



## Multi-Faceted Rasch Model in Practical Tests: Study Subject: Sorayesh Test

SeyyedeH Hoda Naji<sup>1</sup>, Ali Moghadamzadeh<sup>2</sup>, Balal Izanloo<sup>3</sup>, Ebrahim Khodaie<sup>4</sup>

1. Ph.D student, Department of Curriculum Planning and Educational Methods, Faculty of Psychology and Education, University of Tehran, Tehran, Iran. Email: s.hodanaji@ut.ac.ir

2. Associate professor, Department of Curriculum Planning and Educational Methods, Faculty of Psychology and Education, University of Tehran, Tehran, Iran; (Corresponding Author), Email: amoghadamzadeh@ut.ac.ir

3 Associate Professor, University of Kharazmi, Faculty of Psychology and Education, Karaj, Iran. Email: izan.b@khu.ac.ir

4. associate professor. Department of Curriculum Planning and Educational Methods, Faculty of Psychology and Education, University of Tehran, Tehran, Iran. Email: khodaie@ut.ac.ir

### Article Info

**Article Type:**  
**Research Article**

**Received:** 2024/03/24

**Received in revised form:** 2024/08/10

**Accepted:** 2024/09/05

**Published online**  
**2024/10/06**

### ABSTRACT

**Objective:** The purpose of this research is to analyze the data obtained from the practical tests conducted by National Organization of Educational Testing, using the Multi-Faceted Rasch model and comparing it with the results of classical analysis methods.

**Methods:** The research method is quantitative a type of descriptive-analysis method. The participants included all the songwriting candidates taking the Sorayesh practical test. The analyzed data has been obtained from the evaluation forms filled out by four evaluators for each candidate.

**Results:** The findings show that although the correlation coefficient among raters was high (more than 0/90), the agreement of the raters in terms of the Kappa coefficient was average. Furthermore, based on the Multi-Faceted Rasch model, the raters strictness parameter was moderate.

**Conclusion:** The difference between the correlation and the Kappa coefficient shows that these two indicators cannot be used alone to analyze the rater. Also, these indicators present the group status of the raters while the Multi-Faceted Rasch model shows the individual status of each rater. The results of the Multi-Faceted Rasch model indicated that the raters didn't exhibit strictness or leniency errors.

**Keywords:** Multi-Faceted Rasch - Rasch model - Classical Test Theory - Rater - practical test - Kappa

**Cite this article:** Naji, SeyedeH Hoda; Moghadamzadeh, Ali; Izanloo, Balal; Khodaie, Ebrahim (2024). Multi-Faceted Rasch Model in practical tests: Study subject: Sorayesh test. *Educational Measurement and Evaluation Studies*, 14 (47): 07-26 pages. DOI:10.22034/emes.2024.2003752.2479



© The Author(s).

Publisher: National Organization of Educational Testing (NOET)



## مدل راش چند وجهی در آزمون‌های عملی: مورد مطالعه آزمون سُرایش

سیده هدی ناجی<sup>۱</sup>، علی مقدم‌زاده<sup>۲</sup>، بلال ایزانلو<sup>۳</sup>، ابراهیم خدایی<sup>۴</sup>

۱. دانشجوی دکتری، گروه آموزشی روش‌ها و برنامه‌های آموزشی، دانشکده روانشناسی و علوم تربیتی، دانشگاه تهران، ایران. رایانامه: s.hodanaji@ut.ac.ir
۲. گروه آموزشی روش‌ها و برنامه‌های آموزشی، دانشکده روانشناسی و علوم تربیتی، دانشگاه تهران، ایران؛ (نویسنده مسئول)، رایانامه: amoghdamzadeh@ut.ac.ir
۳. گروه آموزشی برنامه‌ریزی درسی، دانشکده روانشناسی و علوم تربیتی، دانشگاه خوارزمی، کرج، ایران. رایانامه: izan.b@khu.ac.ir
۴. گروه آموزشی روش‌ها و برنامه‌های آموزشی، دانشکده روانشناسی و علوم تربیتی، دانشگاه تهران، تهران، ایران. رایانامه: khodaie@ut.ac.ir

اطلاعات مقاله	چکیده
نوع مقاله: مقاله پژوهشی	هدف: هدف این پژوهش تحلیل داده‌های حاصل از آزمون‌های عملی سازمان سنجش با استفاده از مدل راش چندوجهی و مقایسه آن با نتایج حاصل از تحلیل‌های کلاسیک است.
دریافت: ۱۴۰۳/۰۱/۰۵	روش پژوهش: روش پژوهش کمی و از نوع توصیفی-تحلیلی است. مشارکت کنندگان شامل تمامی داوطلبان رشته آهنگسازی حاضر در آزمون عملی سُرایش بودند. داده مورد تحلیل از فرم‌های ارزیابی پر شده توسط چهار ارزیاب برای هر داوطلب بدست آمده است.
اصلاح: ۱۴۰۳/۰۵/۲۰	یافته‌ها: یافته نشان می‌دهد اگر چه ضریب همبستگی میان ارزیاب‌ها بالا بوده (بیش از ۰/۹۰)، توافق ارزیاب‌ها براساس ضریب کاپا، در حد متوسط است. همچنین براساس نتایج راش چند وجهی، پارامتر سخت‌گیری ارزیاب‌ها در حد متوسط بوده است.
پذیرش: ۱۴۰۳/۰۶/۱۵	نتیجه‌گیری: تمایز بین ضریب همبستگی و کاپا، نشان می‌دهد که نمی‌توان از هر یک از این شاخص‌ها به تنهایی برای تحلیل ارزیاب استفاده کرد، همچنین این ضرایب، وضعیت گروهی ارزیاب‌ها را نشان می‌دهند در حالی که راش چند وجهی وضعیت هر یک از ارزیاب‌ها را نشان می‌دهد. نتایج راش چند وجهی نشان داد که ارزیاب‌ها دچار خطای سخت‌گیری و سهل‌گیری نیستند.
انتشار: ۱۴۰۳/۰۷/۱۵	

واژه‌های کلیدی: راش چند وجهی - مدل راش - نظریه کلاسیک اندازه‌گیری - ارزیاب - آزمون عملی - کاپا

استناد: ناجی، سیده هدی؛ مقدم‌زاده، علی؛ ایزانلو، بلال؛ خدایی، ابراهیم (۱۴۰۳). مدل راش چند وجهی در آزمون‌های عملی: مورد مطالعه آزمون سُرایش. *مطالعات*



حق مؤلف © نویسنده‌گان.

اندازه‌گیری و ارزشیابی آموزشی، ۱۴ (شماره ۴۷)، ۲۶-۷ صفحه. DOI: 10.22034/emes.2024.2003752.2479

ناشر: سازمان سنجش آموزش کشور

## مقدمه

در بسیاری از جوامع، سنجش و اندازه‌گیری آموزشی و روانشناختی، سهم بسیار زیادی در علوم شناختی و رفتاری دارد و منابع اطلاعاتی اساسی و قابل توجهی در مورد افراد و گروه‌ها فراهم می‌آورد. استفاده درست از آزمون‌ها می‌تواند زمینه تصمیم‌گیری بهتری در مورد افراد و برنامه‌ها را فراهم کند و راه دسترسی عادلانه‌تری به آموزش در سطوح بالاتر ایجاد کند (انجمن تحقیقات آموزشی آمریکا<sup>۱</sup> و همکاران، ۱۳۹۸/۲۰۱۴).

