

# How much is too much? Item Response Theory procedures to shorten tests

Ottavia M. Epifania<sup>1</sup> and Livio Finos<sup>2</sup>

<sup>1</sup> Department of Psychology and Cognitive Science, University of Trento, IT  
`ottavia.epifania@unitn.it`

<sup>2</sup> Université de Paris-Sud, Laboratoire d'Analyse Numérique, Bâtiment 425,  
F-91405 Orsay Cedex, France

**Abstract.** The abstract should summarize the contents of the paper using at least 70 and at most 150 words. It will be set in 9-point font size and be inset 1.0 cm from the right and left margins. There will be two blank lines before and after the Abstract. ...

**Keywords:** Item response theory, careless error, information functions, short test forms

## 1 Introduction

As a general rule of thumb, the higher the number of items in a test, the better the measurement in terms of validity and reliability. However, there is a trade-off between the number of administered items and the response quality. As such, the trade-off between the number of administered items and the tiredness of the respondents should be kept in mind to obtain reliable and precise measurement tools. Item Response Theory (IRT) provides an ideal framework for shortening existing tests (or for developing tests from item banks) given the detailed information that they provide of the measurement precision of each item with respect to different levels of the latent trait. In this contribution, we present a new IRT-based algorithm for developing short test forms (STFs) from existing tests, denoted as Léon. This algorithm accounts for the tiredness of the respondents during the item inclusion process, such that it attempts at minimizing the number of selected items while accounting for the tiredness of the respondents in order to maximize the measurement precision (as expressed by the test information function, TIF) of the STF. The ability of Léon of developing STFs able to approximate the TIF that would be obtained by administering all the items under the assumption that respondents would never get tired is investigated in a simulation study. Specifically, Léon's ability is compared against that of another algorithm for developing STFs, which does not account for the tiredness of the respondents.

## 2 Item Response Theory and Information Functions

In IRT models for dichotomous responses (e.g., correct vs. incorrect), the probability of observing a correct response on item  $i$  by person  $p$  depends on both the

characteristics of the respondent (as described by their latent trait level,  $\theta_p$ ) and on the characteristics of the item, which can be described by different parameters. IRT models differentiate according to the number of parameters used for describing the characteristics of the items. According to the 4-parameter logistic model (4-PL), the probability of a correct response can be formalized as:

$$eq : 4plP(x_{pi} = 1 | \theta_p, b_i, a_i, c_i, d_i) = c_i + (d_i - c_i) + \frac{\exp[a_i(\theta_p - b_i)]}{1 + \exp[a_i(\theta_p - b_i)]}, \quad (1)$$

where  $\theta_p$  is the latent trait level of person  $p$ ,  $b_i$  is the location of the item on the latent trait (i.e., difficulty parameter, the higher the value, the higher the difficulty of the item),  $a_i$  describes the ability of  $i$  to discriminate between respondents with different latent trait levels (i.e., discrimination parameter, the higher the value, the higher the discrimination ability of the item), and  $c_i$  and  $d_i$  describe the probability of observing a correct response when  $\theta \rightarrow -\infty$  and  $\theta \rightarrow +\infty$ , respectively. When  $\theta \rightarrow -\infty$ , the probability of observing a correct response should tend to 0. Likewise, when  $\theta \rightarrow +\infty$ , the probability of observing a correct response should tend to 1. However, there might be instances where respondents with  $\theta$  levels below the difficulty of the item, whom are hence expected not to respond correctly, might provide the correct response out of luck. The lucky guess parameter  $c_i$  describes this probability, such that the probability of observing a correct response for  $\theta \rightarrow -\infty$  tends to  $c_i$  instead of 0. The same but inverse consideration applies when  $\theta \rightarrow +\infty$ . The  $d_i$  parameter describes the probability of not endorsing the item given that the latent trait is above the location of the item, such that the probability of observing a correct response for  $\theta \rightarrow +\infty$  tends to  $d_i$  instead of 1.

By constraining  $\forall i \in B, d_i = 1$  (where  $B$  is the set of items in a test) in equation ??, the 3-Parameter logistic (3-PL) model is obtained. From the 3-PL model, the 2-Parameter logistic model (2-PL) is obtained by constraining  $\forall i \in B, c_i = 0$ , and the 1-parameter logistic model (1-PL, equivalent to the Rasch model) is obtained by constraining  $\forall i \in B, a_i = 1$ .

