# Analysis Difficulties and Characteristics of Item Test of on Biology National Standard School Examination

Aditya Nugraha Surya Saputra[1,*] Heri Retnawati[2,] Eri Yusron[3]

[1] *Master of Educational Research and Evaluation, Graduate School, Universitas Negeri Yogyakarta, Indonesia*
[2] *Department of Mathematics Education, Faculty of Mathematics and Natural Sciences, Universitas Negeri Yogyakarta, Indonesia*
[3] *Department of Elementary Teacher Education, Cibiru Campus, Universitas Pendidikan Indonesia, Indonesia*
*Corresponding author. Email: adityanugrahass737@gmail.com*

**ABSTRACT**

This research aims to describe the characteristics of biology test items in the test instrument of USBN (National Standard School Examination) through test items analysis to check of reliability estimation, difficulty level, and item discrimination, as well as to identify items that have a high level of difficulty, and to analyze the factors that cause the items these questions are difficult for students. Data collection technique was done through documentation. The collected data were in the form of a set of Biology National Standard School Examination key answers, and students' response of senior high school in Banjarbaru City, South Kalimantan Indonesia. Quantitative data analysis was conducted by using the approaches of Classical Test Theory (CTT) and Item Response Theory (IRT). Qualitative data analysis was carried out to describe students' difficulty in answering test items with difficult category. The results show that reliability estimation is quite reliable. Based on the analysis result of CTT and IRT, the mean of difficulty level is categorized as medium, and the mean of distinguishing ability is categorized good. In qualitative analysis, it is obtained that one item becomes difficult to answer due to the low concept comprehension, unfinished learning materials, and learning difficulties during the class.

*Keywords: Difficult item, Assessment, Examination, Biology education*

## 1. INTRODUCTION

Measurement, assessment, and evaluation are important components interrelated in education system. Each component has its own role in learning process, one of which is assessment. In education, assessment is an important matter in order to identify an educational success. The results of the educational assessment have a major function that will be useful in further educational processes [1]. Assessment is a process of data collection purposing to decision making on students or teachers. School's stakeholder uses information from the assessment process to decide on what has been learned by the students, what and where they should be taught, and several types of related services (e.g. speech and language service, psychological service) that they need [2]. National assessment is arranged to describe students' achievement in the area where that curriculum is applied and later collected to give estimation of students' achievement level in the education system as a whole on particular grade or age level. It provides data for the type of national education audit which is done to inform the policy makers about the key aspect of a system [3].

In Indonesia, educational assessment can be carried out by teachers, education units, or the government. Several assessments carried out by the government through National Education Standards Agency are in the forms of UN or National Examination and USBN or National Standard School Examination. National Standard School Examination is an activity to measure the competency attainment of students carried out by educational units for certain subjects by referring to SKL (Graduate Competency Standards) to obtain recognition for learning achievement [4]. Competencies measured in USBN are in accordance with graduate competency standards derived from basic competencies in the curriculum. Then, the results of the USBN were analyzed by the school with the aim of knowing the achievement of

competencies against standards in one subject or between subjects [5]. By considering various purposes and functions of National Standard School Examination, A test instrument with decent quality is required, so that the test items included can identify the mastery level of learning materials assigned to the students.

During the implementation, the quality of these test items compiled in National Standard School Examination test instrument is still questionable, given that the preparation of secondary education level National Standard School Examination test items carried out by teachers who are members of the Subject Teacher Forum does not have adequate quality, due to human resources which have not been maximally trained, especially the competence of teachers in composing questions [4]. This assessment activity is bound with measurement one to obtain accurate data, thus it requires good measuring tools. Among many types of education assessment, test is the most frequent used measuring tool [6]. In realisation, the measuring tool used in National Standard School Examination is test instrument. Thus, it is crucial to analyse each item to find out its quality, so that in the measurement of test items in National Standard School Examination test instrument is capable to give the needed information in assessment. This information is later able to illustrate the actual students' condition, as in knowledge mastery and skill to achieve learning purposes well.

