

B venom OK!

IMPLICIT PERFORMANCE ESTIMATION FOR SCORE-BASED CLASSIFIERS USING COGNITIVE DIAGNOSIS

A PREPRINT

Nikita Breskanu

Lomonosov Moscow State University
nbreskanu73@gmail.com

Archil Maisuradze

Lomonosov Moscow State University
artchil@mail.ru

← →
но пропату бие норта?

ABSTRACT

Score-based binary classifiers are widely used in machine learning. When it comes to their validation, well-known performance metrics, such as ROC-AUC, F1-score, and Accuracy, are most often used. However, all these metrics have their flaws and represent the performance of the classifier only from a certain angle. This work attempts to aggregate all traditional performance metrics using cognitive diagnosis models. Cognitive diagnosis models are widely researched in smart education and proved to be successful in estimating students' latent knowledge from their exercise solutions. In the context of binary classification, classifiers can be viewed as students and performance metrics as exercises. This reduction represents a novel approach to the validation of binary classifiers and produces latent knowledge attributes, which can be interpreted as new implicit performance attributes.

Keywords Validation · Binary classification · Cognitive diagnosis · Item Response Theory · Machine learning

1 Introduction

The validation is an important step in the machine learning models lifecycle. For this reason, many performance metrics (Accuracy, F1-score, ROC-AUC, etc.) have been developed. However, none of them can entirely characterize the model behavior [9, 6], and there is no clear agreement on which of them to use. It may be impossible to get the true attributes of the model by performing direct aggregation of the answers.

In psychometrics, it is believed that the desired attributes of the subject are only partially manifested in direct measurements. Applying this idea to ML validation, researchers actively try to create a better performance metric by using Item Response Theory (IRT) [23], a classical tool of the psychometrician [14, 17, 3]. However, IRT estimates only one latent attribute, which is not enough to completely describe the model, and the only non-one-dimensional approach approach only estimated the Recall equivalents in multi-class classification [12].

Validating models by their performance metrics can be reduced to the cognitive diagnosis task in smart education, where the goal is to estimate certain predefined attributes of the students by their exercise solutions. For the latter, a wide variety of models have been developed [13, 5, 21, 11]. This work is the first attempt to apply cognitive diagnosis models to create new, potentially better metrics from traditional ones.

