Semantic Web 0 (0) 1 IOS Press

Experts vs. Automata: A Comparative Study of Methods for a Priori Prediction of MCQ Difficulty

Ghader Kurdi ^{a,*}, Jared Leo^a Nicolas Matentzoglu ^a Bijan Parsia^a Uli Sattler^a Sophie Forge^b Gina Donato^b and Will Dowling^b ^a Department of Computer Science, The University of Manchester, Manchester, UK ^b Elsevier, Philadelphia, USA

Abstract. Successful exams require a balance of easy, medium, and difficult questions. Question difficulty is generally either estimated by an expert or determined after an exam is taken. The latter is useless for new questions and the former is expensive. Additionally, it is not known whether expert prediction is indeed a good proxy for difficulty.

In this paper, we compare two ontology-based measures for difficulty prediction with each other and with expert prediction (by 15 experts) against exam performance (of 12 residents) over a corpus of 231 medical case-based questions. We found that one measure (relation strength indicativeness) to be of comparable performance (accuracy = 47%) to the experts (average accuracy = 49%).

Keywords: ontologies, semantic web, automatic question generation, difficulty modeling, difficulty prediction, multiple choice questions, student assessment

1. Introduction

Multiple choice question (MCQ) examinations are widely used to assess the knowledge and skills of students and the quality of the teaching instruments. Using good-quality questions is essential for achieving these purposes. Several criteria exist for measuring question quality, as discussed in [1–3]. Good quality questions need to be, among other things, 1) valid (i.e., they measure what they are supposed to measure); 2) discriminating (i.e. discriminate between high- and low-information students); 3) fair (i.e., their results are not biased in favour of a subgroup within the cohort); and 4) of appropriate difficulty. Difficulty of MCQs is usually¹ defined as the proportion of students solving a question correctly out of the total number of students attempting the question, and is known as *percentage correct*.

The difficulty criterion is of importance, attributed to its effect on the other quality criteria. Knowledge

¹This is based on a recent systematic review we conducted on difficulty prediction and on a survey of studies investigating MCQ examination quality.

^{*}Corresponding author. E-mail: ghader.kurdi@manchester.ac.uk.

about difficulty level and sources of difficulty in ques-1 tions provides insights into whether other quality cri-2 teria are satisfied or not. With regards to validity, be-3 4 ing able to answer the question 'what makes a partic-5 ular question easy or difficult?' is an important step in 6 understanding 'what does the question measure?' For 7 example, questions that are difficult due to their lin-8 guistic complexity are usually not valid in tests other 9 than in language tests. This is because it is not clear 10 whether students failure in answering these questions 11 is due to the language factor or to their lack of the 12 knowledge or skills of interest. In addition, inappropri-13 ately difficult or easy questions are tend to have bad 14 discrimination because, either almost none of the stu-15 dents solve them or all of the students solve them cor-16 rectly. Finally, the difficulty level of the questions is a 17 major determinant of the fairness of exams, especially 18 when different exam forms are used (equally difficult 19 forms are needed), or when question selection is al-20 lowed (equally difficult questions are needed). 21

While information about the difficulty of questions 22 is essential for designing exams, percentage correct 23 can only be retrospectively determined. Traditional 24 means of estimating difficulty are by obtaining it from 25 previous administrations of the questions, if previous 26 statistics are available, or by relying on experts' evalu-27 ation, which is usually the case in small-scale exams.

28 With the recent advances in automated procedures 29 for generating questions [4-10], giving the ability to 30 generate a huge number of new questions, the need 31 for measures that approximate prospective difficulty 32 becomes more vital. These measures can be incorpo-33 rated into the generation process allowing the gener-34 ation of questions with the desired difficulty (satisfy-35 ing the needs of exam developers), or at least, with 36 appropriate difficulty (filtering inappropriately easy 37 and difficult questions). Furthermore, organising auto-38 generated questions by difficulty will reduce experts' 39 efforts in sorting through, and trying to predict the dif-40 ficulty of, a large number of questions. Finally, good 41 predictive measures will allow moving progress to-42 ward the goal of generating exams automatically. 43

The majority of existing automatic difficulty pre-44 diction models are machine-learning based approaches 45 [see, for example, 11–14] that have merely been 46 used for finding correlations in existing data as op-47 48 posed to prediction. Existing cross-validated models that have been developed for prediction [15–17] 49 are highly domain-specific which limit their utility. 50 However, in a prior work [6, 7], we have developed 51

two ontology-derived measures which are based on a domain-independent model of difficulty.

Since the aforementioned ontology-based measures have neither been evaluated thoroughly nor compared to each other in a systematic way, we compare their prediction, along with expert prediction, against a gold standard of student performance. This allows us to validate our measures and determine whether they are suitable for replacing expert estimations when constructing exams.

This paper aims to address the following research questions:

- RQ1: How accurate is expert prediction of difficulty against student performance compared to guessing?
- RQ2: How accurate are automatic difficulty prediction methods against student performance?
 - * compared to guessing?
 - * compared to each other?
 - * compared to domain experts?

We have collected difficulty information for 231 questions through a study involving 15 medical experts and a cohort of 12 residents. We found that MCQ difficulty was moderately predicted by domain experts (average accuracy = 49%). We also found the automated measure developed in [7] to be of comparable performance to experts (accuracy = 47%) and to represent an economical alternative.

The main contributions of this paper are:

- an ontology-based measure for predicting the difficulty of auto-generated MCQs;
- user studies in the medical domain investigating the predictive performance of domain expert and automated ontology-based measures;
- a detailed analysis of the performance of difficulty prediction measures that show, by example, the minimum set of criteria that need to be considered in evaluating the performance of similar measures:
- a fairly large question set (231 questions, of which 92 were answered by at least 10 participants) annotated with percentage correct and expert prediction and can be used for testing the performance of new approaches to difficulty prediction.

1

2

4

5

6

7

8

9

12

13

15

16

17

18

21

22

25

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

2. Background

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

29

30

31

48

49

50

51

2.1. Multiple choice questions

MCQs consist of two components:

 the stem: a textual element that represents a problem to be solved, possibly accompanied by supplementary elements such as tables or graphs;

 the options: a set of alternatives to select from. Standard MCQs, known as single response questions, have one correct option (known as the key) and three or four incorrect options (the distractors) [18]. Another popular type of MCQs, known as multiple response questions, have at least two keys. Distractors are added so that the number of options typically sums to four or five.

18 Writing high-quality MCQs is known to be chal-19 lenging and expensive. The challenges faced by exam 20 developers are apparent from the low quality of MCO 21 examination as indicated by several studies investigat-22 ing their quality. For example, the authors of [19–23] 23 found more than 50% of investigated MCQs to contain at least one item writing flaw². Other studies [24-26]24 25 report that the percentage of MCQs with all their dis-26 tractors being considered functional³ is low (between 27 5% and 23%). 28

2.2. Ontology-based MCQ generation and difficulty prediction

32 Given the challenges faced by test developers in 33 constructing high-quality MCQs, automated approaches 34 for question generation have come into play. Ontolo-35 gies have been increasingly used, in research contexts, 36 as a source for automatic generation of questions [4-37 8]. We attribute their increased use to the following 38 reasons. The first reason is the availability of ontolo-39 gies with potential educational value. These ontolo-40 gies contain exact facts and represent domains of in-41 terest precisely and non-ambiguously in a machine-42 processable way. Besides that, ontologies are sup-43 ported by standard reasoning services and the devel-44 opment of further supporting tools and services is an 45 active research area. Another reason is that, compared 46 to text, the process of finding good distractors is easier. 47

As an example, consider the question 'Which city is located in UK?', generated from a Wikipedia⁴ article about the United Kingdom. The cities mentioned in the articles are most likely to be UK cities. Even if a non-UK city is mentioned, detailed information about it, which is important in deciding whether it serves its intended purpose as a distractor, cannot be found in the same article. For a detailed systematic review of automatic question generation methods, the reader is referred to [32].

