R2DE: a NLP approach to estimating IRT parameters of newly generated questions

Luca Benedetto luca.benedetto@polimi.it Politecnico di Milano Milan, Italy

Roberto Turrin

roberto.turrin@cloudacademy.com Cloud Academy Sagl Mendrisio. Switzerland

ABSTRACT

The main objective of exams consists in performing an assessment of students' expertise on a specific subject. Such expertise, also referred to as skill or knowledge level, can then be leveraged in different ways (e.g., to assign a grade to the students, to understand whether a student might need some support, etc.). Similarly, the questions appearing in the exams have to be assessed in some way before being used to evaluate students. Standard approaches to questions' assessment are either subjective (e.g., assessment by human experts) or introduce a long delay in the process of question generation (e.g., pretesting with real students). In this work we introduce R2DE (which is a Regressor for Difficulty and Discrimination Estimation), a model capable of assessing newly generated multiple-choice questions by looking at the text of the question and the text of the possible choices. In particular, it can estimate the difficulty and the discrimination of each question, as they are defined in Item Response Theory. We also present the results of extensive experiments we carried out on a real world large scale dataset coming from an e-learning platform, showing that our model can be used to perform an initial assessment of newly created questions and ease some of the problems that arise in question generation.

CCS CONCEPTS

• Computing methodologies \rightarrow Natural language processing; Ensemble methods; Information extraction; • Applied computing \rightarrow Education.

KEYWORDS

learning analytics, natural language processing, knowledge tracing, item response theory, latent traits estimation, educational data mining

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

LAK '20, March 23-27, 2020, Frankfurt, Germany

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-7712-6/20/03...\$15.00 https://doi.org/10.1145/3375462.3375517

Andrea Cappelli andrea.cappelli@cloudacademy.com Cloud Academy Sagl Mendrisio, Switzerland

> Paolo Cremonesi paolo.cremonesi@polimi.it Politecnico di Milano Milan, Italy

ACM Reference Format:

Luca Benedetto, Andrea Cappelli, Roberto Turrin, and Paolo Cremonesi. 2020. R2DE: a NLP approach to estimating IRT parameters of newly generated questions. In *Proceedings of the 10th International Conference on Learning Analytics and Knowledge (LAK '20), March 23–27, 2020, Frankfurt, Germany.* ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3375462.3375517

1 INTRODUCTION

Being able to estimate with a low degree of uncertainty the knowledge level of a student is crucial for providing effective material tailored to his expertise - and thus improving the learning experience. The task of estimating the skill level of a student by analysing the results of his interactions with assessment items (i.e., questions) is known as Knowledge Tracing (KT), and it is most commonly addressed with logistic models or neural networks [2].

Although Deep Knowledge Tracing (DKT) [19] - which is KT performed by means of neural networks - generally provides the highest accuracy in predicting the results of future answers [27, 32], logistic models and in particular dynamic Item Response Theory (IRT) models [26] are still used because of their explainability. Indeed, while in DKT the model is considered as a black-box, logistic models estimate latent traits of students and items, which enable a straightforward interpretation. In particular, with IRT models it is possible to estimate the skill level of each student and its evolution over time, as well as the difficulty and the discrimination of each question. The concept of difficulty is straightforward: if a question is more difficult than another, it requires a higher skill level to be answered correctly with the same probability. On the other hand, the discrimination determines how rapidly the odds of a correct answer increase or decrease with the skill level of the student. Therefore, discrimination can be used as a measure of the quality of an item. Indeed, questions with low discrimination provide little or no information about the skill level of the student answering, regardless of the difficulty of the item, because students of all skill levels have similar probability of answering correctly. Because of this, an estimation of the discrimination offers immediate feedback to content creators, who can modify the questions accordingly, in order to produce better (i.e., more discriminative) assessment items.

Having access to a history of exam results, the latent traits of the students and the questions involved in the exams can be estimated via likelihood maximization. In a similar way, when a student takes

an exam composed of calibrated questions (i.e., items whose latent traits are known), it is possible to estimate the skill level of the student from the results of the exam. Therefore, when a new question is created, it cannot be used for assessing students until a reliable estimation of its latent traits is performed. Also, some items might prove unsuited for assessing students (e.g., because of a discrimination which is too low or a difficulty which is too high or too low), thus having to be removed from the set of possible questions; it is important to do this as soon as possible.

A standard solution to the lack of an estimation of the latent traits of newly-created items, which is often referred to as the cold-start problem, consists in pretesting [29]: before using a newly developed item for assessment, it is administered to a certain number of students (usually few hundreds of few thousands) as if it was a regular exam question, but it is not used for scoring. On the contrary, the other questions of the exam are used for assessing the students. Then, the estimated skill level of such students is used together with the answers given to the item under pretesting to estimate its latent traits. Although this procedure leads indeed to an estimation of the latent traits of each item, it causes a long delay between the creation of an item and being able to use it for assessing students, and it also increases the development costs.

Another solution to the cold-start problem consists in using latent traits manually set by human experts: this approach enables the immediate usage of newly created questions in tests for assessing students, but it introduces a high uncertainty in the estimation, due to its nature intrinsically subjective.

In this work we introduce R2DE (a Regressor for Difficulty and Discrimination Estimation), a model that is capable of estimating the difficulty and the discrimination (as defined by Item Response Theory) of multiple-choice questions from the text of the questions and the text of the possible options. We present the results of extensive experiments performed on a real-world large scale dataset coming from an e-learning platform, showing that this model leads to a good estimation of the latent traits of new items, reducing both the importance of pretesting (making it necessary only for a fine-tuning of the initial estimation) and the uncertainty of the estimation of the latent traits (in comparison with latent traits manually set by human experts). Thus, it can be used as a first assessment of the difficulty and the discrimination of newly created questions, enabling an immediate usage in assessment tests and reducing the number of questions that have to be dropped after pretesting because of quality issues.

