

Using G-theory in developing performance assessment of the physical domain of children



M.R. Nor Mashitah ^{1,2,*}, M.N. Mariani ¹

¹Department of Educational Psychology and Counseling, Faculty of Education, University of Malaya, Kuala Lumpur, Malaysia

²Department of Early Childhood and Education, Faculty of Education and Human Development, Sultan Idris Education University, Perak, Malaysia

ARTICLE INFO

Article history:

Received 3 November 2016

Received in revised form

21 January 2017

Accepted 22 January 2017

Keywords:

Physical domain

Performance-based assessment

G-Theory

ABSTRACT

This study investigates potential applications of Generalizability theory (G-theory) in the development of performance-based assessment procedure. 77 kindergarten children were assessed as participants in this study. Analysis of variance showed that nested rater variance component in person and item ($r:\pi$) component accounted for the highest percentage of the total variance, i.e. by $\sigma^2:r:\pi = 0.12208$; 33.1% and the smallest, variance of person $\sigma_p = 0.05879$; 15.9%. Secondly, through analysis in G-study, 74% of the overall variance can be explained by the design. Next, based on optimization analysis in D-study that the overall absolute Coefficient G reading ϕ (Φ) remains at 0.86 which was an acceptable value. Lastly, for reliability test from G-facets analysis, the overall physical domain reliability was recorded at 0.85 as the reliability of the 25 items was ranging from 0.84 to 0.85. This study base on Theory-G had an impact on minimizing the error of measurement and determining the appropriateness use of items in the administration of the assessment.

© 2017 The Authors. Published by IASE. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

Task completion in an actual context includes performance based assessment. The ability to complete a task demonstrates the real capability of the children. Performance assessment or authentic assessment is used to understand how children relate or apply what they have learned; the learning experience provided must be authentic and meaningful as well. When children are related to authentic learning, they are given the opportunity to link new information with the existing information while solving problems.

To clarify the relevance of authentic learning and actual abilities of children, it is appropriate to refer to Kleinert et al. (2002) who stated the objectives of this approach is to allow children to show how they use what they know to represent learning in the form of product or performance. In other words, by authentic learning, it has stimulated children to show their knowledge or true feelings of themselves. According to Wehlage et al. (1996),

authentic learning fosters knowledge construction and focuses on higher-order thinking. The aim is to enhance knowledge level and construct new knowledge. Therefore, for children, opportunities provided through a variety of activities during assessment is intended to observe the level of existing knowledge and new experiences as well as new knowledge when aid granted during activity.

In particular, researchers have looked at the issue of variability in assessment tasks and rater judgments as sources of measurement error in performance testing (Shohamy, 1983; 1984; Pollitt and Hutchinson, 1987; Barnwell, 1989; McNamara and Adams, 1991). This study aimed to investigate potential applications of Generalizability theory (G-theory) in the development of such a performance-based assessment procedure (Cardinet et al., 2009).

2. Literature review

Previous studies of performance-based assessment using instrument is to support children through evidence and proof obtained as well as to identify the strength and weaknesses (Gardner, 1993). When referring to the first purpose of using performance-based assessment on children, it is assure that this assessment is a good tool to assess

* Corresponding Author.

Email Address: normashitah1604@gmail.com (M. R. Nor Mashitah)

<https://doi.org/10.21833/ijaas.2017.03.024>

2313-626X/© 2017 The Authors. Published by IASE.

This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

the progress of development of children as performance-based assessment is designed to measure the actual performance of children or their assignments or activities related to learning. Using observations towards performance is closely related to or directly linked to the development of the achievement rate in children (Harrington et al., 1997). Secondly, performance-based assessment is integrated to teaching. Performances in activities are the natural learning outcome that is parallel to the curriculum and teaching which cannot be separated.

Hills (1993) elaborated while using performance-based assessment, teachers have to know the suitability of its design, the relationship as a mean of testing, interpret the results of assessment to understand the progress of the children, plan further lessons and deliver results to parents and administrators.

Dependability refers to accuracy in generalizing scores obtained from respondents in a test to average score obtained by students in various situations (Shavelson and Webb, 1991). In this research context, dependability is the index obtained in a test analysis based on different individual and item.

