



PROVING CONTENT VALIDITY OF SELF-REGULATED LEARNING SCALE (THE COMPARISON OF AIKEN INDEX AND EXPANDED GREGORY INDEX)

Heri Retnawati

Faculty of Mathematics and Natural Science, Universitas Negeri Yogyakarta, Jl. Colombo No. 1,
Depok, Sleman, 55281, Yogyakarta, Indonesia

Abstract

This study aims to prove the content validity of the self-regulated learning (SRL) scale using Likert model and multiple-choice model with content validity coefficient based on expert assessments with Aiken formula and expanded Gregory formula. In this study, the SRL scale with Likert and multiple-choice model are developed using the same outline/format. There are three experts who assess the items' relevancy using indicators of both scale formats. The results of the expert assessments are then used to calculate the coefficient of the validity with Aiken formula and the expanded Gregory formula. The results showed that the content validity coefficient based on expert assessment on Likert and multiple-choice format with Aiken formula is at 0.9 for each, while using the Aiken formula and expanded Gregory formula, the coefficient is 0.6 for Likert, and 0.8 for multiple-choice.

Keywords: *validity coefficient, Aiken formula, expanded Gregory formula, SRL scale*

How to cite item:

Retnawati, H. (2016). Proving content validity of self-regulated learning scale (The comparison of Aiken index and expanded Gregory index). *Research and Evaluation in Education*, 2(2), 155-164. doi:<http://dx.doi.org/10.21831/reid.v2i2.11029>

***Corresponding Author.**

e-mail: retnawati.heriuny1@gmail.com

Introduction

Successful learning is driven by many factors. One of them is self-regulated learning which is related to independent learning like what college students do. College students are students who study at college and categorized as adults. They are categorized so because of their age, and because of the demands of independent learning in college. For college students, managing themselves to learn is a factor that supports their success in learning at college. The ability to manage themselves in the study is often referred to as self-regulated learning.

Various opinions related to self-regulated learning are presented by experts. Pintrich in Schunk (2005) states that self-regulated learning, or self-regulation, is an active, constructive process whereby learners set goals to review their learning and then attempt to monitor, regulate, and control the reviews of their cognition, motivation, and behavior, guided and constrained by reviewing their goals and the contextual features in the environment. Zimmerman (1989; 1990) writes that self-regulated learning strategies are actions and processes directed at acquiring information or skills that involve agency, purpose, and instrumentality perceptions by learners. It means that a person carries out self-regulated learning in the learning process if he/she controls his/her behavior and cognition systematically by noting the rules made by him/herself, controlling the learning process, integrating the knowledge, practicing to remember the information obtained, and developing also maintaining positive values from his/her learning.

Social cognitive theory of Bandura (Kivinen, 2013) presents the theoretical basis of the self-regulated learning development model in an individual, in which contextual factors and interactional behavior give advantages to students to organize their study and to set themselves at the same time. Social cognitive perspective differs from the standpoint of personal interaction, behavior and his/her environment that is often referred to a triadic process from Bandura, as seen in Figure 1.

Self-regulation is a cyclical process, because the input of the initial capabilities is used to make decisions to repeat the efforts that have been made. The effort of those repetitions is necessary because people, environment, and behavior always change during a learning process that is always observed and monitored.

Discussion on self-regulated learning includes three phases: forethought and planning phase, performance monitoring phase, and reflection on performance phase (Zumbrunn, Tadlock, and Danielle, 2011). In the forethought and planning phase, there are two related things: task analysis, and self-confidence and motivation.

The determination or performance monitoring phase includes self-control and specific observations. Self-reflection phase consists of self-development and self-reaction. These three phases are interrelated and they affect each other, so that they make up a cycle. The cycle is described in Figure 2.

