions from different subjects, outperforming the ASSESS and SACS,

²The original works of ASSESS and SACS did not release their collected datasets. While the original works of BERT and IG_WN have tested their methods on the same Short Answer Scoring dataset from Kaggle, and the results we re-implemented basically match their results in the paper.

which also utilize model answers by an average Acc of 24.85%, 24.94%, and by an average QWK of 46.27%, 53.50%. Compared to the 'best' baseline methods, our ProtSi Network also has a competitive performance, outperforming the IG_WN by an average Acc of 7.84% and an average QWK of 39.59%, and outperforming the BERT by an average QWK of 12.12%.

Even if the BERT method has the highest average accuracy in directly predicting students' scores, our ProtSi Network is better at distinguishing discriminative texts and has the most reliable prediction results, according to its highest average quadratic weighted kappa score. For example, in Q_{10} , there are two answers with different scores.

Score 0 light gray: on cold days the light gray absorbs the hotter temp. and on hot days it does not

	IG_WN [†]		ASSESS		SACS		BERT [†]		ProtSi Network	
	Acc	QWK	Acc	QWK	Acc	QWK	Acc	QWK	Acc	QWK
Q1	0.5075	0.5438	0.4533	0.5815	0.3307	0.4074	0.5938	0.7779	0.6181	0.9054
Q2	0.3086	0.1826	0.3748	0.2907	0.2754	0.1189	0.5312	0.5254	0.4819	0.6691
Q3	0.4116	0.0724	0.2047	0.0191	0.2909	0.0425	0.4062	0.2241	0.4325	0.4000
Q4	0.5776	0.3317	0.2865	0.2262	0.4212	0.1361	0.7500	0.6266	0.7193	0.8197
Q5	0.8273	0.4652	0.2540	0.1513	0.4735	0.2797	0.8750	0.7591	0.7796	0.9067
Q6	0.8556	0.3282	0.3745	0.2192	0.3027	0.0762	0.8438	0.7393	0.7028	0.8611
Q7	0.5889	0.3661	0.3535	0.1676	0.3435	0.1171	0.7812	0.6474	0.6535	0.6429
Q8	0.4583	0.2565	0.5025	0.2927	0.4235	0.2453	0.6562	0.5433	0.6051	0.6667
Q9	0.5528	0.4919	0.4900	0.4101	0.4522	0.2978	0.6875	0.7354	0.6627	0.7931
Q10	0.4939	0.4071	0.5878	0.4194	0.5579	0.3337	0.8438	0.6146	0.7101	0.7407
Average	0.5582	0.3446	0.3881	0.2778	0.3872	0.2055	0.6968	0.6193	0.6366	0.7405

Table 1: Performance comparison on Short Answer Scoring dataset. IG_WN (Sakhapara et al., 2019), ASSESS (Johri et al., 2021), SACS (Girkar et al., 2021), and BERT (Krishnamurthy et al., 2019) results are re-implemented by ourselves. † represents the methods that do not utilize model answer to evaluate the input answer but use labeled data with supervised model to directly predict score.

absorb as much.

Score 2 light gray: Light gray would be a little cool in temperature as the black or dark gray one would get to hot and the white one would get a little cold.

The BERT method considered these two answers having the same context and assigned them to score 1. However, our ProtSi Network is more sensitive to the difference between the student and model answers. It realized that the second answer contained more information and precise expression, since the other colors were also taken into consideration. Therefore, our ProtSi Network marked it with a higher score and matched the correct human-made score.

4.4.2 Ablation Study

We conducted ablation studies on Q10 by removing one of the unsupervised losses. As shown in

\mathcal{L}_1	\mathcal{L}_2	\mathcal{L}_3	Acc	QWK
✓	✓	✓	0.7101	0.7407
✓	✓		0.6809	0.7018
✓		✓	0.6961	0.6557
✓			0.6909	0.6552

Table 2: Ablation study on Q10 by removing one of the unsupervised loss

Table 2, we found that removing \mathcal{L}_2 or \mathcal{L}_3 leads to

performance droppings. Removing \mathcal{L}_3 results in 2.92 point and 3.89 point Acc and QWK decline, reflecting that unsupervised contrastive loss does teach the model to classify discriminative answers and make the results more consistent with the real score distribution. The unsupervised paraphrasing loss \mathcal{L}_2 is intended to improve the robustness of the model, and removing \mathcal{L}_2 leads to 1.40 point and 8.5 point Acc and QWK decrease because the model performs poorly in classifying unseen answers. Removing both \mathcal{L}_2 and \mathcal{L}_3 resulted in a more severe performance reduction than just removing single loss.

