

Reliability of the adapted mathematical knowledge for teaching number concepts and operations measures

Dun Nkhoma Kasoka¹, Arne Jakobsen², and Mercy Kazima¹

¹University of Malawi, Malawi, ²University of Stavanger, Norway

Measuring and understanding teachers' mathematical knowledge for teaching is accepted as an important element for effective mathematics teaching and learning, but mathematics teacher education do not naturally attend to and make sense of these. We describe, in part, a validation study in which 212 pre-service teachers participated in a previously adapted mathematical knowledge for teaching test with an aim of establishing validating and reliability of the adapted measures and items. We discuss the test and items' performance using classical test theory. The results show low overall internal consistence and identify eight very problematic items out of forty-six. We also give possible reasons for the low internal consistency. The study offers hope for defining standardized tools for assessing mathematical knowledge for teaching in Malawi.

Keywords: Mathematical knowledge for teaching, pre-service teacher.

INTRODUCTION

According to the Organization for Economic Cooperation and Development (OECD) education systems should equip all humans in a modern society with knowledge, skills and tools to stay competitive and engaged (OECD, 2016). In terms of Mathematics learning, education systems need to equip all learners with at least knowledge and skills about number concepts and operations. This is one of fundamental topics that is a prerequisite to understanding mathematics in later stages of schooling or appreciating mathematics in the society. To teach mathematics effectively, primary school teachers do not only need to understand the fundamental mathematical concepts, but they are also expected to have a mastery of the art and science of teaching mathematics (Reid, 2011).

The need for mathematics teachers to have a robust mathematical knowledge is undeniable. However, researchers (e.g. Chitera, 2011; Jurdak, 2009) observe that most mathematics teachers are not well prepared to teach mathematics in terms of their subject knowledge. Ball, Thames and Phelps (2008) built on Shulma's work and developed two constructs of subject knowledge namely: content knowledge (CK) and pedagogical content knowledge (PCK). For them, CK comprises teachers' comprehension of the subject matter to be taught and PCK is teachers' comprehension of how to make subject matter learnt by students. In our case, therefore we consider CK as entailing the mathematics content to be taught as prescribed by the curriculum, while PCK as the knowledge that the teacher needs in order to organise and present the mathematics content into forms that can be easily understood by learners, and explain the connections between them.

While there seems to be an obvious distinction between CK and PCK, Hill, Ball and Schilling (2008) contend that the distinction is rather not obvious at the lower level of mathematics teaching and hence they developed the term mathematical knowledge for teaching (MKT). Different schools of thought have emerged about CK, PCK and MKT since the work of Shulman, and Ball and colleagues. However, researchers generally agree that knowledge for teaching is precursor to learning and that this knowledge is affected by, among other factors, teacher development programmes (e.g. Blömeke, Suhl, & Kaiser, 2011). Furthermore, Kleickmann et al. (2013, p. 93) report that the extent to which different teacher education programmes influence knowledge for teaching or the extent to which the knowledge for teaching changes among different populations of teachers remain unclear. In Malawi context, for example, there is an evident need for standardized tools for assessing teachers' understanding of important mathematical concepts, ideas and skills, and the level to which teachers are mathematically prepared for the work of teaching. Research in this direction could assist in informing the educational system to provide relevant professional development for mathematics teachers and teacher-education programs through improved curriculum in terms of content and corresponding pedagogy.

