

# Standardized Measure for Performance Assessment of Athletes in The CrossFit Open: Theoretical Structuring and Item Response Theory

Rafael da Silva Fernandes (✉ [rafasfer2@ufra.edu.br](mailto:rafasfer2@ufra.edu.br))

Federal Rural University of the Amazon <https://orcid.org/0000-0002-3035-8025>

**Bruna Gabriele Biffe**

Centro Universitário Católico Salesiano Auxilium <https://orcid.org/0000-0001-9650-5713>

**Mário Jefferson Quirino Louzada**

Centro Universitário Católico Salesiano Auxilium <https://orcid.org/0000-0002-5744-2235>

**Antônio César Bornia**

Universidade Federal de Santa Catarina <https://orcid.org/0000-0003-3468-7536>

**Dalton Francisco de Andrade**

Universidade Federal de Santa Catarina <https://orcid.org/0000-0002-4403-980X>

---

## Research Article

**Keywords:** continuous response data, graded-response model, growing grouping, latent trait theory, physical conditioning

**Posted Date:** January 31st, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1308148/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# Standardized Measure for Performance Assessment of Athletes in The CrossFit Open: Theoretical Structuring and Item Response Theory

In its competitive form, CrossFit® intends on assessing the performance of athletes in a wide variety of aspects that determine their conditioning. CrossFit Games is an official competition and intends to recognize the best conditioned athlete in the world to each class. Thus, measuring an athlete's Physical Conditioning is, in a sense, assigning a set of performance outputs that can determine the efficiency and efficacy of the athlete, discriminating the performance of one or more athletes. Since the scores obtained by the athletes in the various workouts are directly identifiable, the conditioning can be seen through the performance in a competition, it is, therefore, a result of the interpretation and scope of the workouts in measuring the performance. This work analyzed data form "CrossFit Open" and has as an objective to propose a new theoretical arrangement to the sport discipline while being of the Item Response Theory which is capable of providing data such as, discriminatory capacity and difficulty level of the workouts, as well as, an assessment of the competition. In other words, intends to describe the probability of an athlete performing a workout and obtaining a score, given his/her physical conditioning. Analysis of the main indications that refer to a good quality of the measurement tool indicates it is a high quality competition. Lastly, this work accomplished both objectives proposed, methodological as well as practical, and recognizes the limitations derived from the reduced amount of qualitative data on the topic and the little use of applied probability models.

**Key words:** competition, graded-response model, latent trait theory, physical conditioning, workout.

## 1 Introduction

2 CrossFit®, sport discipline that has been growing worldwide, has become a popular sport with  
3 more than 15.000 members all over the world. Such an increase may be highlighted when compared  
4 to the number of athletes subscribed in its annual competition, called "The CrossFit Open", in which  
5 approximately 26.000 athletes participated in 2011, reaching 572.653 subscribed athletes in 2019.  
6 (CrossFit, LLC. 2021)

7 In its official website, CrossFit®, is defined as:

8 *"CrossFit is a lifestyle characterized by safe, effective exercise and sound nutrition.*  
9 *CrossFit can be used to accomplish any goal, from improved health to weight loss to better*  
10 *performance. The program works for everyone — people who are just starting out and people*  
11 *who have trained for years."*

12 Usually, every sport discipline, specially, CrossFit® – which presents multiple physical  
13 requirements – needs to identify effective techniques to analyze the performance through a smaller  
14 number of influential variables, thus facilitating the analysis and the development of training  
15 programs to enhance relevant physical skills. Due to its practical nature, this performance  
16 enhancement usually happens in an evolutive and adaptative manner. (Gómes-Landero and Frías-  
17 Menacho 2020)

18 Most of the studies dedicated to CrossFit® have been directed towards understanding physiological  
19 and nutritional factors, training strategies, physical and psychological recovery and other aspects that

20 may directly influence the performance of the athletes. (Claudino, et al. 2018, Schlegel 2020,  
21 Mangine, Stratton, et al. 2020, Mangine, Tankersley, et al. 2020, Mangine, et al. 2021)

22 Typically, such studies vary depend on whether one is analyzing beginners, athletes with longer  
23 sport experience, athletes who are focused on maintaining their health, high-level competitors and  
24 several classes relating to age group and sex.

25 Regarding functional limitations, sporting performance is regulated via different factors that go  
26 through the ability of efficiently repeating the contractile motor activity, however it is limited by the  
27 progression of the fatigue – characterized for a decrease in the strength or production of  
28 musculoskeletal energy causing the reduction in the capacity of keeping the intensity of the exercise,  
29 so that greater fatigue leads to better performance. (Taylor, et al. 2016, Potvin and Fuglevand 2017,  
30 Wan, et al. 2017, Hargreaves and Spriet 2020, García-Pinillos, et al. 2019, Khassetarash, et al. 2021)

31 Other important aspect regarding competitions is the reduced variability in the performance of  
32 athletes when they are position in a specific class. Thus, after considering decisive performance  
33 aspects, strategies design a better adaptive response or responses to induce the best gain to the  
34 athlete's output. (Hanin 2009, Silva-Grigoletto, et al. 2013)

35 Currently, the athlete's performance evaluation criteria is provided by the score obtained in the  
36 execution of a workout. So, performance can be conceptualized, in the CrossFit® context, as an output  
37 or score, presented in time, number of repetitions or pounds, originated from the execution of a  
38 workout by an athlete, being possible to distinguish between the efficiency and efficacy of the  
39 execution.