This research is about G-Study to identify various variance resources, which might be in an assessment by estimating the variance component that is contributed by each of it. It is carried out to evaluate the measurement of dependability that is done to the variance which can be considered in the future measurement. This research is focus on D-Study to put forward reliability coefficient as generalizability coefficient covering variance towards error resources. D study is also able to differentiate between the relative decision and the ultimate decision. By using the information which has been collected through G study, D study can design a better and more suitable measurement application for a measurement and assessment suggested (Shavelson and Web, 1991).

G-Theory produced a more integrated approach to assess reliability which has been carried out whether for the purpose of making relative decision (norm-reference test) or actual decision (criteria reference test). Relative decision is based on individual's place in a group compared to actual decision. Actual decision is based on actual score without any comparison with other individuals score in the group (Ary et al., 1996). G-Theory does not make assumption regarding comparison of error resources but estimate simultaneously the variety of error resources including interaction between those errors (Thompson and Crowley, 1994).

3. This study

This study emphasized performance-based assessment towards physical domain in a fun learning environment which involves learning activities with teachers in the playschool. To assess is to collect information. Observation method is used to collect information and evidences. Observation

means children's behavior is under scrutiny. This approach can be used without the consciousness of the children that they are being observed. This study used the role of the Rater, which is the teachers themselves observe the children. Every child will be evaluated by raters.

The broad research questions that guided this study were:

- Contribution of facet towards variance resource according to the Generalizability Theory,
- Score coefficient value of children's performance according to the G-Study,
- Best optimization value towards facet in order to increase the value of coefficient G by using D-study, and
- Reliability score for each item in the performance-based assessment in G-facets analysis.

4. Methodology

Research design of this study is in the form of survey and analyzed data in quantitative method. This study is a descriptive research in order to collect feedback from respondents as well as to survey error resources in measurement. Research design is as in Table 1 and Fig. 1.

Dependability of test score will be used Two facet (r:pi) partially Nested Random Design. Data will be analyzed using EduG software in order to get result for G study and D study. Design model of two facet (r:pi) partially Nested Random Design is shown in Fig. 2.

Fig. 1 shows Venn graph for the research design of this study, that is Two facet (r:pi) partially Nested Random Design.

Fig. 2 is component of variance resources. The p circle represent children (person) being evaluated in the domain of physical development. However, the item circle, i represent item of the physical development domain which is tested on children. This item is made up of item which requires children to show response of their ability in doing it. The person circle, p intersects with the item circle, i produced interaction between people and item, that is the pi interaction. pi interaction shows how children give response towards item which is being tested in the assessment. Following that, in the intersect part between the p and i circle, nested circle is the rater, r. This shows that different rater will evaluates the children's performance, yet item being tested is the same.

In this study, person (p) is the object of measurement. Two facet involved is the nested rater (r) and item (i) in children as well as item p/ri. Observation design is r:pi. All measurement object children and facet are infinite random because the population of inspector and student are infinite, also having variability with universe set.

Table 2 shows variance resources in this study. Based on the research design, two facet (r:pi) partially Nested Random Design, it has produced 4

variance resources, that is person (p), item (i), rater nested in children and item (r:pi) as well as

interaction between person item (pi) and residual (e).

Table 1: Research design

Research Design	Data Collection Method	Respondent	Data Resource
Quantitative research	NOaMA Instrument	18 raters who give score to the performance of 77 children	Performance score