One point worth mentioning is that underlying difficulty models are not part of most existing question generation approaches. According to [32], apart from the similarity-based approach (outlined in Section 3.1), only two question generation approaches [33, 34] take into account generating questions with controlled difficulty but without providing experimental evaluation of the performance of difficulty prediction. The automatic measures compared in this study represent existing domain nonspecific measures of MCQ difficulty. Other measures are either variants of the similarity approach [35], designed for questions with other response formats, or categorised by being domain- or questionspecific [15, 16, 33].

2.2.1. Case-based question generation

One of the limitations of current question generation approaches is the simplicity of the generated questions in terms of both cognitive level,⁵ with the majority of generated questions in [4, 6, 8] testing recall of information, and structure, with generated stems in [4, 6, 8]containing at most two concepts. In a recent study [7], we have tackled the generation of medical case-based questions (see question Q2) using a large medical ontology. What is interesting about these questions is that they are widely used in medical education, and that answering these questions requires more than just recall of information [37-39]. From a computational point of view, complex structure of their stem, consisting of multiple concepts, introduces additional challenges of coordination between these concepts and understanding the role they play in question difficulty. The generation approach was evaluated through expert review of questions generated from four medical specificities. More details on the set of generated questions will be given in Section 4.2.2.

3

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

² Violations of best practices as suggested in MCQ-writing guidelines.

³ Functional distractors are those selected by at least 5% of examinees [25, 27, 28].

⁴Wikipedia has been used as a source for question generation by [29–31].

⁵The mental process involved in question-solving as described in Bloom's taxonomy [36], a popular classification of cognitive levels.

3. Competing measures

The target of difficulty prediction is to assign difficulty levels (easy, medium, difficult), as derived from percentage correct (to be discussed in Section 4.2.4). The two ontology-based measures compared are described in this section.

3.1. Similarity-based measure

A plausible prediction model has been proposed 11 12 in [6], in which the similarity between the key and 13 the distractors was suggested as an indicator of MCQ 14 difficulty. Increasing the similarity between the key 15 and distractors results in increasing the difficulty of 16 MCQs. The rationale is that more knowledge is re-17 quired to differentiate between key and similar distrac-18 tors. As an example, consider the following question 19 (Q1) taken from [32]. The most similar distractor to 20 the key, and the most difficult to eliminate, is the op-21 tion 'the tongue' since this option shares with the key 22 the feature of being a body part. On the other hand, 23 elimination of the options 'disease' and 'glossitis' is 24 easier since they do not have shared features with the 25 key. 26

Q1: Pyorrhoea occurs in:
A. the tongue
B. glossitis
C. the gums (key)
D. a disease

To control the difficulty of questions, [6] has developed a similarity measure that is based on Jaccard similarity [40] and intended to be used with ontologies. The similarity measure is defined as follows:

$$sim(k,d) = \frac{Com(k,d)}{Union(k,d)}$$

where Com(k, d) is the number of common subsumers between the key k and a distractor d, and Union(k, d) is the number of all subsumers of both k and d.⁶ The overall difficulty of the question is then defined as the average similarity between the key and distractors.

Preliminary studies have shown that the similarity measure has a good difficulty prediction [32]. In the

absence of other domain-independent measures that are empirically supported, the similarity measure is considered as the gold standard for automatic difficulty prediction. However, one of the limitations of this measure is that it does not take into account the contribution of the stem to the difficulty of questions. While this did not represent a problem in questions having simple stems (e.g. What is X? where X is a concept name), we believe that the role a stem plays is a major influencer on the difficulty of case-based questions that are characterized by stems that contains multiple terms (i.e. multi-term questions). In addition, the similarity measure is developed based on the assumption that all relational axioms have the same strength (i.e. a disease is either associated or not associated with a clinical finding). However, this is not always the case, especially in the medical domain where relations such as hasClinicalFinding have different degrees of strength (e.g. common clinical finding and rare clinical finding). These limitations motivate us to develop the new difficulty measure described below.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

3.2. Relation Strength Indicativeness

A new measure of question difficulty was introduced in [7] which estimates difficulty by combining several calculations that exploit the relational axioms of an ontology, along with their *strength*. This measure, coined *relation strength indicativeness (RSI)*, requires an ontology to contain existential class axioms, i.e., those axioms of the form $A \sqsubseteq \exists R.B,^7$ where *A* and *B* are classes, that incorporate an associated strength of the relation *R*.

The proposed difficulty measure targets more complex types of questions, such as Q2 below, when compared to simple questions, such as Q1. The two main calculations RSI uses involve *stem indicativeness* and *option entity difference*. The former intuitively represents the degree to which stem entities are indicative of the key, whilst the latter captures the difference between how indicative the stem entities are to the distractors, when compared to the key. The final difficulty measure is based on an average of these two measures.

Consider the following case-based medical MCQ (Q2), similar to those generated in $[7]^8$:

4

1

2

3

4

5

6

7

8

9

10

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

⁶ Different ways of counting subsumers have been defined in [32].

⁷The corresponding Manchester OWL syntax is: A SubClassOf R some B.

 $^{^{8}}$ A simple and modified version of a question generated in [7] is used for the sake of a non complex example.

Q2: A 13-year-old female patient presents with Hemorrhage of urethra and Hematuria. What is the most likely diagnosis? A. Dysmenorrhea B. HIV infection

C. Urethritis (kev)

RSI's primary data source is an OWL ontology representation of Elsevier's Merged Medical Taxonomy (EMMeT), dubbed EMMeT-OWL [7, 41]. RSI uses the EMMeT relation *hasClinicalFinding (hCF)*, which relates Diseases or Symptoms to Diseases, Symptoms or ClinicalFindings, each of which can be used as a question's stem entities (in this case, the patients symptoms). A fragment of the ontology from which the question was generated is listed in Figure 1:

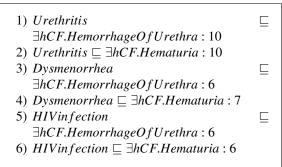


Fig. 1. A snippet of EMMeT-OWL used to provide data for Q2 where the annotations (: n) represent the strength of the hCF relation which range from *most common clinical finding* (10) to *rare clinical finding* (7), including a rank for a known non-relation *not a clinical finding* (6).

Since the question is asking for the *most likely* diagnosis, the option entity that has the strongest relation to the stem entities will be the key.

Definition 3.1 (*stemInd*). Let S be the set of symptoms and **k** be the key. Let *rank* be a function that returns the rank of any annotated axiom and let *min* and *max* be functions that return the minimum and maximum ranks that a given relation can have (usually 7 (rare clinical finding) and 10 (most common clinical finding) respectively). Then *Stem indicativeness* (*stemInd*) is defined as follows⁹:

⁹ Note that hCF relations used in the equations only serve as an example and it can be replaced by any relations associated with strength.

$$stemInd(\mathcal{S}, \mathbf{k}) = 1 - \left(\frac{\sum\limits_{s}^{S} (rank(\mathbf{k} \sqsubseteq \exists hCF.s) - min(hCF))}{|\mathcal{S}| \times (max(hCF) - min(hCF))}\right)$$

The Option entity difference measure (optDiff) is defined in terms of each individual distractor difference (disDiff).

Definition 3.2 (*disDiff*). Let S be the set of symptoms, **d** be a distractor and **k** be the key. Then *disDiff*, is defined as follows:

$$disDiff(\mathcal{S}, \mathbf{k}, \mathbf{d}) = \frac{1}{\left(\frac{\sum\limits_{s}^{\mathcal{S}}(rank(\mathbf{k} \sqsubseteq \exists hCF.s) - \mathbf{d}_{s}) \times \mathbf{d}_{s}}{|\mathcal{S}|}\right)}$$

where 1 is the number of stem components (usually the histories and symptoms, however in this example only symptoms are used) and $\mathbf{d}_s = rank(\mathbf{d} \sqsubseteq \exists h CF.s)$.

