



Title	A comparison of the psychometric properties of three- and four-option multiple-choice questions in nursing assessments
Author(s)	Tarrant, M; Ware, J
Citation	Nurse Education Today, 2010, v. 30 n. 6, p. 539-543
Issued Date	2010
URL	http://hdl.handle.net/10722/134855
Rights	NOTICE: this is the author's version of a work that was accepted for publication in Nurse Education Today. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in Nurse Education Today, 2010, v. 30 n. 6, p. 539-543. DOI: 10.1016/j.nedt.2009.11.002

Manuscript Number: NET-D-09-00044R2

Title: A COMPARISON OF THE PSYCHOMETRIC PROPERTIES OF THREE- AND FOUR-OPTION
MULTIPLE-CHOICE QUESTIONS IN NURSING ASSESSMENTS

Article Type: Full Length Article

Corresponding Author: Dr. Marie Tarrant, RN MPH PhD

Corresponding Author's Institution: University of Hong Kong

First Author: Marie Tarrant, RN MPH PhD

Order of Authors: Marie Tarrant, RN MPH PhD; James Ware, MA MB FRCS DMSc

Manuscript Region of Origin: HONG KONG

Abstract: In multiple-choice tests, four-option items are the standard in nursing education. There are few evidence-based reasons, however, for MCQs to have four or more options as studies have shown that three-option items perform equally as well and the additional options most often do not improve test reliability and validity. The aim of this study was to examine and compare the psychometric properties of four-option items with the same items rewritten as three-option items. Using item analysis data to eliminate the distractor with the lowest response rate, we compared three- and four-option versions of 41 multiple-choice items administered to two student cohorts over two subsequent academic years. Removing the non-functioning distractor resulted in minimal changes in item difficulty and discrimination. Three-option items contained more functioning distractors despite having fewer distractors overall. Existing distractors became more discriminating when infrequently selected distractors were removed from items. Overall, three option-items perform equally as well as four-option items. Since three option-items require less time to develop and administer and additional options provide no psychometric advantage, teachers are encouraged to adopt three-option items as the standard on multiple-choice tests.

A COMPARISON OF THE PSYCHOMETRIC PROPERTIES OF THREE- AND FOUR-OPTION MULTIPLE-CHOICE QUESTIONS IN NURSING ASSESSMENTS

Marie Tarrant RN MPH PhD¹
Assistant Professor
Department of Nursing Studies
4/F, William M. W. Mong Block
Li Ka Shing Faculty of Medicine
21 Sassoon Road, Hong Kong
Tel: +852 2819 2643
Fax: +852 2872 6079
Email: tarrantm@hku.hk

James Ware MA MB FRCS DMSc
Director of Medical Education
Centre of Medical Education, Faculty of Medicine
Health Sciences Centre, Kuwait University
PO Box 24923
Safat, 13110 Kuwait

¹Corresponding author

Word Count: 3146 words

Short Title: Three-option multiple-choice questions

Keywords: multiple-choice questions, multiple-choice tests, distractors, item analysis, test construction, number of choices per item, item discrimination, assessment.

Acknowledgements:

Financial support for this study was provided by the Leung Kau Kui / Run Run Shaw Research and Teaching Endowment Fund, the University of Hong Kong.

1 **Introduction**

2 In nursing education, multiple-choice questions (MCQs) are one of the most
3 popular written assessment formats. Single best-answer MCQs consist of a question,
4 two or more choices from which examinees must choose the correct option (the
5 distractors), and one correct or best response. While MCQs are often criticized for
6 largely assessing factual recall over higher cognitive thinking (Pampllett and Farnhill,
7 1995), MCQs still offer many advantages when compared with other types of written
8 assessment. Despite what many teachers believe, MCQs are adaptable to different,
9 although not all, levels of learning outcomes (Gronlund and Waugh, 2008). High
10 quality MCQs present clinical vignettes to students that mimic actual clinical
11 problems and assess application of knowledge rather than simple factual recall (Case
12 and Swanson, 2003). Therefore, well constructed MCQs can accurately discriminate
13 between high- and low-ability students (Schuwirth and van der Vleuten, 2003).
14 MCQs are objective and they allow teachers to test a wider range of content and
15 educational objectives than many other written assessment methods. Additionally,
16 MCQs allow teachers to efficiently assess large numbers of candidates as they are
17 easy to administer and score (McCoubrie, 2004). Furthermore, because of this
18 broader sampling of content and because MCQ test items can be subjected to post-
19 test review using item analysis procedures, MCQ tests have higher validity than
20 other test methods such as short-answer or essay-style questions (Gronlund and
21 Waugh, 2008).

