

NEPS SURVEY PAPERS

Inga Hahn and Jana Kähler

NEPS TECHNICAL REPORT FOR
SCIENCE: SCALING RESULTS OF
STARTING COHORT 4 IN 11TH GRADE

NEPS Survey Paper No. 6
Bamberg, July 2016

Survey Papers of the German National Educational Panel Study (NEPS)

at the Leibniz Institute for Educational Trajectories (LifBi) at the University of Bamberg

The NEPS Survey Paper Series provides articles with a focus on methodological aspects and data handling issues related to the German National Educational Panel Study (NEPS).

The NEPS Survey Papers are edited by a review board consisting of the scientific management of LifBi and NEPS.

They are of particular relevance for the analysis of NEPS data as they describe data editing and data collection procedures as well as instruments or tests used in the NEPS survey. Papers that appear in this series fall into the category of 'grey literature' and may also appear elsewhere.

The NEPS Survey Papers are available at <https://www.neps-data.de> (see section "Publications").

Editor-in-Chief: Corinna Kleinert, LifBi/University of Bamberg/IAB Nuremberg

Contact: German National Educational Panel Study (NEPS) – Leibniz Institute for Educational Trajectories – Wilhelmsplatz 3 – 96047 Bamberg – Germany – contact@lifbi.de

NEPS Technical Report for Science: Scaling Results of Starting Cohort 4 in 11th Grade

*Inga Hahn & Jana Kähler, Leibniz Institute for Science and Mathematics
Education (IPN), Kiel*

E-mail address of the lead author:

hahn@ipn.uni-kiel.de

Bibliographic data:

Hahn, I. & Kähler, J. (2016): NEPS Technical Report for Science: Scaling Results of Starting Cohort 4 in 11th Grade (NEPS Survey Paper No. 6). Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study. doi:10.5157/NEPS:SP06:1.0

NEPS Technical Report for Science: Scaling Results of Starting Cohort 4 in 11th Grade

Abstract

The National Educational Panel Study (NEPS) aims at investigating the development of competences across the whole life span and designs tests for assessing these different competence domains. In order to evaluate the quality of the competence tests, a wide range of analyses based on item response theory (IRT) have been performed. This paper describes the data on scientific literacy for starting cohort 4 in grade 11. Besides presenting descriptive statistics for the data, the scaling model applied to estimate competence scores and analyses performed to investigate the quality of the scale as well as the results of these analyses are also explained. The science test in grade 11 originally consisted of 29 multiple choice and complex multiple choice items and covers two knowledge domains as well as three different contexts. Five items had to be removed due to insufficient item quality. The test was administered to 4,417 students. A partial credit model was used for scaling the data. Item fit statistics, differential item functioning, Rasch-homogeneity, and the test's dimensionality were evaluated to ensure the quality of the test. The results of the remaining test items illustrate good item fit values and measurement invariance across various subgroups. Moreover, the test showed a moderate reliability. The data shows that the assumption of unidimensionality of scientific literacy measured by this test seems adequate. Among the challenges of this test is the lack of very easy items. But overall, the results emphasize the good psychometric properties of the science test, thus supporting the estimation of reliable scientific literacy scores. In this paper, the data available in the Scientific Use File are described and the ConQuest-Syntax for scaling the data is provided.

Key words:

scientific literacy, 11th grade, differential item functioning item response theory, scaling, scientific use file

Content

Abstract.....	2
1. Introduction.....	4
2. Testing Scientific Literacy	4
3. Data	4
3.1 The design of the study	4
3.2 Sample	6
4. Analyses.....	6
4.1 Missing responses.....	6
4.2 Scaling model.....	6
4.3 Checking the quality of the scale.....	7
5. Results	8
5.1 Exclusion of cases from the analyses	8
5.2 Descriptive statistics of the responses	9
5.3 Missing responses.....	9
5.3.1 Missing responses per person	9
5.3.2 Missing responses per item	11
5.4 Parameter estimates	14
5.4.1 Item parameters	14
5.4.2 Person parameters.....	14
5.4.3 Test targeting and reliability	14
5.5 Quality of the test.....	17
5.5.1 Fit of the subtasks of complex multiple-choice items	17
5.5.2 Item fit.....	17
5.5.3 Differential item functioning	17
5.5.4 Rasch-homogeneity	21
5.5.5 Unidimensionality of the test	21
6. Discussion	21
7. Data in the Scientific Use file	22
References.....	23
Appendix.....	25

1. Introduction

Within the National Educational Panel Study (NEPS), different competences are measured coherently across the life span. Tests have been developed for different domains including scientific literacy. Weinert et al. (2011) give an overview of the competence domains measured in NEPS.

Most of the competence data are scaled using models based on item response theory (IRT). Since most of the competence tests were developed solely for implementation in NEPS, several analyses have been performed to evaluate the quality of the test. The IRT models chosen for scaling, the competence data, and the analyses performed for checking the quality of the scale are described in Pohl and Carstensen (2012a). In this paper the results of these analyses are presented for scientific literacy in the starting cohort 4 for grade 11.

The present report has been modeled along the technical reports of Pohl, Haberkorn, Hardt, and Wiegand (2012) and Haberkorn, Pohl, Hardt, and Wiegand (2012). Note that the analyses of this report are based on preliminary data releases. Due to data protection and data cleaning issues, the data set in the Scientific Use File (SUF) may differ slightly from the data set used for the analyses in this paper. We do, however, not expect severe changes in results.

2. Testing Scientific Literacy

The science test aims at assessing two types of scientific sub-competencies. These are a) knowledge of science (KOS) and b) knowledge about science (KAS). Using the definition by PISA (OECD, 2007; Prenzel et al., 2007), KOS is specified as knowledge of basic scientific concepts and facts whereas KAS can be regarded as the understanding of scientific processes.

KOS is divided into content-related components: matter, system, development, and interaction. KAS is divided in the process-related components scientific enquiry and scientific reasoning. KAS and KOS are implemented in three contexts: health, environment, and technology (see Hahn et al., 2013, and Weinert et al., 2011, for the description of the framework). The test items are organized in units (testlets). Thus, one unit consists of two or three items. Each unit refers to one context-component combination.

There are two types of response formats. These are simple multiple choice (MC) and complex multiple choice (CMC) in the special form of true-false items. In MC items the test taker has to find the correct answer out of four response options. In CMC items the test taker has to decide at each answer option whether the answer is correct or not.

3. Data

3.1 The design of the study

Since scientific literacy was the only competency tested in this study, there was only one testing group who received the science test first and afterwards completed their questionnaire. The test time for the scientific literacy test was 29 minutes, with one additional minute for the procedural metacognition item. There was no multi-matrix design

regarding the choice and order of the items within a test. All students got the same test items in the same order.

The scientific literacy test in grade 11 originally consisted of 29 items. The characteristics of these 29 items are depicted in Table 1. Table 2 is concerned with the response format whereas Table 3 shows how the items cover the different contents and components of the science framework (see Hahn et al., 2013). Five of the 29 items had to be removed from the final analysis presented in this paper due to insufficient item quality.