40 Based on this concept, three points can be highlighted:

- 41 i) Performance, based on the output, it is a tool to measure physical conditioning.
- 42 ii) Performance ascertains the athlete's physical conditioning, in other words, the efficiency and  
43 efficacy of the workout performed by the athlete.
- 44 iii) Performance allows the differentiation or distinction between the conditioning of two or more  
45 athletes.

46 It is important to clarify that the definition of conditioning is not directly identifiable and  
47 observable. What is directly observable and identifiable are the outputs obtained by the athletes in the  
48 various workouts performed. In other words, conditioning is perceived through the performance of  
49 the athletes in the workouts proposed in a competition, therefore, it is a resultant of the interpretation  
50 and scope of the workouts in measuring this execution.

51 For this purpose, evaluating or determining the athlete with the best conditioning based on a single  
52 workout is flawed, since it is insufficient to encompass the wide variety of exercises proposed by  
53 CrossFit® itself.

54 Despite the validation of the “*CrossFit Games*” as a measurement tool of the athlete’s  
55 conditioning, principles of the Classical Test Theory (CTT) are applied to rank athletes on each  
56 workout. Thus, CTT determines the “final score” as a simple rating score, which in CrossFit® is the  
57 sum of the “ranks” obtained. (CrossFit, LLC 2021, Mangine, Tankersley, et al. 2020, Nunnally 1994)

58 Due to the large diversity and number of athletes that can sign up for “CrossFit Games”, it is  
59 expected that when applying CTT, subgroups of athletes are placed on the same rank, thus,  
60 information regarding performance on the different workouts is lost. Hence, Item respond theory  
61 (IRT) has been employed to measure latent traits and characteristics of the measurement. (Baker and  
62 Kim 2004, Hambleton 2000)

63 IRT application contributes to provide information regarding the performance of each athlete in  
64 different workouts, moreover, it becomes possible to obtain a scale for measuring and interpreting  
65 the scores in the CrossFit® setting. (Bonifay 2019, Henninger and Meiser 2020)

66 In this scope, the first objective of this study is methodological: it is proposed the application of a  
67 probabilistic model of the Item respond theory (IRT), called Graded-Response Model proposed by  
68 (Bock, et al. 1997) to describe the probability of an athlete executing a workout obtaining a certain  
69 output, given its physical conditioning, the latent trait being measured. Thus, we have “*CrossFit*  
70 *Games*” as the measurement tool to assess performance (output) of the athletes in different workouts  
71 (items) and determine their latent traits (physical conditioning).

72 The second objective of this work is practical: the application of the model encompasses the  
73 performance analysis, mechanisms to provide additional information that identify execution  
74 characteristics per workout and data regarding the quality of the measurement tool as a criterion for  
75 performance evaluation. In particular, when analyzing the athlete’s performance in various workouts  
76 it is possible to distinguish or differentiate the physical conditioning of two of more athletes.

## 77 **Notation and Framework**

### 78 *Measuring Instrument*

79 “The CrossFit Open” is a qualifying event that, since 2012, is composed by 5 workouts that are  
80 completed by the athletes and mobilizes thousands of athletes around the world to compete in the  
81 biggest participative CrossFit® event, “CrossFit Games”. (CrossFit, LLC 2021)

82 For this purpose, we can describe workout as a group or repetitive series of exercises that require  
83 some combination of strength, cardiopulmonary ability and/or gymnastic, to be performed within a  
84 specific time frame (Time Cap).

85 In most cases, there are two ways of determining the stop criterion: first, after a specified number  
86 of completed repetitions during a predetermined time frame or time cap, which is called truncated by  
87 repetitions, the score is given based on the execution time. In the second case, after a specified time,

88 the athletes complete the maximum number of repetitions, this is called truncated by time, and the  
89 score is given by the number of repetitions, workouts that do not utilize the metrics of the number of  
90 repetitions and/or time may eventually appear, the most common amongst them being the one set by  
91 strength movement, in which the athlete is assessed (score) through the maximum load executed  
92 within a specified time frame. It is worth highlighting the cases in which there is a repetition  
93 sectioning, a case where the athlete could not complete the execution before the time cap the score  
94 is given by the number of completed repetitions.