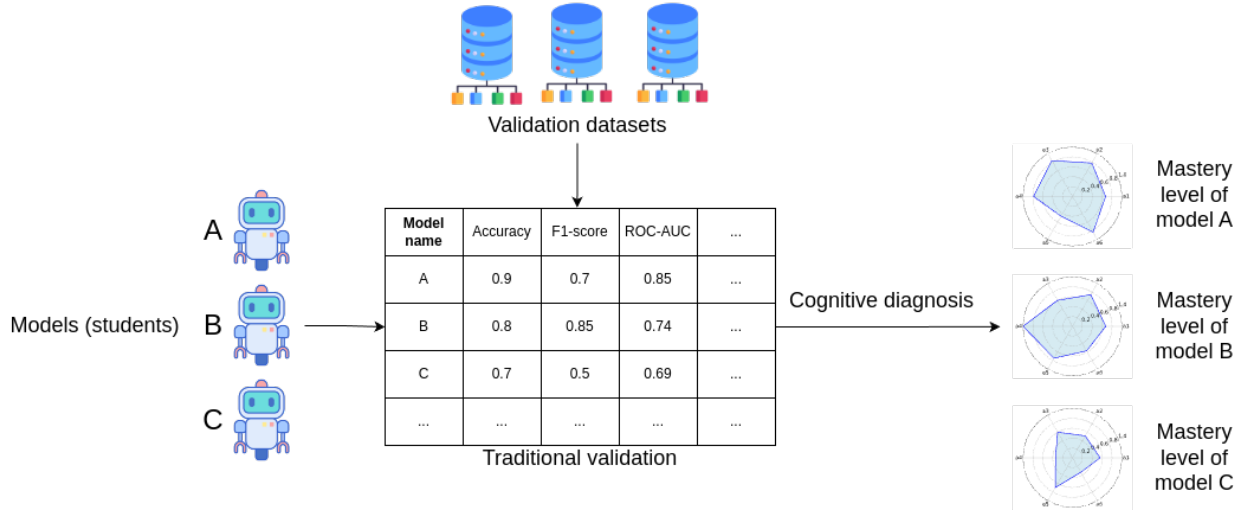


Figure 1: Validation using cognitive diagnosis framework

Our task is to create a framework that would allow estimating mastery of machine learning models for each predefined skill. For that problem, first traditional performance metrics from validation datasets are calculated, and then cognitive assessment is performed, which assigns mastery for each skill to every model. It’s important to mention that this approach only works for a pool of models, not for a single one, due to the nature of cognitive diagnosis models.

We defined specific skills, and then performed experiments on score-based classifiers to obtain their cognitive mastery levels (1). The new metrics (mastery levels) turned out to be competitive with the traditional ones, describing the model from a slightly different angle and considering other models results, and, most importantly, allow using multiple validation datasets, thus capturing model behavior in different learning contexts.

The proposed validation framework can be used to perform validation and comparison of multiple binary classifiers. We also propose a method for adding newly created model to the pool of existing ones, and obtaining its mastery levels. This might open the way to creating multi-skill ML model leaderboards, capturing multiple various datasets.

2 Related work

It is known that traditional performance metrics can’t fully describe the model’s performance, and each one of them has their flaws [6]. For example, in binary classification, Accuracy does not see the difference between errors in the positive and negative classes; Precision and Recall do not know the number of correctly identified negative classes (True Negative); ROC-AUC is sensitive to class imbalance [9].

In 2016, the first attempt of applying psychometric tools for ML validation was made [14], where the author tried to apply IRT [23] for estimated a better version of Accuracy for multi-class classifiers. That study created a great interest for other researchers in applying IRT in machine learning. IRT-based ensembles were proposed [3], where the weights are the IRT scores. IRT-based leaderboard for NLP models validation [17]. Researchers also tried to use IRT to reduce the validation dataset [15], or to make manual validation more efficient [19]. An attempt was made to use IRT for clustering examples in multi-dataset NLP benchmarks [18].

One of the advantages of using IRT for evaluation is that it assigns parameters to every item (object in the dataset), which can later be used to enhance interpretation [17]. The framework for estimating this parameters for newly generated questions was proposed [2].

Another widely researched area is the cognitive diagnosis task, where it is required to estimate students’ mastery levels on every predefined skill by looking at their exercise solutions, and exercise-skill correspondence matrix, which is also known as Q-matrix. Previously, only classical models like DINA [5] or multidimensional IRT (MIRT) [20] models were used. But in 2022, there was a first attempt of using deep cognitive diagnosis model with trainable interaction function, this deep model was called NeuralCD [21]. Later, a lot of extensions of NeuralCD appeared, which were designed to fix some of its flaws, most apparent of which is the lack of knowledge association [21, 11, 13].

To our knowledge, there has been only one attempt of using cognitive diagnosis models for machine learning models validation. In 2023, Camilla framework was proposed for validating deep computer vision multi-class classifiers [12]. The authors estimated the new equivalents of respective Recalls for each class and argued that they describe the performance better by taking into account difficult and easy samples. However, despite their success, we believe that for binary classifiers estimating 2 Recall equivalents is not enough for the full description of the model performance. Our approach is different from the described above in several ways:

- Binary classification task is considered instead of the multi-class.
- Performance metrics are used as exercises instead of the objects.
- Multiple validation datasets are used instead of one, with retraining classifiers for each dataset. This can potentially test the algorithm performance in different learning contexts.