To find out the quality of one instrument, it can be done by observing its characteristic through item analysis. It can be carried out with two approaches, namely Classical Theory Test (CTT) and Item Response Theory (IRT). Measurement model on Classical Theory Test (CTT) can be illustrated as observation scores consisting of the actual scores and measurement errors [7]. IRT is a connection between skill or characteristic measured with instrument and item response [8]. IRT is related to the measurement of latent hypothesis construction and can only be measured indirectly through other manifest variable measurement. This hypothesis construction is a latent variable and often represents one's competence and skills. Latent variable will also be called as competence parameter symbolized as $\theta$ [9]. As more simply, Item Response Theory (IRT) is a relation between competence or characteristics measured with instrument and item responses [8].

This research applied Classical Theory Test and Item Response Theory approaches to describe the characteristics of items in National Standard School Examination test instrument of biology subject

through test items analysis to observe reliability estimation, difficulty level, and discrimination, and to identify test items with high-level difficulty, and to analyse the factors which make those items difficult for students.

## 2. METHOD

This research was a descriptive explorative research by applying both quantitative and qualitative approaches. The data were collected through documenting the National Standard School Examination in 2018/2019 in which it was last held before being changed and returned to School Examination in 2019/2020. Documentation was in the form of Biology USBN test consisting of 35 test items, and because there was a limitation in the documentation, so that the responses obtained were only as many as 160 students during National Standard School Examination in Banjarbaru City, South Kalimantan. Data analysis was carried out quantitatively and qualitatively. The first one was done through observing estimation of instrument reliability, then continued with the analysis of Classical Theory Test (CTT) assisted by Quest software and Item Response Theory (IRT) assisted by Bilog-MG software to calculate difficulty levels and discrimination. The second one was done by describing those items categorized as difficult. The steps in qualitative analysis include identifying test items with high-level of difficulty and analysing the factors which make those items difficult for students.

Acceptable reliability estimation generally has minimum coefficient ranging from 0.80 to 0.90 [10]. In CTT approach, the approximation of difficulty level from 0.3 to 0.7 is the best one [11]. Item discrimination can still be used if the value is not negative [12]. In IRT approach, item difficulty level or b-parameter theoretically ranges from $-\infty$ until $+\infty$ yet generally a good instrument is ranging between -2 to +2 [8].

## 3. RESULT

The analysis results of reliability estimation in Biology National Standard School Examination test items acquired is as many as 0.48. The data of students' response are analysed using Quest to find out difficulty level of test items and discrimination based on CTT approach. The former is shown in Table 1 and the latter is shown in Table 2.

**Table 1.** The results of difficulty level ($\rho$)

| No. Item | $p$ | Category Item |
|---|---|---|
| 1 | 0.377 | moderate |
| 2 | 0.294 | difficult |
| 3 | 0.344 | moderate |
| 4 | 0.792 | easy |
| 5 | 0.742 | easy |
| 6 | 0.331 | moderate |
| 7 | 0.319 | moderate |
| 8 | 0.500 | moderate |
| 9 | 0.331 | moderate |
| 10 | 0.731 | easy |
| 11 | 0.525 | moderate |
| 12 | 0.635 | moderate |
| 13 | 0.570 | moderate |
| 14 | 0.453 | moderate |
| 15 | 0.231 | difficult |
| 16 | 0.575 | moderate |
| 17 | 0.888 | easy |
| 18 | 0.344 | moderate |
| 19 | 0.350 | moderate |
| 20 | 0.400 | moderate |
| 21 | 0.706 | easy |
| 22 | 0.384 | moderate |
| 23 | 0.363 | moderate |
| 24 | 0.575 | moderate |
| 25 | 0.241 | difficult |
| 26 | 0.956 | easy |
| 27 | 0.673 | moderate |
| 28 | 0.644 | moderate |
| 29 | 0.563 | moderate |
| 30 | 0.394 | moderate |
| 31 | 0.719 | easy |
| 32 | 0.313 | moderate |
| 33 | 0.319 | moderate |
| 34 | 0.157 | difficult |
| 35 | 0.342 | moderate |

From the results in Table 1, four items are categorized as difficult, twenty-four items as moderate, and seven items as easy. The mean of

difficulty level from the whole items is 0.479 and is categorized as moderate.