Table 1: Classification of the science test items for grade 11

Knowledge domains	Frequency
Knowledge of Science (KOS)	21
Knowledge about Science (KAS)	8
Total number of items	29

Table 2: Response formats of the science test items for grade 11

Response format	Frequency
Simple Multiple-Choice	14
Complex Multiple-Choice (True-false items)	15
Total number of items	29

Table 3: Number of items for the different contexts of the science test for grade 11

Context	Frequency
Health	10
Environment	9
Technology	10
Total number of items	29

3.2 Sample

Overall, 16,425 students are part of the sample. 4,417 of these students took the science test. The sample was reduced to this size because only persons in high-school (“Gymnasium”) or persons attended an integrated school (“integrative Gesamtschule”) took the test. 3,921 students attended a regular high school while 496 students an integrated school.

All 4,417 persons who took part in the science test are included in the descriptive analyses. The results are presented in the following sections.

4. Analyses

4.1 Missing responses

There are different kinds of missing responses. These are a) invalid responses, b) missing responses due to omitted items, c) missing responses due to items that have not been reached, d) missing responses due to items that have not been administered, and e) multiple kinds of missing responses that occur in an item and are not determined. In this study, all subjects received the same set of items. As a consequence, there are no items that were not administered to a person.

Invalid responses occur, for example, when two response options are selected in simple MC items where just one is required, or when numbers or letters that are not within the range of valid responses are given as a response. Missing responses due to omitted items occur when test persons skip items. Due to time limits, it might happen that not every person finishes the test within the given time. Consequently, missing responses occur due to the fact that items are not reached. As complex multiple choice items are aggregated from several subtasks, different kinds of missing responses or a mixture of valid and missing responses may be found in these items. A CMC item is coded as missing if at least one subtask contained a missing response. When just one kind of missing response occurs, the item is coded according to the corresponding missing response. When the subtasks contain different kinds of missing responses, the item is labeled as a not-determinable missing response.

Missing responses provide information on how well a test works (e.g., time limits, understanding of instructions, handling of different response formats) and they need to be accounted for in the estimation of item and person parameters. We, therefore, thoroughly investigated the occurrence of missing responses in the test. First, we looked at the occurrence of the different types of missing responses per person. This gave an indication of how well the persons were coping with the test. We then examined the occurrence of missing responses per item, in order to get some information on how well the items worked.

4.2 Scaling model

For estimating item and person parameters for scientific literacy, a partial credit model (Masters, 1982) was used and estimated in ConQuest (Wu, Adams, & Wilson, 1997). A detailed description of the scaling model can be found in Pohl and Carstensen (2012a).

CMC items consist of a set of subtasks that were aggregated to a polytomous variable for each CMC item, indicating the number of correctly responded subtasks within that item¹. If at least one of the subtasks contains a missing response, the whole CMC item was scored as missing. When categories of the polytomous variables had less than 200 persons of the sample, the categories were collapsed in order to avoid possible estimation problems. This usually occurred for the lower categories of polytomous items; especially, when the item consisted of many subtasks. In these cases, the lower categories were collapsed to one category.

To estimate item and person parameters, a scoring of 0.5 points for each category of the polytomous items was applied while simple MC items were scored dichotomously as 0 for an incorrect and 1 for the correct response (see Haberkorn, Pohl, Carstensen, & Wiegand, 2012, and Pohl & Carstensen, 2012b, for studies on the scoring of different response formats).

Ability estimates for scientific literacy will be estimated as weighted maximum likelihood estimates (WLEs, Warm, 1989) and later also in form of plausible values (Mislevy, 1991). Person parameter estimation in NEPS is described in Pohl and Carstensen (2012a) while the data available in the SUF are described in section 7. The item parameters were plotted to the ability estimates of the persons was done in order to judge how well the item difficulties are targeted to the ability of the persons. The test targeting gives some information about the precision of the ability estimates at the different levels of ability.

4.3 Checking the quality of the scale

The grade 11 science test was specifically constructed to be implemented in NEPS. In order to ensure appropriate psychometric properties, the quality of the test was evaluated in pilot studies but also checked in several analyses for the data from the main study.

The responses on the subtasks of CMC items are aggregated to a polytomous variable for each CMC item. In order to justify such an aggregation, the fit of the single subtasks is checked in analyses. For this, the single subtasks are separately included in a Rasch model together with the MC items and the fit of the subtasks is evaluated based on the weighted mean square error (WMNSQ), the respective *t*-value, point-biserial correlations of the responses with total correct score, and the item characteristic curve. Only if the subtasks have a satisfactory item fit, they were used to construct polytomous CMC item variables.

MC and CMC items consisted of one correct response and a number of distractors (incorrect response options). We investigated whether the distractors worked well, that is, whether they are chosen by the students with a lower general ability in science more often than by those with a higher general ability in science. For this, we evaluated the point-biserial correlation of giving a certain incorrect response and the total number correct score estimated in the analysis treating all subtasks of CMC items as single items. We judged correlations below zero as very good, correlations below 0.05 as acceptable and correlations above 0.05 as problematic.

¹ As described later, due to collapsing of categories, this interpretation does not necessarily hold for the variables in the SUF.

Item fit was then evaluated for the MC items and the polytomous CMC items based on results of a partial credit model. Again, the weighted mean square error (WMNSQ), the respective t -value, point-biserial correlation of the correct response with the total score, and the item characteristic curve were evaluated for each item. Items with a WMNSQ > 1.15 (t -value > |6|) were considered as having a noticeable item misfit and items with a WMNSQ > 1.2 (t -value > |8|) were judged as a considerable item misfit and their performance was further investigated. Point-biserial correlations of the correct responses with the total score greater than 0.3 were considered as good, greater than 0.2 as acceptable and below 0.2 as problematic. The overall judgment of the fit of an item was based on all fit indicators.

We aimed at constructing a science literacy test that measures the same construct for all students. If there are items that favor certain subgroups (e.g., that are easier for boys than for girls), measurement invariance would be violated and a comparison of literacy scores between the subgroups (e.g., males and females) would be biased and, thus, unfair. Test fairness was investigated for the variables gender, the number of books at home (as a proxy for socio-economic status), and migration background (see Pohl & Carstensen, 2012a, for a description of these variables). In order to test for measurement invariance, differential item functioning (DIF) is estimated using a multi-group IRT model, in which main effects of the subgroups as well as differential effects of the subgroups on item difficulty are estimated. Differences in the estimated item difficulties between the subgroups are evaluated. Based on experiences with preliminary data, we considered absolute differences in estimated difficulties that are greater than 1 logit as very strong DIF, absolute differences between 0.6 and 1 noteworthy to further investigate, and differences smaller than 0.4 as no considerable DIF. Additionally, model fit was investigated by comparing a model including DIF to a model that only includes main effects and no DIF.

The competence data in NEPS are scaled using the partial credit model (1PL), in which Rasch-homogeneity is assumed. The partial credit model was chosen because it preserves the weighting of the different aspects of the framework intended by the test developers (Pohl & Carstensen, 2012a). Nevertheless, Rasch-homogeneity is an assumption that may not hold for empirical data. We, therefore, checked for deviations from a uniform discrimination by estimating item discrimination with the generalized partial credit model (2PL; Muraki, 1992) using the software mdltm (von Davier, 2005), and by comparing model fit indices of the 2PL model to those obtained when applying the partial credit model.