22 Four-option MCQs remain the standard in nursing, both on in-house
23 developed tests (Tarrant et al., 2006) and in test banks and text books used in
24 nursing education (Masters et al., 2001). In other health science disciplines, such as

1 medicine, five-option items are more common (Haladyna and Downing, 1993).
2 Although measurement specialists have long discovered that there are few evidence-
3 based reasons for MCQs to have four or five options, many introductory books on
4 item writing continue to recommend this practice and a majority of teachers
5 continue to follow this recommendation (Owen and Froman, 1987). Three-options
6 items however, have many advantages over four- and five-option items, including
7 less time required to construct items and less testing time required. Alternately, with
8 less time required to complete three-option items, teachers are able to increase the
9 number of items administered on a test and thereby increase the amount of content
10 tested (Haladyna and Downing, 1993). Furthermore, researchers have shown that in
11 both teacher-generated (Tarrant et al., 2009) and professionally-developed
12 (Haladyna and Downing, 1993) four-and five-option MCQs, students rarely select
13 more than two or three of the options.

14 In most nursing programmes, the amount of content that requires
15 assessment can be overwhelming. A substantial proportion of a teacher's time is
16 spent on developing written assessments and since a substantial proportion of those
17 assessments will likely contain MCQs, it is important that teachers are basing those
18 practices on the best available research evidence. Additionally, student numbers are
19 generally increasing to meet workplace shortages, while at the same time the
20 number of available teaching faculty is often getting smaller (Broome, 2009).
21 Therefore, because of their efficiency and ability to assess different learning
22 outcomes, MCQs are likely to continue to remain an important component of written
23 assessment in many nursing programmes for the foreseeable future. Thus if the time
24 required to develop multiple-choice tests can be reduced without reducing the

1 reliability and validity of the assessment, this is an important consideration for
2 nursing faculty.

3 **Background**

4 Numerous research studies have compared three-, four-, and five-option
5 MCQs and most have found that three-option items perform equally as well or better
6 than either four- or five-option items. Sidick, Barrett, & Doverspike (1994) rewrote
7 68 five-option items on public sector employment tests by removing the two least
8 functional distractors. Overall, there was little difference in the psychometric
9 properties between the three- and five-option items. Similarly, Rogers and Harley
10 (1999) rewrote 31 four-option items on a senior secondary school mathematics test
11 into three-option items by eliminating the least functioning distractor. The test with
12 three-option items was less difficult than but equally as discriminating and reliable as
13 the four-option test. As part of pre-college admissions testing, Trevisan et al. (1991)
14 administered the same 45 items in three-, four-, and five-option formats. The three-
15 option test form contained more highly discriminating items and fewer items with
16 non-functioning distractors than the four- or five-option test forms. Owen and
17 Froman (1987) randomly administered 100 items to 114 undergraduate psychology
18 students as either five-option items or three-option items and found no significant
19 differences in either item discrimination or difficulty. Crehan et al. (1993) eliminated
20 the least functional distractor in 12 four-option items and found no differences in
21 discrimination between the three- and four-option formats, although again three-
22 option items were slightly less difficult. In one of the few studies done in a health-
23 science discipline, Cizek and O'Day (1994) reduced 31 five-option items to four-
24 option items on a medical specialty examination by removing a non-functioning

1 distractor. Study findings were consistent with other research in that four-option
2 items were less difficult and equally as discriminating and reliable as five-option
3 items. A recent meta-analysis (Rodriguez, 2005) and a review of research (Haladyna
4 et al., 2002) on the optimal number of options in MCQs both concluded that in most
5 educational settings, three-option items perform best.

6 By using item-analysis data and eliminating non-functioning distractors from
7 MCQs, four or five-option items can easily be reduced to three-option items. Rogers
8 and Harley (1999) have called for additional studies that examine the impact of
9 reducing the number of distracters on item psychometric properties. To date, no
10 such studies have been conducted in nursing and only one has been conducted in
11 medicine (Cizek and O'Day, 1994). Therefore, the purpose of this study was to
12 examine and compare the psychometric properties of four-option items with the
13 exact same items rewritten as three-option items in nursing assessments.

14 **Methods**

15 Data for this study consisted of two tests administered to two cohorts of
16 students in an undergraduate public health nursing course over two subsequent
17 academic years. The first test consisted of 50 four-option items administered to 36
18 students at the end of the fall semester in 2006. The second test consisted of 70
19 three-option items administered to a subsequent cohort of 106 students at the same
20 time the next year. Using item analysis data from the four-option test administered
21 in 2006, the first author reduced the number of options to three by eliminating the
22 least frequently selected distractor. A subset of 41 items was used on both tests.
23 Items for both tests were identical except for the removed option and the course
24 teacher and course content were the same for both 2006 and 2007. Tests were

1 criterion-referenced and absolute passing scores (50%) were used. Ample time was
2 given to complete both tests (three hours) and all students completed the tests
3 within the allotted time.