**Table 2.** The results of discrimination ($r_{pbis}$)

| No. Item | $r_{pbis}$ | Category Item |
|---|---|---|
| 1 | 0.410 | good |
| 2 | -0.050 | not good |
| 3 | 0.150 | good |
| 4 | 0.380 | good |
| 5 | 0.280 | good |
| 6 | 0.250 | good |
| 7 | 0.060 | good |
| 8 | 0.170 | good |
| 9 | 0.280 | good |
| 10 | 0.240 | good |
| 11 | 0.200 | good |
| 12 | 0.170 | good |
| 13 | 0.170 | good |
| 14 | 0.250 | good |
| 15 | 0.150 | good |
| 16 | 0.280 | good |
| 17 | 0.050 | good |
| 18 | 0.180 | good |
| 19 | 0.320 | good |
| 20 | 0.450 | good |
| 21 | 0.220 | good |
| 22 | 0.020 | good |
| 23 | 0.320 | good |
| 24 | 0.410 | good |
| 25 | 0.360 | good |
| 26 | 0.130 | good |
| 27 | 0.200 | good |
| 28 | 0.360 | good |
| 29 | 0.240 | good |
| 30 | 0.130 | good |
| 31 | 0.200 | good |
| 32 | 0.200 | good |
| 33 | 0.350 | good |
| 34 | 0.210 | good |
| 35 | 0.300 | good |

Table 2 shows that item discrimination has thirty-four items categorized as good and one item categorized as not good. The mean of item discrimination from the entire items is as many as 0.23

which is categorized as good. Through these results, the overall items on the test instrument have a good discrimination to differentiate the abilities of students.

Students' response data are analysed with Bilog-MG to discover item difficulty level based on IRT approach. Before going into that, it is necessary to pay attention on PH1 output which is an analysis result with CTT approach on biserial point as a representation of discrimination. Based on the analysis, there are 12 items from total in a set of test items which have negative biserial point value, namely in number 1, 2, 3, 5, 7, 12, 16, 17, 21, 22, 29, and 30. Thus, these items need to be omitted in the next analysis, as it only requires positive biserial point value. It must be done because as there exists the negative value, the items resulted cannot distinguish between upper class (students with high competence) and lower class (students with low competence) [8].

After omitting those 12 items, the analysis using Bilog-MG is continued with the remaining 23 items. It is shown in PH1 outputs that these 23 items do not have negative biserial point anymore, thus it can be continued by looking at PH2 outputs, i.e. analysis results to discover item difficulty level, demonstrated in Table 3.

**Table 3.** Item difficulty level (*b*)

| No. Item | b | Category Item |
|---|---|---|
| 4 | -2.727 | easy |
| 5 | -2.148 | easy |
| 8 | -0.003 | moderate |
| 9 | 1.426 | moderate |
| 10 | -2.039 | easy |
| 11 | -0.208 | moderate |
| 13 | 5.632 | difficult |
| 14 | 0.265 | moderate |
| 15 | 2.440 | difficult |
| 18 | 1.317 | moderate |
| 19 | 1.261 | moderate |
| 20 | 0.573 | moderate |
| 23 | 1.149 | moderate |
| 24 | -0.621 | moderate |
| 25 | -0.333 | moderate |
| 26 | -6.154 | easy |
| 27 | -1.473 | moderate |
| 28 | -1.21 | moderate |
| 31 | -1.913 | moderate |
| 32 | 1.605 | moderate |

| No. Item | b | Category Item |
|---|---|---|
| 33 | 1.546 | moderate |
| 34 | 3.4 | difficult |
| 35 | 1.341 | moderate |

As demonstrated in Table 3, there are three items in difficult category, sixteen items in moderate one, and four items in easy one. The mean of difficulty is by 0.135 categorized as moderate.

## 4. DISCUSSIONS

Test reliability estimation can be relied as it is far from reliability coefficient limit, that is 0.80. The smaller it is, the bigger the measurement errors. There are several factors which influence estimation of reliability, such as group homogeneity, allocated time, and the length of test [10]. Other influential factors is the more items classified as difficult, the lower its estimation of reliability [13].

Based on the analysis, there are 4 items in CTT and 3 items in IRT with difficulty level from total items in National Standard School Examination instrument test. Several items with identical difficulty level are chosen to be examined further and to observe the factors of students' difficulty in answering each item.