95 Therefore, it is possible to consider that a workout has its complexity defined by the number of  
96 repetitions to be executed, the time frame defined for execution, the number of types of exercises,  
97 and the complexity of their execution, the complexity being able to differentiate the various workouts.  
98 Thus, we have the following specifications that specify its complexity:

- 99 • Number of Repetitions: referring to the total amount of repetitions to be completed (for repetition  
100 sectioning) or number of completed repetitions (for time sectioning), usually, it is divided in  
101 amount or repetitions per exercises or number of rounds.
- 102 • Execution time or Time Cap: referring to the time-limit available for the athlete to execute all or  
103 the maximum number of constant repetitions in a workout.
- 104 • Number of types of exercises: it is the number of different exercises proposed in a workout. The  
105 difference among the exercises can be determined by the increase in complexity and/or the  
106 increase of the imposed load.
- 107 • Complexity of execution on each exercise: referring to the categorization of the exercises as to  
108 type, exercises that include gymnastic elements, Olympic weightlifting, or aerobic conditioning.

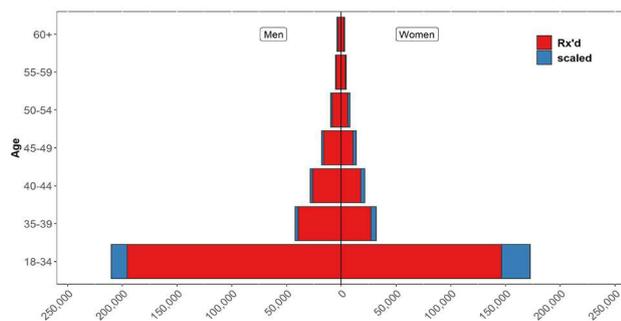
109 In general, workouts are defined to consider the diversity of athletes and could be classified in two  
110 classes: types or division.

111 The first class aims on differentiating beginner athletes to the ones with larger experience in sports  
112 practice and are known as: Rx'd and Scaled. It is worth highlighting the assumption that athletes who  
113 have an extended time of practice tend to be able to execute more complex workouts or have better  
114 skills, whereas beginners need workouts with adapted movements and reduced load.

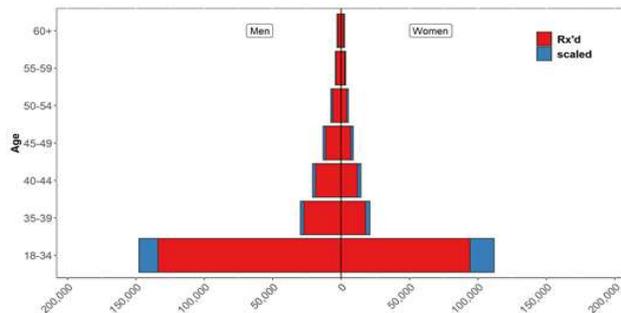
115 The second class is consisted of factors such as sex and age range and takes into consideration  
116 physical and biological characteristics. The workout is specified with varying complexity, according  
117 to these characteristics.

118 As shown in Figure 1 there is a distribution of the athletes in the various divisions and types for  
119 2019 and in Figure 2 for 2020.

## Performance Assessment of The CrossFit



121 **Figure 1.** Number of athletes per division (bar stack) – The Open 2019.



123 **Figure 2.** Number of athletes per division (bar stack) – The Open 2020.

## 124 **Data Set**

125 For this study, data obtained at CrossFit Games (CrossFit, LLC 2021) website referring the years  
126 2019 and 2020 of “The CrossFit Open” were used.

127 Data showed the relevance relative to the variety of the participants that correlates age range and  
128 sex. A greater interest of the athletes is observed in competing in the same “Rx’d” type, the male sex  
129 in larger number compared to the female sex, regarding the 18 to 34 age range. In the “Rx’d” type,  
130 athletes with extended practice time are expected to have longer practice time.

131 Besides, the “Scaled” class is an indication of a smaller amount of practice and/or non-competitive  
132 objectives, in other words, beginner athletes may or may not seek the practice for health purposes and  
133 face competition only as personal challenge.

## 134 **Item Response Theory**

135 In the CrossFit® context, the outputs obtained in a set of workouts have been traditionally used  
136 as an assessment and selection process to find the most conditioned athlete. However, due to the  
137 complexity of each workout, those may benefit athletes with abilities on certain movements. In order  
138 to avoid this situation, and even as a premise of physical conditioning, the set of workouts need to  
139 have a wide variety of requirements.

140 Particularly in “The CrossFit Open”, the score obtained in a specific workout serve as criterion to  
141 rank the athletes, and after the 5 workouts, the general ranking is calculated through the sum of the  
142 rankings in each workout. Consequently, an athlete who is placed in a lower ranking position in each  
143 workout indicates a greater contribution in the final sum, this athlete with lower score being the most  
144 conditioned one.