در زمینه آموزش و روانشناسی، نگرانی در مورد قضاوت ارزیاب‌ها رو به افزایش است. اثرات ارزیاب مانند شدت یا ملایمت و گرایش مرکزی، معمولاً به عنوان منبع واریانس مرتبط با طرح ارزیابی در نظر گرفته می‌شود، یعنی به عنوان منبعی از واریانس نظامدار در رتبه‌بندی‌های مشاهده شده که به ارزیاب‌ها مرتبط است نه با آزمون شوندگان (لی و چا، ۲۰۱۶). آزمون‌های عملکردی، در بسیاری از موارد مستلزم استفاده از چند ارزیاب برای ارزیابی عملکرد آزمودنی‌ها هستند تا قابلیت اطمینان و عینی بودن رتبه‌بندی‌ها را افزایش دهند و خطاهای ارزیاب را به حداقل برسانند، که همین مساله سبب پیچیدگی فرایند ارزیابی می‌شود (همبو<sup>۲</sup> و همکاران، ۲۰۰۱؛ رویتر و استینفلد<sup>۴</sup>، ۲۰۱۸). آزمون شوندگان به سوالات پاسخ می‌دهند و ارزیاب‌ها طبق درک خود از سازه مورد اندازه‌گیری، ابزار اندازه‌گیری و اصول و قواعد ارزیابی مورد استفاده، پاسخ آزمون شوندگان را ارزیابی می‌کنند (اکس<sup>۵</sup>، ۲۰۱۵). بنابراین می‌توان گفت فرایند ارزیابی آزمون شوندگان با استفاده از ارزیاب انسانی، فرآیندی پیچیده و غیر مستقیم است. با توجه به اینکه آزمون شوندگان به طور طبیعی از نظر توانایی مورد اندازه‌گیری متفاوت هستند، این انتظار وجود ندارد که همه آن‌ها امتیاز یکسانی دریافت کنند؛ بلکه انتظار می‌رود در رتبه‌بندی توانایی‌های آزمون شوندگان اختلاف وجود داشته باشد. هر گونه تغییر در رتبه‌های آزمون شوندگان که به دلیل تفاوت قابل اعتماد در توانایی‌های آن‌ها باشد، مطلوب است. با این حال، رتبه‌بندی‌ها تحت تأثیر چندین عامل خارجی مانند نوع سوالات، رتبه‌دهندگان، شرایط و موارد دیگری قرار خواهند گرفت (لی و چا، ۲۰۱۶). اگر هنگام ارزیابی عملکرد توسط ارزیاب یک یا چند سوگیری رخ دهد، میزان سوگیری در پیش‌بینی عملکرد افراد زیاد خواهد بود در نتیجه این پیش‌بینی‌ها باعث اندازه‌گیری‌های غیرقابل اعتماد می‌شوند. سوگیری ارزیاب به واریانس نامربوط سازه نسبت داده می‌شود که مستقیماً اعتبار اندازه‌گیری را تهدید می‌کند. بنابراین، چگونگی معرفی واریانس نامربوط سازه توسط ارزیاب‌ها مورد بحث است (یشیل چنر و ساتا، ۲۰۲۱).

بطور خلاصه اندازه‌گیری پیشرفت تحصیلی، سنجش استعداد تحصیلی و شناخت میزان توانایی علمی و مهارتی افراد، از جمله مهم‌ترین موقعیت‌های تحقیقاتی و بررسی‌های آموزشی است. در بسیاری از اندازه‌گیری‌ها و موقعیت‌های اندازه‌گیری، ناگزیر به استفاده از آزمون‌های تشریحی و عملی خواهیم بود. استفاده از این آزمون‌ها، عموماً نیازمند استفاده از ارزیاب‌های انسانی و نمره‌گذاری‌های ذهنی هستند. در این گونه نمره‌گذاری‌ها همواره اثر ارزیاب بر روی برآورد توانایی افراد وجود دارد. هر چند با آموزش ارزیاب‌ها و تهیه دستورالعمل‌های ارزیابی تلاش می‌شود اینگونه نمره‌گذاری‌ها به نمره‌گذاری عینی نزدیک شود، اما پژوهش‌ها نشان داده که بدلیل ماهیت انسان بودن ارزیاب‌ها و در نتیجه وجود خطاهای انسانی، هرگز نمی‌توان چنین اندازه‌گیری‌هایی را دقیقاً عینی در نظر گرفت، بنابراین نظریه‌های مختلف اندازه‌گیری به بحث اثر ارزیاب و یا خطای ارزیاب پرداخته‌اند و تلاش کرده‌اند در نمره‌گذاری‌هایی که با واسطه ارزیاب انسانی انجام می‌شود، دقت و اعتبار نمرات و توانایی برآورد شده را افزایش دهند.

در ایران برای سنجش داوطلبان ورود به آموزش عالی در برخی از رشته‌های مهارتی مانند موسیقی، علاوه بر کنکور از آزمون‌های عملی یا تشریحی استفاده می‌شود. با توجه به اینکه در این آزمون‌ها از ارزیاب برای ارزیابی داوطلبان و یا تصحیح اوراق استفاده می‌شود، بررسی روانسنجی چنین آزمون‌هایی با استفاده از نظریه‌های جدید اهمیت بسیاری دارد، بنابراین در این پژوهش ویژگی‌های سوال، توانایی افراد و وضعیت ارزیابان با استفاده از نظریه کلاسیک و راش چندوجهی تحلیل شده و مورد مقایسه قرار گرفته است.

1. American Educational Research Association

2. Lee and Cha

3. Hombro

4. Robitzsch and Steinfeld

5. Eckes

## مبانی نظری و پیشینه پژوهش

در نظریه کلاسیک آزمون (CTT)، نمره مشاهده شده آزمون (X) مساوی مجموع دو جزء است: نمره واقعی (T) و نمره خطای تصادفی (E).

$$X = T + E$$

بر مبنای این نظریه نمره مشاهده شده قابل رؤیت و نمرات واقعی و خطا دارای ساخت نظری و غیر قابل مشاهده هستند. برای افزایش قابلیت اطمینان نتایج حاصل از آزمون، باید واریانس ناشی از خطای اندازه‌گیری (واریانس نمره خطا) را کاهش داد تا بتوان واریانس نمرات مشاهده شده را به واریانس نمرات واقعی نسبت داد. برای انجام این کار، معمولاً اندازه‌گیری‌ها در بعضی از شرایط از پیش تعیین شده اندازه‌گیری انجام می‌گیرد و میانگین اندازه‌گیری‌های انجام شده به عنوان یک برآورد از اندازه‌گیری "ایده‌آل" در نظر گرفته می‌شود. باید توجه داشت، تمام اندازه‌گیری‌های انجام شده برای دستیابی به نمره واقعی، باید در شرایط یکسان اندازه‌گیری باشند علاوه بر این، نتایج حاصل را تنها به همان شرایط می‌توان تعمیم داد. از آنجا که برقرار کردن این شرایط یکسان همواره امکان پذیر نیست، استفاده از نتایج اندازه‌گیری محدود می‌شود. از سوی دیگر برای تعیین اعتبار<sup>۲</sup> اندازه‌گیری، روش‌های مختلفی در نظریه کلاسیک وجود دارد از جمله همبستگی میان نمرات دو آزمون موازی و این در حالی است که ایجاد دو آزمون موازی دارای ابهامات و دشواری‌های اساسی است، همچنین روش‌های مختلف برآورد پایایی آزمون، منجر به ضرایب پایایی مختلف و در نتیجه خطای اندازه‌گیری مختلف برای یک آزمون می‌شود و نظریه کلاسیک در تعیین مناسبترین ضریب پایایی، ناتوان است. علاوه بر این، نظریه کلاسیک، هنگامی که طرح اندازه‌گیری پیچیده می‌شود (برای مثال آزمون‌های عملکردی و آزمون‌های غیر عینی)، نمی‌تواند پاسخگو باشد (تیلور<sup>۳</sup>، ۱۳۹۸/۲۰۱۳). در چنین موقعیت‌هایی، بدلیل اینکه منابع اثرگذار بر واریانس ناخواسته مانند ارزیاب به طرح اندازه‌گیری اضافه می‌شود، احتمال خطای اندازه‌گیری افزایش پیدا می‌کند بنابراین باید با استفاده از روش‌هایی اعتبار اندازه‌گیری انجام شده را نشان داد. در نظریه کلاسیک اندازه‌گیری، در چنین موقعیت‌هایی، علاوه بر تحلیل سوال و بررسی پایایی ابزار، ارزیاب نیز مورد بررسی قرار می‌گیرد. برای بررسی ارزیاب در نظریه کلاسیک اندازه‌گیری، روش‌های مختلفی برای بررسی اعتبار ارزیابی‌ها وجود دارد. ضریب توافق، ضریب کاپا<sup>۴</sup>، ضریب همبستگی پیرسون، و مقایسه میانگین ارزیاب‌ها براساس آزمون t یا تحلیل واریانس یک راه<sup>۵</sup>، از جمله روش‌های مورد استفاده در نظریه کلاسیک اندازه‌گیری است (ووغان<sup>۶</sup>، ۱۹۹۱ نقل در پُلِت<sup>۷</sup> و همکاران، ۲۰۲۲). همانگونه که توماس اِکس بیان کرده است، یک مشکل اساسی در رویکردهای سنتی وجود دارد که می‌توان نام آن را پارادوکس توافق-دقت<sup>۸</sup> گذاشت. اِکس در پژوهش خود نشان می‌دهد که پایایی بالا و همچنین توافق بالا در بین ارزیابان، لزوماً به معنای دقت بالا در ارزیابی مهارت آزمون شوندگان نیست، بنابراین پایایی و یا ضریب توافق بالا در بین ارزیابان می‌تواند منجر به نتیجه‌گیری‌های اشتباه شود. از سوی دیگر پایین بودن توافق در بین ارزیابان موجب نتیجه‌گیری نادرست در مورد ارزیابی‌ها می‌شود و همان‌گونه که در بسیاری از موقعیت‌های ارزیابی انجام می‌شود، این گونه ارزیابی به عنوان ارزیابی‌های نامناسب شناسایی شده و ارزیاب‌های ناهمسو با گروه ارزیابی، با وجود داشتن دقت بالا در ارزیابی، بدلیل همسو نبودن با سایرین، از گروه ارزیابی حذف و افراد دیگری جایگزین می‌شوند تا ضریب توافق افزایش یابد (اکس، ۲۰۰۹). اگر چه ممکن است ارزیاب‌ها با یکدیگر موافق باشند، ولی به طور دسته جمعی برداشت اشتباهی از دستورالعمل نمره‌گذاری داشته باشند. این مساله تناقض در رفتار ارزیاب تلقی می‌شود. حتی وقتی به نظر می‌رسد داده‌های مربوط به توافق بین ارزیابان قدرت بالایی دارد، اگر عدم توافق بین ارزیابان به هر نحو با سوگیری همراه باشد، نمرات متناقض می‌تواند بسته به ارزیاب، منجر به نمره کل متفاوتی برای آزمودنی‌ها شود (تیلور، ۱۳۹۸/۲۰۱۳)، بنابراین در نظریه کلاسیک نمی‌توان عملکرد افراد را مستقل از عملکرد ارزیاب محاسبه کرد و همان‌گونه که نمره افراد وابسته به ویژگی‌های روانسنجی سوال و ویژگی‌های سوال وابسته به گروه آزمودنی است، در مورد ارزیاب نیز این وابستگی وجود دارد. اِکس راه حل این تناقض را استفاده از رویکردهای جدید اندازه‌گیری بجای رویکردهای کلاسیک