4 Item analysis data from both tests was generated using Ideal 4.1, an item-
5 analysis software program (Precht et al., 2003) and then imported into Stata 9.2
6 (StataCorp, 2005) for data analysis. The Ideal program generates item difficulty and
7 discrimination statistics, distractor performance statistics, test reliability coefficients,
8 and mean test scores, along with other item and test performance indicators. Item
9 difficulty is the proportion of examinees answering the item correctly, with lower
10 values reflecting more difficult items. Items of moderate difficulty (.40 to .80) are
11 preferable (Osterlind, 1998). Item discrimination is a measure of how effectively an
12 item discriminates between high- and low-ability examinees (Haladyna, 2004).
13 Discrimination is computed using either the point-biserial correlation coefficient, the
14 correlation between the item and total test score (Osterlind, 1998) or the more
15 simple item discrimination index, the difference in the proportion of responses
16 between the upper and lower 27% of examinees (Ebel and Frisbie, 1991). Both the
17 point-biserial and the discrimination index are highly correlated and discrepancies
18 between the two statistics are extremely small or nonexistent (Beuchert and
19 Mendoza, 1979, Oosterhof, 1976). The discrimination index was used in this analysis
20 because it is simple to compute and explain (Ebel and Frisbie, 1991). Items are
21 considered discriminating if the discrimination index for the correct response is
22 positive and the same statistic for the distractors is negative. Items with higher
23 discrimination are more desired, although recommendations for acceptable indices
24 vary. The following categories were used to classify the item discrimination in this

1 study: $<.10$ poor; $.10$ to $.19$ low; $.20$ to $.29$ acceptable; $.30$ to $.39$ good; and $\geq.40$
2 excellent (Ebel and Frisbie, 1991, Trevisan et al., 1991). We evaluated distractor
3 performance using two criteria to define non-functioning distractors: those chosen
4 by fewer than 5% of examinees and those with a positive discrimination index
5 (Rodriguez, 2005).

6 The psychometric properties of both tests and the subsets of 41 items were
7 compared using descriptive statistics. We compared the mean item difficulty and
8 discrimination of the 41 items on the two tests using the paired t-test and product
9 moment correlations (Pearson's r). We also compared item difficulty and
10 discrimination using the previously defined categories with chi-square statistics. For
11 both tests, we evaluated distractor performance by assessing the following distractor
12 characteristics: the proportion of distractors with low selection frequency ($<5\%$) and
13 positive discrimination (≥ 0); the proportion of functioning distractors per test; the
14 proportion of items with 0, 1, 2, and 3 functioning distractors; and the mean number
15 of functioning distractors per item. Finally, we evaluated the effect of removing the
16 least frequently selected option by comparing individual distractor performance on
17 both tests using chi-square statistics.

18 The unit of analysis for this study was the test item and no identifying
19 participant information was used in any part of the analysis. Since the Institutional
20 Review Board of the participating institution approves only human subjects'
21 research, this study was exempted from the ethical review process. During the 2007
22 administration of the test, however, students were given a choice of having either
23 the traditional test with 4-option items or a test with three-option items. All students

1 preferred the three-option test and students were informed that results from the test
2 comparison could potentially be used for future publication.

3 **Results**

4 Table 1 shows the summary characteristics of the two tests and the subsets
5 of 41 items. In total, 142 students were tested. On the original tests, overall mean
6 test scores and the range of test scores were similar for both the 2006 and 2007
7 cohorts. The pass rate for the 2007 cohort was marginally lower than the 2006
8 cohort and the reliability was lower for both subsets of 41 items when compared
9 with the whole test. However, this would be expected with fewer test items. The 41-
10 item subset of three-option items, however, was more reliable than the subset of
11 four-option items (.71 vs. .65).

12 Mean item difficulty values indicate that overall, the 41 three-option items
13 were more difficult than the four-option items ($.70 \pm .15$ vs. $.73 \pm .14$) but the
14 difference was not statistically significant ($t=1.95$; $p=.06$) (data not shown). Figure 1
15 presents a categorical comparison of item difficulty between the two 41-item subsets
16 delivered in 2006 and 2007. Overall, the three-option test contained a greater
17 number of items of moderate difficulty and fewer easy items. However, item
18 difficulty on the two tests was similar and again this difference was not statistically
19 significant. Differences in item discrimination present a similar picture. Mean values
20 show that three-option items were marginally more discriminating than four-option
21 items ($.26 \pm .13$ vs. $.25 \pm .14$), although again the difference was not statistically
22 significant ($t=-0.76$; $p=.45$) (data not shown). When examined categorically, there
23 was no significant difference in the discrimination index of three- and four-option
24 items (Figure 2). Pearson's correlations between item difficulty and item

1 discrimination on both item subsets were $r=.82$ ($p<.001$) and $r=.51$ ($p<.001$)
2 respectively.