3 Problem statement – *что такое зогоры о мкаб.*

Table 1: Definitions

Quantity	Description
$\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_L\}$	Set of datasets
$\mathcal{S} = \{s_1, \dots, s_N\}$	Set of models (students)
$\hat{\mathcal{E}} = \{\hat{e}_1, \dots, \hat{e}_{\hat{M}}\}$	Set of performance metrics (exercises)
$\mathcal{E} = \{e_1, \dots, e_M\}, M = \hat{M} \times D$	Set of performance metrics, taking datasets into account
$\mathcal{K} = \{\mathcal{K}_1, \dots, \mathcal{K}_K\}$	Set of concepts
L	Number of datasets
N	Number of models (students)
\hat{M}	Number of performance metrics (exercises)
K	Number of concepts
T	Number of response logs
$l \in \{1, \dots, L\}$	Index of the dataset
$i \in \{1, \dots, N\}$	Index of the student
$j \in \{1, \dots, \hat{M}\}$	Index of the exercise
$k \in \{1, \dots, K\}$	Index of the concept
$t \in \{1, \dots, T\}$	Index of the response log
$x^s \in \{0, 1\}^N$	One-hot representation of the model (student)
$x^e \in \{0, 1\}^{\hat{M}}$	One-hot representation of the performance metric (exercise)
$Q = \{Q_{jk}\}_{\hat{M} \times K} \in [0, 1]^{\hat{M} \times K}$	Q-matrix
$G \in \{0, 1\}^{K \times K}$	Directed Acyclic Graph (DAG) of concept dependency
$y \in [0, 1]$	Model output
$R = \{(x_t^s, x_t^e, r_t)\}_{t=1}^T$	Response logs
$r \in [0, 1]$	Result of solving the exercise (value of the performance metric)
$\mathcal{M} = \{m_{ik}\}_{N \times K} \in [0, 1]^{N \times K}$	Latent students' mastery levels
\mathcal{L}	Loss function

Task definition Suppose there are L datasets $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_L\}$, N machine learning models (algorithms) $\mathcal{S} = \{s_1, \dots, s_N\}$, \hat{M} traditional performance metrics $\hat{\mathcal{E}} = \{\hat{e}_1, \dots, \hat{e}_{\hat{M}}\}$, and K predefined skills (or concepts) $\mathcal{K} = \{\mathcal{K}_1, \dots, \mathcal{K}_K\}$. Q-matrix $Q \in \mathbb{R}^{\hat{M} \times K}$ is a binary matrix that represents correspondence between performance metrics and concepts: $Q_{jk} = 1 \iff$ knowledge of concept \mathcal{K}_k is required for having a high value of \hat{e}_j . Computing performance metrics $\hat{\mathcal{E}}$ for each dataset forms a set of response logs $R = \{(x^s, x^e, r)\}_{t=1}^T$, where $T = N \times \hat{M}$, $M = \hat{M} \times L$ — a triples consisting of one-hot representation of model (student), performance metric (exercise), taking the index of dataset in the account, and the metric $r \in [0, 1]$, normalized to $[0, 1]$. The desired models' mastery levels for all concepts can be represented as matrix $\mathcal{M} = \{m_{ik}\} \in [0, 1]^{N \times K}$, where m_{ik} is the mastery level of student s_i for the concept \mathcal{K}_k ; $m_{ik} = 1$ represents total knowledge of the concept, and $m_{ik} = 0$ — total ignorance. The task is to infer the mastery matrix \mathcal{M} using cognitive diagnosis model by predicting the responses r_t .

4 Cognitive diagnosis dataset preparation

4.1 Defining datasets

Cognitive models need a large pool of students and exercises, and according to our reduction, students are classifiers and exercises are performance metrics. So, to prepare a dataset for cognitive diagnosis, one needs to evaluate a large number of models on large number of metrics.

In order to test models in different learning scenarios and simultaneously create more exercises, we collected diverse open-source datasets from OpenML. In fact, we used a subset of datasets from paper [1]. The reason for dropping some of the datasets is too large computational complexity when there are more than 1000 features. In total, we have 16 datasets (2).

For each dataset, we transformed it into a pipeline that would be acceptable for every classifier:

- One-hot encoding was performed for all categorical features.
- Standard scaling was performed for all numerical features.

Table 2: Dataset characteristics. All datasets are taken from OpenML. Some of them contain categorical features which will be one-hot encoded. Datasets have diverse class balances: one of them has 99-1.