1. Classical Test Theory
2. reliability
3. Taylor
4. Kappa
5. Analysis of Variance
6. Vaughan
7. Polat
8. agreement-accuracy

می‌داند (اکس، ۲۰۰۹). یکی از رویکردهای جدید در تحلیل ارزیاب‌ها، استفاده از مدل‌های اندازه‌گیری راش<sup>۱</sup> (RM) است. مدل‌های اندازه‌گیری راش از جمله مدل‌های مبتنی بر نظریه صفت مکتون<sup>۲</sup>، است که در سال ۱۹۶۰ توسط روان‌سنج دانمارکی، جُرج راش<sup>۳</sup>، معرفی گردید. راش مدل و رویکرد خود در اندازه‌گیری‌های اجتماعی را به عنوان مدل‌های احتمالی برای برخی از اندازه‌گیری‌های هوش و پیشرفت، مطرح کرد (ویند و هوا<sup>۴</sup>، ۲۰۲۲). این رویکرد در واقع خانواده‌ای از مدل‌های آماری است که بجای محاسبه نمره در نظریه کلاسیک، احتمال پاسخ درست افراد را با توجه به میزان توانایی فرد و دشواری سوال برآورد می‌کند. فرمول عمومی مدل راش بصورت زیر است:

$$P(\theta) = \left( \frac{e^{(\theta_n - \delta_i)}}{1 + e^{(\theta_n - \delta_i)}} \right)$$

که در آن:

$P_\theta$  = احتمال پاسخ درست فرد  $n$  به سوال  $i$ ؛

$P_{\theta-1}$  = احتمال پاسخ نادرست فرد  $n$  به سوال  $i$ ؛

$\theta_n$  = پارامتر توانایی فرد  $n$ ؛

$\delta_i$  = پارامتر دشواری سوال  $i$ ؛ است.

فرمول ارائه شده مدل ساده راش است و برای تحلیل سوالاتی استفاده می‌شود که نمره‌گذاری آن‌ها بصورت صفر و یکی و یا به عبارت دیگر دو ارزشی<sup>۵</sup> است. در کنار مدل‌های دو ارزشی، مدل‌های چند ارزشی<sup>۶</sup> برای تحلیل داده‌های چند ارزشی ایجاد شده‌اند. مدل‌های خانواده راش می‌توانند براساس الگوی پاسخ افراد، با استفاده از روش‌های آماری و ریاضی مانند برآورد بیشینه درست‌نمایی<sup>۷</sup> (ML)، بیشینه پسین<sup>۸</sup> (MAP) و یا پسین مورد انتظار<sup>۹</sup> (EAP)، پارامتر سوال و فرد را مستقل از یکدیگر برآورد کند (امبرستون و رایس<sup>۱۰</sup>، ۱۳۸۸/۲۰۰۰).

خلاصه‌ای از مدل‌های چند ارزشی در جدول یک ارائه شده است (دی‌آیالا<sup>۱۱</sup>، ۲۰۲۲، آندریچ و ماریس<sup>۱۲</sup>، ۲۰۱۹). در مدل‌های راش چند ارزشی، پارامتر شیب برای سوالات در نظر گرفته نمی‌شود و احتمال پاسخ براساس پارامتر توانایی و پارامتر دشواری برآورد می‌شود، اما در مدل‌های غیر راش<sup>۱۳</sup>، پارامتر شیب برای سوالات در نظر گرفته می‌شود.

1. Rasch Measurement
2. Latent Theory
3. Georg Rasch
4. Wind and Hua
5. dicotomous
6. polytomous
7. maximum likelihood
8. maximum a posteriori
9. expected a posteriori
10. Embreston and Reise
11. de Ayala
12. Andrich and Marais
13. Non-Rasch

جدول ۱. مدل‌های چندارزشی

برآورد پارامتر توانایی - برآورد پارامتر دشواری مورد استفاده برای داده‌های گسسته طبقه‌ای (نیازی به یکسان بودن طبقات پاسخ در بخش‌های مختلف آزمون ندارد) مناسب برای پاسخ‌های چند مرحله‌ای	مدل اعتبار جزئی <sup>۱</sup> (PCM)	مدل‌های راش
برآورد پارامتر توانایی - برآورد پارامتر دشواری مورد استفاده برای داده‌های گسسته طبقه‌ای (نیازی به یکسان بودن طبقات پاسخ در بخش‌های مختلف آزمون ندارد) مناسب برای پاسخ‌های طیفی مانند طیف لیکرت	مدل پاسخ مرتب شده <sup>۲</sup> (RSM)	
برآورد پارامتر توانایی - برآورد پارامتر دشواری - برآورد پارامتر تشخیص مورد استفاده برای داده‌های گسسته طبقه‌ای (نیازی به یکسان بودن طبقات پاسخ در بخش‌های مختلف آزمون ندارد) مناسب برای پاسخ‌های چند مرحله‌ای	مدل اعتبار جزئی تعمیم یافته <sup>۳</sup> (GPCM)	مدل‌های غیر راش
برآورد پارامتر توانایی - برآورد پارامتر دشواری - برآورد پارامتر تشخیص مورد استفاده برای داده‌های گسسته طبقه‌ای (نیازی به یکسان بودن طبقات پاسخ در بخش‌های مختلف آزمون ندارد) مناسب برای پاسخ‌های چند مرحله‌ای و طیفی مانند لیکرت	مدل پاسخ مدرج <sup>۴</sup> (GRM)	
برآورد پارامتر توانایی - برآورد پارامتر دشواری - برآورد پارامتر تشخیص مورد استفاده برای داده‌های گسسته که دسته‌های پاسخ در آن ذاتاً مرتب نیست مناسب برای پاسخ‌های اسمی مانند نگرش‌سنج‌ها	مدل پاسخ اسمی <sup>۵</sup> (NR)	

علاوه بر مدل‌های چندارزشی، مدل‌های پیچیده‌تر راش نیز برای تحلیل داده‌های حاصل از اندازه‌گیری‌های دارای ارزیابی ایجاد شدند. بطور کلی در هر اندازه‌گیری، علاوه بر ویژگی‌های فرد و سوال، وجوه دیگری مانند ارزیابی، جنسیت، نوع سوالات، زمان آزمون و ... می‌تواند بر اندازه‌گیری صورت گرفته اثر بگذارد که با استفاده از مدل‌های راش چند وجهی<sup>۶</sup> (MFRM)، می‌توان وجوه مختلف شناسایی شده را تحلیل کرد. بطور خاص تجزیه و تحلیل MFRM نشان می‌دهد که چگونه می‌توان سخت‌گیری یا سهل‌گیری رتبه‌دهنده را برآورد، درجه ثبات رتبه‌دهنده را ارزیابی، نمره هر آزمودنی تحت تاثیر سخت‌گیری رتبه‌دهنده را اصلاح و عملکرد مقیاس رتبه‌بندی را بررسی کرد، همچنین می‌توان تعاملات بالقوه میان وجوه مختلف را تشخیص داد (اکس، ۲۰۰۹، ویند و هوا، ۲۰۲۲).