Using this measure allows *optDiff* to be defined:

Definition 3.3 (*optDiff*). Let \mathcal{D} be the set of distractors. *optDiff* is defined as follows:

$$optDiff(\mathcal{D}, \mathcal{S}, \mathbf{k}) = \sum_{d}^{\mathcal{D}} \left(disDiff(\mathcal{S}, \mathbf{k}, \mathbf{d})^2 \right)$$

The overall question difficulty is simply the average of *optDiff* and *stemInd*.

We demonstrate the use of RSI using Q2. *Stem indicativeness* equates to 0, showing that the stem is indicative of the key, and therefore has a low difficulty score. The more indicative the stem is of the key, the less difficult the question will be, and vice-versa. The *distractor difficulty* for Dysmenorrhea equates to 0.0444 whilst the difficulty of HIVinfection equates to 0.0416, indicating that Dysmenorrhea is more difficult than HIVinfection, or, it would be harder to eliminate Dysmenorrhea as a distractor compared to HIVinfection since the former has stronger relations to the stem entities than the latter. *Option entity difference* then equates to 0.0037, leading to an overall question difficulty of 0.00185. Suppose that instead of axioms 3 and 4 in Figure 1, the following axioms were present:

- 3) Dysmenorrhea ⊑ ∃hCF.HemorrhageOfUrethra : 10
- 4) *Dysmenorrhea* $\sqsubseteq \exists hCF.Hematuria : 9$

Demo

The distractor difficulty for Dysmenorreah would instead equate to 0.2222, and thus the option entity dif-2 ference would change to 0.0511. This demonstrates the 3 effectiveness of RSI: the more similar the distractors 4 5 are to the key, i.e., the more indicative the stem is to 6 the distractors when compared to the key, the more difficult a question is considered, and vice-versa. 7

The questions studied and reviewed in [7] often use more complex stems. These include multiple types of stem entities such as: risk factors (via the hasRiskFactor relation); and patient demographics. The difficulty and similarity calculations are adjusted to account for additional stem entities and relations, where averages are usually taken over each calculation.

4. Method

To evaluate the performance of both experts and automated measures, we conducted two experiments: an expert review and a mock exam. Both experiments are described below.

4.1. Expert review

26 In a previous study [7], we carried out an expert review to evaluate the ontology-based approach we de-27 veloped for generating medical case-based MCQs. As 28 part of the review, experts rated the usefulness of gen-29 erated questions (i.e. whether or not they are ready to 30 use in an exam context) and predicted their difficulty. 31 In what follows, we explain aspects of the review that 32 are centered around expert prediction of difficulty. For 33 a detailed description of the generation approach and 34 the assessment of question usefulness, see [7]. 35

4.1.1. Subjects

Fifteen experts were recruited to review the ques-38 tions and were paid for their participation. Demo-39 graphic information including education level, practical experience, teaching experience and exam construction experience were collected at the start of the review (Table 1).

4.1.2. Questions

The EMMeT-OWL ontology, which contains def-45 initions of concepts such as diseases, clinical find-46 47 ings, drugs, symptoms, and risk factors, was utilised 48 as a source for question generation. Four physician specialties (internal medicine, cardiology, orthopedics, 49 and gastroenterology) were selected and a total of 50 3,407,493 case-based questions were automatically 51

Demographic characteristics	Number	
Speciality		
Internal medicine	5	
Gastroenterology	4	
Cardiology	5	
Orthopedics	1	
Level		
Resident	1	
Generalist	7	
Specialist	7	
Experience as a practitioner		
None	2	
Less than 1 year	0	
1-3 years	4	
3-6 years	3	
More than 6 years	6	
Teaching experience		
None	0	
Less than 1 year	1	
1-3 years	6	
3-6 years	3	
More than 6 years	5	
Exam construction experience		
None	4	
Less than 1 year	6	
1-3 years	2	
3-6 years	1	
More than 6 years	2	

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

generated from these specialties. The generated questions belong to four templates: 'what is the most likely diagnosis?', 'what is the most likely clinical finding?', 'what is the drug of choice?', and 'what is the differential diagnosis?'. A stratified random sample of 435 questions was selected for expert review. Five stratifiers were used: speciality, question template, the number of distractors (key-distractor combinations in the case of differential diagnoses questions), the number of stem entities, and difficulty as predicted by relation strength indicativeness measure. We targeted an equal number of questions from each strata but this was not possible due to the small number of questions in some strata. More details about each type of question and sample selection can be found in [7]. Out of the 435 questions, 375 questions were rated as appropriate by at least one reviewer. We obtained expert predictions for these 375 questions as described next.

6

1

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

36

37

40

41

42

43

4.1.3. Procedure

1

2

3

4

5

6

15

16

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

The expert review was conducted through a webbased questionnaire tool we developed. Each expert reviewed approximately 30 questions belonging to their specialty. To check agreement among experts, questions were reviewed by two experts whenever possible.

Each question was displayed individually and ex-7 perts were asked to solve the displayed question with-8 9 out a time limit. After the experts submitted their se-10 lected answer, they were shown the correct answer, while an explanation was shown only if experts an-11 swered a question incorrectly. The following data 12 about the performance of domain experts were col-13 lected: 14

Selected answer(s)

- Score: Each single response question answered 17 correctly is given one mark while an incorrect an-18 swer is awarded zero marks. With regards to ques-19 tions with multiple responses (i.e. differential di-20 agnosis questions), a mark for each correct an-21 swer is added to the final mark for each ques-22 tion and a mark of zero is given for fully incor-23 rect answers¹⁰. The awarded mark is compared to 24 the full mark of each question, which is equal to 25 the number of correct options, in order to distin-26 guish fully correct answers from partially correct 27 answers. 28
 - Time to solve: The time starts from displaying the question on the screen and ends by the expert clicking the 'submit' button.

After answering each question, experts were instructed to rate different aspects of the questions (e.g. usefulness, difficulty, and correctness of explanation) while keeping in mind that the questions are targeting resident specialists or practising specialists. They started by rating the usefulness of the question. They were then asked to classify the question as belonging to one of the following difficulty levels:

- Easy: More than 70% of examinees would be expected to answer the question correctly;
- Medium: 30% to 70% of examinees would be expected to answer the question correctly;

- Difficult: Less than 30% of examinees would be expected to answer the question correctly.

They were also provided with an optional comment box for any additional information that they may have wanted to add. The main aim of obtaining expert prediction is to compare it with student performance. Therefore, we did not collect their predictions for questions rated as inappropriate to use in an exam context, since these questions would not be used in the mock exam.

4.2. Mock exam

To obtain the empirical difficulty of the selected set of questions (i.e. percentage correct), we administered the questions to a cohort of residents. Details about the cohort, the questions, and the procedure are explained next.

4.2.1. Subjects

Twelve residents, with a mean age of 32 years (standard deviation = 2.3), were recruited to participate in this experiment and were paid for their participation. Participants completed a demographics questionnaire, which asked them to indicate their age, sex, and practical experience (i.e. number of years working as a practitioner). Table 2 summarises their demographic information.

Table 2
Demographic characteristics of residents who took the mock exam.

		55
Demographic characteristics	Number	34
Sex		35
Male	10	36
Female	2	37
Specialty		38
Orthopedics	5	39
Internal medicine	4	40
Gastroenterology	2	41
Cardiology	1	42
Experience as a practitioner		43
None	2	44
Less than 1 year	0	45
1-3 years	3	46
3-6 years	3	47
6-9 years	2	48
More than 9 years	2	49
		50

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

¹⁰As the exam was experimental and no marks were displayed for participants, it made no sense to use negative marking. One could argue that participants could get the full mark on multiple response questions (only differential diagnosis questions in our sample) by selecting all options. We ensured that this was not the case by looking for such a pattern in the responses to differential diagnosis questions.

4.2.2. Questions

We used disproportional stratified random sampling, aiming for equal group proportions whenever possible, to select questions from our sample space which consists of auto-generated questions rated as appropriate by at least one domain expert in the expert study (345 questions). We used this sampling technique to obtain a representative sample of each group in the population which was not possible using other sampling techniques (e.g. random sampling or proportional stratified sampling) due to the large difference in size between groups in the population.