Dataset name	Samples \times features	Numerical \times categorical features	Class balance
Banknote-authentication	1372 \times 5	5 \times 0	55–45%
Blood-transfusion-service-center	748 \times 5	5 \times 0	76–24%
Breast-w	683 \times 10	10 \times 0	65–35%
Climate-model-simulation-crashes	540 \times 21	21 \times 0	99–1%
Cylinder-bands	277 \times 40	25 \times 15	64–36%
Dresses-sales	99 \times 13	2 \times 11	59–41%
Diabetes	768 \times 9	9 \times 0	65–35%
ilpd	583 \times 11	10 \times 1	71–29%
kc1	2109 \times 22	22 \times 0	84–16%
kc2	522 \times 22	22 \times 0	79–21%
pc1	1109 \times 22	22 \times 0	93–7%
pc3	1563 \times 38	38 \times 0	89–11%
Phoneme	5404 \times 6	6 \times 0	70–30%
qsar-biodeg	1055 \times 42	42 \times 0	66–34%
wdbc	569 \times 31	31 \times 0	62–38%
wilt	4839 \times 6	6 \times 0	94–6%

4.2 Defining classifiers

no guess

We defined 295 binary classifiers by varying hyperparameters. In order to create more diversity, similarly to [14], several artificial classifiers were implemented:

- Optimal classifier always predicts the correct class (either 0 or 1).
- Pessimist classifier always predicts the incorrect class (either 0 or 1).
- Majority classifier always predicts the majority class (either 0 or 1).
- Minority classifier always predicts the minority class (either 0 or 1).
- Mean target classifier always predicts the mean target value. For example, in 90-10 class balanced it will always output 0.9.
- Uniform random classifier predicts classes randomly with equal probabilities.
- Balanced random classifier predicts classes randomly with probabilities proportional to their class balance. For example, in 90-10 class balance there will be a 90% probability of class 0.

Table 3: Binary classifiers, replicated by varying hyperparameters.

Classifier	Implementation	Varying parameters	Number of models
Logistic regression	sklearn	C, solver	120
Decision tree	sklearn	max_depth, criterion	60
Random forest	sklearn	max_depth, n_estimators	12
Gradient boosting	sklearn	n_estimators, learning_rate	9
Gradient boosting	LGBM	n_estimators, num_leaves	9
SVM	sklearn	C, kernel	30
K nearest neighbors	sklearn	n_neighbors, weights	40
Multilayer perceptron	sklearn	hidden_layer_sizes, activation	15
Optimal classifier	<manual>	<absent>	1
Pessimal classifier	<manual>	<absent>	1
Majority classifier	<manual>	<absent>	1
Minority classifier	<manual>	<absent>	1
Mean target classifier	<manual>	<absent>	1
Uniform Random classifier	<manual>	<absent>	1
Balanced Random classifier	<manual>	<absent>	1

In total, there are 302 score-based binary classifiers.

4.3 Defining performance metrics

Performance metrics that were used, are defined in table (4). We used Equal Error Threshold (EER, a point where 2 recalls are the same) to obtain label-based metrics in conjunction with default score-based ones.

There was a possibility of using also default 0.5 threshold, but since not all the classifiers are self-calibrated, we decided not to include it. Moreover, we found that EER threshold metrics were very similar to the respective 0.5 threshold ones.