Primary teacher Education in Malawi

The current primary school teacher education programme in Malawi referred to as Initial Primary Teacher Education (IPTE) is a two-year pre-service education programme delivered in two phases of one academic year in length each. Phase one is a full-time residential programme characterised by student teachers attending ten taught modules and about a half module of microteaching and teaching practice experience. The second phase is a one academic year of school-based teaching practice (TP) supervised by school mentors and teacher educators. At the end of the TP, student teachers return to college to sit for national exams set by the Malawi National Examination Board (MANEB). The teacher education curriculum in Malawi addresses both content and pedagogy hence has some aspects of CK and PCK (MIE, 2010). Like in most teacher education programmes, the student teachers' CK mostly develops during the first phase of the IPTE programme. During this phase, student teachers are also formally introduced to mathematics related PCK by their lecturers. Furthermore, the students have more opportunities of developing both CK and PCK through micro-teaching and school experience in demonstration primary schools based within the colleges. However, Kunje, Lewin and Stuart (2003) report that the curriculum is ineffectively delivered with material mainly packaged "as facts to be learned and assessment regimes reinforcing this recall-based orientation to curriculum" (p. 112). They also observe that the school-based phase "is peripheral rather than central to the curriculum". Consequently, the development of knowledge for teaching is greatly compromised among student teachers.

Mathematical literacy is important for developing countries like Malawi (Tsafu, 2013) if she is to alleviate poverty among her people and to strategically position

herself for development. Unfortunately, results from regional studies by the Southern and Eastern Africa Consortium for Monitoring Educational Quality (SACMEQ) show that Malawian learners continue to perform significantly below the regional average in mathematics including on items related to number concepts and operations (SACMEQ, 2011). Influenced, in part, by such comparative studies and the need to have a relevant curricula at all levels of schooling, Malawi has revised and released new curricula recently. Specifically, in September 2015 the Ministry of Education introduced a science based secondary school curricula which puts emphasis on science and mathematics. Despite these changes and new curricula direction, it is unclear how well the IPTE programme prepares teachers to teach mathematics given the persistent poor performance in primary school mathematics.

To what extent do primary school teachers understand mathematics in general and number concepts and operations in particular, and how well do the teachers exhibit mathematical literacy seem important questions to advance in the Malawi context. However, there could be a number of barriers to conducting research to answer these critical questions. One obvious barrier has been and remains the unavailability of appropriate tools with which to assess teachers' knowledge for teaching in Malawi.

Although, there are a number of diagnostic tools available to assess knowledge of mathematics, most of these were developed for specific audiences or are content focussed and do not assess ability to teach mathematics. The Learning Mathematics for Teaching (LMT) measures developed in the US focus on the work of the mathematics teacher (Hill, Ball, & Schilling, 2008). Specifically, the LMT-NCOP items measure conceptual understanding of number concepts and operations, and how a teacher can react to specific situations related to NCOP in a classroom. This renders the LMT-NCOP a potential and meaningful tool for assessing knowledge for teaching numbers concepts and operations of both practicing and pre-service teachers in Malawi. Even though the LMT-NCOP measures have been adapted for use in Africa, Ghana, it is important that we understand how individual NCOP items behave in our context. Therefore, the main purpose of this study was to investigate how the adapted items would function in Malawi.

METHOD

The LMT measures for assessing teachers' mathematical knowledge for teaching number concepts and operations were developed as part of the Learning Mathematics for Teaching (LMT) project at the University of Michigan. The measures assess multiple aspects of teacher knowledge including subject matter knowledge and pedagogical content knowledge (Hill, Ball, & Schilling, 2008). Number concepts and operations as part of the school mathematics curriculum have remained a central part of the LMT measures. We selected LMT items from existing forms for adaptation by aligning the items to the IPTE mathematics curriculum and adapting them to the Malawi context (Kasoka, Kazima, & Jakobsen, 2016).

Kane's (2013) framework for validating tests was employed. Kane remarks that when tests are adopted or adapted, they are used on populations that are generally different from the intended in terms of culture, academic experiences, linguistic abilities, social and economic experiences among others. We were therefore mindful of the need of the adapted measures to be fair and meaningful in the Malawi context while at the same time to provide scores with same meaning as intended in the source context (Kane, 2013). Kane postulates a six component validation argument for exporting tests as follows: (i) domain definition, (ii) evaluation, (iii) generalization, (iv) explanation, (v) extrapolation and (vi) utilization. A rigorous attempt was made to address these components during the adaptation process (see Kasoka, Kazima, & Jakobsen, 2016) as summarised in Table 1. An instrument consisting of forty-six adapted items was administered to pre-service teachers from one teacher education college.