3 Distractor performance is highlighted in Table 2. A substantially higher
4 proportion of items on the three-option test were classified as functioning when
5 compared with the four-option test (74.4% vs. 21.1%). Similarly, 56.1% of items on
6 the three-option test had two functioning distractors compared with only 36.6%
7 (having two or more) on the four-option test. Despite having fewer distractors,
8 three-option items had more functioning distractors per item than four-option items
9 ($1.49 \pm .64$ vs. $1.32 \pm .85$).

10 Changes in distractor performance in three-and four-option items are
11 presented in Tables 3 and 4. The removal of distractors with the lowest response
12 frequency from the four-option items had little impact on the response frequencies
13 of the same distractors in the three-option items. Options that were infrequently
14 selected (<5%) on four-option items were similarly as likely to be infrequently
15 selected on three-option tests (14.6% vs. 17.1%; $p=.76$) and (17.1% vs. 22.0%;
16 $p=.58$) (Table 3). Reducing the number of distractors, however, did have a
17 substantial impact upon distractor discrimination. A greater proportion of distractors
18 were poor discriminators in four-option items when compared with three-option
19 items (34.2% vs. 14.6%; $p=.04$) and (34.2% vs. 17.1%; $p=.08$) (Table 4).

20 **Discussion**

21 To our knowledge, this is the first study in nursing and only the second in a
22 health-science discipline (Cizek and O'Day, 1994) to specifically compare item
23 characteristics of three- and four-option MCQs. Although findings from this study are
24 consistent with other research on this topic, generalizability may be limited by

1 several factors. Since our study examined only two undergraduate nursing
2 examinations, further research in other settings should be done to determine the
3 applicability of our findings. Also, because the number of examinees taking both
4 tests was uneven with substantially fewer taking the four-option test, this may have
5 affected the selection of options that were eliminated from the four-option test.
6 Additionally, because we did not control for examinee ability it is possible that
7 differences in the abilities between the two student cohorts may have accounted for
8 some of the findings of this study.

9 The results of this study, however, do add to the growing body of research
10 supporting three-option items. Overall, the differences in item difficulty and
11 discrimination between four-option items and the same items rewritten as three-
12 option items were small and statistically non-significant. Non-significant results,
13 however, are just as important as significant results. The finding that three-option
14 items perform equally as well as four-option items can have substantial impact upon
15 the practice of item-writing. While there are minimal psychometric differences in
16 item performance characteristics, clearly, three-option items are more efficient to
17 write and administer. Aamodt & McShane (1992) estimate that students can
18 complete an additional 12.4 three-option MCQs in the same time required to
19 complete 100 four-option items. More items also increases test reliability.
20 Furthermore, generating three or four plausible distracters per item is time
21 consuming and if each distractor takes five minutes to generate, writing only three-
22 option items would save over 16 hours of time on a 100-item test (Aamodt and
23 McShane, 1992). Studies have found that students (Owen and Froman, 1987) and
24 teachers (Rogers and Harley, 1999) overwhelmingly prefer items with fewer options.

1 Given the strong empirical and theoretical support for three-option items, Owen and
2 Froman (1987) advise test writers to stop struggling to invent fourth or fifth options
3 when three is almost always sufficient. Furthermore, item analysis data, if available,
4 can be used to effectively eliminate non-functioning distractors from existing MCQs
5 so that testing time can be reduced or content sampling can be increased.

6 Although reducing the number of options had minimal impact on item
7 performance, there were positive effects on distractor performance. First, the
8 proportions of distractors with low selection frequencies and poor discrimination
9 were lower for three-option items. Second, fewer three-option items had 0 or 1
10 functioning distractors. Third, even though the total number of distractors per item
11 was fewer, three-option items had a greater mean number of functioning distractors
12 per item (1.49 vs. 1.32). Finally, existing distractors became more discriminating
13 when infrequently selected distractors were removed from items. These findings
14 illustrate that there is little benefit of including non-functioning distractors in
15 multiple-choice items. In item writing it is challenging to come up with three or more
16 plausible distractors to the correct answer. Consequently, item writers often add
17 superfluous distractors that are so implausible they are selected by only a small
18 proportion of examinees. This study has demonstrated, and others have pointed out,
19 that a three-option item with two functioning distractors is clearly preferable to a
20 four-option item with two or three non-functioning distractors (Schuwirth and van
21 der Vleuten, 2004). Poorly written and clearly implausible distractors may also
22 unintentionally cue test-wise examinees to the correct answer (Owen and Froman,
23 1987). Consequently, implausible distractors can introduce construct-irrelevant
24 variance (CIV) into the assessment of student outcomes. CIV is the introduction of

1 extraneous variables, such as clueing or testwiseness, that are irrelevant to the
2 construct being measured (Downing, 2002) but which can significantly affect
3 examinee test scores (Tarrant and Ware, 2008).