The science test is constructed to measure a unidimensional science literacy score (Hahn et al., 2013). The assumption of unidimensionality was, nevertheless, tested in the data by specifying a two dimensional model with KAS items representing one and KOS the other dimension. The correlation between the subdimensions as well as differences in model fit between the unidimensional model and the two dimensional model were used to evaluate the unidimensionality of the scale.

5. Results

5.1 Exclusion of cases from the analyses

The original data file included 16,425 persons. In an initial step for calculating item parameters, all persons who took part in the test were included ($n=4,417$). For further

analyses, only persons with more than two valid responses were taken into account ($n=4,417$), which means that all persons were included in the analyses. The analyses are based on 24 items which remained after eliminating five items with insufficient item quality.

5.2 Descriptive statistics of the responses

In order to a) get a first rough descriptive measure of item difficulty and b) check for possible estimation problems before performing IRT-analyses, we evaluated the relative frequency of the responses given. The percentage of persons correctly responding to an item (relative to all valid responses) varies over items from 20.1% to 64.4% for the MC items. For the CMC items, the percentage of persons who correctly answered all subtasks varies from 15.1% to 65.0%. From a descriptive point of view, the items cover a relatively wide range of difficulties. However, there are no very easy items as the majority of items show a medium or high difficulty. The mean item difficulty of 0.06 ($SD=0.87$) matches the mean person ability (fixed at zero).

5.3 Missing responses

5.3.1 Missing responses per person

The number of not-valid responses per person is shown in Figure 1. The number of not-valid responses is very small. For 79.6 % of the persons, all answers were valid.

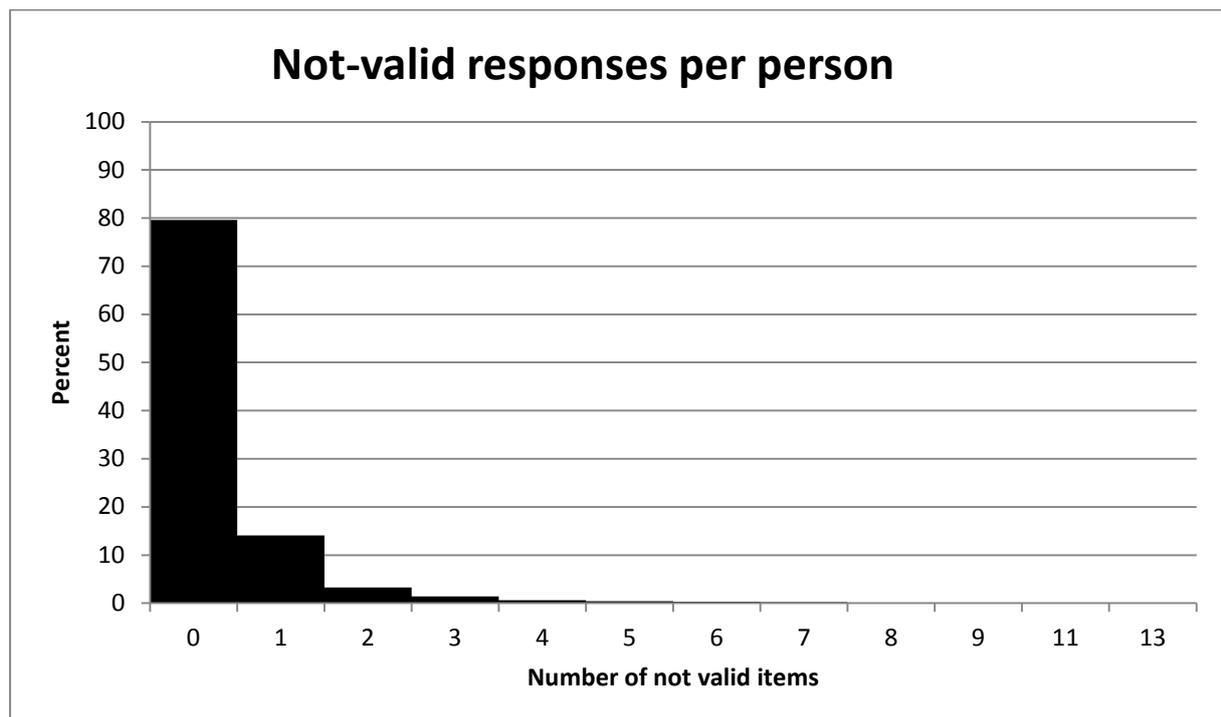


Figure 1: Number of not-valid responses

The number of omitted responses per person is depicted in Figure 2. 82.3 percent of the persons did not omit a single item. Only 4.8% omitted 3 or more than 3 items.

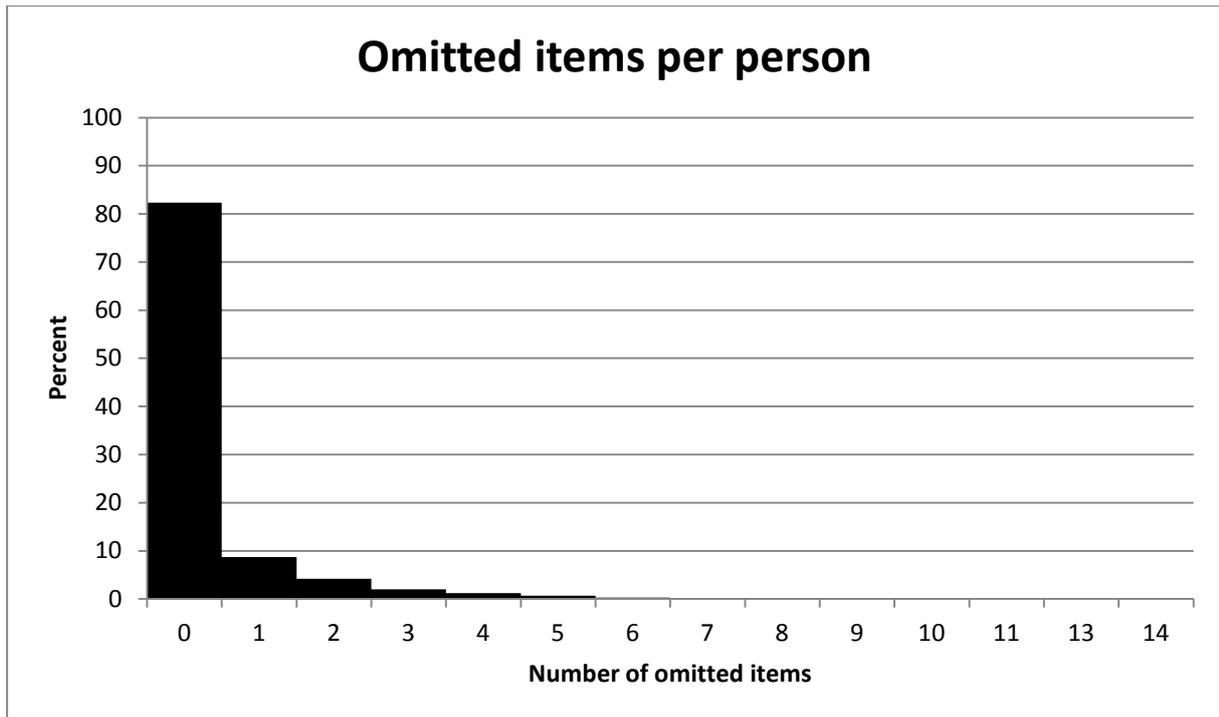


Figure 2: Number of omitted items

Only 39.1% of the students reached the end of the test which might be due to the larger amount of CMC items compared to tests for other cohorts. However, most students managed to finish at least two thirds of the test.