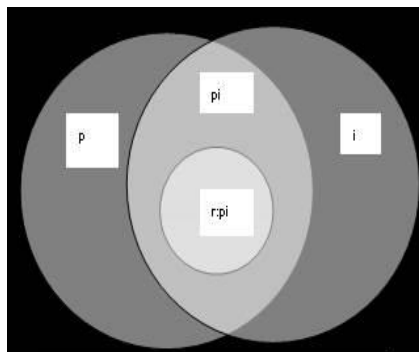


Fig. 1: Variance resources

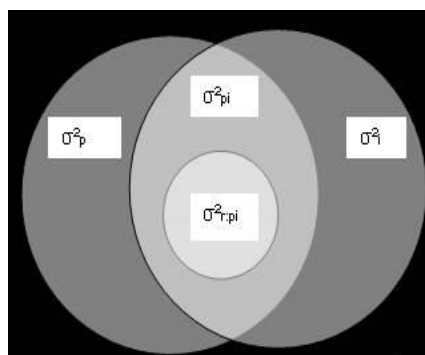


Fig. 2: Variance components

Table 2: Variance resources for two facet partially nested (r:pi) design

Variability Resources	Variance Note
Person (p)	$\sigma^2(p)$
Item (i)	$\sigma^2(i)$
Rater nested in pi	$\sigma^2(r:pi)$
Interaction pi, e	σ^2pi, e

4.1. Sample

Based on the 77 children who were enrolled in the registered playschools, sample selection is based on stratified random method. A total of 77 children as research sample represented the population being studied. Performance-based assessment were carried out among 77 children and were given scores by two different raters (teacher) from 9 playschools,

that is all together 18 raters and identified as rater 1 and 2.

4.2. Research instrument

NOaMA assessment is a learning assessment approach and children development in this study have been re-designed in year 2013 to include scoring procedures in Likert scale (five points). This instrument was re-designed to comply with the assessment concept of National Early Childhood Care and Education Policy.

This instrument reflects the overall skill at the age group which requires children to relate with the learning and development domain. The physical domain contains performance's item that require children to perform a task. Activities prepared will translate such performance items. The physical domain contains a number of 25 items.

The data has been analyzed by using EduG is able to estimate every variance component and determine the dependability score in a test. There are various designs that can be analyzed using EduG according to the desired facet. In the research carried out, researcher used the Two-Facet Partially Nested Design. Analysis outcome of EduG have produced two types of research, that is G-study (Generalizability studies) and D study (Decision Studies). G-study is able to identify variance resources and variance magnitude, while D-study is able to determine coefficient G as well as the design suitable to the number of item in a particular test.

5. Results

5.1. Contribution of facet towards variance resource according to the G-theory

From analysis, variance component which contributed to the dependability of test is shown in [Table 3](#).

Table 3: Variance component of performance based assessment on physical development domain

Facet	Df	Square Total	Square Mean	Variance Component	% Variance
Person, p	76	251.13766	3.30444	0.05879	15.9
r:pi	1925	235.00000	0.12208	0.12208	33.1
Item, i	24	254.23377	10.59307	0.06642	18.0
pi	1824	665.96623	0.36511	0.12152	32.9
Total	3849	1406.33766			100%
Koeffisien G_absolute		= 0.86			

[Table 3](#) shows the variant component of each facet which contributed to the difference of children's score evaluated by the rater in the assessment of physical development domain performance. The variant analysis shows that nested

rater variant component in person and item, (r:pi) shows the highest value of variant component which is 33.1%, followed by the variant component of interaction among person and item (pi) 32.9%. Then, variant component of items (i) is 18% and the

smallest variant component is the person component (p) which is 15.9%.

Through the analysis, it was found that nested rater variant component in person and item, ($r:pi$) shows the highest value of variant component, ($\sigma^2_{r:pi} = 0.12208$; 33.1% from the overall total variant component). This shows differences between raters in giving scores to the children. This is because raters had understood that the scoring based on rubrics and all raters have dissimilar consistency while giving scores for the evaluation of physical development domain. This shows children scores dependability in the test is much influenced by raters.

Based on the analysis, variant component that shows huge reading is the interaction among person and item (pi) which is ($\sigma_{pi} = 0.12152$; 32.9% from the overall total of variant component). This shows that there is significant difference among the children in giving response on the tested items.

Next, the variant component from item shows average reading ($\sigma^2_i = 0.06642$; 18.0% from the overall total of variant component. The average percentage for item component shows that tested items in the evaluation is different in terms of difficulties. The different in level of difficulties influence the performance showed by the children averagely.

