
Diagnostic Questions - The NeurIPS 2020 Education Challenge

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Tasks 1 and 2 of the NeurIPS 2020 Education Challenge are based on student
2 modelling and response prediction. This paper explores classical Item Response
3 Theory (IRT) based approaches to the problem. It can be demonstrated that simple 1
4 parameter IRT model provides reasonable prediction accuracy.

5 1 Task 1: Answer correctness prediction

6 Classical Item Response Theory has been widely used for estimating student ability and predicting
7 student performance [2, 7, 4, 1]. The simplest model, also called the one parameter model fits a
8 difficulty parameter for each question and an ability parameter for each student. The probability of
9 student s answering the question i correctly is the logit of the difference between student ability and
10 question difficulty. The following equation represents the One parameter IRT model

$$11 \quad P(C = 1) = \text{logit}(\sigma_s - \beta_i) \quad (1)$$

12 Where σ_s indicates the ability of the student s and β_i indicates the difficulty of the question i . IRT
13 models can be extended to incorporate multidimensional ability and difficulty parameters and called
14 as Multidimensional IRT (MIRT). [6] is a great resource on MIRT models and their applications. The
15 MIRT equations are the same as IRT except that each parameters is a vector. Furthermore, if the skills
16 (or knowledge components) associated with each question is available, equation 1 can be modified as

$$17 \quad P(C = 1) = \text{logit}((\sigma_s - \beta_i) \cdot q_i) \quad (2)$$

18 where q_i is an n dimensional binary vector with 1s in the columns of the skills associated with
19 question i , and \cdot indicates the dot product. The collection of q_i s is called Q Matrix.

20 The two parameter model adds a second parameter to each question called discrimination denoted
21 by α . The equations for probability of correctness for two parameter IRT and MIRT are given by
22 equations 3 and 4 below.

$$23 \quad P(C = 1) = \text{logit}((\sigma_s - \beta_i) * \alpha_i) \quad (3)$$

$$24 \quad P(C = 1) = \text{logit}(((\sigma_s - \beta_i) * \alpha_i) \cdot q_i) \quad (4)$$

25 The addition of discrimination parameter allows for questions difficulty to be scaled with student
26 ability. As a consequence, students with different abilities might find the relative difficulties of 2
27 questions to be swapped at a threshold value of their ability. For example students with ability ≤ 2.3
28 might find Q1 easier than Q2 whereas students with higher ability find Q2 easier than Q1.

29 Traditionally, IRT parameters are estimated using Estimation Maximisation (EM) or Monte Carlo
30 Markov chain (MCMC) techniques. This limits the size of dataset on which IRT could be applied.
31 Recent work by [5] and [3] showed that Stochastic Variational Inference (SVI) can be used to estimate

32 IRT parameters. This increases scalability and allows IRT with datasets containing millions of records
33 to be compatible with IRT estimation. To the best of our knowledge IRT has not been applied to a
34 dataset with the size comparable to NeurIPS Education challenge.

35 1.1 Models

36 For Task 1 and Task 2, five models were evaluated for response prediction.

- 37 1. One parameter IRT (1pIRT)
- 38 2. Two parameter IRT (2pIRT)
- 39 3. One parameter MIRT with binary Q Matrix (1pMIRT)
- 40 4. Two parameter MIRT with binary Q Matrix (2pMIRT)
- 41 5. Two parameter MIRT with continuous valued Q Matrix (2pMIRTq)

42 1pIRT is the simplest model based on equation 1. 1pMIRT is its multidimensional equivalent based
43 on equation 2. 2pIRT and 2pMIRT are based on equations 3 and 4 respectively. For the construction
44 of Q matrix in 2pMIRT model, only the top 2 level skills were considered. This was done to limit
45 the number of dimensions. This resulted in 9 dimensional vectors (since there are only 9 level 0
46 and level 1 skills). 2pMIRTq is based on equation 4 with the q matrix being estimated or learned
47 rather than provided to the model. Findings in [1] indicate that expert assigned skill / knowledge
48 component tags add little to no value in many datasets. Also, a continuous valued Q matrix can
49 represent temporal information better than binary Q matrix associated with manual tagging. 20
50 dimensions was arbitrarily chosen for 2pMIRTq. Higher dimensions converged quicker and had
51 higher training set accuracy but lost generalizability and had a much lower test set accuracy.

52 1.2 Results

Table 1: Task 1 results

Model Name	Training set accuracy (%)	Public test set accuracy (%)
1pIRT	74.28	73.80
2pIRT	72.77	72.11
1pMIRT	74.65	73.28
2pMIRT	75.35	73.63
2pMIRTq	80.72	73.62

53

...

54 Table 1 shows the performance of the different models. Although 2pMIRTq had the highest training
55 set accuracy, it was prone to overfitting and 1pIRT had the best test set accuracy overall. The number
56 of parameters seem to have little effect on the accuracy. However, 2 parameter models converged
57 faster. The number of dimensions also seems to have a pronounced effect. Higher dimensions tended
58 to overfit and the disparity between training set accuracy and test set accuracy was large (~12%).