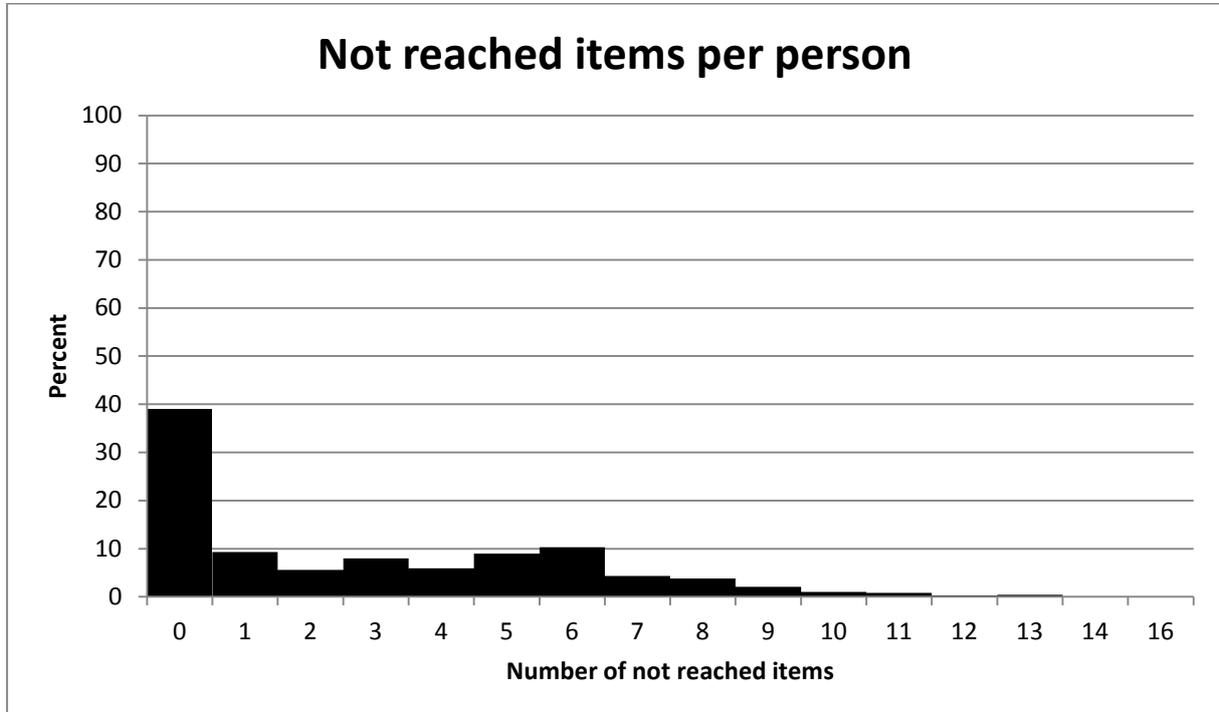


Figure 3: Number of not reached items

Figure 4 shows the total number of missing responses per person. The total number of missing responses is the sum of not-valid, omitted, and not reached missing responses.

26.1% of the students answered all questions and consequently had no missing responses. Only 1.0% of the students have missing responses on more than half of the items which might be due to the comparably large number of CMC items. The amount of missing responses per person can be classified as moderate.

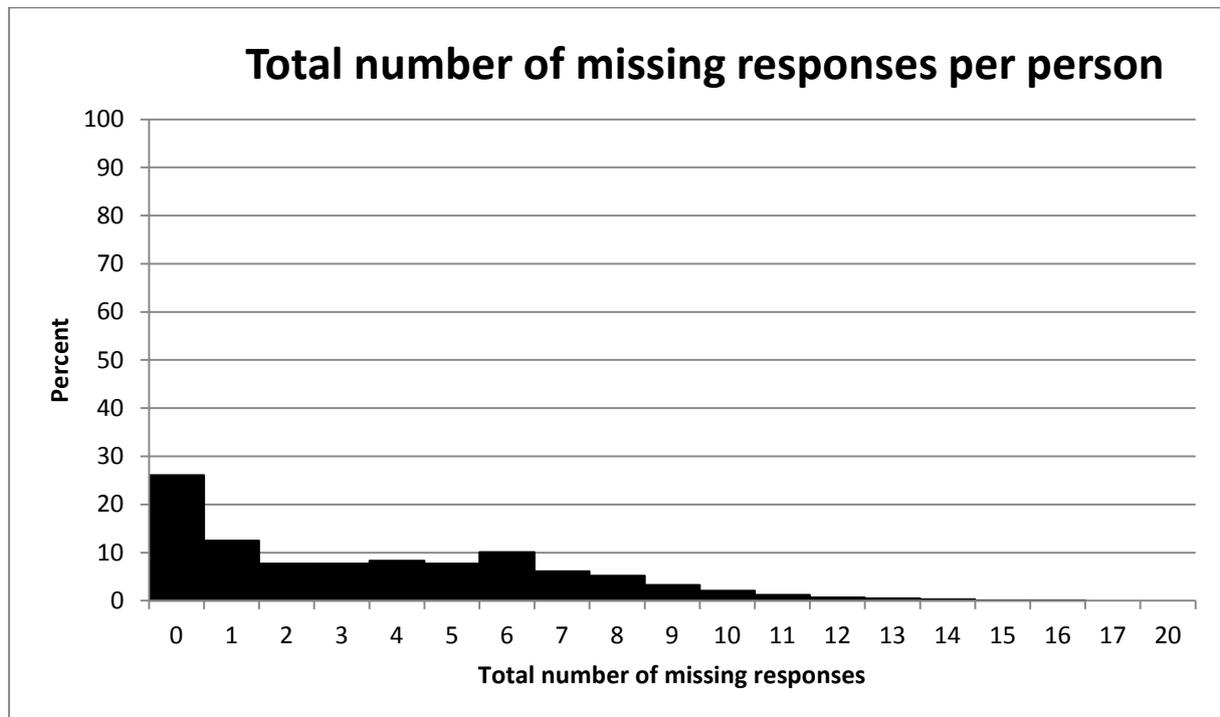


Figure 4: Total number of missing responses

5.3.2 Missing responses per item

Table 4 shows the number of valid responses for each item as well as the number and percentage of missing responses. Overall, the number of persons that omit an item is small. There is no item with an omission rate above 4.6%. The number of missing responses is correlated to .20 with the difficulty of the item. This result indicates that the test takers tend to omit items that are more difficult. The number of invalid responses per item is small. The highest number is 3.8% for item scgb6510_c. The relative frequency of not reached items increases towards the end of the test. Eventually, 60.9% of the students did not reach the last item and thus did not complete the test. The total number of missing responses per item varies between 0.4% and 63.4%. This shows that the test has a slight speed component.