4 Despite years of research supporting fewer options in multiple-choice items,
5 nurse educators have not adopted the shorter items and four-option items remain
6 the norm. Why three-option items have not been widely adopted when they are
7 easier to write and as psychometrically robust as items with more options, is unclear.
8 One possible explanation is that few nurse academics in the health professions have
9 higher education in educational methods such as item construction and are therefore
10 unaware of the literature supporting three-option items. Intuitively, three-option
11 items would appear to be easier for examinees and thus significantly inflate student
12 grades and pass rates. The effect of guessing on multiple-choice tests scores,
13 however, is often overestimated (Rodriguez, 2005). Several studies have shown that
14 reducing items from four or five to three options resulted in a test-score increase of
15 only 1–1.2% (Aamodt and McShane, 1992, Tarrant et al., 2009). If academics are
16 unaware of the research refuting these assumptions about guessing, however, they
17 are unlikely to adopt three-option items. Additionally, policies regarding test format
18 and the number of options in MCQs may not be set by the teacher but by the
19 institutional administrators, who for the same reasons identified above, are reluctant
20 to use fewer than four or five options on summative tests. Finally, the focus of
21 nursing education has traditionally been more on the “what” than the “how.”
22 Although, there is an increasing focus on innovative educational methods and
23 strategies, this is a more recent phenomenon. As universities and academics

1 increasingly look to evidence-based methods to deliver educational programs, item-
2 writing practices may also become more evidence-based.

3 **Conclusion**

4 Results from this study of teacher-generated MCQs lends further support to
5 the conclusion that in most circumstances, three-option items are the more feasible
6 and practical choice when compared with four-option items. Given the time
7 constraints of most nursing faculty today, and the increasing focus on evidence-
8 based education, teachers involved in developing MCQs for nursing assessments are
9 encouraged to use three-option items. Three-option items perform equally as well as
10 the longer four-option items, they require less time to write, and the performance of
11 the remaining distractors improves when implausible options are removed. Time
12 spent writing four and five options is not time well spent and could be used to
13 develop more items rather than more options. Writing tests with more items would
14 increase the amount of content covered in the test, improve the overall reliability
15 and validity of the test, and thus more accurately reflect student achievement.

16

1 **References**

- 2
- 3 Aamodt, M.G., McShane, T., 1992. A meta-analytic investigation of the effect of
4 various test item characteristics on test scores. *Public Personnel Management*
5 21 (2), 151-160.
- 6 Beuchert, A.K., Mendoza, J.L., 1979. A Monte Carlo comparison of ten item
7 discrimination indices. *Journal of Educational Measurement* 16 (2), 109-118.
- 8 Broome, M.E., 2009. The faculty shortage in nursing: global implications. *Nursing*
9 *Outlook* 57 (1), 1-2.
- 10 Case, S.M., Swanson, D.B., 2003. Constructing written test questions for the basic
11 and clinical sciences. National Board of Medical Examiners, Philadelphia, PA.
- 12 Cizek, G.J., O'Day, D.M., 1994. Further investigations of nonfunctioning options in
13 multiple-choice test items. *Educational and Psychological Measurement* 54
14 (4), 861-872.
- 15 Crehan, K.D., Haladyna, T.M., Brewer, B.W., 1993. Use of an inclusive option and
16 the optimal number of options for multiple-choice items. *Educational and*
17 *Psychological Measurement* 53 (1), 241-247.
- 18 Downing, S.M., 2002. Threats to the validity of locally developed multiple-choice
19 tests in medical education: construct-irrelevant variance and construct
20 underrepresentation. *Advances in Health Sciences Education* 7 (3), 235-241.
- 21 Ebel, R.L., Frisbie, D.A., 1991. *Essentials of educational measurement*. Prentice Hall,
22 Englewood Cliffs, N.J.
- 23 Gronlund, N.E., Waugh, C.K., 2008. *Assessment of student achievement*. Pearson,
24 Upper Saddle River, N.J.
- 25 Haladyna, T.M., 2004. *Developing and validating multiple-choice test items*.
26 Lawrence Erlbaum, Mahwah, NJ.
- 27 Haladyna, T.M., Downing, S.M., 1993. How many options is enough for a multiple-
28 choice test item? *Educational and Psychological Measurement* 53 (4), 999-
29 1010.
- 30 Haladyna, T.M., Downing, S.M., Rodriguez, M.C., 2002. A review of multiple-choice
31 item-writing guidelines for classroom assessment. *Applied Measurement in*
32 *Education* 15 (3), 309-334.
- 33 Masters, J.C., Hulsmeyer, B.S., Pike, M.E., Leichty, K., Miller, M.T., Verst, A.L., 2001.
34 *Assessment of multiple-choice questions in selected test banks accompanying*
35 *text books used in nursing education*. *Journal of Nursing Education* 40 (1),
36 25-32.
- 37 McCoubrie, P., 2004. Improving the fairness of multiple-choice questions: a literature
38 review. *Medical Teacher* 26 (8), 709-712.
- 39 Oosterhof, A.C., 1976. Similarity of various item discrimination indices. *Journal of*
40 *Educational Measurement* 13 (2), 145-150.
- 41 Osterlind, S.J., 1998. *Constructing test items: Multiple-choice, constructed-
42 response, performance, and other formats*. Kluwer Academic Publishers,
43 Boston.
- 44 Owen, S.V., Froman, R.D., 1987. What's wrong with three-option multiple choice
45 items? *Educational and Psychological Measurement* 47 (2), 513-522.
- 46 Pampllett, R., Farnhill, D., 1995. Effect of anxiety on performance in multiple-choice
47 examinations. *Medical Education* 29, 298-302.