در راش چند وجهی، برای محاسبه احتمال پاسخ صحیح، علاوه بر ویژگی آزمودنی (پارامتر توانایی) و ویژگی‌های سوال (پارامتر سوال)، پارامترهای مرتبط با رتبه‌دهنده (پارامتر سخت‌گیری) نیز در نظر گرفته می‌شود (دی‌آیالا، ۲۰۲۲، لیناکر، ویند و هوا، ۲۰۲۲). فرمول عمومی راش چند وجهی بصورت زیر است:

1. Partial Ceredit Model
2. Rating Response Model
3. Generaliz Partial Ceredit Model
4. Graded Response Model
5. Nominal Response Model
6. Multi-Faceted Rasch Model

$$\ln \left[ \frac{P_{nij k}}{P_{nij k-1}} \right] = \theta_n - \beta_i - \alpha_j - \tau_k$$

که در آن:

$$P_{nij k} = \text{احتمال دریافت رتبه } k \text{ در سوال } i \text{ توسط ارزیاب } j$$

$$P_{nij k-1} = \text{احتمال دریافت رتبه } k-1 \text{ در سوال } i \text{ توسط ارزیاب } j$$

$$\theta_n = \text{پارامتر توانایی فرد } n$$

$$\beta_i = \text{پارامتر دشواری سوال } i$$

$$\alpha_j = \text{پارامتر سخت‌گیری ارزیاب } j$$

$$\tau_k = \text{دشواری دریافت رتبه } k \text{ نسبت به دریافت رتبه } k-1 \text{ است.}$$

مدل‌های مختلف راش را می‌توان به مدل چند وجهی گسترش داد. با توجه به اینکه در این مدل‌ها پارامتر شیب برای همه سوالات یکسان در نظر گرفته می‌شود، در مدل‌های چندوجهی آن‌ها نیز پارامتر تشخیص برای همه سوالات یکسان است.

پُلْت و همکاران<sup>۱</sup> (۲۰۲۲) در پژوهش خود با عنوان مقایسه نظریه کلاسیک اندازه‌گیری و نظریه راش چند وجهی، به مقایسه نظریه راش چند وجهی با نظریه کلاسیک اندازه‌گیری پرداخته است. در این مقاله پلت علاوه بر معرفی روش‌های کلاسیک تحلیل اثر ارزیاب، به بررسی نتایج تحلیل کلاسیک و راش چند وجهی بر روی داده‌های واقعی آزمون مهارت نوشتاری پرداخته است. پلت برای بررسی کلاسیک ارزیاب‌ها از همبستگی و آزمون  $F$  برای بررسی همسانی درونی و معنی‌داری اختلاف رتبه‌بندی ارزیاب‌ها استفاده کرده است. نتایج تحلیل کلاسیک نشان داد علاوه بر وجود همبستگی بالا (بیشتر از ۰/۷ و یا نزدیک به آن)، اختلاف معنی‌داری بین رتبه‌بندی ارزیاب‌ها وجود دارد. همچنین نتایج راش چند وجهی نشان داد پارامتر شدت ارزیاب از ۰/۴۶- تا ۰/۴۴+ متغیر بوده است. پژوهش ایزانلو و حاجت‌پور (۱۴۰۲) با عنوان بررسی رتبه‌دهی ارزیابان آزمون‌های عملکردی سراسری (طراحی صنعتی، شناخت موسیقی، نمایش عروسکی، طراحی معماری و اسکیس معماری) بر اساس روش‌های کلاسیک و مدل‌های چندوجهی راش، نشان داد که علاوه بر وجود همبستگی قابل قبول بین ارزیاب‌ها، میزان توافق بین آن‌ها قابل قبول نیست. در پژوهش حاضر علاوه بر بررسی وضعیت ارزیابان براساس نظریه کلاسیک و راش چندوجهی، ویژگی روانسنجی سوالات نیز براساس هر دو نظریه مورد بررسی قرار گرفته است.

## روش پژوهش

سازمان سنجش آموزش کشور در فرآیند سنجش داوطلبان مقطع کارشناسی ارشد رشته‌های موسیقی، علاوه بر آزمون تستی از آزمون عملی نیز استفاده می‌کند. با توجه به مطالبی که پیش از این بیان شد و با توجه به استفاده از ارزیاب در آزمون‌های سازمان سنجش، اعتبار نتایج این آزمون‌ها همواره در معرض خطای ارزیاب وجود دارد. با توجه به اهمیت نتایج آزمون‌های عملی سازمان سنجش در پذیرش و یا عدم پذیرش داوطلبان، بررسی و بهبود اعتبار نتایج این آزمون‌ها و نیز تعیین وضعیت ارزیاب‌ها با هدف بهبود ارزیابی‌ها، بسیار ضروری است. بر همین اساس در این پژوهش داده‌های سازمان سنجش آموزش کشور در آزمون عملی سرایش برای داوطلبان رشته آهنگسازی، مورد تحلیل قرار گرفته است.

تحلیل‌های انجام شده با توجه به هدف توصیفی - تحلیلی بوده و از جهت کاربرد، پژوهش کاربردی محسوب می‌شود. داده‌ها براساس نظریه کلاسیک اندازه‌گیری و مدل راش چند وجهی مبتنی بر مدل اعتبار جزئی (PCMFRRM<sup>۱</sup>) و با استفاده از نرم‌افزارهای R نسخه ۴,۰/۳ بسته TAM نسخه ۴-۴/۱، SPSS نسخه ۲۷ و اکسل نسخه ۲۰۱۶، تحلیل شده‌اند.

آزمون عملی سُرایش در سه رشته موسیقی جهانی، موسیقی ایرانی و آهنگسازی وجود دارد. این آزمون شامل سه بخش خواندن ملودی در گام‌های ماژور، خواندن ملودی مینور و وزن خوانی ترکیبی و ساده است که در رشته آهنگسازی هر بخش ۵ نمره دارد. در این آزمون هر داوطلب توسط ۴ ارزیاب مورد ارزیابی قرار می‌گیرد و هر ارزیاب نمرات خود را در فرم جداگانه ثبت می‌نماید. با توجه به اینکه تمام داوطلبان توسط هر چهار ارزیاب مورد ارزیابی قرار می‌گیرند و هر ارزیاب داوطلب را در هر سه بخش ارزیابی می‌کند، طرح ارزیابی مورد استفاده، طرح ضربدری<sup>۲</sup> (کامل) است.

با توجه به اینکه مدل اعتبار جزئی، فقط امکان تحلیل داده‌های چند ارزشی و گسسته را دارد و از سوی دیگر، با توجه به پیوسته بودن نمرات آزمون سُرایش، لازم است پیش از اجرای تحلیل، داده‌ها از حالت پیوسته به گسسته تبدیل شوند، که با توجه به فراوانی نمرات داده شده توسط ارزیاب‌ها، نمره صفر در طبقه ۱ و سایر نمرات به ۱۰ طبقه با فاصله نیم نمره تبدیل شده‌اند (پیوست ۱).

## یافته‌ها

یافته‌های پژوهش در سه بخش ارائه شده است و در هر بخش ابتدا نتایج تحلیل کلاسیک و سپس نتایج راش چندوجهی گزارش شده است.

### تحلیل بخش‌های آزمون

ضرایب دشواری و تشخیص بخش‌های آزمون سُرایش در جدول شماره دو آمده است. براساس نتایج حاصل، دشواری بخش‌های مختلف آزمون از ۰/۳۸ تا ۰/۶۰ متغیر بوده است. در نظریه کلاسیک اندازه‌گیری، ضریب دشواری سوال در دامنه ۰ تا ۱ قرار دارد و به هر میزان که مقدار آن به یک نزدیک‌تر شود، یعنی سوال آسانتر است. در جدول شماره سه، طبقه‌بندی ضرایب دشواری و تشخیص سوالات ارائه شده است. براساس طبقه‌بندی ارائه شده برای ضریب دشواری، می‌توان گفت ضریب دشواری هر سه بخش آزمون سُرایش در سطح متوسطی قرار دارد.