145 This method of evaluation relies on the specific set of workouts that composes the competition;  
146 thus, analysis and interpretation are always associated to the competition as a whole, which is the  
147 main characteristic of the Classical Test Theory. Therefore, it is made unfeasible the comparison  
148 between people that were not subjected to the same competition, or at least to what are called parallel  
149 methods of evaluation. (Mangine, Tankersley, et al. 2020, Andrade, Tavares and Valle 2000)

150 Contextualizing, the Item Response Theory – IRT refers to the set of probabilistic models that  
151 intend to represent the probability of an athlete obtaining a specific score in a workout as a function  
152 of the characteristic parameters of the workout and the athlete's physical conditioning. This relation  
153 is always expressed in a way that the better the physical conditioning higher is the probability of  
154 obtaining a greater score in a workout. (Andrade, Tavares and Valle 2000)

155 From the concept of Performance, the main characteristics of a workout are its complexity and  
156 ability of distinguishing two or more athletes. As a result, from the IRT point of view, the two-  
157 parameters logistic model is an adequate model to this context. (Hori, Fukuhara and Yamada, Item  
158 response theory and its applications in educational measurement Part I: Item response theory and its  
159 implementation in R 2020, Hori, Fukuhara and Yamada, Item response theory and its applications in  
160 educational measurement Part II: Theory and practices of test equating in item response theory 2020,  
161 Chalmers 2012)

### 162 ***Polytomous Model: Categorization and Sorting***

163 In a context of expansion of the dichotomous model, the polytomous models can handle items with  
164 three or more sorted or unsorted classes. Particularly, (Bock, et al. 1997) proposed the Graded-  
165 Response Model – GRM as an extension of the Two-parameter model. Thus, in the context of the  
166 CrossFit, it is intended to describe the probability of an athlete fitting certain group, based on his/her  
167 physical conditioning, hence, it is expected that a better conditioned athlete will have an increased  
168 probability of obtaining improved performances in a set of workouts, thereby obtaining better outputs.  
169 It means that the sectioning of the athletes may be done gradually and orderly. Athletes with better  
170 performance are assigned to the primary groups and, as the score decreases, they are assigned to the  
171 last groups.

172 However, an issue appears when establishing criteria to define the sectioning of the athletes,  
173 mostly due to the continuous property of the scores, regarding the time unit, or the discrete property,  
174 relating to the number of valid repetitions, which results in estimative precision biases. Particularly,  
175 it is possible to make an empiric comparison between the IRT models that encompasses characteristic  
176 parameters of the evaluator when presenting a notation that refers to the data from the performance  
177 evaluation as well as a discussion regarding the common characteristics amongst evaluators. (Uto and  
178 Ueno, Empirical comparison of item response theory models with rater's parameters 2018, Ueno and

179 Okamoto 2008, Uto and Ueno, Item Response Theory for Peer Assessment 2016) In this project, our  
180 aim is to discuss the rater biases on types. Usual rater characteristics on which the accuracy depends  
181 are as follow:

- 182 • Severity: the tendency of ranking with lower positions that what is justifiable by the results.
- 183 • Consistency: the point to which the evaluator classifies similarly the results from similar quality.
- 184 • Range restriction: the tendency to overuse some classes from restricted sections.

185 In practice, the consistency bias is disregarded, once the evaluator does not attribute a result to the  
186 athlete, thus, it does not represent a bias.

187 It will be implied, for the purpose of this project, that only an evaluator is going to determine the  
188 age restriction, also being described as a specialist or professional in the field.

189 Thereunder, the sectioning of the athletes may be completed in one of two ways: based on the  
190 score obtained, called grouping by score or based in the raking of the athlete, called grouping by rank.  
191 Thus, given that the athletes performed a specific workout and that the number of groups and age  
192 restriction are pre-established, we can define it as:

- 193 • Grouping by score refers to the distribution of the athletes through their obtained score. This  
194 grouping has a discriminatory nature and aims to compare athletes in classes, consequently,  
195 results in an inference about athlete's common characteristics and predictor factors and, due to  
196 subjectivity, it is reasonable considering it as an intuitive process that must be done by specialists  
197 or professionals from this field.
- 198 • Grouping by rank refers to the distribution of athletes through their obtained classification. This  
199 grouping has a qualifying nature.

200 Whether by score or rank, the grouping criteria also relies on the type of truncated and the  
201 tiebreaker criterions. Besides, it is still necessary to introduce the premises of growing grouping per  
202 range and that represents the classification of the athlete according to his/her competitive objectives,  
203 presuming that better conditioned athletes will tend to perform a bigger number of repetitions and be  
204 placed in primary groups.

205 The premises is that the grouping per specialist will tend to be rising, given that it is expected that  
206 as the person becomes more competitive, less people would be interested in dedicating time and  
207 effort, and consequently, fitting the primary groups.