Data Collection

The participants of the study were 212 first year students from one primary teacher education college of whom 61 were females and 151 were males. Primary school teachers in Malawi do not specialize to teach specific subjects. Consequently, all first year students at the college participated in the study. All the participants had completed secondary school and held a school certificate of education. The college was purposively selected because we had initially worked with the lecturers at the college during professional development activities. Approval to conduct the study was obtained from the college management and all the participants were briefed about the study and its objectives. We administered the adapted form in one sitting with the help of mathematics lecturers from the college.

Data Analysis

To evaluate the functionality of the adapted measures, we used classical test theory (CTT) and item response theory (IRT). In this paper we focus on CTT results and descriptive statistics. For all the forty-six items, Cronbach's alpha reliability coefficient and descriptive statistics were computed. Furthermore, for each item on the form, we sought four item characteristic indices namely: item difficulty index, upper-lower (U-L) groups discrimination index, point-biserial discrimination index and change in Cronbach's alpha coefficient with the item removed.

The first index we sought was item difficulty index. The item difficulty index measures the proportion of participants that responded to the item correctly. The greater the number of participants that chose the correct response or key, the less difficult the item was. While Mehrens and Lehmann (1991) suggest different acceptable difficulty levels of items based on the number of response options, Allen and Yen (1979) argue that acceptable item difficulty ranges from 0.26 to 0.75 regardless of the number of response options. Since the items have different numbers of response options, we opted to use Allen and Yen range.

Activity	Description								
1. Adopt or adapt	Examined the source instruments and measures in terms of structure and content in line with Malawi context								
2. Selecting measures/items	Aligned the content, tasks and skills in the measures/items with the IPTE mathematics curriculum and syllabus.								
3. Linguistic translation	While Malawi uses English as a language of instruction at all levels of education, we modified some US language terms to terms commonly used in Malawi.								
4. Cultural contextualization	Made changes relating to general cultural context by replacing US names of people, places and objects with familiar names.								
5. School contextualization	Modified the content of the measures to reflect IPTE curriculum, syllabus, and textbooks e.g. <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> <p>US</p> </div> <div style="text-align: center;"> <p>Malawi</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td style="padding: 5px;">9</td> <td style="padding: 5px;">180</td> </tr> <tr> <td style="padding: 5px;">5</td> <td style="padding: 5px;">20</td> </tr> <tr> <td style="padding: 5px;">2</td> <td style="padding: 5px;">4</td> </tr> <tr> <td style="padding: 5px;">2</td> <td style="padding: 5px;">2</td> </tr> </table> </div> </div>	9	180	5	20	2	4	2	2
9	180								
5	20								
2	4								
2	2								

Table 1: Summary of the adaptation process

The second index we used is the U-L discrimination index. According to Kelley (1939), and Hopkins (1998), the U-L coefficient shows the difference between the percent of test takers whose raw scores are above some set limit (i.e. in the upper group) and responded correctly to an item, and the percent of those whose raw scores are below some set limit (i.e. lower group) and responded correctly to an item. For them, the limits are set at the 73rd and 27th percentile points for upper and lower limits respectively. To determine the U-L discrimination index, we therefore divided the participants into upper 27%, middle 46% and lower 27% groups. Although, some researchers (e.g. Taiwo, 1995) have reported U-L discrimination with 25% cut-off point, Kelley argues that the U-L discrimination index stabilises when the 27% cut-off point is used. To calculate the U-L indices, we divided the difference between the number of participants in the upper and lower groups who responded correctly on an item by 57 (i.e. 27% of 212) and used Ebel (1965) to interpret the indices: > 0.4 –

Satisfactory discrimination, > 0.3 – Fair discrimination, > 0.2 – Need some revision and < 0.2 – Item must be removed or totally revised.