جدول ۲. ضرایب دشواری و تشخیص بخش‌های آزمون سُرایش براساس CTT و PCMFRRM

PCMFRRM		CTT		بخش
خطای استاندارد برآورد	پارامتر دشواری	ضریب تشخیص	ضریب دشواری	
۰/۰۴	۰/۵۳	۰/۴۸	۰/۳۸	خواندن ملودی در گام‌های ماژور
۰/۰۳	۰/۰۶	۰/۶۱	۰/۴۱	خواندن ملودی مینور
۰/۰۵	-۰/۵۹	۰/۳۹	۰/۶۰	وزن خوانی ترکیبی و ساده
۰/۰۴	۰	۰/۴۹	۰/۴۶	میانگین

۱. Partial Credit Multi-Faceted Rasch Model

۲. Crossed Desinge

جدول ۳. طبقه‌بندی دشواری و تشخیص براساس CTT و PCMFRM

PCMFRM		CTT			
پارامتر دشواری		ضریب تشخیص		ضریب دشواری	
خیلی آسان	$b \leq -2$	ضعیف	$0 < D \leq 0/02$	بسیار سخت	$P \leq 0/30$
آسان	$-2 < b \leq +0/50$	متوسط	$0/20 < D \leq 0/40$	متوسط	$0/30 < P \leq 0/80$
متوسط	$-0/5 < b \leq +0/5$	خوب	$0/04 < D$	بسیار آسان	$0/80 < P$
سخت	$+0/5 < b \leq +2$				
خیلی سخت	$+2 < b$				

باتوجه به طبقه‌بندی ارائه شده در جدول شماره دو، براساس تحلیل کلاسیک، بخش‌های خواندن ملودی در گام‌های ماژور و خواندن ملودی مینور از نظر ضریب تشخیص در سطح خوب و بخش وزن خوانی ترکیبی و ساده در سطح متوسط قرار دارد. سخت‌ترین بخش آزمون، بخش خواندن ملودی در گام‌های ماژور و آسان‌ترین بخش، بخش وزن خوانی ترکیبی و ساده است. همچنین بخش خواندن ملودی مینور، بیشترین ضریب تشخیص را دارد.

همان‌گونه که پیش از این بیان شد، برای تحلیل راش چند وجهی، از مدل اعتبار جزئی استفاده شده است. در مدل اعتبار جزئی، آنچه تعیین کننده است، پارامتر دشواری است، بنابراین در تحلیل‌های مبتنی بر این مدل، برای بخش‌های آزمون تنها پارامتر دشواری ارائه می‌شود. در نظریه راش، پارامتر دشواری در پیوستاری از ۳- تا ۳+ قرار دارد و هر چقدر از سمت منفی به سمت مثبت حرکت کنیم، میزان دشواری افزایش پیدا می‌کند. مقادیر دشواری نزدیک به ۲- به سوالاتی مربوط می‌شود که خیلی آسان هستند و مقادیر دشواری ۲+ به سوالاتی مربوط می‌شود که خیلی سخت هستند.

همان‌گونه که در جدول شماره دو آمده است، پارامتر دشواری بخش‌های آزمون سُرایش، از ۰/۵۹- تا ۰/۵۳+ متغیر بوده. مطابق با جدول شماره پنج می‌توان گفت، هر سه بخش از نظر پارامتر دشواری متوسط هستند، وزن خوانی ترکیبی و ساده آسان‌ترین بخش و خواندن ملودی در گام‌های ماژور سخت‌ترین بخش در این آزمون بودند.

### تحلیل ارزیاب

در نظریه کلاسیک برای تحلیل ارزیاب‌ها، از شباهت نمره‌گذاری (سازگاری) و میزان توافق (اجماع) میان ارزیاب‌ها استفاده می‌شود. در مبانی نظری، ضرایب همبستگی و تاو بی‌کندال به عنوان شاخص‌های سازگاری و ضرایب کاپا و ضریب توافق به عنوان شاخص‌های اجماع در نظر گرفته می‌شوند. شاخص اجماع نشان‌دهنده اجماع و توافق داوران در ترتیب افراد و شاخص سازگاری نشان‌دهنده رابطه خطی بین نمرات داده‌شده توسط هر ارزیاب است (اکس، ۲۰۰۵). در پژوهش حاضر از ضریب کاپا و ضریب همبستگی پیرسون برای تحلیل ارزیاب براساس نظریه کلاسیک استفاده شده است.

طبق طبقه‌بندی ارائه شده توسط ولف<sup>۱</sup>، مقادیر کاپای کمتر از ۰/۲۰، قابل قبول نیستند (وُلف، ۱۹۹۸). طبقه‌بندی ضریب کاپا در جدول شماره چهار و ضرایب کاپای محاسبه شده برای هر بخش از آزمون سُرایش، در جدول شماره پنج نشان داده شده است.

1. Wolf

جدول ۴. طبقه‌بندی ضریب کاپا

ضعیف	$K \leq 0/0$
کم	$0 < K \leq 0/20$
قابل قبول	$0/20 < b \leq 0/40$
متوسط	$0/40 < b \leq 0/60$
زیاد	$0/60 < b \leq 0/80$
عالی	$0/80 < b$

جدول ۵. ضریب کاپا در هر بخش

بخش	ضریب کاپا	خطای استاندارد	سطح معنی داری
خواندن ملودی در گام‌های ماژور	۰/۵۸	۰/۰۳	$< 0/0001$
خواندن ملودی مینور	۰/۵۷	۰/۰۲	$< 0/0001$
وزن خوانی ترکیبی و ساده	۰/۴۶	۰/۰۲	$< 0/0001$

با توجه به ضرایب کاپای محاسبه شده در هر سه بخش که بیش از ۰/۴۰ بدست آمده است، می‌توان گفت ضریب کاپا برای هر سه بخش قابل قبول و در حد متوسط بوده است. این ضریب برای بخش‌های خواندن ملودی در گام‌های ماژور و خواندن ملودی مینور، بالای ۰/۵۵ است. از دیگر ضرایب مورد استفاده برای تحلیل ارزیاب‌ها در نظریه کلاسیک، ضریب همبستگی میان ارزیاب‌ها است. ضرایب همبستگی پیرسون میان ارزیاب‌ها به تفکیک بخش‌های آزمون، در جدول شماره شش نشان داده شده است.

جدول ۶. ضریب همبستگی ارزیاب‌ها

وزن خوانی ترکیبی و ساده				خواندن ملودی مینور				خواندن ملودی در گام‌های ماژور				
			۱				۱				۱	ارزیاب اول
		۱	۰/۹۴				۱			۱	۰/۹۶	ارزیاب دوم
	۱	۰/۹۷	۰/۹۳			۱	۰/۹۷		۱	۰/۹۵	۰/۹۷	ارزیاب سوم
۱	۰/۹۴	۰/۹۴	۰/۹۸	۱	۰/۹۶	۰/۹۶	۰/۹۸	۱	۰/۹۷	۰/۹۵	۰/۹۸	ارزیاب چهارم

مقدار این ضریب در هر سه بخش آزمون، برای هر چهار ارزیاب بیش از ۰/۹۰ و در سطح ۰/۰۱ معنادار بدست آمده است. نتایج ضریب همبستگی نشان می‌دهد که ارزیاب‌ها در رتبه‌بندی افراد دارای توافق بسیار بالایی هستند. به عبارت دیگر اگر هر ارزیاب بطور جداگانه افراد را از ضعیف به قوی مرتب کند، ترتیب افراد در بین هر چهار ارزیاب بالای ۹۰ درصد شباهت خواهد داشت. در مدل راش چند وجهی، پارامتر ارزیاب مانند پارامتر دشواری و توانایی، در پیوستاری از ۳- تا ۳+ قرار دارد و هر چقدر از سمت منفی به سمت مثبت حرکت کنیم، میزان سخت‌گیری ارزیاب افزایش پیدا می‌کند. مقادیر نزدیک به ۳- به ارزیاب‌های بسیار سهل‌گیر و مقادیر ۳+ به ارزیاب‌های بسیار سخت‌گیر مربوط می‌شود.