208 Regarding the section scores, workout have as a characteristic the repetition grouping, in other  
209 words, after a Time Cap the athlete interrupts the execution and the score is given based on the  
210 execution time, and not for the maximum number of repetitions. Therefore, time may be adjusted as  
211 follows:

$$time_{fit} = (reps_{max} - reps_{exec}) \times \left( \frac{time_{cap}}{reps_{exec}} \right) + time_{cap} \quad (1)$$

212 in which,  $reps_{max}$  represents the number of the repetitions to be done within the  $time_{cap}$  and  
 213  $reps_{exec}$  represents the number of repetitions the athlete was able to perform.

214 Thus, if the athlete has finished all the repetitions within the time cap, their time will be kept.  
 215 Otherwise the time would be the Time Cap plus the average time of execution of the constant  
 216 repetitions.

217 It is worth highlighting that in practice, in the case of time cap being a really big number, the  
 218 athletes would tend to take longer in the completion of the remaining repetitions, given that they are  
 219 spending more body energetic resources.

### 220 **Graded-Response Model (Samejima – 1969)**

221 The Graded-Response Model by (Samejima 1968) assumes that the classification of the response  
 222 to an item may be sorted with each other. This model obtains more information from people's answers  
 223 than simply if they have given yes or no answers. (Bonifay 2019) (Andrade, Tavares and Valle 2000)  
 224 (Uto and Ueno, Empirical comparison of item response theory models with rater's parameters 2018)

225 In the CrossFit context, we assume the grouping (classes) by score, representing the output of the  
 226 athlete's performance in a workout (item), may be sorted amongst each other, thus the Graded-  
 227 Response Model may be applied. Furthermore, the GRM is useful and allows an estimative of the  
 228 probability of an athlete obtaining a score in a workout given his/her physical conditioning.

229 For instance, assuming that the scores of the workout classes are arranged in order, from lowest to  
 230 highest, and denoted by  $k = 0, 1, \dots, m_i$ , where  $(m_i + 1)$  it is the same number of classes of the  $i$ -th  
 231 workout. The probability of and athlete  $j$  being placed in a certain group, or a higher one of the  $i$   
 232 workout is given by the extension of the Two-parameter logistic model

$$P_{i,k}^+ = \frac{1}{1 + e^{a_i(\theta_j - b_{i,k})}} \quad (2)$$

233 with  $i = 1, 2, \dots, I$ ,  $j = 1, 2, \dots, n$  and  $k = 0, 1, \dots, m_i$ , where  $b_{i,k}$  is the parameter of difficulty of the  $k$ -  
 234 th class of the  $i$  workout and  $\theta_j$  represents the physical conditioning (latent trait) of the  $j$ -th athlete.

235 Regarding the models for dichotomous items, the slope parameters  $a_i$  is the item discrimination.  
 236 However, regarding models for non-dichotomous items, the discrimination of a specific class depends  
 237 on the slope parameter, common to all the item classes, as well as the distance from adjacent difficulty  
 238 classes.

239 Thus, the probability of a person  $j$  receiving a score  $k$  in the  $i$  item is given by the expression:

$$P_{i,k}(\theta_j) = \frac{1}{1 + e^{a_i(\theta_j - b_{i,k})}} - \frac{1}{1 + e^{a_i(\theta_j - b_{i,k+1})}} \quad (3)$$

240 Notice that if we have a test with  $i$  items, each one with  $(m_i + 1)$  output classes, then we shall  
 241 have  $[\sum_{i=1}^I m_i + I]$  parameters to be estimated.

## 242 **Results and Discussions**

243 Typically, in world level sports the quality of a competition in determining the best athletes or  
 244 teams is given by their acceptance in recognizing and validating this competition. Another way of  
 245 recognizing a competition is based on the rules and norms that a bigger authority states and in turn,  
 246 exerts validation of the competitions around the world. As an example, the soccer regulation by FIFA.

247 Regarding CrossFit, the validation of the competitions is exerted by the entity itself, which  
 248 sanctions world events and function as qualifying stages for the final competition called *CrossFit*  
 249 *Games*. However, due to the proportions reached by the sport and its growing pace, several  
 250 competitions, classified as amateurs, aim at determining the fittest athletes, those being validated by  
 251 their own competitors.

252 In this section, we present an analysis of the events and an assessment of the quality of the “*The*  
 253 *CrossFit Open*” as a mechanism for measuring physical conditioning.

254 Due to the size of the data set, a preliminary analysis was done regarding the “Rx’d” type, section  
 255 “Men (18-34)”, 2019.

### 256 ***Preliminary Analysis of the Data Set***

257 We initiate our analysis focusing on the workout 19.1, because it is a time grouping, the score is  
 258 given by the number of executed repetitions until  $T_{cap} = 900 \text{ seconds}$ . Therefore, it is necessary to  
 259 establish new values to represent the output of the athletes, thus, if the athlete fits a specific group,  
 260 for example, Group 1, it means that he or she obtained a score 5, Group 2 with a score of four and so  
 261 on. Thus, characterizing and sorting the workout 19.1, the data is summarized in Table 1.