The contributions of this work are: i) the introduction of a novel model which uses natural language for estimating IRT latent traits of assessment items; ii) extensive testing of the model on a real world large scale dataset coming from an e-learning platform; iii) publication of the code used for implementing and testing this model, at https://github.com/lucabenedetto/r2de-nlp-to-estimating-irt-parameters.

The rest of the document is organized as follows: after an introduction to the current state of the art and the related works in Section 2, Section 3 focuses on IRT and on the performance prediction task in order to establish a common ground. R2DE is then presented in Section 4, followed by an introduction to the experimental dataset in Section 5 and a preliminary experiment for model choice and hyperparameter tuning in Section 7. The results of the

experiments are shown is Section 8 and, lastly, Section 9 concludes the paper.

2 RELATED WORK

The related work can be classified in the following categories: i) research about Knowledge Tracing (KT) and ii) research about Natural Language Processing (NLP) approaches for the estimation of questions' latent traits.

2.1 Knowledge Tracing

The concept of Knowledge Tracing (KT) was pioneered many years ago by Atkinson [3], but it is still extensively explored in research. As reported in a recent review by Abyaa et al. [2], the methods that are most commonly used in KT are logistic models and neural networks. Belonging to the family of logistic models are the approaches based on Item Response Theory (IRT) (e.g., dynamic IRT [26]) and the Elo rating system [23].

Recent literature claims that Deep Knowledge Tracing (DKT) - which was first introduced by Piech et al. in [19] and consists in performing KT by means of neural networks - outperforms logistic models in predicting the results of future exams [1, 6, 32, 33], but this advantage is not agreed across the community [7, 18, 28, 31]. Also, DKT does not estimate explicitly the skill level of students nor the latent traits of questions, which makes the interpretation of such models a strenuous task. There have been some attempts to make DKT explainable [14, 30], but they did not reach the same level of explainability as logistic models and are much more computationally expensive. Therefore, because of DKT being hard to explain and more complex from a computational point of view, logistic models and in particular models based on IRT are still widely used in the literature [7, 28].

Item Response Theory [10] estimates latent traits of students and items (i.e., questions) involved in an exam: the simplest model, named "Rasch model" [20], associates a skill level to each student and a difficulty level to each question. More complex models take into consideration additional latent traits [15] (e.g., the probability of correct answer by guessing): in this work we consider the twoparameter model, which associates a discrimination to each item. The discrimination determines how rapidly the odds of a correct answer increase or decrease with the skill level of the student and is a measure of how well an item can discriminate between students whose skill levels are above or below a certain threshold. Given a list of interactions between a set of students and a set of questions, latent factors of both students and questions can be estimated maximizing the likelihood of the observed results [24]. Then, the calibrated items can be used for assessing new students. Given a set of calibrated questions (i.e., questions whose latent traits are known) it is possible to estimate the skill level of the students answering those questions by likelihood maximization. Similarly, it is possible to leverage the answers of students that have already been assessed to estimate the latent traits of newly-created questions. The cold-start problem arises as soon as a new question is generated: since the latent factors of the item are unknown and there is no history of interactions to estimate them, it cannot be used for scoring students. A standard solution to this problem is pretesting: the newly-generated question is given to some students without

being used for assessing them and, looking at the answers provided by these students and at their skill level (assessed using other questions), the latent factors of the new question can be estimated [29]. This procedure introduces a long delay between the time when a question is generated and when it can be used for assessment. A possible solution to this problem consists in using the text of the newly-generated question to estimate its latent factors and remove pretesting from the pipeline (or, at least, reduce the number of students required for pretesting). To the best of our knowledge, in previous works, only Huang et al. in [12] explicitly mentioned that their method for estimating the difficulty of questions from their text could also be useful for targeting the cold-start problem. However, their model differs from the one presented in this paper because it estimates only the difficulty, without focusing on the discrimination; also, no code is available.

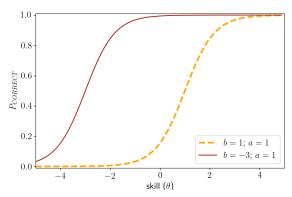
2.2 NLP for Latent Traits Estimation

The idea of using the text of a question to predict its difficulty is not new; however, most of previous works focused on readability estimation [9, 13, 29], which is a concept different from the difficulty defined in IRT. Wang et. al in [25] used textual information together with the interaction with users to estimate the difficulty of questions, but the focus was on Community Question Answering (CQA) services, thus considering a concept of difficulty different from the IRT-estimated difficulty used in this paper.

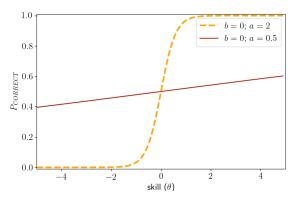
Closer to our paper are some more recent contributions that use NLP techniques to estimate the difficulty of assessment items, but they all define the difficulty as the fraction of wrong answers given to a question (referred to as "wrongness", from now on), which is less accurate than the IRT-estimated difficulty. One of such works is [11], in which the authors propose a neural model to predict the difficulty of reading problems in Standard Tests (i.e., problems whose answer can be found in a text that is given together with the question to the students) given the text of the document, the text of the question and the text of the possible answers. In [12] the authors use a neural network model to extract the Knowledge Components (i.e., the skills) related to each question from its text. A similar approach is adopted by Su et al. in [22], in which the accuracy of difficulty estimation from question's text is measured by looking at the precision in predicting the performance of the students in future exams. The main differences between these works and the present paper can be outlined in the following points: i) we use the IRT estimated difficulty as ground truth, which is more accurate than the "wrongness as difficulty" approach; ii) we estimate both the difficulty and the discrimination of the questions, thus providing a mean to assess the quality of newly-created questions.