13 We based stratification on four stratifiers: specialty, 14 template, difficulty as predicted by our measure, and 15 difficulty as predicted by the domain experts. Strati-16 fying by specialty was necessary to ensure that resi-17 dents from different specialties were tested on ques-18 tions covering areas they are expected to be knowl-19 edgeable about. In addition, using templates as a strat-20 ifier allowed us to investigate the applicability of the 21 measures to different question types, and to investi-22 gate whether differences in difficulty can be attributed 23 to the intrinsic nature of the templates themselves. Fi-24 nally, stratifying based on our difficulty measure and 25 the experts' predictions was used to allow investiga-26 tion of the performance of these measures in predicting 27 empirical difficulty. 28

The sample size for each specialty was determined 29 considering a reasonable duration of testing (60-30 minute exam). This resulted in a sample of 231 ques-31 tions in total to be administered to the residents in-32 volved in the experiment. The distribution of these 33 questions is stated in Table 3. Variation in the number 34 of questions across specialties was due to the unequal 35 36 number of experts in each specialty and therefore, the 37 unequal number of reviewed questions. The selected questions were reviewed for linguistic issues and min-38 39 imal edits were applied where necessary. For example, 40 the stem "A patient with a history of acetaminophen 41 presents with ..." was edited to read: "A patient who 42 has used acetaminophen presents with ...". This step 43 was carried out to eliminate the effect of linguistic am-44 biguity on empirical difficulty.

4.2.3. Procedure

45

46

A web-based system was developed to administer the questions and collect performance data. Residents agreed to complete a 60-minute mock exam using their own machines and were assigned questions belonging to their specialty, in addition to internal medicine questions.¹¹ For example, orthopedic residents were assigned the 29 orthopedic questions in addition to the 92 internal medicine questions. The questions were presented in a random order in order to avoid systematic bias resulting from position effects on difficulty. Residents were not shown feedback indicating whether they answered the questions correctly or not. For each question attempted, the following data were collected:

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

- Selected answer(s)

- Score: the same as in the expert review;
- Time to solve: the same as in the expert review.

4.2.4. Data analysis

A standard test theory analysis [42] was conducted for internal medicine questions that were administered to ten residents or more. The possible values that difficulty (percentage correct) can take and how they are interpreted is as follows:

- Easy: percentage correct >70%
- Medium: $30\% \leq \text{percentage correct} \geq 70\%$
- Difficult: percentage correct <30%

The percentage correct was then compared to difficulty as predicted by the aforementioned measures. However, this type of item analysis was not possible for questions belonging to the other three specialties due to the low number of participants they had been administered to (1 to 5 residents at most).

We designed a new approach for analysing difficulty data for questions answered by less than ten participants. To investigate the relation between expert prediction and empirical difficulty, we grouped the questions based on expert prediction, resulting in three groups: easy, medium, and difficult questions according to the experts. We then computed the percentage correct for each group by dividing the total number of correct responses to all questions in the group by the total number of responses (correct and incorrect) to all questions in the group. One would expect that the number of correct responses to difficult questions, for example, to be low, and therefore the percentage correct for the difficult group to be low. A similar procedure was followed to investigate the relation between automated difficulty measures and percentage correct.

While studies concerned both with investigating expert ability in predicting question difficulty [43, 44] and with building difficulty models [6, 17] use the ac-

8

1

2

3

4

5

6

7

8

9

10

11

¹¹Domain experts indicated that all residents are expected to have knowledge in internal medicine.

1	estion sample per specialty and quest ate 3 = What is the most likely clinic	¥1 1		•	0	
Γ	Specialty	Template 1	Template 2	Template 3	Template 4	Total
	Cardiology	41	7	8	7	63
	Gastroenterology and hepatology	30	10	4	3	47
	Internal medicine	53	14	8	17	92
	Orthopedics	17	9	3	0	29
Γ	Total	141	40	23	27	231

curacy metric (Appendix A) for performance evalua-11 tion, we extend the evaluation by using approaches and 12 metrics borrowed from the information retrieval and 13 machine learning communities. The analysis was ex-14 tended to include other metrics because accuracy does 15 not reflect the performance of prediction when the dis-16 tribution of classes (easy, medium, and difficult ques-17 tions in our case) is not balanced. Another reason is 18 that difficulty is an ordinal variable, and it is therefore 19 important to find how close or far away the prediction 20 is from the empirical difficulty. 21

1

2

3

4

5

6

7

8

9

10

45

The following metrics, which are standard in clas-22 sification problems, were used to compare measures 23 for difficulty prediction: accuracy, precision, recall, F-24 score, and kappa. We also used the evaluation metric, 25 26 'average relative error', which has been used in [17] for evaluating the performance of different machine 27 learning models for predicting the difficulty of reading 28 comprehension questions. We explain how we calcu-29 lated these metrics in Appendix A. 30

31 Since different performance metrics focus on different aspects of the prediction, it is therefore essen-32 tial to consider all of them, prioritising them based 33 on the problem at hand, to allow comparison between 34 35 the performances of the different methods. That is, 36 which metrics do we care about in the case that differ-37 ent metrics give contradictory results? For example, it 38 is usually the case that classification methods have a 39 high precision but low recall, or vice versa. The win-40 ning method depends on the metric that is prioritised, 41 whether it is higher precision or better recall. Our dis-42 cussion of metrics is guided by the following characteristics of the problem of prediction of question diffi-43 44 culty:

 The distribution of difficulty levels is not balanced, with the difficult questions being the minority class. This is apparent from the distribution of difficulty levels in the test set in addition to the literature about MCQ examinations [for example, see: 21, 28, 45, 46]. All of the classes are of importance, with little preference for good performance on difficult questions for two reasons: in addition to them being the minority class, appropriately difficult questions play an important role in discriminating between low- and high-information students.

As we were interested in performance for all difficulty levels, we averaged over the precision for each difficulty level, thereby penalising prediction methods that perform well on some of the difficulty levels. A similar calculation was performed for recall and F-score.

To answer the question of 'whether experts and automated measures do better than random guessing?', we compare their performance with the performance of the following three naive methods as baselines:

- random guesser which assigns difficulty levels arbitrarily;
- weighted guesser which assigns difficulty levels according to their distribution in the test set;
- majority class classifier which assigns the most common difficulty level in the test set (medium) to all questions.

5. Results and Discussion

5.1. Residents' performance

Following the description of the difficulty levels in Section 4.2.4, 39.1% (n=36) of the 92 internal medicine questions were easy, 44.6% (n=41) were medium, and 16.3% (n=15) were difficult. We consider this to be a good indicator of question suitability as a test set, since this distribution of difficulty levels is similar to the distribution of difficulty levels reported in analyses of real exams (as examples, see [21, 45]). Residents' scores range from 58.49 to 77.65 with an average of 67.69 (\pm 5.85) (see Table 4 for details). Comparing these results to the results achieved by domain experts (range: 63.64 to 80.65,

9

1

2

3

4

5

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

Id	No. of questions	Score (out of 100)	% of questions answered correctly
S1	155	77.65	74.19
S2	139	75.47	71.94
S3	92	73.40	69.57
S 4	121	71.02	67.77
S5	92	69.73	65.22
S6	92	66.97	61.96
S7	121	65.94	61.98
S8	121	65.22	60.33
S9	92	64.22	57.61
S10	121	62.32	57.85
S11	103	61.86	56.31
S12	139	58.49	53.96
Average	115.67	67.69	63.22

Table 4										
nts'	performance on the mock exam.	Score is calculated as the	e percentage of the tota	al possibl						

mean: 72.09 ± 5.30) indicates that participants are adequately knowledgeable.