The smallest variant component is the variant component of person (p) which indicates the lowest value of component ($\sigma^2_p = 0.05879$; 15.9% from the overall total variant component). The respondents

show that there are little differences in children abilities; it means that the children who participate as respondents have average abilities.

5.2. Score coefficient value

Relative coefficient G (0.89) and absolute coefficient G (0.86) in Table 4 showed value beyond the accepted conventional value, 0.8. Research design is good to analyze children’s dependability score because coefficient G value beyond conventional value. Absolute coefficient G is considered as this research aimed to evaluate children’s dependability score individually based on the contribution of variant component in different raters. Through analysis, 74% of the findings from children’s score are attributable to the universe score. This means that 74% of the overall research can be explained. However, only 26% of finding score is attributable to random impacts which are not identifiable. This design produced reliability measurement or dependable measurement. It is also can be interpreted as 74% of the factors that contributed to the children’s variance score can be explained, while 26% contributing factors found from error resources which are not identifiable. Findings also show that standard error related to children’s decision score is small while absolute standard error is 0.09979. Standard error shows value that is smaller than the estimated standard deviation 0.24246 for true score dispersion.

Table 4: G-study (random)

Source of Variance	Differentiation Variance	Source of Variance	Relative Error Variance	% Relative	Absolute Error Variance	% Absolute
K	0.05879
	P :KI	0.00244	33.4	0.00244	24.5
	I	0.00266	26.7
	KI	0.00486	66.6	0.00486	48.8
Sum of variances	0.05879		0.01541	100%	0.03133	100%
Standard deviation	0.24246		Relative SE: 0.08545		Absolute SE: 0.09979	

Coef_G relative: 0.89; Coef_G absolute: 0.86

5.3. Best optimization value towards facet

In D-study, the relative coefficient G ($\hat{E}p^2$) displays different level of relative error variance. In D-study, absolute coefficient G phi (Φ) shows degree of difference in absolute error variance. Table 5 shows the difference of reliability value or coefficient G when number of children and rater increase or modified.

In this research, the absolute coefficient G phi (Φ) will only be taken into account because this

research is to examine error variance towards children’s score evaluated by two different raters in the performance based assessment in physical development domain in playschools. This research also compares score given by two raters of different playschools.

Based on Table 5, it is found that number of children that are suitable to be evaluated in the assessment is 77 by taking into account the number of raters remained at 2 person.

Table 5: The variances component of D Study based on the modification number of person and raters

	G-Study	Opt.1	Opt. 2	Opt. 3	Opt. 4	Opt.5
Amount of Children (K)	77	100	120	140	160	180
Amount of Rater (P:KI)	2	4	4	3	2	2
Coef_G relative ($\hat{E}p^2$)	0.89	0.91	0.91	0.90	0.89	0.89
Coef_G absolute (Φ)	0.86	0.87	0.87	0.86	0.85	0.85

With reserves of 77, the Coef_G absolute phi (Φ) remained at 0.86, which is a high value and it is accepted. Coef_G absolute value of phi (Φ) exceeds

the accepted conventional value 0.8. The decision to choose the number of children that are suitable for assessment is based on the consideration of factors

such as time, cost, logistics and others. This means that if the number of children which were maintained at 77 children; it is accepted and sufficient to deal with restrictions on time, cost logistics and others.

Therefore, for this study, researcher suggested number of children to be 77 children and 2 raters in the performance based assessment in the physical development domain is maintained for the value of coef_ G absolute phi (Φ) or high reliability parallel with these findings.

5.4. Reliability score

G-Facets Analysis is carried out to identify the contribution of each item to be tested in the performance-based assessment of the value of the coefficient G or reliability. This analysis estimates the coefficient G adequate for each item tested.

Table 6 shows the relative and absolute value of the coefficient G for each item tested. Generally, all items are functioning well because the value of coefficient G is greater than 0.8. Among these items, item 11 is seen as an item that contributed the largest error in the scoring to children. Item 11 can be said to represent an item which has a high difficulty level or testing children in achieving high level of performance. However, a conclusion can be made that these items are consistent as performance assessment items used to evaluate children. So, these items should be retained and can be used as a test set for children performance-based assessment bank item in physical development domain.