جدول ۷. پارامتر ارزیاب‌ها براساس PCMFRM

خطای استاندارد برآورد	پارامتر سختگیری	ارزیاب
۰/۰۴	۰/۰۷	ارزیاب اول
۰/۰۴	-۰/۰۹	ارزیاب دوم
۰/۰۴	۰/۰۰	ارزیاب سوم
۰/۰۷	۰/۰۲	ارزیاب چهارم

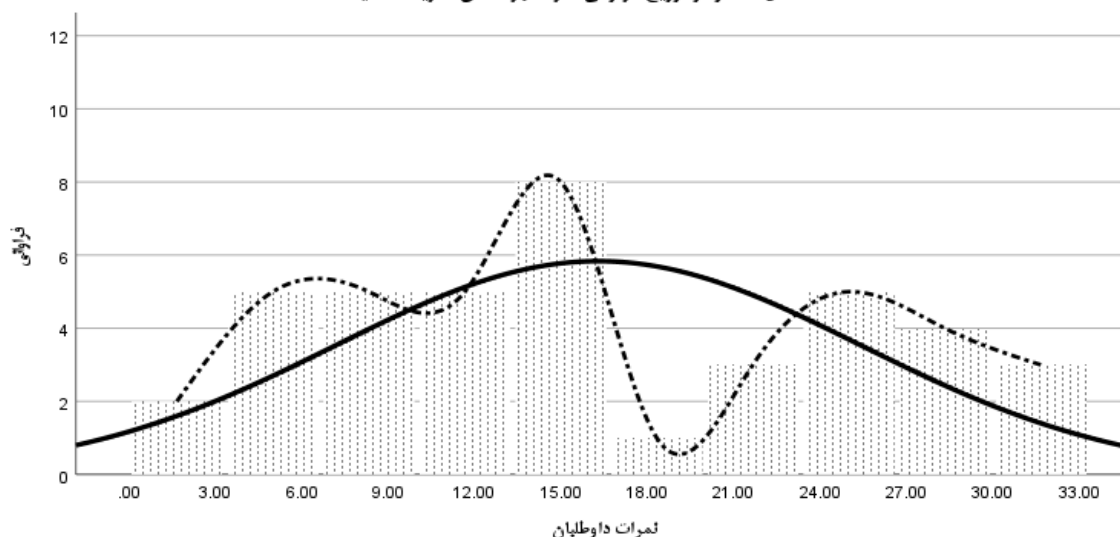
پارامتر ارزیاب برآورد شده برای هر چهار ارزیاب در جدول شماره نه آمده و نشان می‌دهد که شدت سخت‌گیری ارزیاب‌ها از ۰/۰۹- تا ۰/۰۷+ متغیر بوده است. براساس نتایج بدست آمده می‌توان گفت شدت سخت‌گیری ارزیاب‌ها در حد متوسط بوده، ارزیاب اول سخت‌گیرترین و ارزیاب دوم سهل‌گیرترین ارزیاب در بین چهار ارزیاب بوده است.

### تحلیل داوطلبان

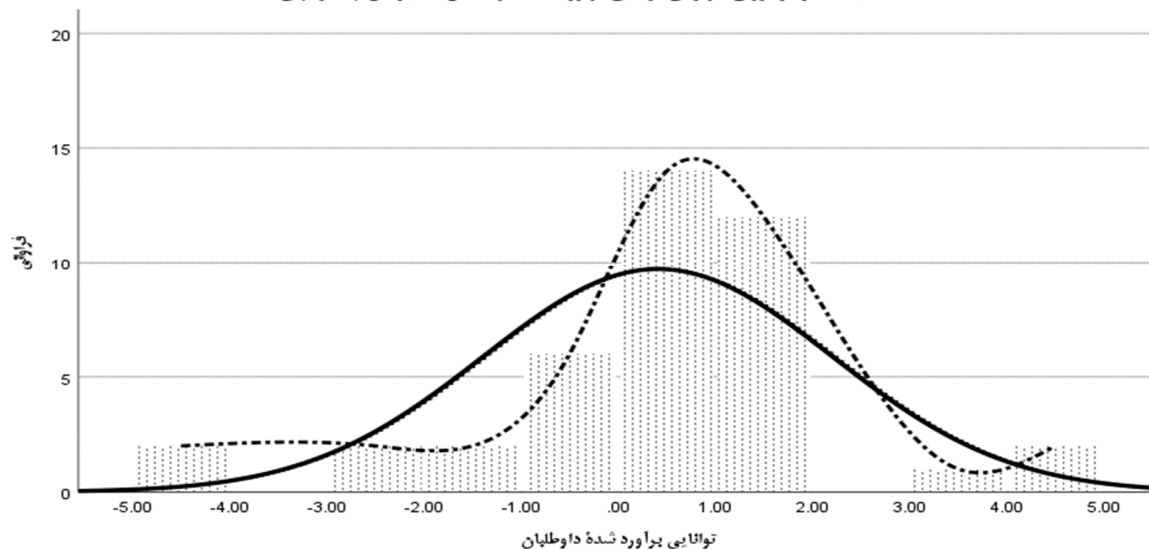
برای مقایسه توانایی برآورد شده از تحلیل راش چند وجهی با نمره‌گذاری کلاسیک از ضریب همبستگی پیرسون و همچنین مقایسه نمودارهای فراوانی استفاده شده است. ضریب همبستگی پیرسون ۰/۹۰+ محاسبه شده است، یعنی نمرات کلاسیک می‌تواند، ۸۱ درصد از واریانس پارامتر توانایی داوطلبان را تبیین کند.

شکل‌های یک و دو، به ترتیب نمودار فراوانی نمرات کلاسیک و توانایی برآورد شده را نشان می‌دهند. همانگونه که نمودار فراوانی نمرات نشان می‌دهد، توزیع نمرات سه‌نمایی بوده و دارای چولگی به چپ هستند. به عبارت دیگر، بیشتر نمرات دریافت شده پائین‌تر از میانگین نمرات است. توزیع پارامتر توانایی تک‌نمایی بوده و به منحنی نرمال نزدیک است، با این وجود نمودار توزیع توانایی نشان می‌دهد که اکثر داوطلبان از نظر توانایی برآورد شده بالاتر از میانگین قرار داشته و از توانایی نسبتاً بالایی برخوردار هستند. نمره و پارامتر توانایی هر یک از داوطلبان در جدول شماره هشت آمده است.

شکل ۱. نمودار توزیع فراوانی نمرات براساس نظریه کلاسیک



شکل ۲. نمودار توزیع فراوانی توانایی برآورد شده براساس مدل رانش چند وجهی



جدول ۸. نمره و توانایی برآورد شده داوطلبان

نمره	توانایی	داوطلب	نمره	توانایی	داوطلب	نمره	توانایی	داوطلب
۱۰/۲۵	۰/۱۳	۱۴۰۱۲۹	۲۰/۷۵	۱/۰۰	۱۴۰۱۱۵	۳۳	۴/۲۴	۱۴۰۱۱
۸/۲۵	-۰/۲۸	۱۴۰۱۳۰	۱۷	۰/۷۴	۱۴۰۱۱۶	۳۳	۴/۲۴	۱۴۰۱۲
۸/۲۵	-۰/۲۸	۱۴۰۱۳۱	۱۶/۵	۰/۷۰	۱۴۰۱۱۷	۳۲/۷۵	۳/۵۱	۱۴۰۱۳
۸	-۰/۳۵	۱۴۰۱۳۲	۱۶/۵	۰/۷۰	۱۴۰۱۱۸	۲۸/۲۵	۱/۹۳	۱۴۰۱۴
۷/۷۵	-۰/۴۲	۱۴۰۱۳۳	۱۶/۲۵	۰/۶۸	۱۴۰۱۱۹	۲۸/۲۵	۱/۹۳	۱۴۰۱۵
۷/۲۵	-۰/۵۹	۱۴۰۱۳۴	۱۵/۷۵	۰/۶۵	۱۴۰۱۲۰	۲۷/۷۵	۱/۸۳	۱۴۰۱۶
۶/۲۵	-۰/۹۹	۱۴۰۱۳۵	۱۵/۷۵	۰/۶۵	۱۴۰۱۲۱	۲۷	۱/۶۹	۱۴۰۱۷
۵/۲۵	-۱/۵۳	۱۴۰۱۳۶	۱۵/۵	۰/۶۳	۱۴۰۱۲۲	۲۵/۷۵	۱/۴۹	۱۴۰۱۸
۴/۷۵	-۱/۸۸	۱۴۰۱۳۷	۱۴/۲۵	۰/۵۴	۱۴۰۱۲۳	۲۵/۵	۱/۴۶	۱۴۰۱۹
۴/۲۵	-۲/۳۱	۱۴۰۱۳۸	۱۴	۰/۵۲	۱۴۰۱۲۴	۲۴/۷۵	۱/۳۶	۱۴۰۱۱۰
۴	-۲/۵۷	۱۴۰۱۳۹	۱۳/۲۵	۰/۴۶	۱۴۰۱۲۵	۲۴/۲۵	۱/۳۰	۱۴۰۱۱۱
۳	-۴/۳۴	۱۴۰۱۴۰	۱۲/۵	۰/۴۰	۱۴۰۱۲۶	۲۳/۵	۱/۲۳	۱۴۰۱۱۲
۳	-۴/۳۴	۱۴۰۱۴۱	۱۰/۷۵	۰/۲۰	۱۴۰۱۲۷	۲۳/۲۵	۱/۲۰	۱۴۰۱۱۳
			۱۰/۵	۰/۱۷	۱۴۰۱۲۸	۲۱/۷۵	۱/۰۷	۱۴۰۱۱۴

در نظریه پرسش-پاسخ چندین روش برای مدل‌سازی اثرات ارزیاب پیشنهاد شده است که بیشتر مبتنی بر مفهوم آیت‌م مجازی است. آیت‌م مجازی مجموعه‌ای از تمام ترکیبات میان ارزیاب‌ها و آیت‌ها است. در این پژوهش، ۳ بخش و چهار ارزیاب وجود دارد و با توجه به اینکه از طرح ارزیابی متقاطع استفاده شده است، بنابراین ۱۲ بخش مجازی (۳ بخش  $\times$  ۴ ارزیاب) در تحلیل وجود دارد. در جدول شماره نه، شاخص‌های برازش مدل در دوازده بخش مجازی آمده است.