Table 4: Valid responses and missing values

Serial No.	Variable name	Number of valid responses	Position in the test	Relative frequency of not reached items %	Relative frequency of omitted items %	Relative frequency of invalid responses %
1	scgb6420_c	4,347	2	0.0	0.7	0.8
2	scgb0620_c	4,199	3	0.0	4.6	0.4
3	scgb0630_c	4,300	4	0.0	1.8	0.9
4	scgb012s_c	4,262	6	0.0	1.2	2.3
5	scgb083s_c	4,303	7	0.0	0.9	1.7
6	scgb0720_c	4,399	8	0.0	0.3	0.1
7	scgb032s_c	4,349	9	0.0	0.3	1.2
8	scgb0330_c	4,306	10	0.0	0.6	1.9
9	scgb6510_c	4,238	11	0.0	0.2	3.8
10	scgb652s_c	4,312	12	0.0	0.2	2.1
11	scgb602s_c	4,371	13	0.2	0.3	0.6
12	scgb0510_c	4,248	14	0.5	1.9	1.4
13	scgb0520_c	4,300	15	0.9	1.2	0.6
14	scgb0540_c	4,215	16	1.6	2.1	0.8
15	scgb123s_c	4,138	17	2.6	0.5	3.3
16	scgb102s_c	4,065	18	4.7	1.3	2.0
17	scgb021s_c	3,839	19	8.5	2.8	1.8
18	scgb022s_c	3,647	20	12.9	3.3	1.2
19	scgb112s_c	3,264	22	23.1	1.1	1.9
20	scgb6210_c	2,783	23	32.1	3.6	1.3
21	scgb622s_c	2,543	24	38.0	2.4	2.0
22	scgb6320_c	2,233	25	46.1	3.0	0.3
23	scgb0930_c	1,998	26	51.6	2.9	0.2
24	scs3131s_c	1,616	28	60.9	1.2	1.3

Remark. The numbers left out in the column *position in the test* show the positions of the test where the eliminated items were located.

Table 5: Item parameters

Serial No.	Item	Difficulty/location parameter	SE (difficulty/location parameter)	Weighted MNSQ	Weighted t-value	Pt.bis of correct response	Discrimination (2PL)
1	scgb6420_c	-0.391	0.032	1.03	3.0	0.31	0.77
2	scgb0620_c	0.087	0.033	0.97	-2.9	0.41	1.34
3	scgb0630_c	0.171	0.032	0.98	-2.0	0.41	1.30
4	scgb012s_c	0.390	0.033	0.94	-5.6	0.48	1.83
5	scgb083s_c	-0.879	0.031	0.94	-3.7	0.38	0.92
6	scgb0720_c	-0.528	0.033	0.99	-0.8	0.38	1.18
7	scgb032s_c	-1.170	0.041	1.01	0.5	0.26	0.49
8	scgb0330_c	-0.273	0.032	0.98	-2.4	0.43	1.51
9	scgb6510_c	-0.608	0.033	1.01	1.2	0.36	1.02
10	scgb652s_c	-0.147	0.035	0.98	-1.2	0.30	0.72
11	scgb602s_c	-0.661	0.035	1.06	3.5	0.24	0.29
12	scgb0510_c	-0.407	0.033	1.00	0.1	0.38	1.15
13	scgb0520_c	-0.654	0.033	1.01	1.0	0.35	1.04
14	scgb0540_c	1.516	0.040	1.02	1.1	0.27	0.78
15	scgb123s_c	-0.406	0.029	1.06	2.9	0.25	0.36
16	scgb102s_c	1.403	0.039	0.98	-0.9	0.36	1.40
17	scgb021s_c	1.329	0.040	0.99	-0.6	0.33	1.12
18	scgb022s_c	1.895	0.048	1.01	0.3	0.25	0.86
19	scgb112s_c	-0.678	0.039	1.06	3.9	0.26	0.56
20	scgb6210_c	0.024	0.040	0.98	-1.9	0.42	1.30
21	scgb622s_c	1.458	0.051	1.03	1.1	0.23	0.67
22	scgb6320_c	0.619	0.046	0.99	-0.6	0.39	1.25
23	scgb0930_c	-0.618	0.049	0.97	-1.8	0.42	1.46
24	scs3131s_c	0.031	0.068	0.98	-0.7	0.32	0.65

5.4 Parameter estimates

5.4.1 Item parameters

In the end, 24 of the original 29 items (including all subtasks for the polytomous items) were included in the analyses. Five items had to be removed due to insufficient item qualities. The estimated item difficulties for polytomous variables (CMC items) and location parameters for dichotomous variables (MC items) are listed in Table 5. The step parameters (for polytomous variables) are depicted in Table 6. For three of the twelve CMC items (items scgb083s_c, scgb652s_c and scgb602s_c), the two lowest categories were collapsed. Furthermore, for two of the twelve CMC items (scgb032s_c and scs3131s_c), the three lowest categories were collapsed. As these items were CMC items with a maximum score of 2, these items were scaled using the following intervals 0, 0.5, 1 and 1.5 or 0, 0.5, and 1 respectively. CMC-items 14, 16, 17, 18, 19 and 21 were reduced to a 0 and 1 scoring since they showed a decrease in one or two of their step parameters instead of an increase.

Table 6: Step parameters for the CMC items

Item	Step 1 (SE)	Step 2 (SE)	Step 3 (SE)	Step 4 (SE)
scgb083s_c	-0.400 (0.031)	0.277 (0.035)	0.124	
scgb032s_c	-0.059 (0.033)	0.059		
scgb652s_c	-1.230 (0.034)	0.520 (0.034)	0.710	
scgb602s_c	-1.307 (0.033)	0.420 (0.032)	0.888	
scgb123s_c	-0.794 (0.037)	-0.665 (0.032)	0.888 (0.039)	0.571
scs3131s_c	-0.461 (0.050)	0.461		

For estimating item difficulties, the mean of the ability distribution was constrained to be zero. The estimated item difficulties (or location parameters for polytomous variables) vary between -1.17 (scgb032s_c) and 1.90 (scgb022s_c) with a mean of 0.06. Due to the large sample size, the standard error of the estimated item difficulties is very small, $SE(\beta) \leq 0.07$. Overall, the test items are rather difficult. Except for item scgb032s_c, the test lacks items with < -1 logits.

5.4.2 Person parameters

Person parameters are estimated as WLEs and PVs (Pohl & Carstensen, 2012a). WLEs will be provided in the first release of the SUF. PVs will be provided in later analyses. A description of the data in the SUF can be found in section 7. An overview of how to work with competence data is given in Pohl and Carstensen (2012a).

5.4.3 Test targeting and reliability

Test targeting was investigated in order to evaluate the measurement precision of the estimated ability scores and to judge the appropriateness of the test for the specific target population. In the analyses, the mean of the ability is constrained to be zero. The variance was estimated to be 0.475 and is relatively small since most of the students in the sample

were in 11th grade of high school (“Gymnasium”) or an integrative school “integrative Gesamtschule”), respectively, when they were tested.

The reliability of the test (WLE reliability = .635) was moderate. The amount to which the item difficulties and location parameters are targeted to the ability of the persons is shown in Figure 5. The figure shows that the items cover a great range of the ability distribution of the persons. However, there is a lack of items covering the medium to above medium person abilities and there are not enough items available for persons with low science ability. Instead, the majority of items are easy or of medium difficulty. As a consequence, persons with a medium and rather high ability will be measured relatively precisely with a low standard error while ability estimates for students with medium to high science ability or low science ability will have a larger standard error.