Table 6: G-Facets analysis towards item (i)

Level	Coef_G Relative	Coef_G Absolute
1	0.88510	0.85134
2	0.88392	0.85602
3	0.88836	0.85477
4	0.88310	0.84584
5	0.88789	0.85168
6	0.88863	0.85504
7	0.89072	0.85795
8	0.89025	0.85677
9	0.88895	0.85544
10	0.89072	0.85795
11	0.89211	0.85828
12	0.88167	0.84432
13	0.87916	0.84027
14	0.88076	0.84375
15	0.88258	0.84641
16	0.88031	0.84235
17	0.88258	0.84601
18	0.88741	0.84849
19	0.87934	0.84122
20	0.88400	0.84625
21	0.88927	0.85533
22	0.88114	0.84339
23	0.88540	0.85146
24	0.88073	0.84224
25	0.88744	0.85145

6. Discussion

Model design of this study is Two facet (r:pi) partially Nested Random Design, it has 4 variance resources, that is person (p), nested rater in person and item (r:pi), item (i), and interaction between

person-item (ki) and residual (e). The analysis of nested rater variance in person and item (r:pi) shows variant component contributed the highest percentage from overall variance total, that is 33.1%, followed by variant component among person and item (pi) which is 32.9%. Subsequently, the variant component of item (i) which is 18.0% and the smallest variant component is the person (p) which is 15.9%. From this analysis, it is found that nested rater variant component in children and item (r:pi) shows highest value of variant component, ($\sigma^2_{r:pi} = 0.12208$; 33.1% from the overall total of variant component). This shows differences between raters in giving scores to the children. This is because raters had understood that the scoring based on rubrics and all raters have dissimilar consistency while giving scores for the evaluation of physical development domain. This shows that the dependability in children scores in the test is influenced by the raters. Next, variant component that shows huge reading is the interaction among person and item (pi) which is ($\sigma_{pi} = 0.12152$; 32.9% from the overall total of variant component). This shows that there are significant differences among the children in responding towards the tested items. Besides, the item variant component (i) shows average reading ($\sigma^2_i = 0.06642$; 18.0% from the overall total of variant component). The average percentage for item component means that the item tested in the evaluation is different in terms of difficulties. The different of difficulty in all items influence the performance showed by the children averagely. The smallest variant is the person component (p) which is the lowest value of variant component, ($\sigma^2_p = 0.05879$; 15.9% from the overall total variant component). Through these respondents, it shows that there are little dissimilarities in children abilities, this shows that children who participate in the study have average abilities.

Analysis based on Generalizability Theory by using EduG software is able to show variant component of every facet that contributed to the difference of children's score. G coefficient worth 0.86 is interpreted as 74% of the factors contributed to the children's score variance, while 26% of the contributing factors found from the error resources are not identifiable. Nested rater variant component in person and item contributed the highest percentage of the total variances, that is ($\sigma^2_{r:pi} = 0.12208$; 33.1%). This shows differences between raters in giving scores to the children. This is because raters had understood that the scoring based on rubrics and all raters have dissimilar consistency while giving scores for the evaluation of physical development domain. This shows that dependability of children scores in the test is influenced by raters. The variant component among person and item, (pi) shows the highest value of variant component, ($\sigma^2_{pi} = 0.12152$; 32.9% from the overall total variant component). This shows significant differences between children in responding to the tested items

Based on optimization analysis, it is suggested to remain the 77 children with absolute Coef_G phi (Φ) which maintained at 0.86, that is a high value and accepted. This absolute Coef_G phi (Φ) value is beyond the accepted conventional value; that is 0.8. The decision to choose the number of children which is the most suitable for the assessment is made by consideration of factors such as time, cost, logistics and other. This means that if the number of children which were assessed remains at 77 children, it is accepted and sufficient to cope with the constraint of time, cost, logistics and others. Therefore, in this study, the researcher suggests the number of children to be remained at 77 children and rater 2 persons in the performance based assessment in the language development domain in order to obtained high absolute Coef_G phi (Φ) value or high reliability value which parallel with the research findings.