To be able to answer correctly, it requires a strong knowledge about the concept of urine formation which is part of the excretory system, including its chemical substance and kidney parts which play an important role in urine formation process. There three phases in this process, namely Filtration, Reabsorption, and Augmentation. Filtration is formed when liquid streams from bloodstream to Bowman lumen capsule. Glomerular capillaries and particular cell of Bowman



**Figure 1** Item 15

capsule retain the blood cells and large molecules, such as plasma protein, in which it can be penetrated by water and other small solutes. Hence, the produced filtration in bowman capsule contains salt, glucose, amino acids, vitamin, nitrogen waste, and other small molecules. Since such a molecule passes freely between the glomerular capillaries and Bowman's capsule, the concentration of this substance at the initial filtration equals the concentration in blood plasma. Meanwhile, reabsorption in the proximal tubule is essential to recapture ions, water, and valuable nutrients from the large volume of the initial filtrate. Glucose, amino acids, potassium ions (K+), and other essential substances are also actively or passively transported from the filtrate to the interstitial fluid and then to the peritubular capillaries in this reabsorption phase [14]. The answer key for the item above is C.

The difficulty level in this item is 0.261 based on CTT analysis and 2.974 from IRT analysis where both are categorized as difficult. The students chose D in Moderate by 36.88% compared to those who chose C by 23.12%. It is classified as Higher Order Thinking Skills (HOTS) type, where it requires students' high intelligence. A discussion with a biology teacher explains that during the class, students are often given Lower Order Thinking Skills (LOTS) in school textbooks or student's worksheets. It can be said that the form of questions that do not provide a stimulus for students to think at high levels will cause them to have difficulty answering questions with the HOTS type in the future, including this USBN item. It also points out the students' weakness in applying urine formation concept within a test item. Even though they are able to describe verbally, they are still unable to maximizes it during problem solving [15]. Excretion system is also a difficult concept or topic for middle school students to learn it [16].
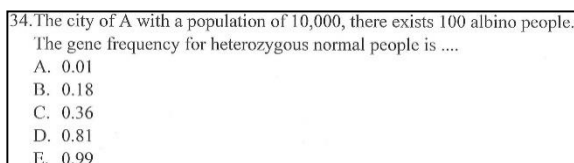

**Figure 2** Item 34

It requires a strong knowledge regarding Hardy-Weinberg's law in Evolution material in order to answer the item above correctly. In an undeveloped population, the allele and genotype frequencies will remain constant from generation to generation, provided only Mendelian segregation and allele recombination are in effect. Such a population is said to be in Hardy-Weinberg equilibrium. The Hardy-Weinberg equilibrium equation for the three

genotypes will appear with the proportion $p2 + 2pq + q2$. Consequently, the amount of frequency from all three genotypes must equal 1 (100%) in any population apart from whether it is situated in Hardy-Weinberg equilibrium. The main point is that one point is in this equilibrium only if the genotype frequency observed from one homozygote is $p2$, another frequency observed from other homozygote is $q2$, and the heterozygote frequency observed is $2pq$ [14]. The answer key for the item above is B.

The difficulty level is by 0.156 from CTT analysis and 2.874 from IRT one, in which both are categorized as difficult. The students chose A in Moderate by 36.88% compared to those who chose B by 15.67%. A discussion with a Biology teacher explains that the factors which influence students in answering this item is that they do not know the pattern or they know it but do not know to implement it during the test. Among the main difficulties in evolutionary biology teaching and learning, is that there are many alternative concepts [17]. In addition, the material on Evolution in biology learning is located at the end of discussion or the last semester in 12th grade. In this phase, special treatment is needed as at the end of the semester, the 12th grade students are busied with final exam simulation, such as practical exams, National Standard School Examination, and National Examination. It causes the teacher occasionally unable to deliver the material completely due to collisions with those activities. In the end, learning in class that is not optimal makes it difficult for students to interpret what they have learned into questions, so that what makes the items difficult for students. Thus, an effective teaching is needed in order to produce effective learning.