The forethought phases can be classified into two points, namely the task analysis (covering self-regulation purpose and strategic

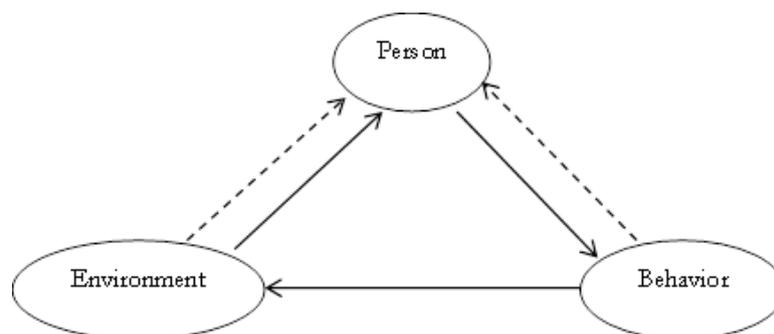


Figure 1. Self-regulation Triadic Form from Zimmerman (Kivinen, 2013)

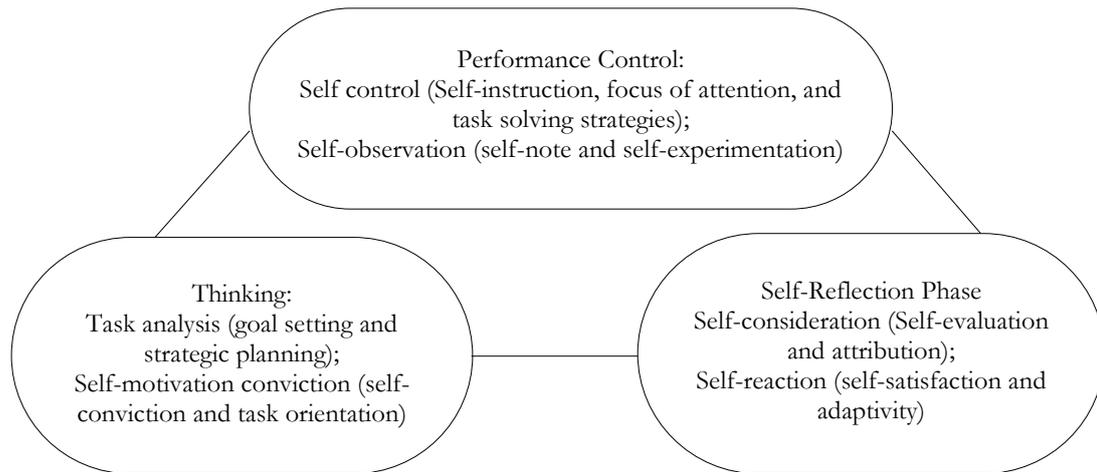


Figure 2. SRL Phase (Zumbrunn, Tadlock, & Danielle, 2011)

planning) and self-motivation (self-confident and task-oriented). The forethought phases can be classified into two points: the task analysis (covering self-regulation purpose and strategic planning) and self-motivation (self-confident and task-oriented). The performance monitoring phase includes self-control (covering self-instruction, focus of attention, task-solving strategies). Self-reflection consists of self-consideration (self-evaluation and attribution) and self-reaction (self-satisfaction and adapt-ability). To determine the SRL scale, Wolkers, Pintrich, and Karabenick (2003) write that developing items is essential to do first to measure the cognition arrangements, followed by regulation, motivation, and behavior. These three things need to be measured in the academic context.

Some researches show that the SRL is strongly associated with motivation (Vrieling, Bastiaens, and Stijnend, 2012). The SRL can be reinforced by educators in the learning process by preparing tasks that support the improvement of SRL (Zumbrunn, Tadlock, and Danielle, 2011). SRL is recognized as an important predictor of student academic motivation and achievement (Zumbrunn, Tadlock, and Danielle, 2011). Related to the importance of the SRL contribution to the success of college education, the SRL of students need to be measured. The result of the measurement can be interpreted to be followed up as an effort to maintain or improve the SRL. Therefore, the valid SRL measurement instrument is needed to develop based on the instrument development steps

each of which can be accounted. SRL measuring instrument development steps consist of several stages, including constructing a format based on the proper construction theory, preparing items, proving the content validity, trying out instruments on the correlating respondents, estimating the reliability, understanding the characteristics of the items, and reassembling the decent items into the instrument that is ready for use.