3 THEORETICAL BACKGROUND

The objective of the model introduced in this work consists in using an NLP approach to estimate the difficulty and the discrimination of assessment items, as they are defined in IRT. Such estimations are then evaluated considering i) the accuracy with respect to ground truth values of the latent traits and ii) the accuracy in the performance prediction task. In order to establish common ground this



(a) Same discrimination a, different difficulty b.



(b) Same difficulty b, different discrimination a.

Figure 1: Example showing the effects of different difficulties and discrimination on the item response functions.

section presents an introduction i) to IRT, providing the mathematical explanation of the concepts of difficulty and discrimination, and ii) to the performance prediction task.

3.1 Item Response Theory

We use a two-parameter IRT model [10], which is characterized by three latent traits (the "two" refers to the number of items' latent traits): i) a skill level θ associated to each student, ii) a difficulty level b, and iii) a discrimination a associated to each assessment item. These latent traits are then used to compute the probability that student i correctly answers question j with the item response function:

$$P_{\text{CORRECT}} = \frac{1}{1 + e^{-a_j \cdot (\theta_i - b_j)}}$$

An example of the item response functions of two questions with equal discrimination and different difficulties is displayed in Figure 1a. According to the intuition, students with the same skill level (represented on the x-axis) have a lower probability of answering correctly the question with higher difficulty. The discrimination *a*, on the other hand, affects the steepness of the logistic curve, and that is the reason why it can be used as a measure of how well an item can discriminate between students whose skill level is above or

below a certain threshold (i.e., the difficulty of the question). Figure 1b shows the item response function of two questions with equal difficulty and different discrimination: the plot for the question with low discrimination is almost flat, showing that students with very different skill levels have similar probabilities of correctly answering the question. Thus, the information that can be gathered from that item is very limited.

Given the correctness of the answers that a student gave to a set of calibrated assessment items $Q = \{q_1, q_2, ..., q_{N_q}\}$, it is possible to estimate the knowledge level $\tilde{\theta}$ of the student. This is done by maximizing the result of the multiplication between the item response functions of the questions that were answered correctly and the inverse of the item response functions of the questions that were answered erroneously, with the following formula:

$$\tilde{\theta} = \max_{\theta} \left[\prod_{q_j \in Q_C} \frac{1}{1 + e^{-a_j \cdot (\theta - b_j)}} \cdot \prod_{q_j \in Q_W} \left(1 - \frac{1}{1 + e^{-a_j \cdot (\theta - b_j)}} \right) \right]$$

In the equation above Q_C is the set of questions correctly answered and Q_W is the set of questions that were answered erroneously.

3.2 The Performance Prediction Task

The latent traits estimated with IRT are non-observable by definition. Therefore, even though they can be considered as ground truth and are commonly used to evaluate the accuracy of a model (e.g., [29]), they have to be carefully dealt with. For this reason, in this work we validate our model not only by measuring the accuracy in predicting the latent traits of assessment items, but also by measuring its effects on the performance prediction task, which is the only way to work with an observable ground truth: the correctness of students' answers.

The performance prediction task consists in predicting the correctness of the answers given by a student to a sequence of assessment items. It can be used to measure how well the latent traits of the items are estimated by our model since these latent traits are a key element in predicting the correctness of students' answers. Given the ordered sequence of questions that a student interacted with, the correctness of each answer and an estimation of the latent traits of the items, it is performed as follows: i) given the latent traits of the first item and the estimated skill level of the student at that time (possibly unknown, in case of the first item), the correctness of his answer is predicted; ii) the actual answer is observed and compared to the predicted value in order to measure the accuracy of the prediction; iii) the actual answer is used to update the estimation of the skill level of the student; iv) this sequence of steps is repeated for all the assessment items the student interacted with.

In this work, the performance prediction task is used to evaluate the accuracy of our model in estimating the latent traits of assessment items. This is done by comparing the accuracy obtained in predicting the correctness of the answer using different algorithms to estimate the latent traits of the questions.

4 R2DE

This section introduces R2DE, which is a *Regressor for Difficulty* and *Discrimination Estimation*. The structure of R2DE, from the input question to the estimated latent traits, is shown in Figure

2. This section will focus on presenting the building blocks of the model and the steps that lead to the estimation of difficulty and discrimination from text.

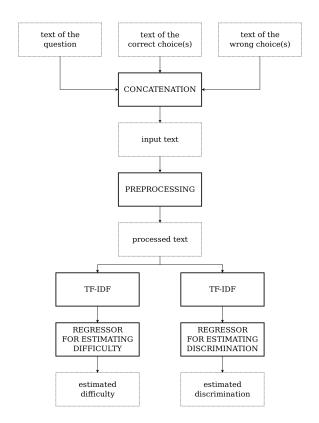


Figure 2: Structure of R2DE, from the input question to the estimated latent traits.

4.1 Data Model

The model introduced in this paper requires two different types of information. First of all, since R2DE works on text of multiple-choice questions, it needs the text of all the questions and the text of all the possible choices, as well as an indication of which choice contains the correct answer to each question. The model can also deal with the scenario where a question has multiple correct choices. This text information is used to generate the feature arrays that are used as input to the model.

Secondly, it requires the history of interactions between a set of students and a set of assessment items. For each interaction (i.e., the answer given by a student to a question), four fields of information are needed: i) a unique identifier of the student, ii) a unique identifier of the assessment item, iii) the correctness of the answer, and iv) the timestamp of the interaction. This data is used both to perform the IRT estimation of the latent traits of each question, which are used as ground truth for training R2DE, and to perform validation with the performance prediction task.