## 2.1 Information Functions

The measurement precision of each item with respect to different levels of the latent trait can be expressed by the so-called *item information functions* (IIFs), which for the 4-PL is formalized as:

$$eq : iifIIF_i = \frac{a_i^2 [P(\theta) - c_i]^2 [d_i - P(\theta)]^2}{(d_i - c_i)^2 P(\theta) Q(\theta)}. \quad (2)$$

The informativeness of each item is strongly influenced by the location of the item on the latent trait with respect to a specific latent trait level  $\theta$ , its discriminativity, and the probability of lucky guess and careless error. In absence of lucky guess and careless error, the IIF reaches its maximum when the location of the item  $b_i$  matches the latent trait level  $\theta$ , and it decreases as the distance

between  $b_i$  and  $\theta$  increases. Moreover, the higher the discrimination of the item  $a_i$ , the more informative the item. However, when the lucky guess and careless error are taken into account, the informativeness of the item decreases, and, most importantly, the maximum of the IIF does not corresponds to the item location on the latent trait and it is generally lower.

The *test information function* (TIF) is the sum of the IIFs, such that its shape (i.e., its informativeness with respect to different levels of the latent trait) and its height (i.e., the amount of information for different levels of the latent trait) depend on the distribution of the items along the latent trait. The more the items with high discrimination and low lucky guess and careless error parameters are spread throughout the latent trait, the more the test would be informative of different regions. The more the items have lucky guesses and careless errors, the less the TIF.

### 3 item Selection Algorithms

#### 3.1 Frank

Frank considers the entire latent trait for the item selection, in that it selects the item whose IIF is best able to reduce the distance from the  $TIF^*$  along the entire latent trait, as follows:

At  $k = 0$ :  $\forall i \in B$ ,  $IIF_i$ ,  $TIF^0(\theta) = 0 \forall \theta$ ,  $Q^0 = \emptyset$ . For  $k \geq 0$ ,

1.  $A^k = B \setminus Q^k$
2.  $\forall i \in A^k$ ,  $pTIF_i^k = \frac{TIF^k + IIF_i}{|Q^k| + 1}$
3.  $i^* = \arg \min_{i \in A^k} |TIF^* - pTIF_i|$
4. Termination criterion:  $|TIF^* - pTIF_{i^*}| \geq |TIF^* - TIF^k|$ :
  - FALSE:  $Q^{k+1} = Q^k \cup \{i^*\}$ ,  $TIF^{k+1} = pTIF_{i^*}$ , iterates 1-4
  - TRUE: Stop,  $Q_{Frank} = Q^k$

The first operation done by Frank at the beginning ( $k = 0$ ) is to compute the IIFs of all the items in  $B$ , considering the available item parameters. At  $k = 0$ , the subset of items  $Q^0$  is empty and the  $TIF^0$  is 0 for all the  $\theta$  levels. At each iteration  $k$ : (1.) a set of available items is generated as the items in the item bank that have not been included in the STF yet,  $A^k = B \setminus Q^k$ ; (2.) An average provisional TIF,  $pTIF$ , is computed by adding the IIF of each of the items in the set of the available items  $A^k$ , one at the time, to the TIF obtained from the items in  $Q^k$  (The denominator is obtained by adding 1 to the cardinality of  $Q^k$ ); (3.) Among all the items in  $A^k$ , the one that allows for minimizing the distance between  $pTIF$  and  $TIF^*$  is included in  $i^*$ ; (4.) The termination criterion is tested. If the distance between the  $TIF^*$  and the  $pTIF_{i^*}$  is greater than or equal to the distance between the  $TIF^*$  and the  $TIF^k$  (i.e., the TIF obtained from the items in the subset  $Q^k$ , without item  $i^*$ ) (TRUE), then the item  $i^*$  does not contribute in the reduction of the distance from the  $TIF^*$ , the algorithm stops, and the final item selection is the one without the item in  $i^*$ ,  $Q_{Frank} = Q^k$ . Conversely (FALSE), the item in  $i^*$  does contribute in the reduction of the distance from the  $TIF^*$ , hence it is included in the set of items and a new iteration starts,  $Q^{k+1} = Q^k \cup \{i^*\}$ .

## 4 Simulation study

### 4.1 Simulation design

In this study, we operationalized the tiredness of the individuals as the probability of committing careless error ( $1 - d_i$  QUESTA è LA PROBABILITÀ DI NON CARELESS ERROR DEVO SPIEGARE MEGLIO), such that the more the items in a test, the higher the probability of careless error. The careless error probability is intrinsically related to the rank of the item in the test, and it increases as the position of the item in the test increases. It increases exponentially as a function of the number of administered items, such that the first administered item has  $d_i = 1$  (i.e., does not have a probability of observing an incorrect response given the latent trait level  $\theta$ ).