Based on G-facets analysis, a conclusion can be made that these items are consistent as performance assessment items used to evaluate children. So, these items should be retained and can be used as a test set for children performance-based assessment bank item in physical development domain.

7. Conclusion

These findings lead to a number of implications in the construction of early learning standard instrument in early childhood development. Practically, it is difficult to build a truly fair and equitable item for all students who have different abilities. G-study and D-study according G-Theory that have been carried out gives impacts in efforts to minimize the measurement error besides making wise decisions in number of item that is the most suitable to be administered in this assessment in the future. Items that functioned well can be included into the assessment item bank of physical development domain. Analysis of children's abilities by using rater assessment based on G- Theory gives a different dimension. By Generalizability Theory analysis, the contribution of each error in the measurement can be identified separately, making analysis of Generalizability Theory a more precise and detailed. In assessing the ability of children, the set of assessment need to be implemented carefully after taking into account various factors that contribute to the result scores in the assessment. The constructor of the assessment item is responsible to ensure the constructed items show continuing consistency if tested on other children and validated according to the needs and purpose the instrument is constructed. The existence of internal and external factors that may contribute to the variance of score should be controlled so that the reliability of findings and validity of the instrument can be improved. G-Theory may explain the error components which become the contributing factor to

the difference of assessment score. Analysis of physical development domain items based on the above theories has clarified directly or indirectly on the quality of the test and the improvements that need to be implemented to ensure that the instrument is truly able to meet the objectives of the measure.

Acknowledgment

Special appreciation to the Institute of Research Management and Monitoring (IPPP), University of Malaya, Kuala Lumpur in allowing and giving postgraduate grant for us to conduct this study. Similarly, the cooperation of respondents and teachers from all kindergartens involved.

References

- Ary D, Jacobs LC, and Razavieh A (1996). Introduction to research in education. Harcourt Brace College Publishers, Florida, USA.
- Barnwell D (1989). 'Naive' native speakers and judgements of oral proficiency in Spanish. *Language Testing*, 6(2): 152-163.
- Cardinet J, Johnson S, and Pini G (2009). Applying generalizability theory using EduG: Quantitative methodology series. Routledge, New York, USA.
- Gardner H (1993). Multiple intelligences: The theory in practice. Basic Book, New York, USA.
- Harrington HL, Meisels SJ, McMahon P, Dichtelmiller ML, and Jablon JR (1997). Observing, documenting, and assessing learning: The work sampling system handbook for teacher educators. Rebus, Michigan, USA.
- Hills TW (1993). Assessment in context: Teachers and children at work. *Young Children*, 48(5): 20-28.
- Kleinert H, Greene P, and Harte M (2002). Creating and using meaningful alternative assessments. *Teaching Exceptional Children*, 34(4): 40-47.
- McNamara TF and Adams RJ (1991). Exploring rater behaviour with Rasch techniques. In the Annual Language Testing Research Colloquium, Princeton, USA. Available online at: <http://files.eric.ed.gov/fulltext/ED345498.pdf>
- Pollitt A and Hutchinson C (1987). Calibrating graded assessments: Rasch partial credit analysis of performance in writing. *Language Testing*, 4(1): 72-92.
- Shavelson R and Webb N (1991). Generalizability theory: A primer. SAGE, California, USA.
- Shohamy E (1983). The stability of oral proficiency assessment on the oral interview testing procedures. *Language Learning*, 33(4): 527-540.
- Shohamy E (1984). Does the testing method make a difference? The case of reading comprehension. *Language Testing*, 1(2): 147-170.
- Thompson B and Crowley S (1994). When classical measurement theory is insufficient and generalizability theory is essential. In the Annual Meeting of the Western Psychological Association, Kailaukona, Hawaii. Available online at: <http://files.eric.ed.gov/fulltext/ED377218.pdf>
- Wehlage GG, Newmann FM and Secada WG (1996). Standards for authentic achievement and pedagogy. In: Fred MN (Ed.), *Authentic Achievement: Restructuring Schools for Intellectual Quality*: 21-48. Jossey-Bass, San Francisco, USA.