4.2 Features Engineering and Target Labels

This subsection shows the steps necessary to obtain the target labels and the feature arrays from data structured as presented in Subsection 4.1.

4.2.1 Input Features. The first step consists in creating an input text for each question. In this work three approaches (later referred to as encodings) were tested:

- *question_only*: considering only the text of the questions;
- question_correct: concatenating the text of the correct options (possibly more than one option is correct) to the text of the question;
- question_full: concatenating the text of all the possible options (both correct and wrong) to the text of the question.

Considering a fictitious example, let us assume that the student is shown the question "Which is the capital city of Germany?" and the possible answers are "London", "Berlin", "Madrid" and "Paris". Then, the body of text to represent the question would become, in the three cases: i) "Which is the capital city of Germany", ii) "Which is the capital city of Germany Berlin", and iii) "Which is the capital city of Germany London Berlin Madrid Paris". The outcome of this first step consists - in all three cases - in an input text characterizing each question.

The second step consists in preprocessing the corpus made of all the input texts using standard steps of NLP: removal of stop words, removal of punctuation, and stemming [17].

The third step consists in creating arrays of features from the input text of each item using a technique from Information Retrieval: TF-IDF (i.e.,Term Frequency-Inverse Document Frequency) [16]. The TF-IDF weight represents how important is a word (or a set of words) to a document in a corpus. The importance grows with the number of occurrences of the word in the document but it is limited by its frequency in the whole corpus: intuitively, words that are very frequent in all the documents of the corpus are not important to any of them. In particular, the formula used to compute the TF-IDF weight of word w in document d belonging to the corpus $C = \{d_1, d_2, ..., d_{N_d}\}$ is the following

$$TFIDF(w, d, C) = count(w, d) \cdot \left(\log_e \frac{N_d + 1}{count(w, C) + 1} + 1 \right)$$

where count(w, d) is the number of occurrence of w in document d, N_d is the number of documents in the corpus C and count(w, C) is the number of documents in the corpus C where w appears.

Lastly, since the feature set produced as outcome of TF-IDF is too large to be directly used as input of R2DE (i.e., one feature for each stemmed word generated from the original corpus) and standard approaches for dimensionality reduction such as Principal Component Analysis (PCA) would heavily affect the possibility of understanding the impact of specific concepts on the latent traits of the questions, a simpler method was used to reduce the size of the feature set. Specifically, only the top N_W features are considered: this is done by sorting the features (i.e., the stemmed words as obtained with preprocessing) according to their number of occurrences in the corpus and keeping only the N_W most frequent ones. This threshold (N_W) is considered as one of the parameters of the model and therefore can be chosen with cross-validation in order to be tuned for a specific dataset (as will be shown in Section 6).

4.2.2 Target Labels. The target labels are the latent traits (specifically, difficulty and discrimination as defined in IRT) of the items in the dataset. The latent traits can be estimated from the history of interactions between students and questions with a two-parameter IRT model.

4.3 The Regression Algorithm

R2DE contains two regressors that work in parallel to estimate i) the difficulty and ii) the discrimination of multi-choice questions. For both the elements a set of different algorithms should be tested in order to choose the ones that perform better on a specific dataset. Specifically, we tested Random Forests, Decision Trees, Support Vector Regression, and Linear Regression.

Using the same approach as in [29], model choice and tuning of the parameters are performed with 5-fold cross validation. Differently from previous works, we also use cross validation to choose which one of the three encodings described in Subsection 4.2 - i.e., i) question's text only, ii) question's text and correct answer's text, iii) question's text plus text of all the possible answers - to use and N_W , the number of most frequent keywords that should be used for the estimation. The fact that we can tune all these parameters makes this model more flexible and likely to perform comparably well on several datasets.

5 EXPERIMENTAL DATASET

To the best of our knowledge there are no public datasets containing both the text of the questions and the results of the answers, and all previous works experimented on private data collections. Our model as well is evaluated on a private database, which is a sample of actual data provided by the e-learning provider *Cloud Academy*¹. In particular, two data collections are used in this work, according to the required data format described in Subsection 4.1: i) one contains the information about the assessment items, ii) the other contains the information about the interactions between the students and the questions (i.e., the answers given by the students to the assessment items).

5.1 Question Dataset

In total, there are about 10K questions and the average number of possible choices is 4.2; the distribution of questions per number of choices is displayed in Table 1. In any case, regardless of the number of possible choices, when a question is prompted to a student only four of the possible choices are shown, among which there are always the correct ones.

Table 1: Distribution of questions per number of possible choices.

# of choices	percentage of items
4	86%
5	8%
6	5%
> 6	< 1%

¹https://cloudacademy.com/

The average length of the questions is 26.75 words and the answers are on average 6.83 words long, but the length of both the questions and the possible choices varies considerably. Table 2 presents the distribution of questions and Table 3 presents the distribution of choices per length.

Table 2: Distribution of questions per length.

length (# of words)	percentage of items
len <= 5	1%
5 < len <= 10	10%
10 < len <= 20	39%
20 < len <= 50	37%
len > 50	13%

Table 3: Distribution of possible choices per length.

