



addition to quality and discrimination, construct validity, and appropriate difficulty should be taken into account.

We have sequentially developed three AQG methods and conducted evaluations in terms of the proposed viewpoints. First, we have built a semi-automatic question generator that allows the test author to compose questions just with some clicks on the screen. Secondly, we evaluate randomly-generated questions with a group of students. In terms of discrimination power, our AQG method for grammar questions is as efficient as workbook questions. Evaluation on construct validity shows some evidence that the pattern-generated questions measure intended grammar targets. Finally, we present a CAT (Computer Adaptive Testing) system that administrates automatically generated questions. We have developed a difficulty predictor using machine learning techniques, which can be used for newly generated questions. Evaluation on difficulty adaptivity shows that the predicted difficulty value has more information on the actual correct response rate than the baseline index (sentence length) alone.

Keywords: Language testing, Grammar and vocabulary test, Multiple-Choice Fill-In-the-Blank (MC-FIB), IRT (Item Response Theory), NLP (Natural Language Processing), Machine learning