جدول ۹. شاخص برازش آیت‌های مجازی

Infit_p	Infit_t	Infit	شماره بخش مجازی	بخش مجازی
۰/۲۳	-۱/۲	۰/۶۷	۱	بخش ۱-ارزیاب ۱
۰/۰۵	-۱/۹۷	۰/۵۴	۲	بخش ۱-ارزیاب ۲
۰/۰۱	-۲/۵۶	۰/۴۲	۳	بخش ۱-ارزیاب ۳
۰/۰۸	-۱/۷۵	۰/۵۶	۴	بخش ۱-ارزیاب ۴
۰/۷۵	-۰/۳۲	۰/۹۰	۵	بخش ۲-ارزیاب ۱
۰/۹۱	۰/۱۱	۱/۰۱	۶	بخش ۲-ارزیاب ۲
۰/۳۸	۰/۸۸	۱/۲۱	۷	بخش ۲-ارزیاب ۳
۰/۱۷	-۱/۳۸	۰/۶۷	۸	بخش ۲-ارزیاب ۴
۰/۱۳	۱/۵۱	۱/۴۰	۹	بخش ۳-ارزیاب ۱
۰/۱۰	۱/۶۶	۱/۴۳	۱۰	بخش ۳-ارزیاب ۲
۰/۰۱	۲/۸۴	۱/۸۱	۱۱	بخش ۳-ارزیاب ۳
۰/۲۴	۱/۱۸	۱/۲۸	۱۲	بخش ۳-ارزیاب ۴
۰/۹۹				میانگین
۰/۴۳				انحراف استاندارد

برای بررسی برازش مدل با داده‌ها از شاخص‌های برازش Infit استفاده شده است. آماره Infit تفاوت بین نمرات مشاهده شده و نمرات مورد انتظار است. چنانچه تفاوت بین نمرات مشاهده شده داوطلبان و نمرات مورد انتظار آن‌ها زیاد باشد، می‌تواند دقت اندازه‌گیری را تهدید کند. مقدار قابل قبول برای شاخص Infit در فاصله دو انحراف استاندارد از میانگین است، بنابراین دامنه قابل قبول برای این شاخص در محدوده ۱/۸۵ تا ۰/۱۳ است. براساس مقادیر Infit در جدول شماره یازده می‌توان گفت برازش مدل با داده در تمام آیت‌های مجازی قابل قبول بوده است.

در جدول شماره ده، شاخص‌های برازش داوطلبان آمده است. مقادیر Infit که بیش از دو انحراف استاندارد با میانگین فاصله داشته باشد، نشان‌دهنده توافق کم ارزیاب‌ها در مورد داوطلب است. بنابراین در خصوص داوطلبانی که مقدار Infit آن‌ها بیش از ۲/۱۰ است، توافق کمی بین ارزیاب‌ها وجود دارد. همانگونه که در جدول مشخص است، ارزیاب‌ها در ارزیابی داوطلب شماره ۱ و ۱۰ توافق کمی داشته‌اند.

جدول ۱۰. شاخص برازش توانایی داوطلبان

infitPerson	داوطلب	infitPerson	داوطلب
۰/۱۷	۲۲	۲/۷۴	۱
۰/۳۶	۲۳	۱/۶۵	۲
۰/۰۵	۲۴	۰/۲۶	۳
۰/۵۷	۲۵	۰/۸۷	۴
۰/۳۰	۲۶	۱/۱۷	۵
۰/۲۳	۲۷	۱/۷۶	۶

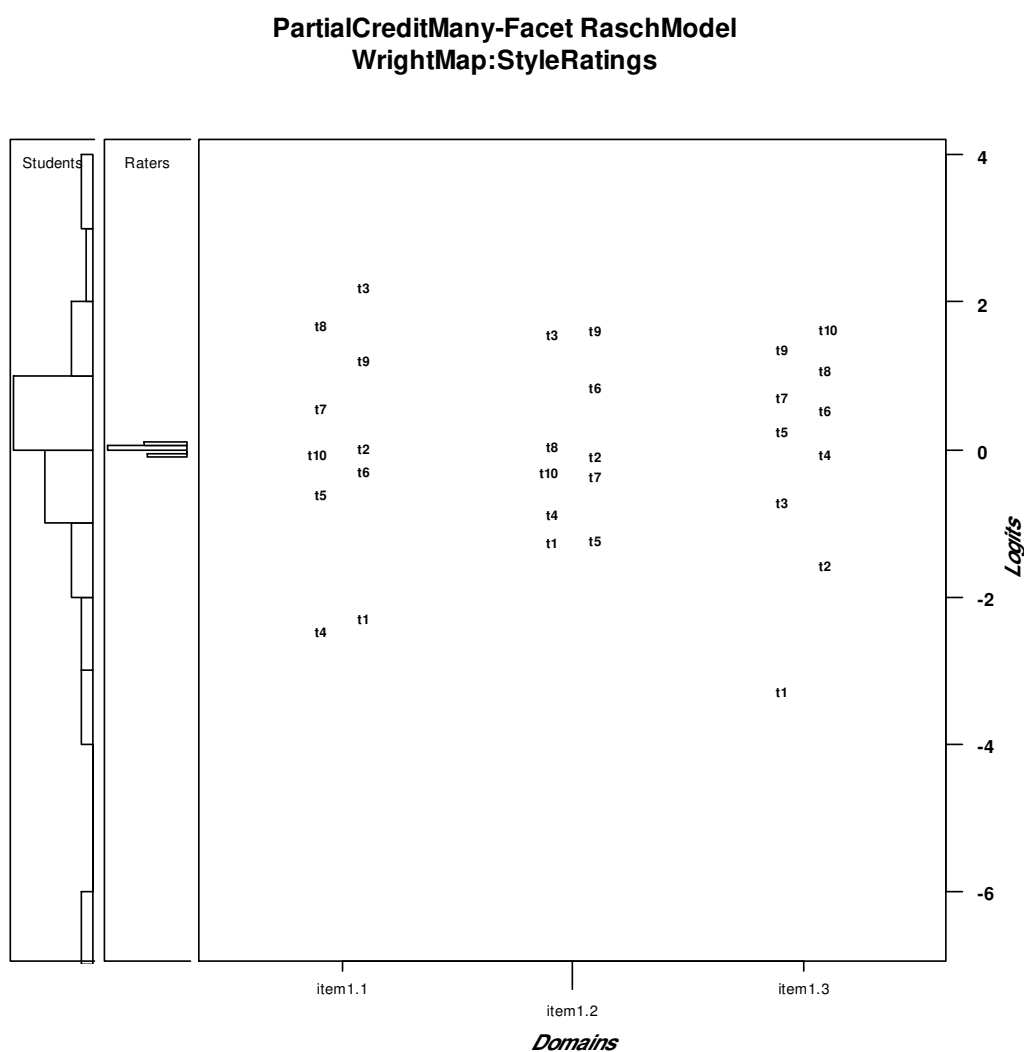
۰/۱۴	۲۸	۰/۹۴	۷
۱/۱۹	۲۹	۰/۰۵	۸
۰/۰۷	۳۰	۰/۵۴	۹
۱/۱۵	۳۱	۲/۸۷	۱۰
۰/۴۷	۳۲	۱/۱۷	۱۱
۰/۲۶	۳۳	۱/۳۰	۱۲
۰/۷۲	۳۴	۰/۷۷	۱۳
۰/۵۱	۳۵	۰/۱۷	۱۴
۰/۴۶	۳۶	۰/۷۶	۱۵
۰/۰۹	۳۷	۰/۲۵	۱۶
۰/۳۸	۳۸	۰/۹۰	۱۷
۰/۰۹	۳۹	۱/۷۹	۱۸
۱/۲۷	۴۰	۰/۷۶	۱۹
۰/۳۵	۴۱	۰/۵۱	۲۰
		۰/۲۸	۲۱
۰/۷۴		میانگین	
۰/۶۸		انحراف استاندارد	

خلاصه‌ای از وضعیت پارامتر توانایی داوطلبان، پارامتر شدت ارزیاب‌ها و دشواری بخش‌های مختلف آزمون در نقشه رایت<sup>۱</sup> قابل مشاهده است. همانگونه که در تصویر شماره دو مشاهده می‌شود، دامنه توانایی داوطلبان از بسیار ضعیف تا بسیار قوی بوده و اکثر داوطلبان از نظر توانایی در سطح متوسط بودند. همچنین هر چهار ارزیاب از نظر شدت سختگیری، در حد متوسط بودند.