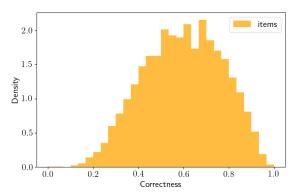
length (# of words)	percentage of choices
len = 1	21%
1 < len <= 5	35%
5 < len <= 10	20%
len > 10	24%

5.2 Interactions Dataset

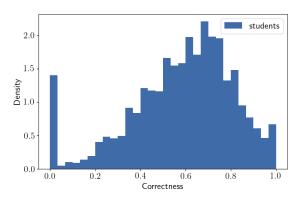
The interaction dataset used for the experiments contains about 2.3M interactions, collected over two years and involving a total of about 17K distinct students and 8K distinct assessment items. Overall, the interactions with a correct answer are the 64.69%, but this varies considerably depending on the specific students and items. This dataset was built from the original sample of data provided by Cloud Academy in order i) to contain only "first timers" (i.e., for each student-item pair, only the answer given during the first attempt is considered), and ii) to contain only questions that were answered by at least 100 distinct students. This was done for two different reasons: i) students that answer several times a specific question are more likely to answer correctly, thus impacting the accuracy of the latent traits estimated with IRT; ii) since the objective of this work consists in using textual information for estimating the latent traits of questions and not in analyzing the effectiveness of IRT, working on items with low support would affect the IRT estimation of questions' latent traits. Having the most possible accurate IRT estimation of difficulty and discrimination is crucial: in fact, the IRT-estimated latent traits are considered as ground truth while training R2DE and, in case of a bad estimation, we would train our model with noisy samples.

Figure 3 displays the distribution of students and questions per correctness, showing that both present a Gaussian-shaped distribution, although - in the case of students - the distribution shows two peaks for values of correctness close to 0 and 1. This is probably due to the fact that some students have low support (i.e., they are involved in few interactions). However, differently from what was done for the items, students with low support were not removed in order not to reduce too much the size of the dataset.

Some additional statistics about the dataset are presented in Table 4.



(a) Distribution of questions per correctness.



(b) Distribution of students per correctness.

Figure 3: Distribution of students and items per correctness after filtering questions answered by less than 100 students.

Table 4: Statistics about the Cloud Academy dataset, after data cleaning.

value	mean	std. dev.
# interactions per student	130.54	193.19
# interactions per item	283.03	407.47
correctness per student	58.25%	22.86%
correctness per item	59.46%	17.30%

6 EXPERIMENTAL SETUP

Training of our model is performed in two steps: i) an IRT model is trained in order to estimate the "true" difficulty and discrimination of each question, then ii) these IRT-estimated latent traits are used as target labels (i.e., ground truth) to train R2DE, which gets the text as input. Therefore, in order to avoid any leaks of information between the training data and the test data, two different splits are performed on the dataset presented in Section 5. Figure 4 displays how the four datasets are generated from the interaction dataset and the question dataset with the two splitting operations mentioned above.

INTERACTIONS DATASET USER ID TIMESTAMP QUESTION ID CORRECT IRT (estimated difficulties, $\mathrm{DS}_{\mathrm{GTE}}$ 70% estimated discriminations) **ESTIMATION** ${\rm DS_{V\!AL}}$ 30% QUESTIONS DATASET CHOICES TARGETS EXPERIMENTS ON PERFORMANCE PREDICTION MODEL trained 80% SELECTION regressor EXPERIMENTS ON 20% LATENT TRAITS TEST **ESTIMATION**

Figure 4: Experimental setup.

First of all, a 70:30 split stratified on the questions is performed on the interaction dataset. This leads to the generation of two smaller datasets: the interactions were randomly split, but we imposed the constraint of having at least one entry for each question; let us call the larger one $\mathrm{DS}_{\mathrm{GTE}}$, since it is used for the Ground Truth Estimation (i.e., the initial estimation of the IRT latent traits), and the smaller one $\mathrm{DS}_{\mathrm{VAL}}$, since it will be used to validate the latent traits estimated with of R2DE on the performance prediction task.

The second split is performed on the questions of the question dataset with a 80:20 rate, thus generating two smaller datasets containing non-overlapping sets of questions. The larger of the two question datasets generated from this split (DS $_{TRAIN}$) is used to train R2DE, while the smaller (DS $_{TEST}$) is used to test the estimation of the latent traits from the input text.

Thanks to this split performed in two different steps, it is possible to i) perform the ground truth estimation of the IRT latent traits using the interactions stored in $\mathrm{DS}_{\mathrm{GTE}}$; ii) train R2DE on $\mathrm{DS}_{\mathrm{TRAIN}}$; iii) test its capability of estimating latent traits form text on $\mathrm{DS}_{\mathrm{TEST}}$; and iv) validate it on $\mathrm{DS}_{\mathrm{VAL}}$ by measuring its accuracy on the performance prediction task.

7 MODEL CHOICE

The ground truth latent traits of each question are estimated, using the interactions stored in $\mathrm{DS}_{\mathrm{GTE}}$, with a two-parameter IRT model implemented with $pyirt^2$. Then, several regression models are tested with five-fold cross validation in order to find the best configuration for each of them. This is done using the input text of the questions in $\mathrm{DS}_{\mathrm{TRAIN}}$. The following models are tested, and for all of them the scikit-learn³ implementation is used:

- Random Forest (RF) Regressor [4], with the following hyperparameters:
 - $n_{estimators} = [10, 25, 50, 100, 150, 200, 250]$
 - $max_depth = [2, 5, 10, 15, 25, 50]$

- Decision Tree (DT) Regressor [5], with the following hyperparameters:
 - *max_features* = [1, 2, 3, 4, 5, None]
 - $max_depth = [2, 5, 10, 20, 50]$
- Linear Regression (LR)[21], with the following hyperparameters:
 - normalize = [True, False]
- Support Vector Regression (SVR)[8], with the following hyperparameters:
 - kernel = ['linear', 'poly', 'rbf']
 - gamma = ['auto', 'scale']
 - shrinking = [True, False]
 - degree = [1, 2, 3, 4]

Also, each configuration is tested on the three encodings presented in Subsection 4.2: i) $question_only$, ii) $question_correct$, iii) $question_full$. Lastly, each configuration is tested after performing dimensionality reduction with different values of N_W : in particular, values in the range [100, 2000] are tested. This preliminary experiment is performed twice, in order to choose the model configuration for both difficulty and discrimination estimation.