5 Conclusion

In this paper, we presented the first study on subjective answer evaluation using few-shot learning. We proposed a novel semi-supervised architecture ProtSi Network following the teacher’s real scoring ways. Our ProtSi Network used data augmentation with unsupervised paraphrasing loss to solve the overfitting problem and used contrastive learning to alleviate the discriminative text problem. Experimental results on the Short Answer Scoring dataset showed that ProtSi Network outperformed the current baseline methods under the limitation of labeled data and demonstrated its robustness across different subjects.

Limitations

While the ProtSi Network seems suitable for real-world applications with a small number of labeled answers and outperforms the baseline methods, it still has limitations regarding its architecture. It uses an untrainable BERT layer to convert words to vectors without fine-tuning and might lose some performance. In ProtSi Network, we use a neural network to learn answer similarities and a simple Euclidean distance function to find the distances. While these yield better performance, they do not guarantee that the similarity of two high-dimensional vectors is correctly calculated, because the two vectors are far apart by the Euclidean distance but they may still be oriented closer together according to the Cosine similarity or other similarity measures in information retrievals. Another limitation of the ProtSi Network is that it does not consider the complexity of questions. There are multiple sub-questions under one primary prompt, and each sub-question may have its scoring rubric and rubric range. That requires the model to have a better learning ability to evaluate multiple sub-questions simultaneously.

Ethics Statement

we claim that our solution beats the performance of previous methods on a common open-source dataset. To demonstrate the improvement, We re-implement the baseline methods published in the recent three years by strictly following the papers. Our model has not achieved about 100% testing accuracy on the public datasets, so when educational systems want to deploy our model, please test with their private dataset to ensure high testing accuracy because subjective answer evaluation requires exceptionally high accuracy to ensure the fairness of scoring.

References

- Maruan Al-Shedivat, Liam Li, Eric Xing, and Ameet Talwalkar. 2021. [On data efficiency of meta-learning](#). In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 1369–1377. PMLR.
- Muhammad Farrukh Bashir, Hamza Arshad, Abdul Rehman Javed, Natalia Kryvinska, and Shahab S Band. 2021. Subjective answers evaluation using machine learning and natural language processing. *IEEE Access*, 9:158972–158983.
- Vishal Bhonsle, Priya Sapkal, Dipesh Mukadam, and Vinit Raut. 2019. An adaptive approach for subjective answer evaluation. *VIVA-Tech International Journal for Research and Innovation*, 1(2):1–6.
- Junfan Chen, Richong Zhang, Yongyi Mao, and Jie Xue. 2022. Contrastnet: A contrastive learning framework for few-shot text classification.
- Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J Passonneau, and Rui Zhang. 2021. Container: Few-shot named entity recognition via contrastive learning. *arXiv preprint arXiv:2109.07589*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina N. Toutanova. 2018. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- Thomas Dopierre, Christophe Gravier, and Wilfried Logerais. 2021a. [A neural few-shot text classification reality check](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 935–943, Online. Association for Computational Linguistics.
- Thomas Dopierre, Christophe Gravier, and Wilfried Logerais. 2021b. Protaugment: Unsupervised diverse short-texts paraphrasing for intent detection meta-learning. *arXiv preprint arXiv:2105.12995*.
- Abhishek Girkar, Mohit khambayat, Ajay Waghmare, and Supriya Chaudhary. 2021. Subjective answer evaluation using natural language processing and machine learning. *International Research Journal of Engineering and Technology*.
- Gaël Guibon, Matthieu Labeau, H el ene Flamein, Luce Lefeuvre, and Chlo e Clavel. 2021. Few-shot emotion recognition in conversation with sequential prototypical networks. In *The 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP 2021)*.
- Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*.
- Bei Hui, Liang Liu, Jia Chen, Xue Zhou, and Yuhui Nian. 2020. Few-shot relation classification by context attention-based prototypical networks with bert. *EURASIP Journal on Wireless Communications and Networking*, 2020(1):1–17.
- Daphne Ippolito, Reno Kriz, Jo ao Sedoc, Maria Kustikova, and Chris Callison-Burch. 2019. [Comparison of diverse decoding methods from conditional language models](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3752–3762, Florence, Italy. Association for Computational Linguistics.
- Tomoharu Iwata and Atsutoshi Kumagai. 2020. Meta-learning from tasks with heterogeneous attribute spaces. *Advances in Neural Information Processing Systems*, 33:6053–6063.