One of the instruments that can be used to measure the SRL is a questionnaire. The questions in the questionnaire have various forms, including dichotomy questions, multiple-choice questions, rank ordering, rating scale, and also open-ended questions (Cohen, Manion, and Morrison, 2011). Each of these forms has its own characteristic. Dichotomy questions in the questionnaire contain only two answer choices. These questions are used if the researcher wants to ask the respondents questions related to variable containing two answers only, for example, gender (male or female, yes or no, true or false). The multiple-choice questionnaire questions are basically like multiple choice questions in description question. In the multiple-choice, respondents are usually allowed to choose one answer only. The scoring can be done by only right or wrong option, or stratified alternatives. If scoring is done differently, an ideal condition needs to be thought by a questionnaire maker. The questionnaire model that is most often used in Indonesia is rating scale or better known as Likert model.

From the interviews with practitioners in the educational fields, some practitioners question the validity of the questionnaire with Likert model in multiple choice models. Each practitioner has its own reasonable arguments. The Likert questionnaire model is easy to make and easy to read by the respondents, but the data obtained contain desirability bias. The multiple-choice questionnaire model is difficult to make and the respondents need time to read, but more valid data can be obtained from it. Related to this problem, this study describes the proof of the content validity from the questionnaire in Likert and multiple-choice model with stratified scoring.

There are various opinions on the validity of the instruments used for the measurement, both in education and psychology. According to American Educational Research Association (AERA), American Psychological Association (APA), and the National Council on Measurement in Education (NCME) in the Standards for Educational and Psychological Testing, validity refers to the degree of facts and theories that support the interpretation of instrument scoring, and the most important consideration in the development of an instrument (1999). Other experts point out that the validity of a measuring instrument is to what extent the measuring instrument able to measure what should be measured (Nunnally, 1978; Allen and Yen, 1979, p.97; Kerlinger, 1986).

Meanwhile, Linn and Gronlund (1995) explain that validity refers to the adequacy and interpretation appropriateness made of assessment, related to a specific use. This opinion is reinforced by Messick (1989) who writes that validity is an integrated evaluative policy concerning what extent of empirical facts and theoretical reasons support the adequacy and appropriateness of inferences and actions based on test scores or scores of an instrument. Based on those opinions, it can be concluded that validity will show supports to empirical facts and theoretical reasons for the interpretation of test scores or score of an instrument, and it is associated with the measurement precision.

There are three types of validity, namely: (1) criterion validity (criterion-related

validity), (2) content validity, and (3) construct validity (Nunnally, 1978; Allen and Yen, 1979; Fernandes, 1984; Woolfolk and McCane, 1984; Kerlinger, 1986; and Lawrence, 1994). This can be known through validity existence facts. Sources of validity facts can be grouped into content validity, response process, internal structure, relations with other variables, and the consequences of the implementation of data collection (AERA, APA, and NCME, 1999; Cizek, Rosenberg, and Koons, 2008). The validity existence of an instrument can be identified through content analysis and empirical analysis from instrument score of item response data (Lissitz and Samuelsen, 2007).

The criteria of validity are divided into two, namely the predictive validity and concurrent validity. Fernandes (1984) writes that the validity based on criteria is intended to answer the question about the extent to which an instrument can predict the participants' ability in the future (predictive validity) or estimate the ability of other measuring devices in almost the same deadline (concurrent validity). A similar opinion is also expressed by Lawrence (1994) who says that the instrument is said to have predictive validity if it is able to predict capability in the future. In the analysis of the predictive validity, performances to be predicted are called criteria. The size of the estimated predictive validity value of an instrument is described by the correlation coefficient between the predictors of those criteria.

The content validity of an instrument is the extent to which the items in the instrument represents the components in the overall area of the contents of the object to be measured and the extent to which the items reflects behavioral traits that will be measured (Nunnally, 1978; Fernandes, 1984). Meanwhile, Lawrence (1994) explains that content validity is the questionable representation of special abilities that must be measured. Based on this opinion, it can be concluded that the content validity is related to the rational analysis of the domain to be measured to determine the representation of the instrument with the ability to be measured.

Construct validity is the validity which shows to what extent the instruments reveal

the ability or particular theoretical construct to be measured (Nunnally, 1978; Fernandes, 1984). A construct validation procedure starts from an identification and restriction regarding the variables to be measured and is expressed in terms of a logical construct based on the theory of those variables. From this theory, a practical consequence of the results of measurements on certain conditions is drawn, and this consequence will be tested. If the result is in line with expectations, the instrument is considered to have good construct validity.