- 1 Precht, D., Hazlett, C., Yip, S., Nicholls, J., 2003. International Database for
2 Enhanced Assessments and Learning (IDEAL-HK): Item analysis users' guide.
3 IDEAL-HK, Hong Kong.
- 4 Rodriguez, M.C., 2005. Three options are optimal for multiple-choice items: A meta-
5 analysis of 80 years of research. *Educational Measurement: Issues and*
6 *Practice* 24 (2), 3-13.
- 7 Rogers, W.T., Harley, D., 1999. An empirical comparison of three- and four-choice
8 items and tests: susceptibility to testwiseness and internal consistency
9 reliability. *Educational and Psychological Measurement* 59 (2), 234-247.
- 10 Schuwirth, L.W.T., van der Vleuten, C.P.M., 2004. Different written assessment
11 methods: what can be said about their strengths and weaknesses? *Medical*
12 *Education* 38 (9), 974-979.
- 13 Sidick, J.T., Barrett, G.V., Doverspike, D., 1994. Three-alternative multiple choice
14 tests: An attractive option. *Personnel Psychology* 47 (4), 829-835.
- 15 StataCorp, 2005. *Stata Statistical Software: Release 9.2*. StataCorp LP, College
16 Station, Tx.
- 17 Tarrant, M., Knierim, A., Hayes, S.K., Ware, J., 2006. The frequency of item writing
18 flaws in multiple-choice questions used in high stakes nursing assessments.
19 *Nurse Education Today* 26 (8), 662-671.
- 20 Tarrant, M., Ware, J., 2008. Impact of item-writing flaws in multiple-choice questions
21 on student achievement in high-stakes nursing assessments. *Medical*
22 *Education* 42 (2), 198-206.
- 23 Tarrant, M., Ware, J., Mohammed, A.M., 2009. An assessment of functioning and
24 non-functioning distractors in multiple-choice questions: a descriptive
25 analysis. *BMC Medical Education*.
- 26 Trevisan, M.S., Sax, G., Michael, W.B., 1991. The effects of the number of options
27 per item and student ability on test validity and reliability. *Educational and*
28 *Psychological Measurement* 51 (4), 829-837.

Table 1 Characteristics of the tests

	2006 Test		2007 Test	
	4-option items 36 Examinees		3-option items 106 Examinees	
	Original Test 50 items	Subset of 41 items	Original Test 70 items	Subset of 41 items
Mean test score % (SD)	70.3 (11.61)	--	69.7 (9.84)	--
Range of test scores (%)	38.0 – 94.0	--	41.4 – 94.3	--
Pass Rate	97.2%	--	94.4%	--
KR20 Reliability	.74	.65 ^a	.75	.71 ^a

SD = standard deviation; KR-20 = Kuder-Richardson 20

^a Spearman-Brown formula used to estimate reliability for length of original test

Table 2 Distractor performance

	2006 Test 41 Items 4-option items	2007 Test 41 Items 3-option items
No. of distractors	123	82
Distractors with: n (%)		
Frequency <5%	46 (37.4)	16 (19.5)
Discrimination ≥ 0	60 (48.8)	13 (15.9)
Both	9 (8.5%)	8 (27.6)
Functioning distractors n (%)	26 (21.1)	61 (74.4)
Functioning distractors per item n (%)		
None	6 (14.6)	3 (7.3)
One	20 (48.8)	15 (36.6)
Two	11 (26.8)	23 (56.1)
Three	4 (9.8)	--
Functioning distractors per item M (SD)	1.32 (.85)	1.49 (.64)