59 1.3 Challenges and Limitations

60 The most obvious limitation of using IRT for prediction in this challenge is the characteristic of
61 IRT student ability estimates. The estimates learned this way must conform to the assumption that
62 answering the questions does not affect their ability. While this may be a reasonable assumption for
63 tests like GRE, student ability continuously improves as they progress through Intelligent Tutoring
64 Systems (ITSs). This almost always leads to overestimation of student abilities. Also ITS student
65 responses are usually skewed towards correct responses (since many ITSs require students to practice
66 skills to mastery). To validate the theory of overestimation, we performed a 90-10 train test split
67 based on timestamp with training set data occurring before the test data stratified on students. This
68 improved the discrepancy between training and test set accuracies by a small margin.

69 The next biggest limitation is the compensatory nature of MIRT. If a student has a high ability
70 in one of the dimensions associated with the question, their response is predicted to be correct
71 despite very low abilities in other dimensions. To mitigate this, we tried to threshold student ability
72 contribution in each individual dimension while predicting correct response. The results were mixed,
73 with improvement in prediction of certain students and degradation in prediction of others. We also
74 tried implementing a possible alternative model called partial compensatory model presented in [6]
75 which requires students to master all skills associated with a question to predict their response as
76 correct. However, we were unsuccessful in estimating parameters with good fit using SVI.

77 At this stage, we are unable to explain the overfitting of 2pMIRTq model. Further study and discussion
78 with experts is required to better understand the factors affecting this. We believe 2pMIRTq to have
79 tremendous potential as it removes the need to manually tag each question. IRT models also have
80 very high interpretability and transferability. For example, the question difficulties estimated from 1
81 batch can be retained and only the student ability be estimated for a consecutive batch. Similarly,
82 student abilities from one subject can be used to estimate question difficulty from another subject as
83 long as the skills overlap.

84 **2 Task 2: Predicting student response**

85 Our approach to Task 2 was to borrow the predictions from Task 1 and substitute the correct answers
86 for predictions which indicated correct response. For the predictions involving incorrect responses,
87 we clustered students based on their abilities and selected the most frequent answer for each question
88 from the clusters the students belonged to. We arbitrarily chose 1000 clusters and bin values for
89 each cluster from the range of student abilities. Student ability is normally distributed and hence
90 the number of students in each cluster varies. A possible alternative is to cluster students by having
91 varying bin lengths such that the number of students in each group is approximately equal. Due to
92 time constraints, we were unable to try this. Since 1pIRT model had the best test set accuracy, we
93 used the student ability from that model for prediction. The best accuracy achieved was 63.77%

94 **2.1 Challenges and Limitation**

95 The biggest limitation of the approach is tied to skewed predictions resulting from IRT which are
96 used to derive student responses. This is explained in detail in section 1.3. The second challenge is
97 computational complexity. Since we need to compute the most frequent answer for each question
98 for each cluster, the time complexity is $O(N*I*n)$. Where N = number of students, I = number of
99 questions and n = number of clusters. This severely limited our ability to repeat the experiment for
100 different number of clusters.

101 **References**

- 102 [1] Theophile Gervet, Ken Koedinger, Jeff Schneider, Tom Mitchell, et al. When is Deep Learning
103 the Best Approach to Knowledge Tracing? *JEDM| Journal of Educational Data Mining*,
104 12(3):31–54, 2020.
- 105 [2] Jeff Johns, Sridhar Mahadevan, and Beverly Woolf. Estimating student proficiency using an
106 item response theory model. In *International conference on intelligent tutoring systems*, pages
107 473–480. Springer.
- 108 [3] John P Lalor, Hao Wu, and Hong Yu. Learning latent parameters without human response
109 patterns: Item response theory with artificial crowds. In *Proceedings of the Conference on*
110 *Empirical Methods in Natural Language Processing. Conference on Empirical Methods in*
111 *Natural Language Processing*, volume 2019, page 4240. NIH Public Access.
- 112 [4] Youngjin Lee. Estimating student ability and problem difficulty using item response theory (irt)
113 and trueskill. *Information Discovery and Delivery*, 2019.
- 114 [5] Prathiba Natesan, Ratna Nandakumar, Tom Minka, and Jonathan D Rubright. Bayesian prior
115 choice in irt estimation using mcmc and variational bayes. *Frontiers in psychology*, 7:1422, 2016.
- 116 [6] Mark D Reckase. Multidimensional item response theory models. In *Multidimensional item*
117 *response theory*, pages 79–112. Springer, 2009.

- 118 [7] Wim J van der Linden and Ronald K Hambleton. *Handbook of modern item response theory*.
119 Springer Science & Business Media, 2013.