262 Table 1  
 263 Characterization and sorting of workout 19.1.

workout	Group	Score	Truncated	Frequency	Inferior Limit	Upper limit
19.1	01	5	by times	2224	342	418
19.1	02	4	by times	13006	304	342
19.1	03	3	by times	20194	281	304
19.1	04	2	by times	27549	258	281
19.1	05	1	by times	42950	228	258
19.1	06	0	by times	58140	0	228
19.1	n/a	0	n/a	19193	n/a	n/a

264 n/a = not assessed by the measurement tool.

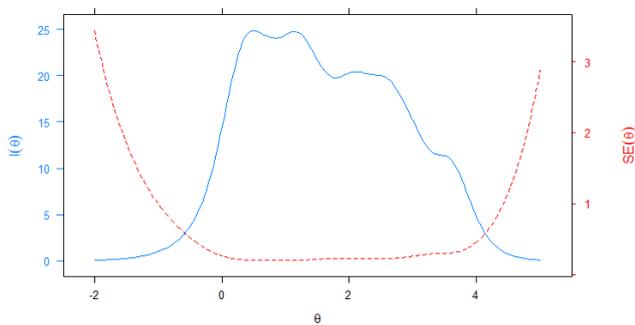
265 The parameters to each of the workouts are estimated assuming that the distribution of  $\theta_j$  follows  
 266 a normal distribution with  $\mu = 0$  e  $\sigma = 1$ . Values for  $a < 1$  indicate the item has little discrimination  
 267 capacity. Values for  $a \geq 1$  mean the item discriminates well. It is possible to observe on Table 2 that  
 268 all the workouts present values for  $a \geq 1$ .

269  
270

Table 2  
Estimates and standard error (SE) of the workout parameters of Graded-Response Mode on scale (0,1).

workout	$\hat{a}$	$\hat{b}_1$	$\hat{b}_2$	$\hat{b}_3$	$\hat{b}_4$	$\hat{b}_5$
19.1	2.26 (0.009)	0.25 (0.004)	0.51 (0.004)	1.11 (0.005)	1.80 (0.006)	3.05 (0.013)
19.2	4.60 (0.025)	0.46 (0.003)	1.23 (0.004)	2.23 (0.007)	3.13 (0.014)	3.76 (0.027)
19.3	3.43 (0.015)	0.05 (0.003)	0.61 (0.003)	1.30 (0.004)	2.35 (0.008)	3.66 (0.023)
19.4	4.41 (0.021)	0.26 (0.003)	0.88 (0.003)	1.44 (0.004)	2.04 (0.006)	2.64 (0.009)
19.5	4.48 (0.022)	0.31 (0.003)	0.73 (0.003)	1.15 (0.004)	1.70 (0.005)	2.76 (0.01)

271 Finally, we may still evaluate the Test Information Function (TIF) and the Standard Error of  
272 Measurement (SEM) presented on Figure 3. So, we can verify the degree of precision of the workout  
273 set to several scale ranges (0,1), and as can be seen, SEM presents lower values, better precision, in  
274 the interval [0,4].

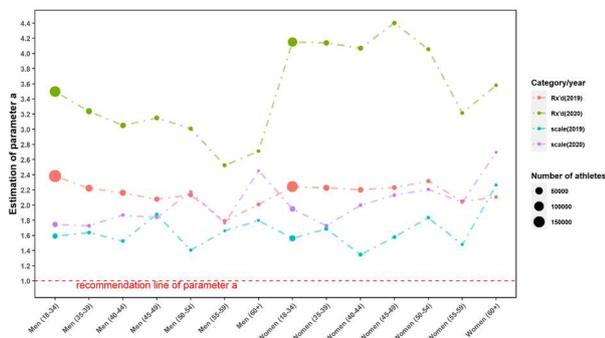


276 **Figure 3.** Test Information Function (TIF – continuous blue line) and Standard Error of Measurement (SEM – dotted red line).

277 **General Analysis of the Data Set**

278 The general analysis consists of an analysis of the main indicators that relate to a measurement  
279 tool of good quality, those indicators being: *a* parameter analysis to know the discrimination capacity  
280 of the workouts; the respondents frequency, which in this context is about the number of athletes  
281 distributed on the sections; and the analysis of the TIF and SEM curves.

282 Initially, all the estimates of the *a* parameters were analyzed and presented on Figure 4. It is  
283 possible to observe that the values of interest, typically  $a \geq 1$ , were presented in all the scenarios.



285 **Figure 4.** Estimates of the discrimination parameters.

**286 Final Considerations**

287 We focused our work in presenting a conceptual arrangement of the main definitions, terms and  
288 expressions applied in the CrossFit context and that were pointed towards Samejima's Graded-  
289 Response Model of the Item Response Theory (Samejima 1968). In this sense, we were able to  
290 accomplish the primary objective of the production by applying GRM to the respective context, and  
291 as a result, describing the probability of an athlete performing a Workout and obtaining a score, given  
292 his/her physical conditioning.