The last two characteristics we considered are correlation based. In this regards, we first assessed the item discrimination using point-biserial correlation. The correlation coefficients were obtained to understand the relationships between the score on an item and the total score from all items. As pointed earlier, the internal consistency reliability of the items collectively was determined by the Cronbach's alpha. We therefore sought changes in the value of the Cronbach's alpha if an item was deleted from the form to decide whether or not the deletion would increase the internal consistency. IBM SPSS Statistics was used for these analyses.

RESULTS AND DISCUSSION

If the adapted measures are to be used for inference about teachers' knowledge, then it is important that they must be both reliable and valid. Our focus in this paper is reliability as a step to determining the validity of the measures. Internal consistency reliability was evaluated by Cronbach's alpha. Analysis of all items produced Cronbach's alpha of 0.581(N=111). The value is below the generally acceptable level of 0.7 (Pallant, 2007). This suggests that the internal consistency of the items collectively is poor and the items appear not to be measuring the knowledge for teaching number concepts and operations. When we examined the changes in alpha (see Table 2) after deleting an item, seventeen items produced positive changes. The maximum change recorded was 0.028 from item 21c. Therefore, the deletion of any items does not significantly improve the internal consistency. While we understand the implication of the low value of Cronbach's alpha, Pedhazur and Schmelkin (1991) argue that reliability values ought to be evaluated by taking into consideration the specific circumstances of a study before condemning the measures due to lack of reliability. For instance, they argue that if test takers are not knowledgeable, they tend to respond rather randomly hence affecting the reliability of the measures. All the participants involved in our study were pre-service teachers in their first year of teacher education. Consequently, most of them did not have any teaching experience to draw on as they responded to the items. This is evident as Cronbach's alpha was calculated based on $N = 111$ out of 212 participants. This shows that almost half (101) of the participants did not respond to at least one item. While the analysis fails to show good reliability of the items based on the alpha value, we are of the view that this is mostly attributable to the participants' characteristics than the items.

The difficulty level of each item was compared to the range proposed by Allen and Yen (1975). Out of the forty-six items in the form, seventeen were below 0.26 and only one was above 0.75. Item 17 was the easiest with a difficulty index of 0.82. This item was attempted by all the participants and 82% (174) of the participants provided a correct response. The most difficult item was item 24 with a difficulty index of 0.01. Only 1% (2) of the participants responded correctly to this item. With an

average difficulty of 0.31, the items were generally hard for the participating pre-service teachers. Given that the adaptation related the items to the IPTE curriculum, student text books, and instructors' guides (see Kasoka, Kazima, & Jakobsen, 2016), the overall low difficulty level would be as a result of the characteristics of the participants rather than a mismatch between the items' content and curriculum.

Item discrimination analysis involved establishing the U-L discrimination index and point-biserial correlation coefficient for each item. U-L indices for twenty-two items were outside Ebel's (1965) recommendation minimum of 0.2. Only one item (item 1) had a negative U-L discrimination index (-0.12) which shows that more low performing pre-service teachers (in lower group) responded to the item correctly than their high performing counterparts (in upper group). The least discriminating items were items 6, 7, 8d, 10, 12a, 24 and 25. The difficulty indices for these items were all below 0.26 (see Table 2). This shows that very difficult items discriminated poorly and hence problematic. Moreover, removal of most of them negatively affected the alpha value. Item 5c was the most discriminating with a U-L index of 0.68.

All the items were analysed to examine the tendency of the participants responding correctly to an item and scoring highly overall. This was done through the point-biserial correlation. The results showed that some items did not correlate well with the total score. Twenty items did not show expected correlation with the total score. These had point-biserials of less than the desirable 0.3 for 'good items' (Pallant, 2007). However, Pallant advises researchers to use a minimum threshold value for point-biserial and recommends 0.15. Considering Hopkins' (1998) observations that discrimination indices above 0.1 indicate "fair discrimination" and that any positive discrimination indices above this value may be accepted because they show that there is a higher probability of a high performing examinee select a correct response than a low performing examinee. Hopkins is, nevertheless, concerned with negative point-biserial indices because they show that low performing examinees are more likely to select a correct response than high performing examinees as in items 1, 6 and 11a (see Table 2). This may suggest that the items have some underlying errors which are preventing high performing participants from responding correctly. The three items are therefore obviously problematic and must be removed.