During learning process other than teaching, assessment is one of the vital components which must be mastered by a teacher. Even though the assessment report is carried out at the end of the cycle, the assessment must be designed since the beginning to ensure that the desired information type can be done when delivering the results [18]. Hence, a good assessment system is required in order to result a beneficial information in decision making or a proper responsibility during evaluation. The needed information during assessment is gained through measurement. Measurement is a process to obtain numerical description regarding how far the particular characteristic one owns [19]. In other words, measurement model, that is a framework employed to communicate with others about the evidence gained from observation, can be used to draw a conclusion regarding students' characteristics contained in construct variable [18]. It requires a measuring tool

during the test in the form of a set of test items to fulfil information necessities.

One of the biggest factors in preparing test items is a teacher [20]. They must pay attention on the characteristics of test items. The effect of disregarding item's characteristics in the test causes the measurement results to be unsuitable with the students' actual competence as a consequence to the items which are too difficult or too easy. In addition, two other influential factors are the competence and readiness of the students. The higher they are, the more the students are able to answer the questions [13]. It can be said that a difficult item finally cannot illustrate the test participant's actual condition or competence. It will make an impact on the follow-up assessment results. Therefore, it can be said that in the end, the difficult test items are unable to explain participants' actual condition or competence. It will certainly bring an impact to the follow-up of assessment results which causes the teachers unable to comprehend which student or groups have high and low competence.

Within the discussion, it is mentioned that the items categorized as difficult is caused by a low concept comprehension and incomplete learning materials and learning difficulties during the class. Furthermore, the difficulty factors are that the materials tend to be abstract, relying on memorizing ability, using Latin, containing complexed knowledge, inconvenient practice activity, non-related materials with human's daily basis [21]. Other factors such as teacher competence and professionalism, teacher experience, and the choice of learning model used by teachers in the classroom are also very influential. Thus, it requires an evaluation towards those matters.

Evaluation can be done by reviewing or changing any forms which become in input in Biology teaching. National Education Academy Panel recommends giving special attention to students' cognition aspect such as problem representation, strategy use and self-regulation skill, and formulation of explanation and interpretation as consideration of required aspects for National Assessment of Educational Progress to assess a complete and accurate achievement in the subject area [18]. Students' concept comprehension is also crucial given that there are a large number of concepts in Biology materials in every phase. This can be done through learning activities by equipping students with problem models, especially those related to real contexts [1]. Then, the students who have mastered conceptual comprehension are not only able to describe the concept properly, but also to explore the concepts in different situation and utilize them in appropriate area within problem solving [15]. Next,

those materials are closely related with abstract, hence teachers must comprehend their surroundings to find learning sources which can be connected with the daily basis. It can support the students learning activities [22]. The schedules collision with final exam simulations or stabilization causes the materials are not entirely delivered. Therefore, it is very important to implement educational training programs that support the professionalism of teacher performance.

## 5. CONCLUSION

Based on the results, it is discovered that the reliability estimation in the Biology National Standard School Examination test items is quite reliable as it is far from 0.80. From CTT analysis, the Moderate difficulty level is categorized as Moderate, and the discrimination is categorized as good. From IRT analysis, the Moderate of it is categorized as Moderate. Several test items with equal difficulty level of items on CTT and IRT results are to be examined further and the factors causing students' difficulty in doing each item are to be checked as well. It is stated in the discussion results that an item becomes difficult to answer is caused by low concept understanding, incomplete learning materials, and difficulties during the class. It needs an evaluation to solve those problems. Firstly, an evaluation on Biology learning by reviewing or changing all forms as inputs in it, such as an emphasis on students' understanding, and selecting learning resources and materials mastery to minimise the abstraction during learning process. Secondly, an evaluation on teachers' competence in managing learning process, hence a support from the government and related departments is needed to improve their competence through professional training programs. It is necessary to do more in-depth research on the students' difficulty in answering the questions in Biology subject, especially in the final exam based on student perceptions.

## REFERENCES

[1] H. Retnawati, B. Kartowagiran, J. Arlinwibowo, E. Sulistyaningsih, Why are the Mathematics National Examination Items Difficult and What Is Teachers' Strategy to Overcome It?, International Journal of nstruction 10(3) (2017) 257–276. DOI: https://doi.org/10.12973/iji.2017.10317a

[2] J. Salvia, J. E. Ysseldyke, S. Bolt, Assessment In Special and Inclusive Education, Eleventh E. Wadsworth, Cengage Learning, 2010.