Validity is an indispensable term required in an instrument's development. According to Sireci supported by Lissitz and Samuelsen (2007), the validation of instruments used in education should involve the content analysis and empirical analysis of the scores obtained from the instrument and the respondents' response to the items. Content analysis of an instrument is associated with content analysis that later, also needs an empirical analysis to prove the construct validity. Both of these analyses are intended to make instruments in the world of education qualified as a standard measurement instrument.

Content validity is determined using expert agreement. Expert agreement, also called as measured domain determines the content validity stratification (content-related). This happens because of the measuring instruments, for example a test or questionnaire is proved to be valid if the expert believes that the instrument measures the mastery abilities defined in the domain or the measured psychological constructs. For understanding this agreement, a validity index can be used, including the index proposed by Aiken (1980; 1985). The item validity index proposed by Aiken is formulated as follows:

$$V = \frac{\sum s}{n(c-1)} \tag{1}$$

where V is the item validity index; s is scores assigned by each rater minus the lowest score in the used category ($s = r - l_0$, with r = rater category selection score and l_0 the lowest scores in the scoring category); n is the number of raters; and c is the number of categories that raters can choose.

Based on the afore-mentioned opinion, V is the rater's deal index of items' suitability

with indicators that need to be measured using the items. If it is applied to the measurement instrument, according to a rater, then n can be replaced by m (the number of items in an instrument). The V index value ranges from 0 to 1. The closer an item to 1, the better it is, because it is more relevant to the indicator.

Another way to prove the content validity with expert agreement is using expert index agreement suggested by Gregory (2007). The index also ranges from 0 to 1. It is done by making contingency tables on two experts, with the first category that is not relevant and less relevant become the weak relevancy category, and the second category which is for quite relevant and very relevant that is created in a new strong relevant category. The expert agreement index for content validity is a comparison of the number of items of the two experts with strong relevance category of overall items.

The expert agreement index for content validity is a comparison of the numbers of items from two experts as validators with strong relevance to the overall items category (Gregory, 2007). While the results of the relevancy tabulation (contingency tables) are presented in Table 1, the validity coefficient is presented in Formula 2.

Table 1. The relevance category scoring with two validators

		Validator 1	
		Weak	Strong
Validator 2	Weak	A	B
	Strong	C	D

$$\text{Content validity coefficient} = \frac{D}{(A+B+C+D)} \tag{2}$$

If the validators are three experts, the size of contingency tables with the number of cells $2 \times 2 \times 2 = 8$ cells is presented in Table 2. The content validity coefficient is an expansion coefficient of Formula 2. The coefficient calculation with the Formula 2 expansion is presented in Formula 3.

This coefficient also ranges from 0 to 1. The coefficient close to 0 means the validators' agreement index on the instrument item relevance with their indicators is getting

lower. Conversely, if the validity coefficient is closer to 1, the validators' agreement index about the instrument items relevance with their indicator becomes greater.

Method

Sub-indicators are compiled by using SRL components and indicators (adapted from Zimmerman, 2000). The results of the indicator development and the item numbers are presented in Table 3.

Instrument items which are the SRL scales are arranged by using the outline above. The scale is set in two forms: a questionnaire in Likert model and a multiple choice model.

For example, item 1 in Table 4 for items with Likert model and Table 5 for items with the multiple choice model.

Two forms of the outline/format and items of the instrument for measuring SRL were then given to three validators. The validators consisted of two educational psychologists and one educational measurement expert. The three validators assessed the items' relevancy with indicators, on both scale forms. Based on the results of the assessment of the three validators, then the validity index and validity coefficient were calculated using Aiken scale (Formula 1), on both of the scale models.

Table 2. Table of contingency to calculate the validity coefficient with Gregory formula involving three validators

Expert 1	Weak	Weak	Weak	Weak	Strong	Strong	Strong	Strong
Expert 2	Weak	Weak	Strong	Strong	Weak	Weak	Strong	Strong
Expert 3	Weak	Strong	Weak	Strong	Weak	Strong	Weak	Strong
Total	A	B	C	D	E	F	G	H

$$\text{Content validity coefficient} = \frac{H}{(A+B+C+D+E+F+G+H)} \tag{3}$$

Table 3. SRL Components and Indicators (adapted from Zimmerman (2000))