None of the statistics we have computed so far show the performance of individual response options to understand their contribution to the overall quality of the items. All problematic items were therefore, subjected to further analysis to examine the performance of the response options in terms of their popularity among the participating pre-service teachers. We use item 6 to exemplify this analysis. Item 6 was intended to measure teachers' ability to explain the divisibility rule for 4. This item had a difficulty index below the overall mean of 0.3 and its U-L index was 0.00 which shows that it could not discriminate between the upper and lower groups. In terms of the usefulness of the response options for item 6, response option a) had positive point-biserial which suggests that high performing participants were more

likely to select this as a correct response than low performing participants. This response option was an ineffective distracter. Eighty-four participants selected this option, 25 of whom were among the 27% most knowledgeable and 16 were among the 27% less knowledgeable. The behaviour depicted by this distracter is expected of a key (correct) response. The key for this item was b) and had a negative point-biserial which is uncharacteristic of a key response. The other two response options c) and d) appeared to be effective distracters given that they both returned negative point-biserials (see Table 3). This suggests that the participants who picked these response options were mainly among those who exhibited low overall performance.

Item No.	Change in alpha	Difficulty	U-L Discrimination	Point-Biserial
1	0.027	0.72	-0.12	-0.049
5c	-0.048	0.50	0.68	0.521**
6	0.016	0.11	0.00	-0.036
7	0.024	0.24	0.05	0.238
8d	-0.012	0.13	0.07	0.209
10	-0.001	0.25	0.09	0.204
11a	0.023	0.38	0.14	-0.043
12a	0.020	0.16	0.07	0.213
17	0.003	0.82	0.16	0.324
21c	0.028	0.25	0.18	0.189
24	-0.001	0.01	0.04	0.233
25	-0.002	0.11	0.05	0.244

*Significant at the 0.05 level. **Significant at the 0.01 level

Table 2: CTT results for selected items

Option	Code	N			Proportion	Point-biserial
		Total	Upper	Lower		
a)	0	84	25	16	0.4	0.145
b) <i>Key</i>	<i>1</i>	<i>23</i>	<i>7</i>	<i>7</i>	<i>0.11</i>	<i>-0.036</i>
C)	0	36	6	13	0.17	-0.055
d)	0	60	13	22	0.28	-0.066

Table 3: Further analysis of item 6

In summary, item 6 requires further examination in terms of how it was phrased and its content in relation to Malawi curricula. It might also be important to examine the

knowledge background of the participants. The key and response option a) fail to make it a quality item.

CONCLUSION

In conclusion, we note that although the study has some limitations mainly to do with the homogenous nature of the sample used, helpful evidence was shown that the adapted measures are a promising tool for use in Malawi albeit the low internal consistency. Furthermore, the current results and performance of the adapted measures are a good basis to speculate that the participants lacked knowledge for teaching number concepts and operations, and to suggest a methodological change in investigating MKT among pre-service teachers in Malawi. The pre-service teachers seem to have responded to the items rather randomly than influenced by their knowledge. It sounds logical therefore to suggest that there is need for further evidence from a less homogeneously population using a longitudinal study. The study has found that there are low levels MKT for NCOP among pre-service teachers and this has possible implications on teacher education programmes. We have also found that not all adapted items (e.g. item 6) function well and may need to be improved or removed. While there is evidence to support the use of the measures, our analysis showed some inconsistencies between item difficulty, discrimination and point-biserials. However, we noted that most items that were problematic based on difficulty index, they also fell out of range on the other two indices.