5.2. Performance of the measures

5.2.1. Is expert prediction a good proxy for difficulty?

Overall, the accuracy of expert prediction ranges be-tween 46% and 53%. As one can see from Table 6, the accuracy of experts is close (less than 10% varia-tion in accuracy between experts). However, looking at other metrics, more variation in performance between-and within-experts can be seen. Of interest are the low values for precision, recall, and thus F-score on dif-ficult questions compared to easy and medium questions,¹² which suggests that domain experts are less precise and complete in classifying difficult questions as compared to easy or medium questions. Given that domain experts who are involved in the experiment have teaching and exam construction experience, it is expected that they have more self-training (compar-ing one's own prediction with student performance) in predicting the difficulty of easy and medium questions since these represent a majority. The amount of self-training is a possible explanation of the difference in performance.

A point of interest is whether or not there are consistent patterns characterising expert prediction. An example of a pattern is experts having a tendency to underestimate or overestimate the difficulty of questions.

¹²We performed a one way repeated measure ANOVA to compare the effect of actual difficulty of questions on F-scores achieved by experts. The F-score differed significantly between the different difficulty levels (F(2,8) = 10.96, p < 0.05).

Looking at the data, we found 44 questions for which experts overestimated the difficulty compared to 21 questions for which experts underestimated the difficulty. This suggests that experts tend to overestimate difficulty as opposed to underestimating it. We ran a further analysis of the relation between experts' performance on questions (getting the question right or wrong) and their prediction. The analysis aimed to answer two questions: 1) Is there a relation between experts' performance and their prediction accuracy? 2) Is there a relation between experts' performance and overestimation or underestimation of difficulty? Regarding the first question, the data suggest that experts were more accurate in their prediction when they answered the questions correctly. The prediction of 51%of questions solved correctly was accurate compared to 36% of question solved incorrectly. Concerning the second question, experts overestimated the difficulty of 63% of the questions they solved correctly, compared to 81% of the questions they solved incorrectly, which hints at an increase in the percentage of overestimation when questions are solved incorrectly. However, the small number of observations, especially the observation about questions solved incorrectly, precludes making a strong conclusion about expert performance and prediction.

Given that expert prediction is considered as a major component of the evaluation framework for difficulty measures, which is apparent from relying heavily on expert prediction as a source of validation in multiple studies [6, 47, 48], the performance of domain experts was lower than anticipated. However, all experts outperform the three baseline classifiers in each of the prioritised metrics (i.e. accuracy, kappa, average pre-

cision, average recall, and average F-score) except for
the relative error metric, which is outperformed by the
majority classifier. However, this is due to the majority
of the questions in the test set belonging to the medium
level and therefore the distance between any misclassified level and the actual difficulty level is minimal.

7 With regards to questions belonging to other specialties, a Fisher's exact test¹³ was performed, compar-8 9 ing the frequency of responses to questions belonging 10 to the three difficulty levels (see Table 5), as predicted 11 by domain experts. Since the P-value of the test (0.003)12 is less than the significance level (0.05), we can con-13 clude that a dependency exists between expert predic-14 tion and student performance. As can be seen in Ta-15 ble 5, easy questions have a higher percentage of cor-16 rect responses and a lower percentage of incorrect re-17 sponses as compared to medium questions. However, 18 this was not the case for difficult questions. This result, 19 along with the results obtained from internal medicine 20 questions, indicates that expert precision was worst on 21 difficult questions. 22

To summarise, the results indicate that experts moderately predicted question difficulty. The results are suggestive of an adverse effect of expert's performance on their accuracy and of experts' tendency to overestimate question difficulty.

23

24

25

26

27

28

29

30

49

50

51

5.2.2. How well did the automated measures perform in comparison with guessing and in comparison with each other?

31 While preliminary evaluations of the similarity mea-32 sure [6] show that it has potential for predicting ques-33 tion difficulty, the current evaluation shows that the ac-34 curacy of this measure on its own is lower than two of 35 the baseline classifiers (Table 6). However, it is impor-36 tant to note that the similarity measure has been eval-37 uated in questions that have simple stems (i.e. consist 38 of two concepts at most). Most of the questions in our 39 dataset have more complex stems that contain two to 40 five concepts. It is expected that these complex stems 41 contribute to the difficulty of the questions which is 42 not captured by the similarity measure. This seems a 43 plausible justification for its low performance. Com-44 bining the similarity measure with stem indicativeness, 45 as explained in Section 3.2, increases the performance 46 on all metrics except for recall on difficult questions as 47 can be seen in Table 6. The performance of the rela-48

¹³ The Fisher's exact test was selected because of the low frequencies observed in some cells (Table 5).

tion strength indicativeness measure is also better than random and weighted guessers.

Another observation we have made is that the similarity measure tends to overestimate the difficulty of questions. The predicted difficulty of 45 questions (48.91%) was higher than the empirical difficulty. On the other hand, the predicted difficulty of 14 questions (15.22%) was lower than the empirical difficulty. We observed a similar pattern for the relation strength indicativeness measure. We expect that cohort exposure to examined concepts, particularly when reviewing previous or sample exam papers, to moderate the effect of difficulty factors captured by the automated measures. Investigating the relation between cohort characteristics and difficulty remains an area for future research.

Performing Fisher's exact test on questions belonging to other specialties did not reveal a significant difference between the frequencies of correct and incorrect responses to questions belonging to different difficulty levels (as predicted by relation strength indicativeness measure). Results obtained from internal medicine indicate that the distance between predicted difficulty and empirical difficulty is higher in automatic prediction than in expert prediction. Classifying easy questions as difficult, and vice versa, is expected to have a strong impact on the frequency of correct and incorrect responses in each group (see Table 5). Therefore, we attribute the failure in detecting a significant relation to the high value of the average relative error (Table 6).

5.2.3. How well did the automated measures perform in comparison to domain experts?

The performance of our measure was competitive 35 compared with the performance of domain experts. 36 Looking at Table 6, it can be seen that the relation 37 strength indicativeness measure ranks higher than low-38 performing experts on all prioritised metrics except for 39 the relative error metric. This indicates that difficulty 40 levels assigned by domain experts are closer to the ac-41 tual difficulty levels than the difficulty levels assigned 42 by the automated measure. This can be explained by 43 the ability of domain experts to recognise other fea-44 tures (e.g. linguistic features) that play a role in the dif-45 ficulty of the questions. For example, while the rela-46 tion strength indicativeness measure predicts questions 47 with indicative stems and low-similarity distractors to 48 be easy, the language complexity of the questions or 49 the use of rare concepts increases the difficulty of the 50 question. In addition, experts have pedagogical con-51

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

Table 5

Frequency of responses to questions belonging to difficulty levels predicted by: a) domain experts; b) relation strength indicativeness measure. Raw numbers are presented between parentheses.

			С	orrectness of respon	ises			
		Incor	rect	Partially correct	Cor	rect	Percentage correct	Total responses
	Easy	17.37	(33)	3.68 (7)	78.95	(150)	79	(190)
a)	Medium	34.07	(46)	5.19 (7)	60.74	(82)	61	(135)
	Difficult	11.11	(2)	0 (0)	88.89	(16)	89	(18)
	Easy	23.78	(39)	4.27 (7)	71.43	(118)	72	(164)
b)	Medium	23.23	(23)	0 (0)	76.77	(76)	78	(99)
	Difficult	28.57	(10)	0 (0)	71.43	(25)	71	(35)

Table 6

Performance of different methods on difficulty prediction of internal medicine questions. Rank of each method among others is enclosed in parentheses and boldface indicates the method with the best performance in each metric. Acc. = accuracy; rel. error= relative error; E = easy; M = medium and D = difficult.