Components	Indicators	Sub Indicators	Items
Thought	Task Analysis	Goals Setting	1
		Strategic Planning	2
	Confidence	Self-Capability	3
		Task-Oriented	4
Performance Control	Self-control	Self-instruction	5
		Study Focus Effort	6
		Task-finishing Strategy	7
	Sufficient Observation	Metacognitive Observation	8
		Self-note	9
		Self-experimentation	10
		Self-reflection	Self-consideration
Self-reflection	Self-consideration	Causal Attribution	12
		Self-Reaction	Self-satisfaction (Reward)
	Self-Reaction	Self-satisfaction (Punishment)	14
		Adaptive/defensive	15

Table 4. Items with Likert model

No	Statements	Never	Seldom	Often	Always
1	I frame my study/course goals before the activity begins	1	2	3	4
8	I make maps of activities that I have done				

Table 5. Items with multiple-choice model

No.	Items
1.	At the beginning of the lecture (semester 1), a statement that is the most suitable with your condition is. . . . A. I frame my purposes clearly after I graduate. (4) B. I just know the best college for me, and my dream after graduate is not important. (2) C. I have a principle that life is just flowing, including the lecture. (1) D. I know what I will do after I graduate, but I am not sure with that. (3)
8.	About the efforts that you have done, which statement describes your condition. . . . A. I record my failure, so it motivates me to correct it. (3) B. Failure, success, and effort that I have been done or will do, I draw them only in my mind. (2) C. I do not map my efforts, success, and failures that I think I fail to correct. (1) D. I make a map or diagram of the efforts that I have done and their results, as a success or failure. (4).

Table 6. Experts final results of items compatibility with indicators data

Likert				Multiple-Choice			
Items	Rater1	Rater2	Rater3	Items	Rater1	Rater2	Rater3
1	4	4	2	1	4	4	2
2	4	4	4	2	4	4	4
3	4	4	4	3	4	4	4
4	4	4	4	4	4	3	3
5	4	2	4	5	4	3	2
6	4	4	4	6	4	4	4
7	4	2	3	7	4	2	3
8	4	2	4	8	4	3	4
9	4	2	4	9	4	3	4
10	4	4	4	10	3	3	4
11	4	4	4	11	4	4	4
12	4	4	4	12	4	4	4
13	4	2	4	13	4	4	4
14	4	4	3	14	4	4	4
15	4	4	4	15	4	4	4

Notes:

(4= Very Relevant, 3= Adequate Relevant, 2= Less Relevant, 4= Irrelevant)

By using the same data, a new category was created for relevancy, weak and strong classifications, with which a contingency table as shown in Table 2 was made. Furthermore, the validity coefficient was calculated using the extended Gregory formula (formula 2) in both scale models.

Findings and Discussion

The results of the assessment of the validators are inserted into Table 6. In addition to providing quantitative assessments, the validators also provide qualitative inputs, which include (1) statement improvement in Likert items, (2) stem items and the multiple-choice option improvement, and (3) according to the validators, Indonesian respondents are not familiar yet with the multiple-choice

questionnaire, because its reading takes longer time than a questionnaire with Likert model.

Furthermore, the results of the quantitative assessment, the items of validity index, and the scale of validity coefficient using Aiken formula are calculated on Likert model or scale with multiple-choice model. The results are presented in Table 7. Comparison on each item of the two models is presented in Figure 3.

Table 7 and Figure 3 show that the calculation results of the item validity index using Likert model and inventory model are not much different. Similarly, the scale using Likert model and Inventory model obtained are exactly the same in the result of validity coefficient calculation.

Table 7. The results of validity calculation using Aiken formula

Items	Likert	Multiple-choice
1	0.78	0.78
2	1.00	1.00
3	1.00	1.00
4	1.00	0.78
5	0.78	0.67
6	1.00	1.00
7	0.67	0.67
8	0.78	0.89
9	0.78	0.89
10	1.00	0.78
11	1.00	1.00
12	1.00	1.00
13	0.78	1.00
14	0.89	1.00
15	1.00	1.00
Scale	0.90	0.90

Based on the same data, the item relevance category that becomes only weak

and strong is created. Furthermore, each category is calculated on Likert questionnaire models presented in Table 8.