- Era Johri, Nidhi Dedhia, Kunal Bohra, Prem Chandak, and Hunain Adhikari. 2021. Assess-automated subjective answer evaluation using semantic learning. Available at SSRN 3861851.
- Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2, page 0. Lille.
- Surya Krishnamurthy, Ekansh Gayakwad, and Nallakuruppan Kailasanathan. 2019. Deep learning for short answer scoring. *International Journal of Recent Technology and Engineering*, 7:1712–1715.
- S Mehri, M Eric, and D Dialoglue Hakkani-Tur. 2020. A natural language understanding benchmark for task-oriented dialogue [j]. *arXiv preprint arXiv:2009.13570*.
- Iaroslav Melekhov, Juho Kannala, and Esa Rahtu. 2016. Siamese network features for image matching. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 378–383.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR.
- Piyush S. Patil, Vaibhav Miniyaar, and Amol. 2018. Subjective answer evaluation using machine learning.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training.
- Avani M. Sakhapara, Dipti Pawade, Bhakti Chaudhari, Rishabh Gada, Aakash Mishra, and Shweta Bhanushali. 2019. Subjective answer grader system based on machine learning. *Advances in Intelligent Systems and Computing*.
- Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30.
- Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. 2018. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. 2016. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Changlong Yu, Jialong Han, Haisong Zhang, and Wilfred Ng. 2020. Hypernymy detection for low-resource languages via meta learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3656. Online. Association for Computational Linguistics.

A More Implementation Details

The concrete model answers of questions from English and English Language Art subjects are not given in the dataset, so we randomly choose one fully marked student answer as the model answer to test ASSESS, SACS, and our ProtSi Network. We first generated the word embeddings of few-shot datasets and saved them as TFRecordDataset files. The generated dataset for each question occupies about 100G of storage, and it takes about 40 minutes to generate using multi-processing. We manually tune the hyperparameters and find that the learning rate as 10^{-3} and EPOCHS as 5 having the highest accuracy score. The monotonically decreasing function t is designed as:

$$t = \frac{EPOCHS - epoch}{EPOCHS + n} \quad (11)$$

where $epoch$ is the current epoch index and n is 1 to avoid $t = 1$. For each question, we use 80% of the answers to generate training tasks and 20% for validation tasks.

B Hardware and Software Configuration

All experiments are performed on Ubuntu Linux machine with 60-core AMD EPYC 7543 32-Core Processor with 360 GB RAM and 4 Nvidia A40 GPUs with 48 GM Memory. We use Python 3.8 and Tensorflow 2.5.0 with CUDA 11.2 as the deep learning framework.

C More Experimental Results

Figure 3 presents the distributions of the predicted score and the real score. We observe that:

- The distribution of predicted grades roughly matches the true distribution, especially for those full-mark answers.
- Some 0-point answers are misclassified as 1-point, illustrating that our model is not strict on scoring and the scoring strategy needs to be further defined and adjusted.

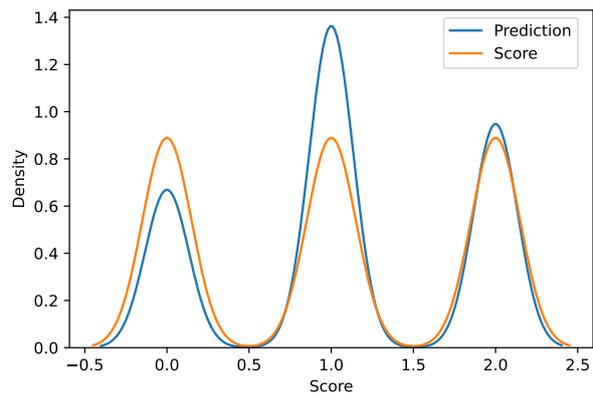


Figure 3: Predicted Score and real score distributions of Q_{10}