REFERENCES

- Allen, M. J., & Yen, W. M. (1979). *Introduction to measurement theory*. California: Brooks & Cole.
- Ball, D. L., Thames, M. H. & Phelps, G. (2008). Content knowledge for teaching: What makes it special? *Journal of Teacher Education*, 59(5), 389-407.
- Blömeke, S., Suhl, U., & Kaiser, G. (2011). Teacher education effectiveness: Quality and equity of future primary teachers' mathematics and mathematics pedagogical content knowledge. *Journal of Teacher Education*, 62(2), 154-171.
- Chitera, N. (2011). Modeling school mathematics teaching in initial teacher training colleges for multilingual classrooms. *US-China Education Review*, A(7), 1005-1015.
- Ebel, R. L. (1965). *Measuring educational achievement*. Englewood Cliff, NJ: Prentice-Hall.
- Hill, H., Ball, D. L., & Schilling, S. (2008). Unpacking "pedagogical knowledge": Conceptualizing and measuring teachers' knowledge of students. *Journal for Research in Mathematics Education*, 39(4), 372-400.
- Hopkins, K. D. (1998). *Educational and psychological measurement and evaluation* (8th ed.). Boston, MA: Allyn & Bacon.

- Jurdak, M. (2009). *Towards equity in quality in mathematics education*. New York: Springer.
- Kane, M. T. (2013). Validating the interpretations and use of test score. *Journal of Educational Measurement*, 50, 1-73.
- Kasoka, D. N., Kazima, M., & Jakobsen, A. (2016). Adaptation of mathematical knowledge for teaching number concepts and operations measures for use in Malawi. In W. Mwakapenda, T. Sedumedi, & M. Makgato (Eds.), *Proceedings of the 24th Annual Conference of the Southern African Association for Research in Mathematics, Science and Technology* (pp. 100-110). Pretoria, South Africa: SAARMSTE.
- Kelly, T. L. (1939). Selection of upper and lower groups for the validation of test items. *The Journal of Education Psychology*, 17-24. Retrieved March 30, 2016 from http://www.inf.ufsc.br/~cezar/tri_material/Artigo27%25.pdf
- Kleickmann, T., Richter, D., Kunter, M., Elsner, J., Besser, M., Krauss, S., & Baumert, J. (2013). Teachers' content knowledge and pedagogical content knowledge: The role of structural differences in teacher education. *Journal of Teacher Education* 64(1), 90-106.
- Kunje, D., Lewin, K., & Stuart, J. (2003). *Primary teacher education in Malawi: Insights into practice and policy*. Sevenoaks: DFID Education Publications.
- Mehrens, W. A., & Lehmann, I. J. (1991). *Measurement and evaluation in education and psychology* (4th ed.). London, UK: Holt, Rinehart & Winston.
- OECD (2016). *Building a high-quality teaching profession: Lessons from around the world*. Retrieved April 12, 2016 from <http://www.oecd.org/edu/Directorate-for-education-and-skills-brochure.pdf>
- Pallant, J. (2007). *SPSS Survival Manual*. New York: McGraw-Hill.
- Pedhazur, E. J., & Schmelkin, L. (1991). *Measurement, design, and analysis: An integrated approach*. Hillsdale, NJ: Lawrence Erlbaum.
- Reid, C. (2011). The art and science of teaching mathematics: One teacher's creative journey. In P. Gouzouasis (Ed.). *Pedagogy in a new tonality: Teacher inquiries on creative tactics, strategies, graphics organizers, and visual journals in the K-12 classroom* (pp. 275-284). Boston: Sense.
- SACMEQ (2011). Trends in achievement levels of grade 6 learners in Malawi. *Policy Brief* (pp. 1-7). Retrieved January 10, 2016 from www.sacmeq.org.
- Taiwo, A. A. (1995). *Fundamentals of classroom testing*. New Delhi: Vikas.
- Tsafe, A. K. (2013). Mathematics literacy: An agent of poverty alleviation and national development. *Journal of Education and Practice*, 4(25), 51-54.