The procedure is replicated 100 times. At each replica, a test of  $B$  of 50 items with difficulty ( $b_i \sim \mathcal{U}(-3, 3)$ ) discrimination ( $a_i \sim \mathcal{U}(.90, 2)$ ) parameters drawn from uniform distributions is generated. Lucky guess and careless error parameters are constant,  $c_i = 0, \forall i \in B$  and  $d_i = 0, \forall i \in B$ . The  $TIF^*$  is generated as the mean TIF of all items in  $B$ , which describes the ideal TIF that one would obtain if respondents are administered with all the items without getting tired.

A new test  $B'$  is generated by keeping the item parameters equal to those in  $B$ , with the only exception of the careless error parameters, which are included as related to the rank of the item in the item bank with an exponential function. The assumption is that the first administered item does not have a careless error probability ( $d_{1st} = 1$ ), but it increases as the administration goes on. As such, the careless error probabilities associated to the second and third administered items are, respectively,  $d_{2nd} = .99$  and  $d_{3rd} = .98$ . In this application, the maximum careless error probability associated to the last administered item is  $d_{50th} = .61$ .

Frank and Léon are applied to generate a STF able to recreate the  $TIF^*$ . Both algorithm are applied considering the items in  $B'$ . While both algorithms are able to consider the careless error parameters in the computation of the IIFs, only Léon actually accounts for the number of items included in the STF.

La differenza è che Frank calcola tutte le iif degli item all'inizio, considerando la careless error legata alla posizione dell'item così come viene generata inizialmente dalla funzione. Léon invece calcola la iif ad ogni iterazione (per questo è più lento) mettendo dentro la careless error che è legata alla posizione dell'item. Nel senso, se Léon per ora ha selezionato 3 item e deve valutare se metterne un quarto, le iif degli item che sono rimasti verrà calcolata considerando la careless error in  $B'$  dell'item somministrato per quarto.

### 4.2 Comparison

The ability of Frank and Léon in reducing the distance from the  $TIF^*$  (i.e., the TIF obtained from all the items in  $B$ ) has been considered. Moreover, the TIF of the items in  $B'$  has been computed as well and compared against those of the STFs resulting from the application of Frank and Léon. The rationale is as

follows. The more the item in a test, the higher the information for different regions of the latent trait. However, administering too many items might be an highly demanding task, such that the respondent might get tired and their response accuracy might decrease during the administration, such that the last administered items might include error variance not related to the construct under investigation. In this light, it should be more convenient to administer less but highly informative items, able to approximate a specific TIF target, than to administer the entire test to obtain precise measurements of the latent trait.

## 5 Results

On average, Léon included more items in the STF than Frank did,  $M_{Leon} = 6.75 \pm 3.06$ ,  $M_{Frank} = 5.29 \pm 2.92$ ,  $t = 3.45$ ,  $df = 197.57$ ,  $p < .001$ . *scrivere meglio* Given that leon account for the tiredness of the respondents during the item selection, it also considers the measurement precision of the STF given that the respondents do get tired during the administration, while Frank does not. As such, while preventing the over inclusion of items in the STF, Léon is also able to account for the measurement prevision, in the attempt of balancing the number of items with the measurement precision.

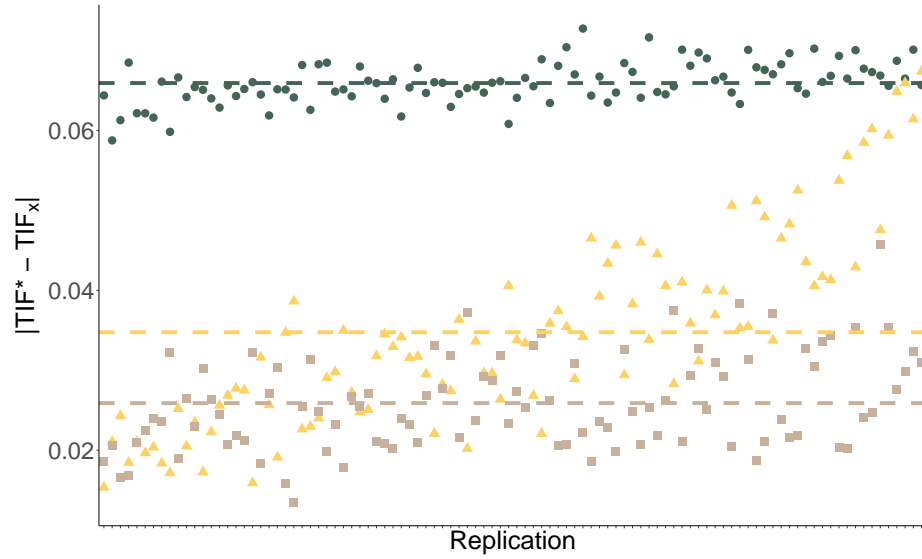
Figure 1 illustrates the distance from  $TIF_B$  of  $TIF_{B'}$  (green dots),  $TIF_{Frank}$  (yellow triangles), and  $TIF_{Leon}$  (brown squares) in each of the 100 replications, along with the average distance across the 100 replications (the horizontal line of the corresponding color).