								Prec	ision			Re	call			F-se	core	
Method	#q	Acc.	Rel. error	Kappa	E	М	D	Avg.	Е	М	D	Avg.	E	М	D	Avg.		
							J	Baseline										
Random	-	.33 (9)	.44 (7)	0 (8)	.33 (8)	.33 (8)	.33 (2)	.33 (7)	.33 (6)	.33 (8)	.33 (5)	.33 (6)	.36 (7)	.39 (7)	.22 (6)	.32 (8)		
Weighted	-	.38 (7)	.38 (4)	0 (8)	.39 (7)	.45 (4)	.16 (7)	.33 (7)	.39 (5)	.45 (6)	.16 (6)	.33 (6)	.39 (6)	.44 (6)	.19 (7)	.34 (6)		
Majority	-	.45 (6)	.28 (1)	0 (8)	Na	.45 (4)	Na	.15 (8)	0 (8)	1 (1)	0(7)	.33 (6)	Na	.62 (1)	Na	Na		
]	Experts										
Expert 1	22	.46 (5)	.39 (5)	.19 (2)	.80 (2)	.40 (7)	.29 (4)	.50 (3)	.40 (4)	.50 (4)	.50 (2)	.47 (2)	.53 (4)	.44 (6)	.36 (4)	.45 (3		
Expert 2	35	.46 (5)	.36 (3)	.16 (5)	1 (1)	.44 (5)	.25 (6)	.56 (1)	.42 (3)	.47 (5)	.50 (2)	.46 (3)	.59 (3)	.46 (5)	.33 (5)	.46 (2		
Expert 3	20	.50 (3)	.36 (3)	.18 (3)	.63 (4)	.63 (1)	0 (8)	.42 (5)	.71 (1)	.42 (7)	0(7)	.38 (5)	.67 (1)	.50 (4)	0 (8)	.39 (4		
Expert 4	23	.52 (2)	.36 (3)	.05 (7)	.63 (4)	.40 (7)	0 (8)	.34 (6)	.71 (1)	.29 (9)	0(7)	.33 (6)	.67 (1)	.33 (9)	0 (8)	.33(7)		
Expert 5	30	.53 (1)	.30 (2)	.24 (1)	.67 (3)	.50 (3)	.40 (1)	.52 (2)	.55 (2)	.57 (2)	.40 (4)	.51 (1)	.60 (2)	.53 (3)	.40 (1)	.51 (1		
							А	utomatic										
[7]	92	.47 (4)	.42 (6)	.17 (4)	.48 (5)	.54 (2)	.32 (3)	.45 (4)	.39 (5)	.54 (3)	.47 (3)	.47 (2)	.43 (5)	.54 (2)	.38 (3)	.45 (3		
[6]	92	.36 (8)	.50 (8)	.08 (6)	.46 (6)	.41 (6)	.27 (5)	.33 (7)	.28 (7)	.29 (9)	.73 (1)	.43 (4)	.35 (8)	.34 (8)	.39 (2)	.36 (5		

tent knowledge (i.e knowledge about challenging concepts that students find difficult to understand or have misconceptions about) which gives them an advantage over automated measures.

6. Methodological Reflection

While we have investigated expert performance on question difficulty prediction, our investigation was fo-cused on medical questions and therefore the general-isability of these results to other domains is unknown. It is possible that other domains are more mature in the sense that pedagogical content knowledge is well-known. This, in turn, would improve expert predic-tion which would provide different results. In addition, we find it worthwhile and interesting to look at do-main experts' characteristics (e.g. teaching experience, and exam construction experience) and how these contribute to their predictive performance. However, the amount of data that we have was limited for conduct-ing such an analysis. Another factor that is expected to

improve expert prediction, and that requires additional studies, is interaction and familiarity with the cohort to be tested.

Automatic measures for difficulty prediction are de-veloped for the purpose of controlling the difficulty of automatically generated questions. This does not pre-clude the use of these measures for predicting the dif-ficulty of hand-written questions (after parsing these questions). One of the limitations of the current study is that our test set consists of automatically gener-ated questions only. These questions are very similar in terms of their linguistic structure. Difficulty predic-tion measures might perform worse on hand-crafted questions that are expected to be inherently more di-verse in their linguistic structure. Another difference between auto-generated and hand-crafted questions is that, as mentioned earlier, the percentage of flawed questions is high among the latter type of questions. This is another expected source of performance vari-ation between different measures for the two sets of questions. However, obtaining hand-crafted examina-tion questions annotated with student performance was

difficult because of exam security issues. Further stud ies that investigate the consistency of the results for
hand-crafted questions are in high demand.

Another point that needs to be emphasised here is 4 5 that, although the questions in the test set belong to 6 four templates, these templates have different characteristics (e.g. number of concepts in the stem, number 7 of keys). In addition, we varied the questions' charac-8 9 teristics within questions belonging to the same template. If the questions had been similar, we would have 10 had no confidence in the generality of the test set and 11 the generalisability of the results. However, at least the 12 different characteristics of the question set increased 13 our confidence in generalizing the results. 14

Finally, it is worth mentioning that the performance 15 16 of both automatic measures investigated in this paper is heavily dependent on the completeness and correct-17 ness of the used ontology. Thus, an interesting next 18 step would be investigating the variation of perfor-19 mance when ontologies with different characteristics 20 21 (e.g., size and expressivity) are used. Taking a different perspective, the performance of these measures can 22 also be used as an indication of ontology quality. 23

7. Conclusion

24

25

26

27

To the best of our knowledge, this study is the first to 28 compare the performance of domain experts and naive 29 and automated methods for MCQ difficulty prediction. 30 With respect to RQ1, experts moderately predicted the 31 difficulty of questions and were more accurate in pre-32 dicting easy and medium questions compared to diffi-33 cult questions. Regarding RQ2, the comparison shows 34 that the relation strength indicativeness measure out-35 36 performs the similarity-based measure. Moreover, the 37 new ontology-based difficulty measure is of comparable performance to that of domain experts, who are 38 heavily relied on in practice. We consider this as a 39 major success since it can be used as an economical 40 alternative. We believe that the ability of our model 41 to explain its decisions (why a particular question is 42 classified as belonging to a particular difficulty level), 43 whether the decision is correct or not, is another point 44 of strength. These justified decisions can make exam 45 designers consider new aspects of questions, which in 46 47 turn provide new insights about the difficulty and va-48 lidity of questions.

However, investigating additional factors that can be
used to predict the difficulty of both automatically gen erated questions and hand-written questions is still a

subject of ongoing research. While doing this, the criteria presented in this study need to be considered as the minimum set of evaluation criteria.

Finally, while we made an attempt at creating an annotated question set that can be used for testing the performance of prediction methods, a larger question set is needed to cross-validate the results and gain more confidence in their consistency, as well as to provide statistical significance. In addition, a larger question set will allow the use of standard machine learning algorithms for building prediction models and investigating whether these models outperform the ontologybased measures compared in this study.

Acknowledgement

We would like to thank all participants of our experiments for their valuable contributions to our work.

Funding

This work was funded by an EPSRC grant (ref: EP/P511250/1) under an Institutional Sponsorship (2016) for The University of Manchester, along with a partial contribution from Elsevier. The funding acts as a secondment to an initial EPSRC grant (ref: EP/K503782/1) awarded as an Impact Acceleration Account (2016) for The University of Manchester.

Appendix A. Calculation of the evaluation metrics

Let $D = \{e, m, d\}$ be a set of difficulties (e = easy, m = medium, and d = difficult) and let Q be a set of questions $\{q_1, ..., q_n\}$. Let $actDif : Q \to D$ be a function over Q and D that returns the actual difficulty of a question (as derived from percentage correct) and let $preDif : Q \to D$ be a function over Q and Dthat returns the predicted difficulty of a question. Let $Q_{pc} \subseteq Q$ be the set of correctly classified questions, i.e. $q \in Q_{pc}$ if actDif(q) = preDif(q). We can define accuracy as follows:

$$Accuracy = \frac{|Q_{pc}|}{|Q|}.$$

Possible values are between 0 and 1 with 1 indicating that all questions are correctly classified.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

For $x \in D$, let $Q_x \subseteq Q$ be the set of questions with the difficulty level x s.t $q \in Q_x$ if actDif(q) = x and let $Q_{px} \subseteq Q$ be the set of questions predicted as being x s.t $q \in Q_{px}$ if preDif(q) = x. Precision for Q_x is defined as follows:

$$Precision_{Q_x} = \frac{|Q_x \cap Q_{pc}|}{|Q_{px}|}.$$

14

The value ranges from 0 to 1 with higher values indicating that the classifier is less likely to identify questions as being x while they are actually not. Next, we define the recall on Q_x as:

$$Recall_{Q_x} = \frac{|Q_x \cap Q_{pc}|}{|Q_x|}$$

The value ranges from 0 to 1 with a value of 1 indicating that the classifier has identified all questions in Q_x and a value of 0 indicating that it has missed all questions in Q_x . In what follow, we define the F - score on Q_x :

$$F - score_{Q_x} = 2 * \frac{Precision_{Q_x} * Recall_{Q_x}}{Precision_{Q_x} + Recall_{Q_x}}$$

 $F - score_{Q_x}$ ranges between 0 and 1. The closer the $precision_{Q_x}$ and $recall_{Q_x}$ to each other, the greater the value.