Table 4: Performance metrics used for score-based classifiers

Performance metric	Definition	Description
ROC-AUC	AUC	Area under the ROC curve
PR-AUC for class 0	PRAUC0	Area under the PR curve, where class 0 is the positive class
PR-AUC for class 1	PRAUC1	Area under the PR curve, where class 1 is the positive class
Gain chart AUC for class 0	GCAUC0	Area under the gain chart, where class 0 is the positive class
Gain chart AUC for class 1	GCAUC1	Area under the gain chart, where class 1 is the positive class
KS statistic	KS	Kolmogorov-Smirnov statistic
Kendall's tau	KTAU	Kendall's correlation between predicted and actual results
Accuracy (EER)	ACC	Accuracy at EER threshold
Precision for class 0 (EER)	PR0	Precision with 0 as positive class at EER threshold
Precision for class 1 (EER)	PR1	Precision with 1 as positive class at EER threshold
Recall (EER)	REC	Recall at EER threshold
Balanced accuracy (EER)	BA	Balanced accuracy at EER threshold
F1-score for class 0 (EER)	FS0	F1-score with 0 as positive class at EER threshold
F1-score for class 1 (EER)	FS1	F1-score with 1 as positive class at EER threshold
Average F1-score (EER)	AVGFS	Average F1-score at EER threshold
Fowlkes-Mallows index for class 0 (EER)	FM0	Fowlkes-Mallows index with 0 as positive class at EER threshold
Fowlkes-Mallows index for class 1 (EER)	FM1	Fowlkes-Mallows index with 1 as positive class at EER threshold
Markedness (EER)	MKNS	Markeness at EER threshold
MCC (EER)	MCC	Matthews correlation coefficient at EER threshold
Jaccard index (EER)	JAC	Jaccard index at EER threshold
Cohen's kappa (EER)	KAPPA	Cohen's kappa at EER threshold

After training all classifiers on all datasets and obtaining performance metrics, there appeared to be 54 classifier duplicates, mostly logistic regressions. After removing them, we obtained the final 248 classifiers \times 353 metrics dataset, ready for cognitive assessment.

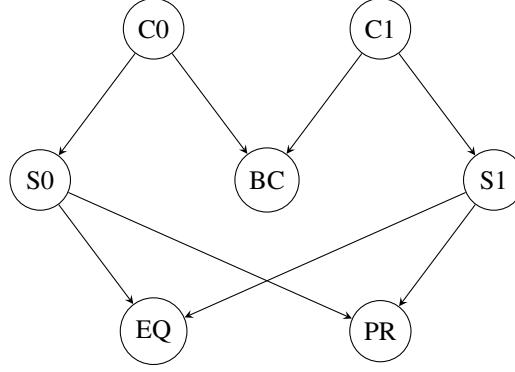


Figure 2: Attribute dependency graph. Attributes C0 and C1 are the root nodes, and all other attributes are assumed to be dependant on them.

4.4 Defining Attributes

Coming up with meaningful attributes for classifiers is a challenging task, and the proposed attributes may be reconsidered in future researches. Let there be two classes in a binary classification problem that need to be distinguished: class 0 and class 1. The following knowledge attributes $\mathcal{K}_n = \{\mathcal{K}_1, \dots, \mathcal{K}_K\}$ will be used to describe the performance of the model:

- **C0**: *Only class 0* - the performance of the model only on class 0. If $m = 1$, then the model makes no mistakes on objects of class 0.
- **C1**: *Only class 1* - the performance of the model only on class 1. If $m = 1$, then the model makes no mistakes on objects of class 1.
- **BC**: *Both classes* - the performance of the model aggregated over both classes, independent of the balance of classes in the validation sample. If $m = 1$, then the model makes no mistakes on any objects.
- **S0**: *Class 0 is superior* - the performance under the condition that the ratio between the type 1 and type 2 errors (Precision and Recall) for class 0 is important. If $m = 1$, then the model makes neither type 1 or type 2 errors on class 0.
- **S1**: *Class 1 is superior* - the performance under the condition that the ratio between type 1 and type 2 errors for class 1 is important. If $m = 1$, then the model makes neither type 1 or type 2 errors on class 1.
- **EQ**: *Equivalent classes* - the performance under the condition that the classes are equivalent. If $m = 1$, then the model always predicts the correct class.
- **PR**: *Prevalent class* - model performance biased towards the prevailing class. In this concept, models that more often predict the prevailing class will have a higher m value. In particular, if the objects of the predominant class occupy almost the entire sample, then the degenerate model predicting only this class regardless of the object will have $m \approx 1$. If $m = 1$, then the model always predicts the correct class.

Some CDMs use a attribute dependency graph, while others consider all concepts independent. Therefore, the graph (Fig. 2) was constructed. In that graph, “Both classes” (BC) is dependent on C0 and C1; “Class i is superior” is dependent on “Only class i” attribute, which is a reasonable requirement, as the former looks at performance on class i together with objects from another class. Finally, EQ and PR are the leaf nodes, and both require mastery of “Class i is Superior” attributes.