Overall,  $TIF_{B'}$ s are the most distant from  $TIF_B$ , while  $TIF_{Leon}$ s are the closest ones. Interestingly, the latter ones present more variability around the mean ( $\bar{\Delta}_{Leon} = 0.03 \pm 0.02$ ,  $CV = 0.67$ ) than the former ones ( $\bar{\Delta}_{B'} = 0.07 \pm 0.03$ ,  $CV = 0.43$ ). Between  $TIF_{B'}$  and  $TIF_{Leon}$ ,  $TIF_{Frank}$  is the less consistent in terms of distance from  $TIF_B$ , with ( $\bar{\Delta}_{Frank} = 0.03 \pm 0.03$ ,  $CV = 1.00$ ).

E che cazzo devo dire ancora?

## 6 Final Remarks

This manuscript presented a first attempt at the development of an IRT-based algorithm, denoted as Léon, for the generation of informative and static STF's able to account for the tiredness of the respondents. In the item selection for inclusion in the STF, Léon considers the number of items included up to that iteration and adds a penalty in terms of higher probability of careless error, which increases as the number of items included in the STF increases. Its performance in approximating a TIF target, which is here conceptualized as the TIF that one would obtain considering the entire administration of the test if respondents would never get tired (i.e., without careless error), is compared against that of another IRT-based algorithm, denoted as Frank. Differently from Léon, Frank does not add any penalization for the number of items included in the STF.



**Fig. 1.** Distance from  $TIF_B$  of  $TIF_{B'}$  (green dots),  $TIF_{Frank}$  (yellow triangles), and  $TIF_{Leon}$  (brown squares) in each of the 100 replications. The horizontal green, yellow, and brown lines are the average distance from  $TIF_B$  of  $TIF_{B'}$ ,  $TIF_{Frank}$ , and  $TIF_{Leon}$ , respectively

Finally, the distance between the TIF obtained from the administration of the entire test with and without tiredness has been considered as well.

The results of a simulation study suggest that the administration of fewer items selected also considering the tiredness of the respondents might provide better measurement tools than administering the entire test and tiring out the respondents.

Although the results are promising, there are several limitations that should be acknowledged. Firstly, this study is focused on the approximation of the TIF target. However, the final aim with which tests are administered is to estimate the latent trait of the respondents. The lack of the precision of estimation of the latent trait of the respondents represents the main limitation and future studies should focus on this issue. Secondly, the operationalization of the tiredness of the respondents as an increase of the probability of committing careless error as the administration goes on is an arbitrary choice.

## References

1. Smith, T.F., Waterman, M.S.: Identification of common molecular subsequences. J. Mol. Biol. 147, 195-197 (1981). doi:10.1016/0022-2836(81)90087-5
2. May, P., Ehrlich, H.-C., Steinke, T.: ZIB structure prediction pipeline: composing a complex biological workflow through web services. In: Nagel, W.E., Walter, W.V.,

- Lehner, W. (eds.) Euro-Par 2006. LNCS, vol. 4128, pp. 1148?1158. Springer, Heidelberg (2006). doi:10.1007/11823285\_121
3. Foster, I., Kesselman, C.: The Grid: Blueprint for a New Computing Infrastructure. Morgan Kaufmann, San Francisco (1999)
  4. Czajkowski, K., Fitzgerald, S., Foster, I., Kesselman, C.: Grid information services for distributed resource sharing. In: 10th IEEE International Symposium on High Performance Distributed Computing, pp. 181?184. IEEE Press, New York (2001). doi:10.1109/HPDC.2001.945188
  5. Foster, I., Kesselman, C., Nick, J., Tuecke, S.: The physiology of the grid: an open grid services architecture for distributed systems integration. Technical report, Global Grid Forum (2002)
  6. National Center for Biotechnology Information. <http://www.ncbi.nlm.nih.gov>