Let max be a function that returns the maximum possible error where each $x \in D$ is associated with numerical values between [1,3], and maximum possible error is the difference between the maximum and minimum values associated with x (in this case, 3-1=2).

Average relative error =
$$\frac{\sum_{n=1}^{|Q|} preDif(q) - actDif(q)}{|Q| * max}$$

The value ranges from 0 to 1. The closer the value to 0, the fewer errors are made by the classifier.

Finally, to define kappa, let p_o be the observed agreement and p_e be the agreement by chance. Then,

$$Kappa(Q, p_o, p_e) = \frac{p_o - p_e}{1 - p_e}$$

The value is less than or equal to 1 with a value of 1 indicating a perfect agreement.

References

434

2 [1] J. Collins, Writing Multiple-Choice Questions for 3 Continuing Medical Education Activities and Self-4 Assessment Modules, RadioGraphics 26(2) (2006), 543-5 551. doi:10.1148/rg.262055145. https://doi.org/10.1148/rg. 6 262055145. 7 [2] J. Considine, M. Botti and S. Thomas, Design, format, validity and reliability of multiple choice questions for use 8 in nursing research and education, Collegian 12(1) (2005), 9 19-24, ISSN 1322-7696. doi:https://doi.org/10.1016/S1322-10 7696(08)60478-3. http://www.sciencedirect.com/science/ 11 article/pii/S1322769608604783. 12 [3] J.D. Wasserman and B.A. Bracken, Psychometric Characteristics of Assessment Procedures, in: Handbook of Psychol-13 ogy, John Wiley and Sons, Inc., 2003. ISBN 9780471264385. 14 http://dx.doi.org/10.1002/0471264385.wei1003. 15 [4] A. Papasalouros, K. Kanaris and K. Kotis, Automatic Genera-16 tion Of Multiple Choice Questions From Domain Ontologies., 17 in: IADIS International Conference e-Learning, 2008, pp. 427-18 [5] M. Cubric and M. Tosic, Towards automatic generation of 19 e-assessment using semantic web technologies, International 20 Journal of e-Assessment $\mathbf{1}(1)$ (2011). 21 [6] T. Alsubait, B. Parsia and U. Sattler, Generating Multi-22 ple Choice Questions From Ontologies: Lessons Learnt., in:

1

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

- OWLED, 2014, pp. 73-84. [7] J. Leo, N. Matentzoglu, G. Kurdi, B. Parsia, S. Forege, G. Donato and W. Dowling, Ontology-based generation of multiterm, exam-ready MCQs, paper submitted to Knowledge-Based Systems (2018).
- [8] M. Al-yahya, Ontology-based multiple choice question generation, The Scientific World Journal 2014 (2014).
- [9] N. Afzal, Automatic Generation of Multiple Choice Ouestions using Surface-based Semantic Relations, International Journal of Computational Linguistics (IJCL) 6(3) (2015), 26-44.
- [10] M. Heilman, Automatic factual question generation from text, PhD thesis, Carnegie Mellon University, 2011.
- [11] V. Crisp and R. Grayson, Modelling question difficulty in an A level physics examination, Research Papers in Education 28(3) (2013), 346-372. https://doi.org/10.1080/02671522. 2012.673005.
- [12] J.D. Scheuneman, Y.V. Fan and S.G. Clyman, An investigation of the difficulty of computer-based case simulations, Medical Education 32(2) (1998), 150-158, ISSN 1365-2923. http://dx. doi.org/10.1046/j.1365-2923.1998.00193.x.
- [13] V. Mesic and H. Muratovic, Identifying predictors of physics item difficulty: A linear regression approach, Phys. Rev. ST Phys. Educ. Res. 7 010110. doi:10.1103/PhysRevSTPER.7.010110. (2011),https://link.aps.org/doi/10.1103/PhysRevSTPER.7.010110.
- [14] J. Stiller, S. Hartmann, S. Mathesius, P. Straube, R. Tiemann, V. Nordmeier, D. KrÄijger and A.U. zu Belzen, Assessing scientific reasoning: a comprehensive evaluation of item features that affect item difficulty, Assessment & Evaluation in Higher Education 41(5) (2016), 721-732. https://doi.org/10. 1080/02602938.2016.1164830.
- [15] R.F. Boldt, GRE Analytical Reasoning Item Statistics Prediction Study, ETS Research Report Series doi:10.1002/j.2333-8504.1998.tb01786.x. 1998(2), 23.

https://onlinelibrary.wiley.com/doi/abs/10.1002/j.2333-8504. 1998.tb01786.x.

[16] M.K. Enright and K. Sheehan, Modeling the difficulty of quantitative reasoning items: Implications for item generation, in: *Item Generation for Test Development*, S.H. Irvine and P.C. Kyllonen, eds, Routledge, 2002, pp. 129–157.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

- [17] D. Hutzler, E. David, M. Avigal and R. Azoulay, Learning Methods for Rating the Difficulty of Reading Comprehension Questions, in: 2014 IEEE International Conference on Software Science, Technology and Engineering, 2014, pp. 54–62. doi:10.1109/SWSTE.2014.16.
- [18] S.M. Downing, Threats to the Validity of Locally Developed Multiple-Choice Tests in Medical Education: Construct-Irrelevant Variance and Construct Underrepresentation, Advances in Health Sciences Education 7(3) (2002), 235–241, ISSN 1573-1677. https://doi.org/10.1023/A:1021112514626.
- [19] J.C. Masters, B.S. Hulsmeyer, M.E. Pike, K. Leichty, M.T. Miller and A.L. Verst, Assessment of multiple-choice questions in selected test banks accompanying text books used in nursing education, *Journal of Nursing Education* 40(1) (2001), 25–32.
- [20] R. Nedeau-Cayo, D. Laughlin, L. Rus and J. Hall, Assessment of item-writing flaws in multiple-choice questions, *Journal for nurses in professional development* 29(2) (2013), 52–57.
- [21] B.R. Rush, D.C. Rankin and B.J. White, The impact of itemwriting flaws and item complexity on examination item difficulty and discrimination value, *BMC Medical Education* 16(1) (2016), 250, ISSN 1472-6920. doi:10.1186/s12909-016-0773-3. https://doi.org/10.1186/s12909-016-0773-3.
- [22] B.M. Hijji, Flaws of Multiple Choice Questions in Teacher-Constructed Nursing Examinations: A Pilot Descriptive Study, *Journal of Nursing Education* 56(8) (2017), 490–496.
- [23] J. Pais, A. Silva, B. Guimarães, A. Povo, E. Coelho, F. Silva-Pereira, I. Lourinho, M.A. Ferreira and M. Severo, Do item-writing flaws reduce examinations psychometric quality?, *BMC Research Notes* 9(1) (2016), 399, ISSN 1756-0500. doi:10.1186/s13104-016-2202-4. https://doi.org/ 10.1186/s13104-016-2202-4.
- [24] Y. Sarin, M. Khurana, M. Natu, A.G. Thomas and T. Singh, Item analysis of published MCQs, *Indian pediatrics* 35(11) (1998), 1103–1105.
- [25] M. Tarrant, J. Ware and A.M. Mohammed, An assessment of functioning and non-functioning distractors in multiplechoice questions: a descriptive analysis, *BMC Medical Education* 9(1) (2009), 40, ISSN 1472-6920. doi:10.1186/1472-6920-9-40. https://doi.org/10.1186/1472-6920-9-40.
- [26] J. Ware and T. Vik, Quality assurance of item writing: During the introduction of multiple choice questions in medicine for high stakes examinations, *Medical Teacher* 31(3) (2009), 238–243. doi:10.1080/01421590802155597. https:// doi.org/10.1080/01421590802155597.
- [27] M.R. Hingorjo and F. Jaleel, Analysis of one-best MCQs: the difficulty index, discrimination index and distractor efficiency, JPMA-Journal of the Pakistan Medical Association 62(2) (2012), 142–147.
- [28] S.K. Namdeo and S.D. Rout, Assessment of Functional and Nonfunctional Distracter in an Item Analysis.
- [29] K. Gautam, I. Gupta and K. Chandramouli, Conceptual Extraction of Questions from Wikipedia.