Based on Table 8, from a 15-item scale, there are nine strong items that have strong relevance according to the three validators' assessment. This shows that with Formula 3, the instrument reliability coefficient SRL measurement using Likert model obtains 0.60. Using the same technique, the relevant category of the validity coefficient in multiple-choice models is also created. The results are presented in Table 9.

Based on Table 9, from 15 items of the scale, there are 12 strongly relevant items according to the three validators' assessment. This shows that with Formula 3, reliability coefficient instrument of SRL measurement with multiple-choice models gains 0.80.

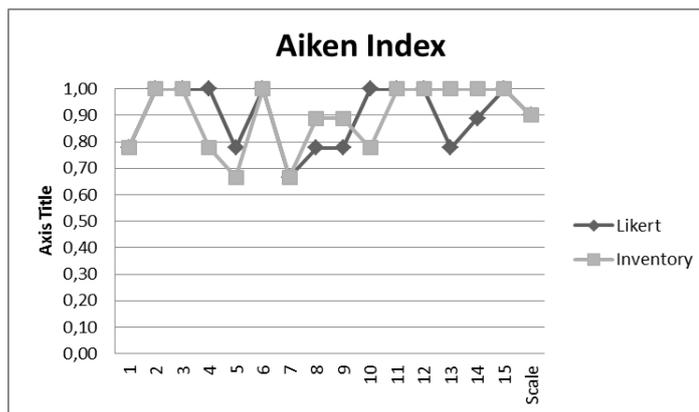


Figure 3. Aiken index on scale of Likert and Multiple-choice model

Table 8. Likert relevance category

Expert 1	Weak	Weak	Weak	Weak	Strong	Strong	Strong	Strong
Expert 2	Weak	Weak	Strong	Strong	Weak	Weak	Strong	Strong
Expert 3	Weak	Strong	Weak	Strong	Weak	Strong	Weak	Strong
Total	0	0	0	0	0	5	1	9

Table 9. Multiple-choice Relevancy Category

Expert 1	Weak	Weak	Weak	Weak	Strong	Strong	Strong	Strong
Expert 2	Weak	Weak	Strong	Strong	Weak	Weak	Strong	Strong
Expert 3	Weak	Strong	Weak	Strong	Weak	Strong	Weak	Strong
Total	0	0	0	0	0	1	2	12

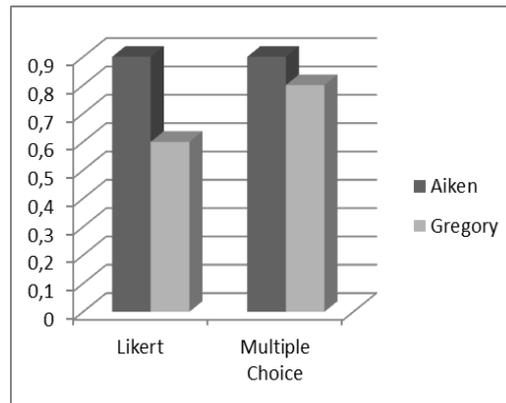


Figure 4. Comparison of validity coefficient using Aiken formula and Gregory formula

The comparison of calculation results of SRL validity coefficient scale' if it is compared based on its forms and formulas, is presented in Figure 4. Based on the image, it can be obtained that the result of the validity coefficient calculation using Aiken formula is more stable compared with it is using Gregory formula. From their shape, these results indicate that the validity coefficient calculated using Gregory formula on SRL scale of the multiple-choice model is lower than the validity coefficient scale calculated on the Likert model.

Conclusion

In this study, two instruments of SRL measurement on Likert model and multiple-choice model using the same format are developed. The formats and the two instrument models are then given to three validators to be assessed for the relevance of the items with indicators. The results of the expert assessment are used to prove the content validity using Aiken formula and expanded Gregory formula. The results of the study show that the content validity coefficients, based on expert assessment on Likert format and multiple choice with Aiken formula, are at 0.9 for each, with the index for each item being almost the same, and with the Aiken formula and expanded Gregory formula being 0.6 for Likert and 0.8 for multiple choice.

These results show that the acquisition of the index and the validity coefficient using Aiken formula on Likert model and multiple-choice model are almost the same. This

happens because both models are developed using the same format. However, when the validity verification is done by using the Gregory formula, the results are different. Coefficient acquisition using the Gregory formula is less than that using the Aiken formula, because in the Gregory formula, the probability to obtain the combination of all three validators on assessing a strong relevant item is very small.