Table 5 and Table 6 display the results obtained with this initial experiment, showing the Mean Squared Error (MSE) for difficulty estimation and discrimination estimation respectively. All the results presented in the table were obtained with the best performing configuration of hyperparameters and N_W for each "family" of models and for each encoding.

For both latent traits the best performing model was the Random Forest (RF) regressor, with the input text encoded using the *question_only* encoding. The configurations of hyperparameters leading to the lowest error, instead, were different for the two latent traits: i) for difficulty estimation, the RF is made of 250 estimators, each with a maximum depth of 50 and N_W is 1000; ii) for discrimination estimation the RF is composed of 100 estimators with 25 maximum depth; the considered N_W is 800.

We also analyzed the effects of varying the value of N_W and observed that it does not have a significant impact on the error

²https://github.com/17zuoye/pyirt

³https://scikit-learn.org/

Table 5: Mean Squared Error for difficulty estimation, results of the preliminary experiment for model choice.

	question_only	question_correct	question_full
Model	MSE	MSE	MSE
RF	0.359	0.306	0.306
DT	0.790	0.767	0.757
LR	0.719	0.691	0.669
SVR	0.438	0.391	0.399

Table 6: Mean Squared Error for discrimination estimation, results of the preliminary experiment for model choice.

	question_only	question_correct	question_full
Model	MSE	MSE	MSE
RF	0.178	0.188	0.123
DT	0.196	0.195	0.189
LR	0.190	0.190	0.185
SVR	0.187	0.188	0.200

of the Random Forest. Figure 5a and Figure 5b show the MSE on difficulty estimation and discrimination estimation respectively, obtained by the best configurations of the RF for varying value of N_W .

8 RESULTS

This section presents and discusses the results of the experiments carried on on the data.

8.1 Latent Traits Estimation

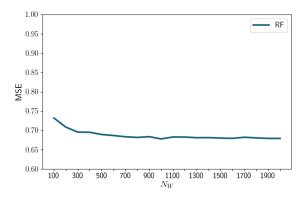
First of all, we test the capability of the selected model configuration to estimate latent traits of new questions given the text of the question and the text of the possible answers. This is done by estimating with R2DE the latent traits of the questions in DS_{TEST} and comparing them with the IRT estimated values. The Root Mean Squared Error (RMSE) and the Mean Absolute Error (MAE) for difficulty estimation and discrimination estimation are presented respectively in Table 7 and Table 8. The tables also show the relative errors (i.e., Relative RMSE and Relative MAE), which represent the errors measured relatively to the range of possible values of difficulty and discrimination. Specifically, the Relative RMSE for the difficulty is computed as:

Relative RMSE =
$$\frac{RMSE}{max_difficulty - min_difficulty}$$

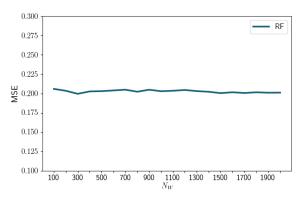
Similarly, it can be computed the Relative RMSE for the discrimination and the Relative MAE for both latent traits. For our experimental dataset, in particular, the IRT model was trained in order to have difficulties in the range [-5; 5] and discriminations in the range [-1; 2.5].

Table 7: Test of difficulty estimation.

RMSE	Relative RMSE	MAE	Relative MAE
0.823	8.23%	0.639	6.39%



(a) Mean Squared Error on difficulty estimation for varying N_W .



(b) Mean Squared Error on discrimination estimation for varying $N_{W'}$.

Figure 5: Effects of varying N_W on the Mean Square Error while estimating difficulty and discrimination.

Table 8: Test of discrimination estimation.

RMSE	Relative RMSE	MAE	Relative MAE
0.447	12.8%	0.329	9.4%

Figure 6 shows, as example, the comparison between the item response function obtained with the IRT estimation and the one obtained with the latent traits estimated with R2DE. The vertical lines represent the difficulty of the questions, implying the value estimated with IRT and the one estimated with R2DE.

Even though it is not possible to perform an actual comparison between the results of this work and previous ones, due to the fact that each research focused on a different (private) dataset, the analysis of the errors of the different approaches in difficulty estimation can still provide useful insight. This cannot be done for discrimination estimation, since R2DE is the first model that is capable of estimating the discrimination as well as the difficulty of assessment items from the input text. Table 9 compares the relative errors obtained in recent works. The table shows that the error obtained in this work is smaller than the errors obtained in

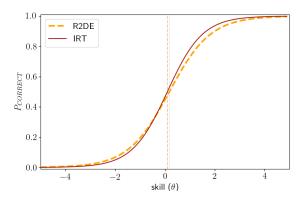


Figure 6: Comparison between the item response function obtained with the latent traits estimated with IRT and with R2DE.

previous works and, although this does not assure that this model is better performing than the others on any datasets, it suggests that simple regression models as ours could perform as well as - maybe even better than - more complex models (e.g., the Convolutional Neural Network with attention mechanism proposed in [11]).

Table 9: Comparison with state of the art.

Paper	Difficulty range	RMSE	Relative RMSE
R2DE	[-5; 5]	0.823	8.23%
Huang et al. [11]	[0; 1]	0.21	21%
Yaneva et al. [29]	[0; 100]	22.45	$\boldsymbol{22.4\%}$

8.2 Performance Prediction

In Subsection 8.1 we presented a comparison between the latent traits estimated with R2DE and the latent traits estimated with IRT, which are considered as ground truth. However, this is not an observable ground truth and, for this reason, we also validate R2DE measuring its efficacy in the performance prediction task, which offers - as presented in Section 3 - the only observable ground truth. In particular, the interactions stored in DSVAL are used to validate the model. The baselines that we use to test our model against are i) the "ground truth" latent traits estimated with IRT (which is a upper threshold), and ii) majority prediction.