4.5 Q-matrix

We define Q-matrix for connecting performance metrics and attributes in Table 5. Some of the metrics, like Precision, Recall and F1-score are asymmetric, i.e. they change when swapping positive and negative class. That’s why they are applied for both cases: when class 0 is positive (for class 0 in the table), and when class 1 is positive (for class 1). Apart from well-known metrics, like Accuracy, Precision, Recall, Balanced accuracy, F1-score [4], we include several other for greater diversity: Average F1-score which is a mean of 2 F1-scores, resulting in a symmetric metric; Fowlkes-Mallows index [8] is a geometric mean of Precision and Recall; Markedness [16], after normalization, is a mean of 2 Precisions. Matthews correlation coefficient [16] is a Pearson correlation coefficient between predicted and

Table 5: Manually created Q-matrix. Attributes *Both classes* and *Prevalent class* weakly connected to performance metrics, as they have only 1 and 2 corresponding performance metrics.

Exercise (performance metric)	C0	C1	BC	S0	S1	EQ	PR
ROC-AUC	0	0	1	0	0	1	0
PR-AUC for class 0	0	0	0	1	0	0	0
PR-AUC for class 1	0	0	0	0	1	0	0
Gain chart AUC for class 0	0	0	0	1	0	0	0
Gain chart AUC for class 1	0	0	0	0	1	0	0
KS statistic	0	0	0	0	0	1	0
Kendall’s tau	0	0	1	0	0	1	0
Accuracy (EER)	0	0	0	0	0	0	1
Precision for class 0 (EER)	0	0	0	1	0	0	0
Recall for class 0 (EER)	1	0	0	1	0	0	0
Precision for class 1 (EER)	0	0	0	0	1	0	0
Recall for class 1 (EER)	0	1	0	0	1	0	0
Balanced accuracy (EER)	0	0	1	0	0	1	0
F1-score for class 0 (EER)	0	0	0	1	0	0	0
F1-score for class 1 (EER)	0	0	0	0	1	0	0
Average F1-score (EER)	0	0	0	0	0	1	0
FM-score for class 0 (EER)	0	0	0	1	0	0	0
FM-score for class 1 (EER)	0	0	0	0	1	0	0
Markedness (EER)	0	0	0	0	0	1	0
Matthews coefficient (EER)	0	0	0	0	0	1	0
Jaccard index (EER)	0	0	0	0	0	0	1
Cohen’s kappa (EER)	0	0	0	0	0	1	0

actual labels; Jaccard index is Jaccard similarity score between predicted and actual labels; Cohen’s kappa [22] is an agreement between predicted and actual labels.

5 Requirements for mastery levels

We propose several interpretability requirements for the mastery levels (new metrics generated by the cognitive model):

- Mastery levels (new metrics) are expected fully cover the old metrics, i.e. traditional metrics are expected to be derived from mastery levels by using some formula.
- If one classifier has higher mastery level on some concept \mathcal{K}_k than the other classifier, it is expected to have higher corresponding performance metrics $e_j : Q_{jk} = 1$ than the other.
- If one classifier has higher performance metric e_j than the other classifier, it is expected to have higher corresponding mastery levels $\mathcal{K}_k : Q_{jk} = 1$ than the other.

6 Experiments

6.1 Cognitive Diagnosis Models Evaluation Metrics

CDMs will be compared by how well their trained knowledge levels satisfy proposed interpretability requirements. To evaluate them, the following metrics will be used:

$$R^2 = 1 - \frac{\text{MSE}(\text{CDM})}{\text{MSE}(\bar{x})} \quad (1)$$

Coefficient of determination R^2 (Equation 1) is a traditional regression metric that shows how well the CDM is fit, i.e. how well traditional performance metrics are derived from the attribute knowledge levels.

Table 6: Comparison of cognitive diagnosis models, trained for 100 epochs with BCE loss. R2 score, Degree of agreement and Degree of consistency are shown as 95% confidence intervals on 3 different-seed runs.