Some future research projects that can be done are the stability of the number of validators. Further research is needed on the number of validators, so the acquisition of the index or the coefficient is maximized. It is better done on both Aiken formula and Gregory formula.

References

- Aiken, L, R. (1980). Content validity and reliability of single items or questionnaires. *Educational and Psychological Measurement*, 40, 955-967.
- Aiken, L, R. (1985). Three coefficients for analyzing the reliability and validity of ratings. *Educational and Psychological Measurement*, 45, 131-142.
- Allen, M. J. & Yen, W. M. (1979). *Introduction to measurement theory*. Monterey, CA: Brooks/Cole.
- American Educational Research Association, American Psychological Association, and National Council on Measurement in Education. (1999). *Standards for educational and psychological testing*.

- Washington, DC: American Psychological Association.
- Cizek, G.J., Rosenberg, S.L. & Koons, H.H. (2008). Source of validity evidence for educational and psychological test. *Educational and Psychological Measurement*, Vol. 68, pp. 397-412.
- Cohen, L., Manion, L., & Morrison, K. (2011). *Research methods in education*. New York, NY: Routledge.
- Fernandes, H. J. X. (1984). *Evaluation of educational program*. Jakarta: National Education Planning, Evaluating and Curriculum Development.
- Gregory, R.J. (2007). *Psychological testing: history, principles, and applications*. Boston, MA: Pearson.
- Kerlinger, F.N. (1986). *Asas-asas penelitian behavioral [Behavioral research principles]* (L.R. Simatupang, trans.). Yogyakarta: Gajahmada University Press.
- Kivinen, K. (2013). *Assessing motivation and the use of learning strategies by secondary school students in three international schools* (Unpublished doctoral dissertation). University of Tampere, Finland.
- Lawrence, M.R. (1994). Question to ask when evaluating test. *Eric Digest*. Retrieved from <http://www.ericfacility.net/ericdigest/edu.385607.html>.
- Linn, R.L. & Gronlund, N.E. (1995). *Measurement and assessment in teaching* (7th ed.). Englewood Cliffs, NJ: Prentice-Hall.
- Lissitz, W. & Samuelsen, K. (2007). Further clarification regarding validity and education. *Educational Researcher*, Vol. 36, No. 8, pp. 482-484.
- Messick, S. (1989). Validity in R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 13-103). New York, NY: Macmillan.
- Nunnally, J. (1978). *Psychometric theory* (2nd ed.). New York, NY: McGraw Hill.
- Schunk, D. H. (2005). Self-regulated learning. The educational legacy of Paul R. Pintrich. *Educational Psychologist*, 40, 85-94.
- Vrieling, E., Bastiaens, T., & Stijnend, S., (2012). Effects of increased Self-Regulated Learning opportunities on student teachers' motivation and use of metacognitive skills. *Australian Journal of Teacher Education*, Vol 37, 6, August 2012.
- Wolkers, C.A., Pintrich, P.R., & Karabenick, S.A., (2003). *Assesing academic Self-Regulated Learning*. Paper presented on the Conference on Indicators of Positive Development: Definitions, Measures, and Prospective Validity, Washington, DC.
- Woolfolk, A. E. & McCune, L. N. (1984). *Educational psychology for teachers*. Englewood Cliffs, NJ: Prentice Hall.
- Zimmerman, B.J. (1989). A social cognitive view of Self-Regulated Academic Learning. *Journal of Educational Psychology*, 81(3).
- Zimmerman, B.J. (1990). Self-Regulated Learning and academic achievement: An overview. *Education Psychologist*, 25(1), 3-17.
- Zimmerman, B. J. (2000). Attaining self-regulation: A social cognitive perspective. In Boekaerts, M., Pintrich, P. R., and Zeidner, M. (eds.), *Handbook of Self-Regulation: Theory, research, and applications*, Academic Press, San Diego, CA, pp. 13–39.
- Zumbrunn, S., Tadlock, J., & Danielle, E. (2011). Encouraging Self-Regulated Learning in the classroom: A review of the literature. *Metropolitan Educational Research Consortium (MERC)*. Virginia Commonwealth University.