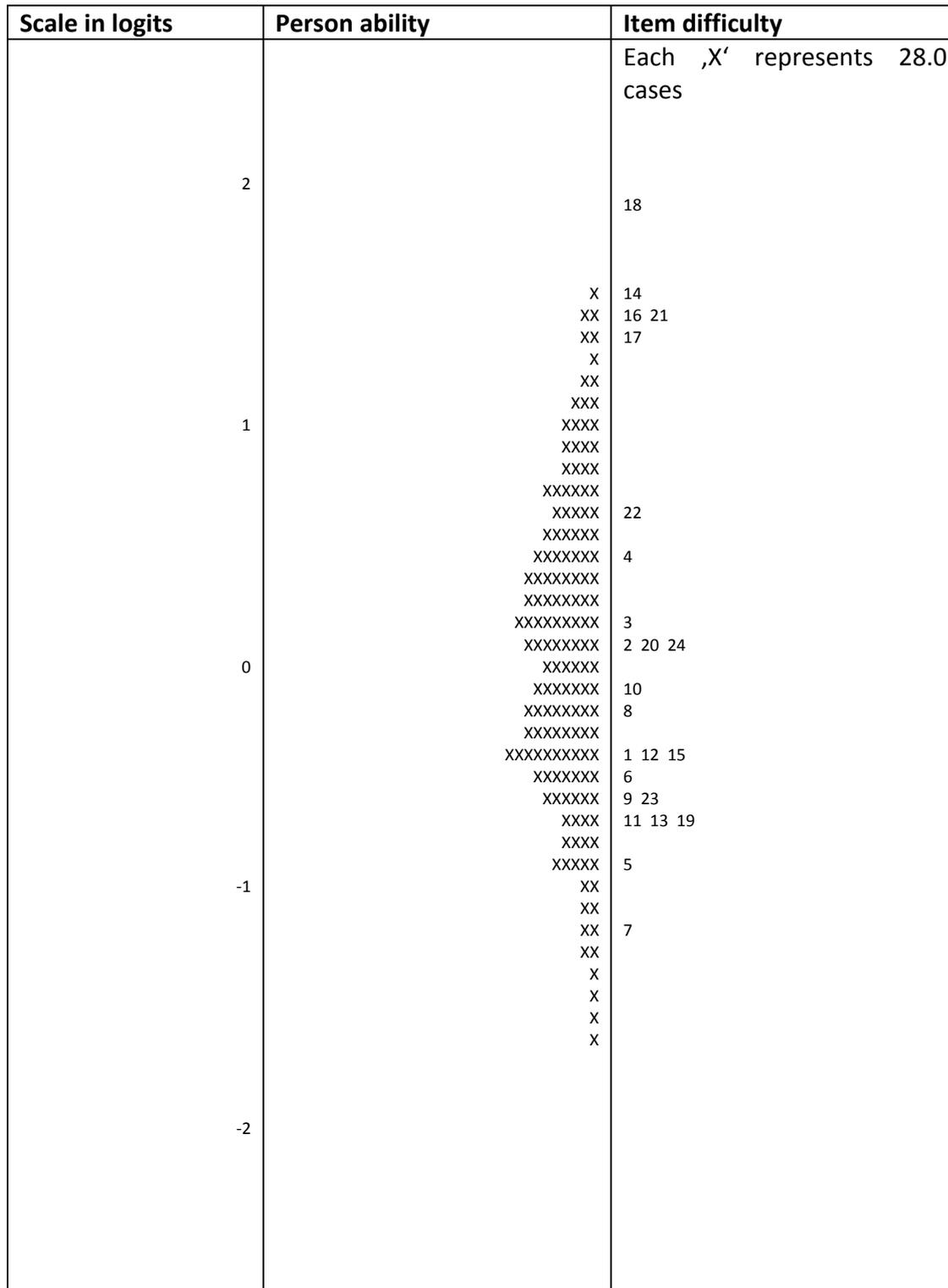


Figure 5: Test targeting. The distribution of person ability in the sample is depicted on the left side of the graph. Each 'X' represents 28.0 cases. The difficulty of the items is depicted on the right side of the graph. Each number represents an item.

5.5 Quality of the test

5.5.1 Fit of the subtasks of complex multiple-choice items

Before the responses on the subtasks of CMC items were aggregated and analyzed via a partial credit model, the fit of the subtasks was checked by analyzing the single subtasks together with the simple MC items in a Rasch model.

No estimation problems occurred and all subtasks showed a satisfactory item fit. The WMNSQs ranged from 0.94 to 1.06, the respective t -values from -5.6 to 3.9. There were no unacceptable deviations of the empirical estimated probabilities from the model-implied item characteristic curves. Hence, an aggregation of polytomous variables seemed to be justified.

In addition to the overall item fit, we specifically investigated how well the distractors performed in the test by evaluating the point-biserial correlation between each incorrect response (distractor) and the students' total scores. All distractors had point-biserial correlations with the total score below zero. The results indicate that the distractors worked well.

5.5.2 Item fit

Regarding the MC and the aggregated CMC items, the fit was very good. WMNSQs were close to 1 with the lowest value being 0.94 (items scgb012s_c and scgb083s_c) and the highest being 1.06 (items scgb602s_c, scgb123s_c and scgb112s_c). Overall, there were no items with a WMNSQ above 1.1. However, there was one item with a t -value above 4.0. But, the item characteristic curve of this item showed a reasonable fit. Hence, no indications for a heavy misfit of the item could be detected and, therefore, it was kept in the analysis for estimating the scientific literacy scores.

5.5.3 Differential item functioning

We checked for test fairness for different groups (i.e., measurement invariance) by estimating the amount of DIF. DIF was investigated for the variables gender, the number of books at home (as a proxy for socio-economic status), migration background, and school type (see Pohl & Carstensen, 2012a, for a description of these variables). Table 7 shows the difference between the estimated item difficulties in different groups. Male vs. female, for example, indicates the difference in difficulty $\beta(\text{male}) - \beta(\text{female})$. A positive value indicates a higher difficulty for males, a negative value a lower difficulty for males as opposed to females. Since the science test was the only competency test administered in grade 11, DIF was not and could not be checked for test rotation.

DIF was investigated for gender. 2,416 (54.7%) of the test takers were female and 2,001 (45.3%) were male. On average, male students showed slightly higher scores in scientific literacy than female students (main effect = 0.362 logits, Cohen's $d = -0.545$). There is no item with a considerable gender DIF. The highest difference in difficulties between the two groups is 0.284 logits.