Two different tests are carried out for performance prediction. First, we filter the validation dataset (DS_{VAL}) in order to keep only the test questions (i.e., the ones stored in DS_{TEST}). Then, only the test questions are used both for measuring the accuracy of the prediction and for updating the estimation of the skill level at each step. The results obtained with this experiment are presented in Table 10 displaying for each approach the accuracy (Acc.), the precision (Prec.) and the recall (Rec.) on the correct interactions, and the precision and the recall on the wrong interactions. The table shows that the latent traits estimated with R2DE lead, for most of the metrics, to values that are close to the ones obtained with ground truth latent traits and generally outperform majority

prediction. This means that R2DE is able to compute the latent traits from the question text with a performance similar to IRT that, instead, is based on hundreds of interactions.

Table 10: Results of the test on performance prediction, using only the questions in DS_{TEST} .

		Correct Answers		Correct Answers Wrong Answer		Answers
Approach	Acc.	Prec.	Rec.	Prec.	Rec.	
R2DE	0.662	0.704	0.794	0.562	0.442	
IRT	0.666	0.713	0.781	0.565	0.475	
Majority	0.625	0.625	1.0	-	0.0	

The second test explicitly reproduces the scenario that occurs in the wild. In real assessment only some of the items are newly generated (thus requiring an estimation of the latent traits from text), most of them are already calibrated (i.e., with known latent traits). Therefore, in this second experiment on performance prediction both the test questions and the train questions of DS_{VAL} are used. However, only the test questions are considered for evaluating the performance of the model on performance prediction; the train questions are used exclusively to update the estimated skill level of the student during the experiment. Anyway, the test questions are used also for updating the skill level estimation, as it was the case in the previous test. Table 11 displays the results obtained with this second experiment, using the same metrics as above. Majority is not reported here since it is computed as in the previous test and performed much worse than the other two approaches. Again, the latent traits estimated with R2DE proved good estimations for newly generated items. Indeed, the accuracy obtained with our model is only 1.57% lower than the accuracy obtained with IRT-estimated latent traits, which is the upper threshold.

Table 11: Results of the test on performance prediction on test interactions, skill estimated using all the questions.

		Correct Answers		Correct Answers Wron		Wrong	Answers
Approach	Acc.	Prec.	Rec.	Prec.	Rec.		
R2DE	0.689	0.744	0.767	0.590	0.559		
IRT	0.701	0.757	0.768	0.603	0.589		

9 CONCLUSIONS

In this work, we introduced R2DE, a model which is capable of estimating the latent traits (i.e., the difficulty and the discrimination) of newly generated multiple-choice questions by looking at the text of the question and the text of the possible choices.

Extensive experiments carried out on a large scale real world dataset provided by *Cloud Academy* showed that the model is capable of estimating with a low uncertainty both the difficulty and the discrimination. Specifically, it reached a MAE of 0.639 for difficulty estimation (6.39% of the whole difficulty range) and a MAE of 0.329 for discrimination estimation (9.4% of the overall discrimination range). R2DE is the first model estimating the discrimination as well as the difficulty, thus a comparison with the errors of other models was not possible for discrimination estimation. On the other hand,

a comparison with recent literature was performed for difficulty estimation; such comparison suggests that this model might be capable of performing at least as well as previously existing models. However, an extensive comparison on the same datasets was not possible due to the unavailability of code from previous research and all the datasets being private.

We showed that this model improves the accuracy on the task of exam results prediction with respect to using a simple estimation such as majority estimation and reaches an accuracy comparable to the upper threshold obtained with the IRT-estimated latent traits. Therefore, it is fair to say that R2DE reduces the importance of pretesting newly generated questions. Indeed, the estimated latent traits only require some fine-tuning, which is much faster than performing pretesting from scratch involving few hundreds or few thousands students. Moreover, having an estimation of the difficulty and the discrimination of the items is also useful for content creators at the time of writing the questions: they can have immediate feedback and modify the question accordingly, before deploying it. Thus, this model enables a reduction in the number of questions that have to be removed from the set of assessment items due either to a too low discriminative power or to a too high (or too low) difficulty.

Future work will continue to explore this research direction, focusing in particular on the following aspects: i) using advanced embeddings for encoding the text of the assessment items; ii) analyzing the importance of specific keywords, in order to understand whether the very structure of questions and items - and not few keywords - is the reason why some questions are more difficult and more discriminative than others; iii) exploring the possibilities of explaining the predictions of the model. Also, we aim at making possible a comprehensive comparison between this model and all the other models that use NLP approaches to estimate the latent traits of questions; for this reason, we make the code available for future research⁴, so that it might be run and tested on other datasets.