293 On the other hand, given the adjustment of the IRT to the context presented, it was possible to  
294 achieve the second objective and incorporate to the athletes' performance analyses useful data that  
295 seek to identify performance characteristics and positively evaluate the quality of the CrossFit Open.

296 Commonly, the evaluation process of a measurement tool requires an interactive process of  
297 regrouping the results, that in a sense, intend to improve the quality and validation of the measurement  
298 instrument. In the Item Response Theory this process means that the classes of responses are to  
299 assessing well enough the analyzed item, and it is necessary a grouping, in other words, in the context  
300 of the CrossFit and this work, it means that the grouping of the athletes, in face of their results, cannot  
301 contain all the data that the Model would be able to collect, and with the regrouping, it could offer  
302 more data about the measurement tool.

303 Considering this and the 56 studied scenarios, individually evaluation each presented situation  
304 could be a costly work. Besides the regrouping analysis, it is still necessary to answer the following  
305 questions: Did the score grouping, designed by a specialist, allocated an adequate number of athletes  
306 or did it in the best way? Is the number of specified groups sufficiently adequate? Given that the  
307 grouping elaborated by the specialist and done per Workout is based on a single scenario, could the  
308 remaining scenarios be better regrouped?

309 Those are the complex answers that, besides considering the premises and characteristics of the  
310 referred context, require a greater discussion on how the optimization process must be completed.  
311 However, this work did not focus on this optimization process.

312 In addition, the Item Response Theory provides us other tools that make possible to infer and  
313 measure, qualitatively, predictor factors of performance, for example, a subject of great relevance in  
314 the field of Sports Science and Exercise Physiology.

315 Studies that merge the Item Response Theory and the Sports Science and Exercise Physiology  
316 context in the presented manner were not found in literature. There are, on the other hand, studies that  
317 merge Sports Science and Psychology and encompasses psychometric assessments, among those,  
318 ones that employ the Item Response Theory.

319 In conclusion, this work accomplished both objectives proposed, methodological as well as  
320 practical, and recognizes the limitations derived from the reduced amount of qualitative data on the  
321 topic and the little use of applied probability models.

322 Future research is needed to build a standardized performance measurement scale that allows for  
323 a contextual and practical interpretation of the performance metrics obtained.

324

## 325 **Acknowledgments**

326 This study was fulfilled employing data collected form CrossFit® official website, which describes  
327 its privacy policy that the publicly available information may contain personal information. However,  
328 this work is restricted to the generic use of the data, not using or disclosing athlete's personal data.  
329 The authors thank CrossFit® for making the data available.

## 330 **Authors' contributions**

331 All the authors contributed equally to this assignment.

## 332 **Competing Interests**

333 The authors declare that they have no competing interests.

## 334 **References**

- 335 Andrade, Dalton Francisco de, Heliton Ribeiro Tavares, e Raquel da Cunha Valle. *Teoria da Resposta*  
336 *ao Item: Conceitos e Aplicações*. São Paulo: Associação Brasileira de Estatística, 2000.
- 337 Baker, Frank B., e Seiock-Ho Kim. *Item Response Theory: parameter estimation techniques*. 2. New  
338 York: Marcel Dekker, 2004.
- 339 Bock, R. Darrell, et al. *In Handbook of Modern Item Response Theory*. Edição: Wim J. Van der  
340 Linden e Ronald K. Hambleton. New York: Springer-Verlag, 1997.
- 341 Bonifay, W. *Multidimensional item response theory*. SAGE Publications, 2019.
- 342 Chalmers, R. Philip. "mirt: A multidimensional item response theory package for the R environment."  
343 *Journal of statistical Software* 48, n° 1 (05 2012): 1-29.
- 344 Claudino, João Gustavo, et al. "CrossFit Overview: Systematic Review and Meta-analysis." *Sports*  
345 *medicine-open* 4, n° 1 (2018): 1-14.
- 346 CrossFit, LLC. *Games CrossFit*. CrossFit. 2021. <https://games.crossfit.com/> (acesso em 22 de 06 de  
347 2021).
- 348 CrossFit, LLC. *CrossFit*. 2021. <https://www.crossfit.com> (acesso em 22 de 06 de 2021).
- 349 García-Pinillos, Felipe, Alejandro Molina-Molina, Juan A. Párraga-Montilla, e Pedro A. Latorre-  
350 Román. "Kinematic alterations after two high-intensity intermittent training protocols in  
351 endurance runners." *Journal of Sport and Health Science* 8, n° 5 (09 2019): 442-449.