Model	# parameters	R2	DOA	DOC
MIRT	4552	-0.011 ± 0.000	0.619 ± 0.001	0.619 ± 0.001
NeuralCD	7177	0.887 ± 0.003	0.586 ± 0.009	0.586 ± 0.009
KaNCD	22467	0.885 ± 0.001	0.584 ± 0.006	0.584 ± 0.006
HierMIRT	9544	0.848 ± 0.028	0.577 ± 0.007	0.577 ± 0.007
HierNCD	11657	0.892 ± 0.001	0.600 ± 0.006	0.600 ± 0.006
QCCDM (small)	9882	0.940 ± 0.038	0.539 ± 0.005	0.539 ± 0.005
QCCDM	144282	0.955 ± 0.054	0.542 ± 0.018	0.542 ± 0.018

$$DOA_k = \frac{\sum_{a,b \in S} [m_{ak} > m_{bk}] \frac{\sum_{j=1}^M Q_{jk} [x_{aj} > x_{bj}]}{\sum_{j=1}^M Q_{jk} [x_{aj} \neq x_{bj}]}}{\sum_{a,b \in S} [m_{ak} > m_{bk}]} \quad (2)$$

$$DOA = \frac{1}{K} \sum_{k=1}^K DOA_k$$

Degree of agreement (Equation 2) is a widely used interpretability metric in cognitive diagnosis [21, 7], a variant of Cohen’s kappa [22] for cognitive diagnosis. It shows how well attribute-exercise monotonicity is maintained, i.e. if student a has higher attribute than student b , the first is expected to solve corresponding to the attribute exercises better than the latter.

$$DOC_j = \frac{\sum_{a,b \in S} [x_{aj} > x_{bj}] \frac{\sum_{k=1}^K Q_{jk} [m_{ak} > m_{bk}]}{\sum_{k=1}^K Q_{jk} [m_{ak} \neq m_{bk}]}}{\sum_{a,b \in S} [x_{aj} > x_{bj}]} \quad (3)$$

$$DOC = \frac{1}{M} \sum_{j=1}^M DOC_j$$

Degree of consistency (Equation 3) is another interpretability metric [10], that shows the inverse to DOA exercise-student monotonicity: if student a solves an exercise better than student b , first is expected to have higher attributes corresponding to this exercise than the latter.

6.2 Results

Table 6 shows number of parameters and 95% confidence intervals for evaluation metrics 1, 2, 3. All CDMs except MIRT achieve acceptable R2-score. However, DOA and DOC metrics are very low for all the models, with the highest value of 0.619. Ideally, it is required that they are at least 0.8. It is worth noting that random assignment of mastery levels achieves 0.5 DOA and 0.5 DOC, meaning that 0.619 is close to being random in terms of interpretability.

Current results show that resulting mastery levels are not interpretable enough to perform further visualizations on mastery levels and exercise parameters.

Most likely, the reason for such small interpretability is that currently defined attributes, taken from label-based experiments, appear to be bad at describing score-based classifiers, which are much more complex than label-based ones.

7 Conclusion

The experiments suggest that the 7 attributes C0, C1, BC, S0, S1, EQ and PR are bad at describing score-based classifiers.