The number of books at home was used as a proxy for socio-economic status. There were 879 (19.9%) test takers with 0 to 100 books at home, 3,305 (74.8%) test takers with more than 100 books at home, and 233 (5.3%) test takers did not give a valid response. DIF was investigated using these three groups. There are considerable average differences between the three groups. Participants with 100 or less books at home on average show a 0.362 logits (Cohen's $d = -0.537$) lower scientific literacy score than participants with more than 100 books. Participants without a valid response on the variable 'books at home' performed 0.144 logits (Cohen's $d = -0.203$) higher than participants with up to 100 and 0.220 logits (Cohen's $d = 0.331$) lower than participants with more than 100 books at home, respectively. There is no considerable DIF comparing participants with many or fewer books (highest DIF = 0.152). Comparing the group without valid responses to the two groups with valid responses, DIF occurs up to 0.271 logits. This is a rather large difference, which may, however, also be the result of the uncertainty in estimation due to the small number of persons with missing responses.

There were 3,359 (76.0%) participants without a migration background, whereas 720 (16.3%) of the participants had a migration background (for 1.6% students neither their mother, father, nor they, themselves, were born in Germany, for 6.7% only the participants were born in Germany and both of their parents were born abroad, for 8.0% of the students only one of their parents was born abroad). 338 (7.7%) students could not be allocated to either group. These three groups were used for investigating DIF for migration. There is a considerable difference in the average performance of participants with or without migration background (main effect = -0.328 logits, Cohen's $d = 0.482$). Participants without a migration background have a higher scientific literacy than participants with a migration background. Also, students without a migration background differ from those with an unknown background on migration (main effect = -0.184 logits, Cohen's $d = 0.274$). However, there was no considerable difference between students with a migration background and those with an unknown background on migration (main effect = 0.144 logits, Cohen's $d = -0.197$). There is no considerable DIF concerning the migration status.

DIF was also investigated for school type. 3,921 (88.8%) of the test takers were high-school students and 496 (11.2%) visited an integrative school. On average, high-school students have a higher scientific literacy score than students who attend an integrated school (main effect = 0.574 logits, Cohen's $d = -0.856$). There is no considerable DIF concerning the school type.

Besides investigating DIF for each single item, an overall test for DIF was performed by comparing models which allow for DIF with those that allow only for main effects. In Table 8, the models including only main effects are compared with those that, additionally, estimate DIF. Akaike's (1974) information criterion (AIC) and the Bayesian information criterion (BIC, Schwarz, 1978) were used for assessing the models. Using the AIC, the models estimating DIF are favored for all four DIF variables. The BIC takes the number of estimated parameters into account and, thus, prevents from overparameterization of models. Using BIC, the more parsimonious model including only the main effect is preferred over the more complex DIF model for the number of books at home, the migration background, and the school type. For the DIF variable gender the more complex DIF model has slightly better information criteria. There is slight DIF in favor of male persons.

Table 7: Differential item functioning (absolute differences between difficulties)

Item	Gender Male vs. female	School type High school vs. Others	Books			Migration status		
			<100 vs. >100	<100 vs. Missing	>100 vs. Missing	Without vs. With	Without vs. Missing	With vs. Missing
scgb6420_c	-0.209	-0.047	-0.026	-0.060	-0.033	-0.139	-0.016	0.125
scgb0620_c	0.140	-0.173	0.000	-0.006	-0.006	0.019	0.044	0.026
scgb0630_c	-0.062	0.149	-0.010	0.024	0.034	0.041	0.056	0.016
scgb012s_c	0.148	0.034	0.038	-0.019	-0.057	-0.069	-0.086	-0.016
scgb083s_c	0.284	0.057	0.040	0.002	-0.040	-0.045	-0.002	0.044
scgb0720_c	-0.026	-0.001	0.096	0.271	0.174	-0.079	0.029	0.110
scgb032s_c	-0.106	0.096	0.053	-0.025	-0.077	0.026	-0.052	-0.075
scgb0330_c	0.223	0.108	0.152	0.146	-0.006	-0.038	-0.036	0.004
scgb6510_c	-0.255	0.138	0.005	-0.056	-0.061	-0.001	-0.083	-0.081
scgb652s_c	0.146	0.106	0.107	0.184	0.085	0.021	-0.032	-0.052
scgb602s_c	-0.238	0.036	-0.049	-0.119	-0.064	0.097	-0.113	-0.221
scgb0510_c	-0.145	-0.064	-0.088	-0.203	-0.114	0.077	-0.031	-0.109
scgb0520_c	0.079	0.004	-0.035	-0.003	0.032	0.020	0.063	0.044
scgb0540_c	-0.043	-0.055	-0.145	-0.094	0.051	0.082	0.010	-0.071
scgb123s_c	-0.029	-0.156	-0.105	-0.064	0.041	0.170	0.133	-0.033
scgb102s_c	0.102	0.112	-0.005	-0.058	-0.051	-0.019	0.035	0.055
scgb021s_c	0.132	-0.052	0.023	0.065	0.043	-0.143	-0.076	0.068
scgb022s_c	0.059	-0.087	-0.041	-0.016	0.025	0.092	0.076	-0.015
scgb112s_c	-0.203	-0.259	-0.140	-0.158	-0.018	-0.038	-0.088	-0.050
scgb6210_c	-0.203	-0.038	0.039	-0.032	-0.071	-0.059	0.006	0.067
scgb622s_c	-0.108	0.130	-0.074	0.067	0.141	0.004	-0.008	-0.011
scgb6320_c	0.185	-0.189	0.013	0.111	0.097	0.016	0.174	0.158
scgb0930_c	0.045	0.106	0.093	0.067	-0.027	-0.013	0.080	0.093
scs3131s_c	0.120	0.075	0.052	0.126	0.069	-0.111	0.022	0.164

Table 8: Comparison of models with and without DIF

DIF variable	Model	Deviance	Number of parameters	AIC	BIC
Gender	main effect	138,332.799	37	138,504.67	138,467.67
	DIF	137,875.563	61	138,158.92	138,097.92
Books	main effect	138,401.555	38	138,578.07	138,540.07
	DIF	138,318.013	86	138,717.49	138,631.49
Migration	main effect	138,440.770	38	138,617.28	138,579.28
	DIF	138,362.572	86	138,762.05	138,676.05
School type	main effect	138,336.660	37	138,508.53	138,471.53
	DIF	138,258.993	61	138,542.35	138,481.35

5.5.4 Rasch-homogeneity

In order to test for the assumption of Rasch-homogeneity, 24 of the original 29 items entered the analysis with the generalized partial credit model (2PL) to test for Rasch-homogeneity. The estimated discrimination parameters are depicted in the last column in Table 5. They range from 0.29 (item scgb602s_c) to 1.83 (item scgb012s_c). The discriminations differ considerably among the items. However, the 2PL model (BIC = 138,680.79, number of parameters = 70) fits the data slightly worse than the 1PL (BIC = 138,535.01, number of parameters = 36).

5.5.5 Unidimensionality of the test

The unidimensionality of the test was investigated by specifying an onedimensional and a twodimensional model.

The first model is based on the assumption that scientific literacy is a onedimensional construct that measures one distinct competence whereas the second model distinguishes between the two sub-competencies, knowledge about science and knowledge of science (for more details see Hahn et al., 2013). For estimating, a twodimensional model based on the Gauss Hermite quadrature estimation implemented in ConQuest was used (n=30 nodes were chosen so that stable parameter estimations could be obtained). The twodimensional model (BIC= 154,095.14, number of parameters = 38) fits the data less well than the unidimensional model (BIC= 138,535.01, number of parameters = 36; correlations of the two dimensions: 0.90). Consequently, scientific literacy as measured by this test can be regarded as unidimensional and, therefore, this simpler model was used for estimating competence scores.