- 352 G3mes-Landero, Luis Arturo, e Juan Miguel Fr3as-Menacho. "Analysis of Morphofunctional  
353 Variables Associated with Performance in Crossfit® Competitors." *Journal of Human*  
354 *Kinetics* 73, n3 1 (07 2020): 83-91.
- 355 Hambleton, Ronald K. "Emergence of Item Response Modeling in Instrument Development and Data  
356 Analysis." *Mecial Care* 38, n3 9 (09 2000): 1160-1165.
- 357 Hanin, Yuri Hanin and Muza. "Optimization of performance in top-level athletes: an action-focused  
358 coping approach." *International Journal of Sports Science & Coaching* 4, n3 1 (03 2009): 47-  
359 91.
- 360 Hargreaves, Mark, e Lawrence L. Spriet. "Skeletal muscle energy metabolism during exercise."  
361 *Nature Metabolism* 2, n3 1 (08 2020): 817-828.
- 362 Henninger, Mirka, e Thorsten Meiser. "Different Approaches to Modeling Response Styles in Divide-  
363 by-Total Item Response Theory Models (Part 2): Applications and Novel Extensions."  
364 *American Psychological Association* 25, n3 5 (2020): 577-595.
- 365 Hori, Kazuki, Hirotaka Fukuhara, e Tsuyoshi Yamada. "Item response theory and its applications in  
366 educational measurement Part I: Item response theory and its implementation in R." *WIRES*  
367 *Computational Statistics*, 2020: e1531.
- 368 Hori, Kazuki, Hirotaka Fukuhara, e Tsuyoshi Yamada. "Item response theory and its applications in  
369 educational measurement Part II: Theory and practices of test equating in item response  
370 theory." *WIRES Computational Statistics* e1531 (2020).
- 371 Khassetarash, Arash, Gianluca Vernillo, Renata L. Kryger, W. Brent Edwards, e Guillaume Y. Millet.  
372 "Neuromuscular, biomechanical, and energetic adjustments following repeated bouts of  
373 downhill running." *Journal of Sport and Health Science*, 06 2021.
- 374 Mangine, G. T., Y Feito, J. E. Tankersley, J. M. McDougle, e B. M. Kliszczewicz. "Workout Pacing  
375 Predictors of Crossfit® Open Performance: A Pilot Study." *Journal of Human Kinetics* 78 (03  
376 2021): 89-100.
- 377 Mangine, Gerald T., et al. "Physiological differences between advanced CrossFit athletes,  
378 recreational CrossFit participants, and physically-active adults." *PLOS ONE* 15, n3 4 (04  
379 2020): e0223548.
- 380 Mangine, Gerald T., et al. "Predictors of CrossFit Open Performance." *Sports* 8, n3 7 (07 2020): 102.
- 381 Nunnally, JUM C. *Psychometric Theory*. 3. New York: Tata McGraw-hill education, 1994.
- 382 Potvin, Jim R., e Andrew J. Fuglevand. "A motor unit-based model of muscle fatigue." *PLoS*  
383 *Computational Biology* 13, n3 6 (06 2017): e1005581.
- 384 Samejima, Fumi. *Estimation of Latent Ability Using a Response Pattern of Graded Scores*. New  
385 Jersey: Psychometric Monograph, 1968.

- 386 Schlegel, Petr. “CrossFit® Training Strategies from the Perspective of Concurrent Training: A  
387 Systematic Review.” *Journal of Sports Science and Medicine* 19, nº 4 (2020): 670-680.
- 388 Silva-Grigoletto, Marzo Edir Da, Teresa Valverde-Esteve, Ciro José Brito, e Juan Manuel García-  
389 Manso. “Ability to repeat strength: effects of recovery between repetitions.” *Revista*  
390 *Brasileira de Educação Física e Esporte* 27, nº 4 (12 2013): 689-705.
- 391 Taylor, Janet L. , Markus Amann, Jacques Duchateau, Romain Meeusen, e Charles L. Rice. “Neural  
392 Contributions to Muscle Fatigue: From the Brain to the Muscle and Back Again.” *Medicine*  
393 *and science in sports and exercise* 48, nº 11 (2016): 2294-2306.
- 394 Ueno, M., e T. Okamoto. “Item response theory for peer assessment.” *Proc. IEEE International*  
395 *Conference on Advanced Learning (IEEE)*, 2008: 554-558.
- 396 Uto, Masaki, e Maomi Ueno. “Empirical comparison of item response theory models with rater's  
397 parameters.” *Heliyon* 4, nº 5 (05 2018): e00622.
- 398 Uto, Masaki, e Maomi Ueno. “Item Response Theory for Peer Assessment.” *IEEE TRANSACTIONS*  
399 *ON LEARNING TECHNOLOGIES* 9, nº 2 (04-06 2016): 157-170.
- 400 Wan, Jing-jing, Zhen Qin, Peng-yuan Wang, Yang Sun, e Xia Liu. “Muscle fatigue: general  
401 understanding and treatment.” *Experimental & Molecular Medicine* 49, nº 1 (10 2017): e384.  
402  
403