References

- [1] Vitor Cirilo Araujo Santos, Lucas Cardoso, and Ronnie Alves. The quest for the reliability of machine learning models in binary classification on tabular data. *Scientific Reports*, 13(1):18464, 2023.
- [2] Luca Benedetto, Andrea Cappelli, Roberto Turrin, and Paolo Cremonesi. R2de: a nlp approach to estimating irt parameters of newly generated questions. In *Proceedings of the tenth international conference on learning analytics & knowledge*, pages 412–421, 2020.
- [3] Ziheng Chen and Hongshik Ahn. Item response theory based ensemble in machine learning. *International Journal of Automation and Computing*, 17:621–636, 2020.
- [4] Peter Christen, David J Hand, and Nishadi Kirielle. A review of the f-measure: its history, properties, criticism, and alternatives. *ACM Computing Surveys*, 56(3):1–24, 2023.
- [5] Jimmy De La Torre. Dina model and parameter estimation: A didactic. *Journal of educational and behavioral statistics*, 34(1):115–130, 2009.
- [6] Peter Flach. Performance evaluation in machine learning: the good, the bad, the ugly, and the way forward. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9808–9814, 2019.
- [7] Francois Fouss, Alain Pirotte, Jean-Michel Renders, and Marco Saerens. Random-walk computation of similarities between nodes of a graph with application to collaborative recommendation. *IEEE Transactions on knowledge and data engineering*, 19(3):355–369, 2007.
- [8] Maria Halkidi, Yannis Batistakis, and Michalis Vazirgiannis. On clustering validation techniques. *Journal of intelligent information systems*, 17:107–145, 2001.
- [9] Nathalie Japkowicz. Why question machine learning evaluation methods. In *AAAI workshop on evaluation methods for machine learning*, volume 6. University of Ottawa, 2006.
- [10] Jiatong Li, Qi Liu, Fei Wang, Jiayu Liu, Zhenya Huang, Fangzhou Yao, Linbo Zhu, and Yu Su. Towards the identifiability and explainability for personalized learner modeling: An inductive paradigm. In *Proceedings of the ACM on Web Conference 2024*, pages 3420–3431, 2024.
- [11] Jiatong Li, Fei Wang, Qi Liu, Mengxiao Zhu, Wei Huang, Zhenya Huang, Enhong Chen, Yu Su, and Shijin Wang. Hiercdf: A bayesian network-based hierarchical cognitive diagnosis framework. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 904–913, 2022.
- [12] Qi Liu, Zheng Gong, Zhenya Huang, Chuanren Liu, Hengshu Zhu, Zhi Li, Enhong Chen, and Hui Xiong. Multi-dimensional ability diagnosis for machine learning algorithms. *arXiv preprint arXiv:2307.07134*, 2023.
- [13] Shuo Liu, Hong Qian, Mingjia Li, and Aimin Zhou. Qccdm: A q-augmented causal cognitive diagnosis model for student learning. In *ECAI 2023*, pages 1536–1543. IOS Press, 2023.
- [14] Fernando Martínez-Plumed, Ricardo BC Prudêncio, Adolfo Martínez-Usó, and José Hernández-Orallo. Making sense of item response theory in machine learning. In *ECAI 2016*, pages 1140–1148. IOS Press, 2016.
- [15] Felipe Maia Polo, Lucas Weber, Leshem Choshen, Yuekai Sun, Gongjun Xu, and Mikhail Yurochkin. tinybenchmarks: evaluating llms with fewer examples. *arXiv preprint arXiv:2402.14992*, 2024.
- [16] David MW Powers. Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation. *arXiv preprint arXiv:2010.16061*, 2020.
- [17] Pedro Rodriguez, Joe Barrow, Alexander Miserlis Hoyle, John P Lalor, Robin Jia, and Jordan Boyd-Graber. Evaluation examples are not equally informative: How should that change nlp leaderboards? In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4486–4503, 2021.
- [18] Pedro Rodriguez, Phu Mon Htut, John P Lalor, and João Sedoc. Clustering examples in multi-dataset benchmarks with item response theory. In *Proceedings of the Third Workshop on Insights from Negative Results in NLP*, pages 100–112, 2022.
- [19] João Sedoc and Lyle Ungar. Item response theory for efficient human evaluation of chatbots. In *Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems*, pages 21–33, 2020.
- [20] Yanyan Sheng and Christopher K Wikle. Comparing multiunidimensional and unidimensional item response theory models. *Educational and Psychological Measurement*, 67(6):899–919, 2007.
- [21] Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang, Yu Yin, Shijin Wang, and Yu Su. Neuralcd: a general framework for cognitive diagnosis. *IEEE Transactions on Knowledge and Data Engineering*, 35(8):8312–8327, 2022.
- [22] Matthijs J Warrens. Five ways to look at cohen’s kappa. *Journal of Psychology & Psychotherapy*, 5, 2015.

- [23] Frances M Yang and Solon T Kao. Item response theory for measurement validity. *Shanghai archives of Psychiatry*, 26(3), 2014.