- [30] A. Singh Bhatia, M. Kirti and S.K. Saha, Automatic Generation of Multiple Choice Questions Using Wikipedia, in: *Pattern Recognition and Machine Intelligence*, P. Maji, A. Ghosh, M.N. Murty, K. Ghosh and S.K. Pal, eds, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp. 733–738. ISBN 978-3-642-45062-4.
- [31] M. Liu, R.A. Calvo, A. Aditomo and L.A. Pizzato, Using Wikipedia and Conceptual Graph Structures to Generate Questions for Academic Writing Support, *IEEE Transactions on Learning Technologies* 5(3) (2012), 251–263, ISSN 1939-1382. doi:10.1109/TLT.2012.5.
- [32] T. Alsubait, Ontology-based question generation, PhD thesis, University of Manchester, 2015.
- [33] S. Williams, Generating Mathematical Word Problems, in: The Association for the Advancement of Artificial Intelligence AAAI Fall Symposium: Question Generation, 2011, pp. 61–64.
- [34] Complexity-based generation of multi-choice tests in AQG systems, in: 2013 IEEE 4th International Conference on Cognitive Infocommunications (CogInfoCom), 2013, pp. 399–402. doi:10.1109/CogInfoCom.2013.6719278.
- [35] V. E.V. and S. Kumar, A novel approach to generate MCQs from domain ontology: Considering DL semantics and openworld assumption, *Web Semantics: Science, Services and Agents on the World Wide Web* 34(Supplement C) (2015), 40– 54, ISSN 1570-8268. http://www.sciencedirect.com/science/ article/pii/S1570826815000475.
- [36] B.S. Bloom, M.D. Engelhart, E.J. Furst, W.H. Hill and D.R. Krathwohl, *Taxonomy of educational objectives, handbook 1: The cognitive domain*, Vol. 19, New York: David McKay Co Inc, 1956.
- [37] J.P.W. Cunnington, G.R. Norman, J.M. Blake, W.D. Dauphinee and D.E. Blackmore, Applying Learning Taxonomies to Test Items: Is a Fact an Artifact?, in: *Advances in Medical Education*, A.J.J.A. Scherpbier, C.P.M. van der Vleuten, J.J. Rethans and A.F.W. van der Steeg, eds, Springer Netherlands, Dordrecht, 1997, pp. 139–142. ISBN 978-94-011-4886-3.
- [38] M.E. Abdalla, A.M. Gaffar and R.A. Suliman, Constructing A-Type Multiple Choice Questions (MCQs): Step By Step Manual, 2011.
- [39] L.W.T. Schuwirth, M.M. Verheggen, C.P.M. Van Der Vleuten, H.P.A. Boshuizen and G.J. Dinant, Do short cases elicit different thinking processes than factual knowledge questions do?, *Medical Education* 35(4) (2001), 348–356, ISSN 1365-2923. doi:10.1046/j.1365-2923.2001.00771.x. http://dx. doi.org/10.1046/j.1365-2923.2001.00771.x.
- [40] P. Jaccard, Étude comparative de la distribution florale dans une portion des Alpes et des Jura, *Bull Soc Vaudoise Sci Nat* 37 (1901), 547–579.
- [41] B. Parsia, T. Alsubait, J. Leo, V. Malaisé, S. Forge, M. Gregory and A. Allen, Lifting EMMeT to OWL Getting the Most from SKOS, in: Ontology Engineering: 12th International Experiences and Directions Workshop on OWL, OWLED 2015, co-located with ISWC 2015, Bethlehem, PA, USA, October 9-10, 2015, Revised Selected Papers, V. Tamma, M. Dragoni, R. Gonçalves and A. Ławrynowicz, eds, Springer International Publishing, Cham, 2016, pp. 69–80. ISBN 978-3-319-33245-1. https://doi.org/10.1007/978-3-319-33245-1_7.
- [42] L. Crocker and J. Algina, *Introduction to classical and modern test theory.*, ERIC, 1986.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

G. Kurdi et al. / Experts vs. Automata: A Comparative Study of Methods for a Priori Prediction of MCQ Difficulty

[43] J.D. Kibble and T. Johnson, Are faculty predictions or item taxonomies useful for estimating the outcome of multiplechoice examinations?, *Advances in Physiology Education* **35**(4) (2011), 396–401. doi:10.1152/advan.00062.2011. https: //doi.org/10.1152/advan.00062.2011.

- [44] G. van de Watering and J. van der Rijt, TeachersâĂŹ and studentsâĂŹ perceptions of assessments: A review and a study into the ability and accuracy of estimating the difficulty levels of assessment items, *Educational Research Review* 1(2) (2006), 133–147, ISSN 1747-938X. doi:https://doi.org/10.1016/j.edurev.2006.05.001. http://www. sciencedirect.com/science/article/pii/S1747938X06000236.
- [45] M. Mukhopadhyay, K. Bhowmick, S. Chakraborty, D. Roy, P.K. Sen and I. Chakraborty, Evaluation of MCQs for Judgement of Higher Levels of Cognitive Learning, *Gomal Journal* of Medical Sciences 8(2) (2010).
- [46] B.S. Malau-Aduli and C. Zimitat, Peer review improves the quality of MCQ examinations, Assessment & Evaluation in Higher Education 37(8) (2012), 919–931. doi:10.1080/02602938.2011.586991. https://doi.org/10.1080/02602938.2011.586991.
 - [47] F.-l. Lee and R. Heyworth, Problem complexity: A measure of problem difficulty in algebra by using computer, *Education Journal* 28(1) (2000), 85–108.

- [48] S. Banerjee, N.J. Rao and C. Ramanathan, Rubrics for assessment item difficulty in engineering courses, in: 2015 IEEE Frontiers in Education Conference (FIE), 2015, pp. 1–8. doi:10.1109/FIE.2015.7344299.
- [49] S.M. Case and D.B. Swanson, Extended matching items: A practical alternative to free response questions, *Teaching and Learning in Medicine* 5(2) (1993), 107–115. doi:10.1080/10401339309539601. https://doi.org/10.1080/10401339309539601.
- [50] E. Ibrahem, Automated MCQ generation: An ontology-based approach for Java knowledge assessment and ontology validation, Master's thesis, The University of Manchester, 2016.
- [51] V. E.V and P.S. Kumar, Automated generation of assessment tests from domain ontologies, *Semantic Web* 8(6) (2017), 1023–1047. doi:"10.3233/SW-170252". "https:// content.iospress.com/articles/semantic-web/sw252".
- [52] S. Gajjar, R. Sharma, P. Kumar and M. Rana, Item and test analysis to identify quality multiple choice questions (MCQs) from an assessment of medical students of Ahmedabad, Gujarat, *Indian journal of community medicine: official publication of Indian Association of Preventive & Social Medicine* **39**(1) (2014), 17–20. doi:10.4103/0970-0218.126347.