[3] V. Greaney, T. Kellaghan, Assessing National Achievement Levels in Education, World Bank, 2008.

[4] I. Ulumudin, Management of National Standard School Examination, Jurnal Penelitian Kebijakan Pendidikan 11(3) (2018). DOI: https://doi.org/10.24832/jpkp.v11i3.207

[5] U. Rosidin, Herpratiwi, W. Suana, R. Firdaos, Evaluation of National Examination (UN) and National-Based School Examination (USBN) in Indonesia, European Journal Educational Research 8(3) (2019) 827–837. DOI: https://doi.org/10.12973/eu-jer.8.3.827

[6] R. Sumner, The Role of Assessment in Schools, A New Edit, The NFER-NELSON Publishing Company Ltd, 2004.

[7] D. N. M. Gruitjer, L. J. T. Kamp, Statistical Test Theory for The Behavioral Sciences, Chapman & Hall/CRC, 2008.

[8] C. DeMars, Item Response Theory, Oxford University Press, Inc, 2010.

[9] J.-P. Fox, Bayesian Item Response Modeling: Theory and Applications, Springer, 2010.

[10] L. Crocker, J. Algina, Introduction to Classical and Modern Test Theory, Cengage Learning, 2008.

[11] J. A. Marry, M. Y. Wendy, Introduction to Measurement Theory. Monterey, Brooks/Cole Publishing Company, 1979.

[12] D. A. Frisbie, Measurement 101: Some Fundamentals Revisited, Educational Measurement: Issues and Practice 24(3) (2005) 21–28. DOI: https://doi.org/10.1111/j.1745-3992.2005.00016.x

[13] F. Sampouw, H. Retnawati, Characteristics of non-compulsory mathematics test items on nationally standardized school examination in Kaimana Regency, West Papua Indonesia, in: Proceedings of 3th International Seminar on Innovation in Mathematics and Mathematics Education, vol. 1581, IOP Publishing, Bristol, 2020, DOI: https://doi.org/10.1088/1742-6596/1581/1/012034

[14] J. B. Reece, L. A. Urry, N. A. Campbell, Campbell Biology, Eleventh e, Pearson Higher Education, 2016.

[15] M. Karagöz, M. Çakir, Problem Solving in Genetics: Conceptual and Procedural Difficulties, Educational Sciences: Theory & Practice 11(3) (2011) 1668–1674.

[16] A. Çimer, What Makes Biology Learning Difficult and Effective: Students' Views, Educational Research and Reviews 7(3) (2012) 61–71. DOI: https://doi.org/10.5897/ERR11.205

[17] R. Tidon, R. C. Lewontin, Teaching Evolutionary Biology, Genetics and Molecular Biology 27(1) (2004) 124–131. DOI: https://doi.org/10.1590/S1415-47572004000100021

[18] National Research Council, Knowing What Students Know: The Science and Design of Educational Assessment, National Academy Press, 2001.

[19] M. D. Millner, R. L. Linn, N. E. Gronlund, Measurement and Assessment in Teaching, Tenth Edit, Pearson Education Inc., 2009.

[20] A. Friatma, A. Anhar, Analysis of validity, reliability, discrimination, difficulty and distraction effectiveness in learning assessment, in: Proceeding of The 3rd International Conference on Education, Science, and Technology, vol. 1387, IOP Publishing, Bristol, 2019, DOI: https://doi.org/10.1088/1742-6596/1387/1/012063

[21] G. Hadiprayitno, Muhlis, Kusmiyati, Problems in learning biology for senior high schools in Lombok Island, in: Proceeding of The International Seminar on Bioscience and Biological Education, vol. 1241, IOP Publishing, Bristol, 2019, DOI: https://doi.org/10.1088/1742-6596/1241/1/012054

[22] T. Ozcan, S. Ozgur, A. Kat, S. Elgun, Identifiying and comparing the degree of difficulties biology subjects by adjusting it is reasons in elemantary and secondary education, in: Procedia - Social and Behavioral Sciences, vol. 116, Elsevier, Amsterdam, 2014, pp. 113–122. DOI: https://doi.org/10.1016/j.sbspro.2014.01.177