6. Discussion

The analyses in the previous sections aimed at providing information on the quality of the science test in grade 11 and at describing how the scientific literacy score is estimated.

The amount of invalid responses and not-reached items is moderate. However, some items show higher omission rates, although, in general, the amount of omitted items is acceptable.

The test has a moderate reliability (WLE reliability = 0.635) and distinguishes well between test takers of average and high scientific literacy, but not as well for medium high and low performers. Very easy items are missing; hence, test targeting is somewhat suboptimal and the test measures scientific literacy of low- and medium high -performing students less accurately. This is depicted by the test's variance (= 0.478) which, ideally, should be higher, but which is also due to the limited variance of the sample itself which mainly consists of high school students.

Indicated by various fit criteria – WMNSQ, t-value of the WMNSQ, ICC – the items exhibit a good item fit. Also, discrimination values of the items (either estimated in a 2PL model or as a correlation of the item score with total score) are acceptable. Different variables were used for testing measurement invariance. No considerable DIF became evident for any of these variables, indicating that the test is fair to the considered subgroups.

A unidimensional partial credit model yielded a better model fit than a twodimensional partial credit model (between-item-multidimensionality, the dimensions being the content areas). Hence, the unidimensional model was used for estimating scientific literacy scores.

Summarizing the results, the test has good psychometric properties that facilitate the estimation of a unidimensional scientific literacy score.

7. Data in the Scientific Use file

There are 24 items in the data set that are either scored as dichotomous variables (MC or SCR items) with 0 indicating an incorrect response and 1 indicating a correct response, or scored as a polytomous variable (CMC items) indicating the (partial) credit. The dichotomous variables are marked with a '0_c' at the end of the variable name, the CMC items are marked with a 's_c' at the end of the variable name. Note that the value of the polytomous variable does not necessarily indicate the number of correctly responded subtasks (see section 4.2 aggregation of CMC items). In the scaling model, each category of CMC items is scored with 0.5 points. Manifest scale scores are provided in form of WLE estimates (scg11_sc1u) including the respective standard error (scg11_sc2u). WLE estimates for longitudinal analyses are not yet provided but will be prepared for a later data release.

Please note that when categories of the polytomous variables had less than 2% of the sample, the categories were collapsed. For the science test, this concerned the two or three lowest categories of two polytomous items (see section 5.4.) on the aggregation of the CMC items. In the scaling model, the collapsed polytomous items are scored in steps of 0, 0.5, 1.0 and 1.5 or 0, 0.5 and 1.0 (denoting the highest). The ConQuest Syntax for estimating the WLE scores from the items is provided in Appendix A. Students that did not take part in the test or those that do not have enough valid responses to estimate a scale score will have a non-determinable missing value on the WLE score for scientific literacy.

Plausible values that allow investigating latent relationships of competence scores with other variables will be provided in later data releases. User interested in investigating latent relationships may alternatively either include the measurement model in their analyses or estimate plausible values themselves. A description of these approaches can be found in Pohl and Carstensen (2012a).

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19, 716-722.
- Haberkorn, K., Pohl, S., Carstensen, C., & Wiegand, E. (2012). Incorporating different response formats in the IRT-scaling model for competence data. Manuscript submitted for publication.
- Haberkorn, K., Pohl, S., Hardt, K., & Wiegand, E. (2012). Technical Report of Reading– Scaling Results of Starting Cohort 4 in Ninth Grade (NEPS Working Paper No. 16). Bamberg: Otto-Friedrich-Universität, Nationales Bildungspanel.
- Hahn, I., Schöps, K., Rönnebeck, S., Martensen, M., Hansen, S., Saß, S., Dalehefte, I.M., & Prenzel, M. (2013). Assessing scientific literacy over the lifespan—A description of the NEPS science framework and the test development. *Journal for Educational Research Online/Journal für Bildungsforschung Online*, 5(2), 110-138.
- Masters (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47, 149-174.
- Mislevy, R.J. (1991). Randomization-based inference about latent variables from complex samples. *Psychometrika*, 56, 177-196.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, 16, 159-176.
- OECD (2007). PISA 2006: Science Competencies for Tomorrow's World: Volume 1: Analysis, PISA, OECD Publishing.
- Pohl, S. & Carstensen, C. H. (2012a). NEPS technical report – Scaling the data of the competence tests. (NEPS Working Paper No. 14). Bamberg: Otto-Friedrich-Universität, Nationales Bildungspanel.
- Pohl, S. & Carstensen, C. H. (2013). Scaling the competence tests in the National Educational Panel Study – Many questions, some answers, and further challenges. *Journal for Educational Research Online/Journal für Bildungsforschung Online*, 5(2), 189-216.
- Pohl, S., Haberkorn, K., Hardt, K., & Wiegand, E. (2012). Technical Report of Reading– Scaling Results of Starting Cohort 3 in Fifth Grade (NEPS Working Paper No. 15). Bamberg: Otto-Friedrich-Universität, Nationales Bildungspanel.
- Prenzel, M., Schöps, K., Rönnebeck, S., Senkbeil, M., Walter, O., Carstensen, C., & Hamann, M. (2007). Naturwissenschaftliche Kompetenz im internationalen Vergleich. In PISA-Konsortium Deutschland (Hrsg.), PISA 2006 - Die Ergebnisse der dritten internationalen Vergleichsstudie (S. 63-105). Münster, Waxmann.
- Schwarz, G. E. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6 (2), 461–464.
- von Davier, M. (2005). A general diagnostic model applied to language testing data (ETS Research Rep. No. RR-05-16). Princeton, NJ: ETS.

- Warm T.A. (1989). Weighted Likelihood Estimation of Ability in Item Response Theory. *Psychometrika*, 54, 427-450.
- Weinert, S., Artelt, C., Prenzel, M., Senkbeil, M., Ehmke, T., & Carstensen C.H. (2011) Development of Competencies Across the Life Span. In H. P. Blossfeld, H. G. Roßbach & J. v. Maurice & (Eds.). Education as a Lifelong Process: The German National Educational Panel Study (NEPS). *Zeitschrift für Erziehungswissenschaft, Sonderheft 14* . Wiesbaden: VS Verlag für Sozialwissenschaften.
- Wu, M.L., Adams, R.J., & Wilson, M.R. (1997). ACER Conquest: Generalised item response modelling software. Melbourne: ACER Press.

Appendix

Appendix A: ConQuest-Syntax for estimating WLE estimates in starting cohort 4

Title Starting Cohort 4, SCIENCE: Partial Credit Model;

data filename.dat;

format pid 1-7 responses * /* insert number of columns with data*/

labels << filename_with_labels.txt;

codes 0,1,2,3,4;

score (0,1) (0,1) !item (1-4,6,8,9,12-14,16-23);

score (0,1,2,3,4) (0,0.5,1,1.5,2) !item (15);

score (0,1,2,3) (0,0.5,1,1.5) !item (5,10,11);

score (0,1,2) (0,0.5,1) !item (7,24);

set constraint=cases;

model item + item*step;

estimate;

show !estimates=latent >> filename.shw;

itanal >> filename.ita;

show cases !estimates=wle >> filename.wle;