REFERENCES

- [1] Ghodai Abdelrahman and Qing Wang. 2019. Knowledge Tracing with Sequential Key-Value Memory Networks. (2019).
- [2] Abir Abyaa, Mohammed Khalidi Idrissi, and Samir Bennani. 2019. Learner modelling: systematic review of the literature from the last 5 years. Educational Technology Research and Development (2019), 1–39.
- [3] Richard C Atkinson. 1972. Ingredients for a theory of instruction. American Psychologist 27, 10 (1972), 921.
- [4] Leo Breiman. 2001. Random forests. Machine learning 45, 1 (2001), 5–32.
- [5] Leo Breiman, Jerome Friedman, Charles J Stone, and RA Olshen. 1984. Classification and Regression Trees. CRC Press.
- [6] Penghe Chen, Yu Lu, Vincent W Zheng, and Yang Pian. 2018. Prerequisite-Driven Deep Knowledge Tracing. In 2018 IEEE International Conference on Data Mining (ICDM). IEEE, 39–48.
- [7] Xinyi Ding and Eric Larson. 2019. Why Deep Knowledge Tracing has less Depth than Anticipated. In The 12th International Conference on Educational Data Mining.
- [8] Harris Drucker, Christopher JC Burges, Linda Kaufman, Alex J Smola, and Vladimir Vapnik. 1997. Support vector regression machines. In Advances in neural information processing systems. 155–161.
- [9] William H DuBay. 2004. The Principles of Readability. Online Submission (2004).
- [10] Ronald K Hambleton, Hariharan Swaminathan, and H Jane Rogers. 1991. Fundamentals of item response theory. Sage.
- [11] Zhenya Huang, Qi Liu, Enhong Chen, Hongke Zhao, Mingyong Gao, Si Wei, Yu Su, and Guoping Hu. 2017. Question Difficulty Prediction for READING Problems in Standard Tests. In Thirty-First AAAI Conference on Artificial Intelligence.
- $^4 https://github.com/lucabenedetto/r2de-nlp-to-estimating-irt-parameters\\$

- [12] Zhenya Huang, Yu Yin, Enhong Chen, Hui Xiong, Yu Su, Guoping Hu, et al. 2019. EKT: Exercise-aware Knowledge Tracing for Student Performance Prediction. IEEE Transactions on Knowledge and Data Engineering (2019).
- [13] Walter Kintsch and Douglas Vipond. 2014. Reading comprehension and readability in educational practice and psychological theory. Perspectives on learning and memory (2014), 329–365.
- [14] Jinseok Lee and Dit-Yan Yeung. 2019. Knowledge Query Network for Knowledge Tracing: How Knowledge Interacts with Skills. In Proceedings of the 9th International Conference on Learning Analytics & Knowledge. ACM, 491–500.
- [15] Eric Loken and Kelly L Rulison. 2010. Estimation of a four-parameter item response theory model. Brit. J. Math. Statist. Psych. 63, 3 (2010), 509–525.
- [16] Christopher Manning, Prabhakar Raghavan, and Hinrich Schütze. 2010. Introduction to information retrieval. Natural Language Engineering 16, 1 (2010), 100–103
- [17] Christopher D Manning, Christopher D Manning, and Hinrich Schütze. 1999. Foundations of statistical natural language processing. MIT press.
- [18] Ye Mao, Chen Lin, and Min Chi. 2018. Deep learning vs. bayesian knowledge tracing: Student models for interventions. JEDM Journal of Educational Data Mining 10, 2 (2018), 28-54.
- [19] Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein. 2015. Deep knowledge tracing. In Advances in neural information processing systems. 505–513.
- [20] Georg Rasch. 1960. Studies in mathematical psychology: I. Probabilistic models for some intelligence and attainment tests. (1960).
- [21] George AF Seber and Alan J Lee. 2012. Linear regression analysis. Vol. 329. John Wiley & Sons.
- [22] Yu Su, Qingwen Liu, Qi Liu, Zhenya Huang, Yu Yin, Enhong Chen, Chris Ding, Si Wei, and Guoping Hu. 2018. Exercise-enhanced sequential modeling for student performance prediction. In Thirty-Second AAAI Conference on Artificial Intelligence.
- [23] Josine Verhagen, David Hatfield, and Dylan Arena. 2019. Toward a Scalable Learning Analytics Solution. In International Conference on Artificial Intelligence in Education. Springer, 404–408.
- [24] Hua Wang, Cuiqin Ma, and Ningning Chen. 2010. A brief review on Item Response Theory models-based parameter estimation methods. In 2010 5th International Conference on Computer Science & Education. IEEE, 19–22.
- [25] Quan Wang, Jing Liu, Bin Wang, and Li Guo. 2014. A regularized competition model for question difficulty estimation in community question answering services. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP). 1115–1126.
- [26] Xiaojing Wang, James O Berger, Donald S Burdick, et al. 2013. Bayesian analysis of dynamic item response models in educational testing. The Annals of Applied Statistics 7, 1 (2013), 126–153.
- [27] Zhiwei Wang, Xiaoqin Feng, Jiliang Tang, Gale Yan Huang, and Zitao Liu. 2019. Deep Knowledge Tracing with Side Information. In *International Conference on Artificial Intelligence in Education*. Springer, 303–308.
- [28] Kevin H Wilson, Yan Karklin, Bojian Han, and Chaitanya Ekanadham. 2016. Back to the basics: Bayesian extensions of IRT outperform neural networks for proficiency estimation. (2016).
- [29] Victoria Yaneva, Peter Baldwin, Janet Mee, et al. 2019. Predicting the Difficulty of Multiple Choice Questions in a High-stakes Medical Exam. In Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications. 11–20.
- [30] Chun-Kit Yeung. 2019. Deep-IRT: Make Deep Learning Based Knowledge Tracing Explainable Using Item Response Theory. arXiv preprint arXiv:1904.11738 (2019).
- [31] Chun-Kit Yeung and Dit-Yan Yeung. 2018. Addressing two problems in deep knowledge tracing via prediction-consistent regularization. In Proceedings of the Fifth Annual ACM Conference on Learning at Scale. ACM, 5.
- [32] Jiani Zhang, Xingjian Shi, Irwin King, and Dit-Yan Yeung. 2017. Dynamic keyvalue memory networks for knowledge tracing. In Proceedings of the 26th international conference on World Wide Web. International World Wide Web Conferences Steering Committee, 765–774.
- [33] Liang Zhang, Xiaolu Xiong, Siyuan Zhao, Anthony Botelho, and Neil T Heffernan. 2017. Incorporating rich features into deep knowledge tracing. In Proceedings of the Fourth (2017) ACM Conference on Learning@ Scale. ACM, 169–172.