Table 3 Cross-tabulation of choice frequency of distractors in three-option and four-option items

Choice Frequency	<u>Distractor 1</u>		<u>Distractor 2</u>	
	4-option items	3-option items	4-option items	3-option items
<5%	6 (14.6)	7 (17.1)	7 (17.1)	9 (22.0)
≥5%	35 (85.4)	34 (82.9)	34 (82.9)	32 (78.0)
	$\chi^2(1, N=82) = .091, p = 0.76$		$\chi^2(1, N=82) = .311, p = 0.58$	

Table 4 Cross-tabulation of distractors discrimination in three-option and four-option items

Distractor Discrimination	<u>Distractor 1</u>		<u>Distractor 2</u>	
	4-option items	3-option items	4-option items	3-option items
≥0 (poor)	14 (34.2)	6 (14.6)	14 (34.2)	7 (17.1)
<0 (good)	27 (65.8)	35 (85.4)	27 (65.8)	34 (82.9)
	$\chi^2(1, N=82) = 4.23, p = 0.04$		$\chi^2(1, N=82) = 3.14, p = 0.08$	

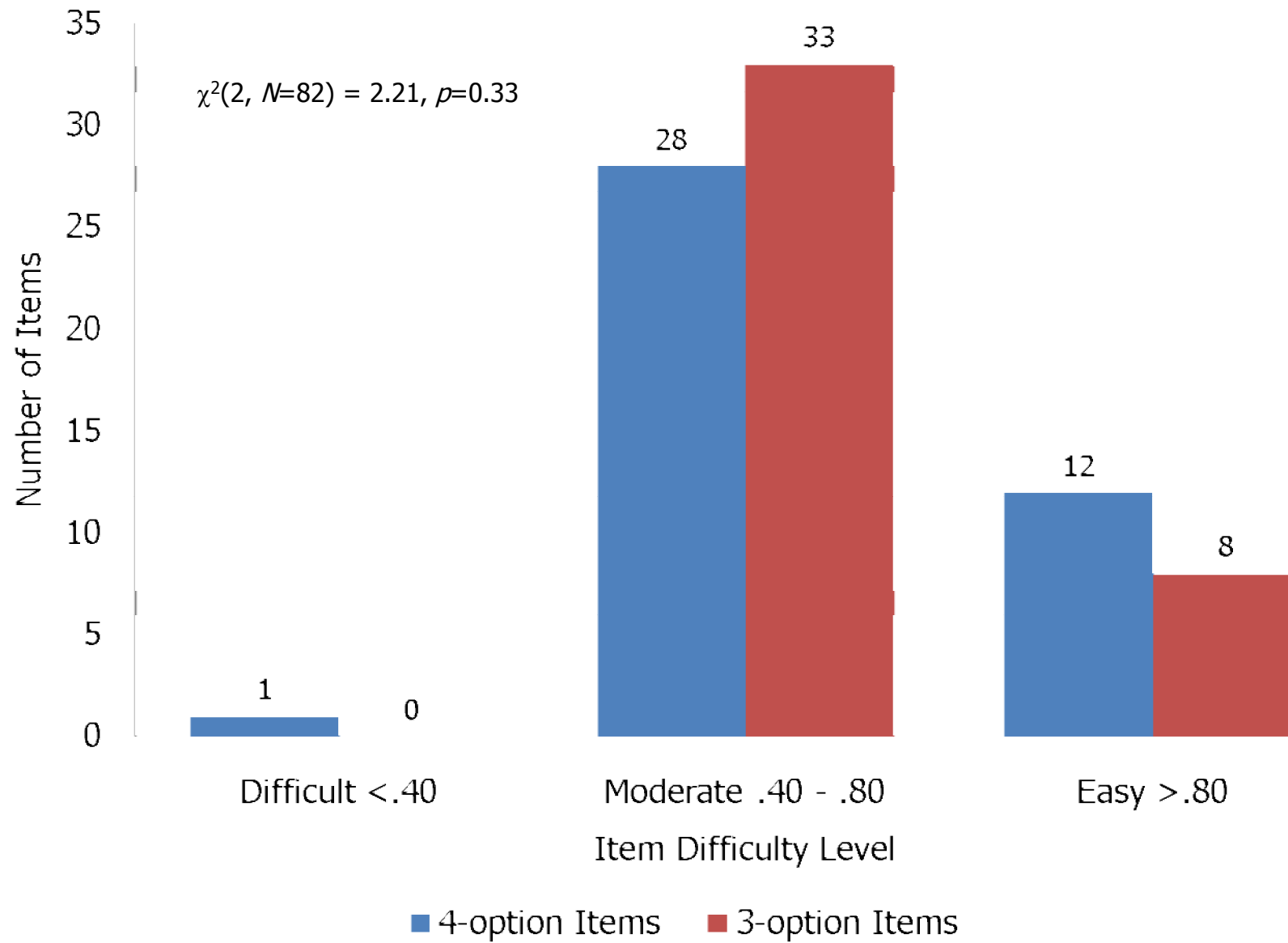


Figure 1 Comparison of item difficulty between three- and four-option items

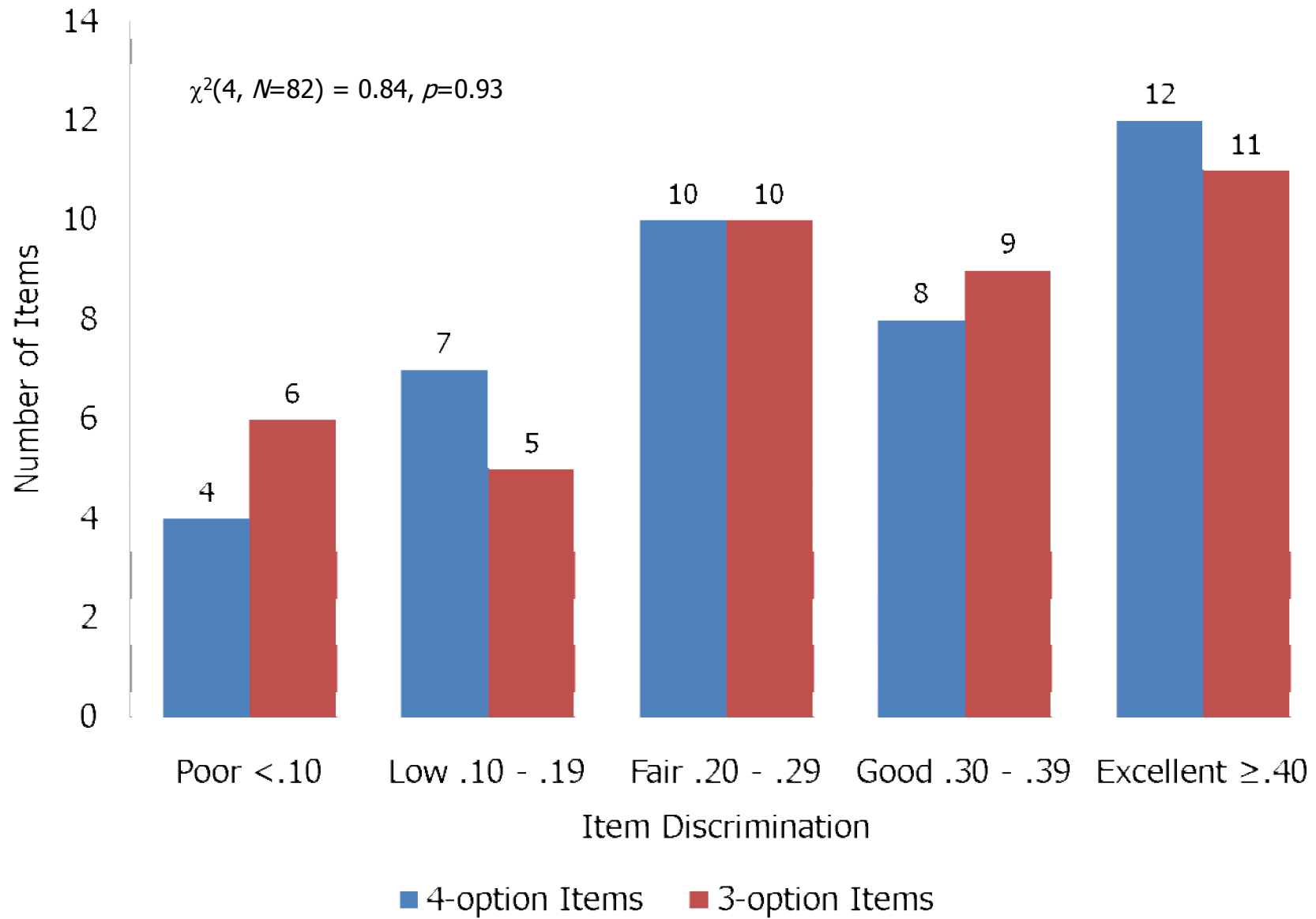


Figure 2 Comparison of item discrimination between three- and four-option items