یکی از خروجی‌های قابل توجه در تحلیل‌های چند ارزشی‌راش، بررسی میزان دشواری دستیابی به هر یک از نمرات است. همانگونه که در نقشه رایت مشاهده می‌شود، در بخش اول و دوم آزمون سُرایش، دشواری دریافت نمرات مختلف نامرتب است اما در بخش سوم، دشواری دریافت نمرات مرتب است به این صورت که با افزایش نمره، دشواری دریافت آن نیز افزایش یافته است.

1 Wright Map

شکل ۳. نقشه رایت



### بحث و نتیجه‌گیری

در این پژوهش، آزمون عملی سُرایش برای داوطلبان رشته آهنگسازی مورد تحلیل قرار گرفت. این آزمون شامل ۳ بخش ۵ نمره‌ای است و از چهار ارزیاب برای ارزیابی داوطلبان استفاده شده است. تعداد کل داوطلبان حاضر در این آزمون ۴۱ نفر بوده که فرم ارزیابی همه داوطلبان مورد تحلیل قرار گرفت. یافته‌های حاصل از تحلیل کلاسیک نشان می‌دهد که بطور کلی آزمون از لحاظ دشواری نسبتاً دشوار بوده است. سخت‌ترین بخش آزمون خواندن ملودی در گام‌های ماژور و آسان‌ترین بخش، وزن خوانی ترکیبی و ساده بود. در آزمون‌های توانایی، کمترین میزان ضریب تشخیص قابل قبول در برخی منابع ۰/۲۰ و در برخی دیگر از منابع ۰/۳۰ گزارش شده است. با توجه به اینکه ضریب تشخیص هر سه بخش آزمون بالای ۰/۳۰ است، می‌توان گفت هر سه بخش از نظر ضریب تشخیص در سطح قابل قبولی قرار دارند.

برای تحلیل ارزیاب در نظریه کلاسیک اندازه‌گیری از روش‌های مختلفی همچون ضریب همبستگی، ضریب توافق، ضریب کاپا، آزمون تی یا تحلیل واریانس یک راهه، استفاده می‌شود. در این پژوهش از دو روش ضریب کاپا برای بررسی پایایی توافق ارزیابان و ضریب همبستگی برای

بررسی میزان همسویی آن‌ها در ارزیابی انجام شده، استفاده شده است. ضریب کاپای بدست آمده برای بخش خواندن ملودی در گام‌های ماژور برابر با  $0/58$ ، برای بخش خواندن ملودی مینور برابر با  $0/57$  و برای بخش وزن خوانی ترکیبی و ساده برابر با  $0/46$  بدست آمده است. مقدار ضریب کاپا عددی از  $0$  تا  $1$  است و مقادیر بیشتر نشان دهنده پایایی بالاتر ارزیاب‌ها است، یعنی ارزیاب‌ها در رتبه‌بندی افراد توافق بیشتری دارند. براساس طبقه‌بندی کیفی ضریب کاپا، مقادیر بین  $0/41$  تا  $0/60$  نشان‌دهنده توافق متوسط بین ارزیابان است. بر این اساس می‌توان گفت در آزمون عملی سُرّایش، توافق ارزیاب‌ها در سطح متوسط است، از سوی دیگر ضریب همبستگی حاصل بین ارزیاب‌ها در هر سه بخش آزمون بیش از  $0/90$  است که نشان‌دهنده همبستگی بسیار بالا بین آن‌ها است. اختلاف بین ضریب کاپا و ضریب همبستگی بین ارزیاب‌ها را می‌توان چنین تحلیل کرد؛ ضریب کاپا نشان‌دهنده توافق ارزیاب‌ها در نمره اختصاص داده شده به افراد است، در حالی که ضریب همبستگی نشان‌دهنده ترتیب افراد از نظر ارزیاب‌ها است. به عبارت دیگر اگر داوطلبان را براساس ارزیابی هر یک از ارزیاب‌ها بصورت جداگانه رتبه‌بندی کنیم، رتبه‌بندی هر چهار ارزیاب یکسان خواهد بود، اما نمره‌ای که هر یک از ارزیاب‌ها داده‌اند متفاوت است، برای مثال داوطلب  $11-14$  نمرات متفاوتی از هر ارزیاب دریافت کرده است اما از نظر هر چهار ارزیاب در رتبه اول قرار می‌گیرد. پژوهش‌های پیشین نیز اختلاف بین همبستگی و ضریب کاپا را گزارش کرده‌اند (پلت و همکاران،  $2022$  و ایزانلو و حاجت‌پور،  $1402$ ). واقع می‌توان گفت اگرچه ترتیب داوطلبان از نظر توانایی در بین ارزیاب‌ها یکسان است، اما الزاماً نمره‌ایی که به داوطلب‌ها داده‌اند، یکسان نیست. در این پژوهش‌ها از ضریب همبستگی به عنوان توافق نسبی و از ضریب کاپا به عنوان توافق مطلق نام برده شده است. براساس نتایج تحلیل راش چند وجهی، بخش‌های مختلف آزمون سُرّایش از نظر دشواری در سطح متوسطی بودند، همچنین پارامتر ارزیاب نشان می‌دهد که هر چهار ارزیاب از نظر سخت‌گیری و سهل‌گیری در میانه قرار داشتند. پارامتر ارزیاب در محدوده  $3-$  تا  $3+$  قرار دارد که هر چه از  $3-$  به  $3+$  نزدیکتر شود، سخت‌گیری ارزیاب بیشتر خواهد بود. با توجه به اینکه پارامتر ارزیاب برای هر چهار ارزیاب در محدود  $0/09-$  تا  $0/07+$  برآورد شده است، باید گفت که ارزیاب‌ها میانه‌رو بوده و دارای خطای سهل‌گیری و یا سخت‌گیری نیستند که این مسئله را می‌توان در همبستگی بالای نمرات کلاسیک و پارامتر توانایی نیز مشاهده کرد.

بطور کلی نتایج این پژوهش نشان می‌دهد که اگر چه همبستگی ارزیاب‌ها می‌تواند بالا باشد، اما ضریب همبستگی به تنهایی نمی‌تواند برای بررسی وضعیت ارزیاب‌ها شاخص مناسبی باشد، از سوی دیگر ضرایب کاپا، همبستگی و سایر شاخص‌های مبتنی بر نظریه کلاسیک برای بررسی ارزیاب، تنها می‌تواند وضعیت گروهی ارزیاب‌ها را از نظر همسویی و همخوانی نشان دهد، در حالی که مدل راش چند وجهی می‌تواند وضعیت هر یک از ارزیاب‌ها را از نظر میزان سخت‌گیری و یا سهل‌گیری نشان دهد. همچنین در مدل راش چند وجهی، پارامتر دشواری سوال، مستقل از عملکرد ارزیاب برآورد می‌شود، به عبارت دیگر عملکرد ارزیاب تاثیری در میزان دشواری سوال ندارد. بطور کلی اگر چه ترتیب دشواری بخش‌های آزمون در هر دو روش تحلیل یکسان است (خواندن ملودی در گام‌های ماژور، خواندن ملودی مینور، وزن خوانی ترکیبی و ساده)، اما میزان دشواری آن‌ها متفاوت بدست آمده است. براساس نتایج حاصل از نظریه کلاسیک، بخش خواندن ملودی در گام‌های ماژور سخت، خواندن ملودی مینور متوسط و وزن خوانی ترکیبی و ساده آسان محاسبه شده‌اند، در حالی که براساس نتایج راش چند وجهی، دشواری هر سه بخش در حد متوسط برآورد شده است. از سوی دیگر براساس نظریه کلاسیک، برای هر بخش ضریب تشخیص محاسبه شده است، در حالی که در تحلیل راش چند وجهی با توجه به استفاده از مدل اعتبار جزئی، پارامتر تشخیص برای همه بخش‌ها یکسان و برابر با یک در نظر گرفته می‌شود. بنابراین می‌توان گفت در تحلیل‌هایی که میزان تشخیص بخش‌ها و یا سوالات آزمون اهمیت تعیین کننده دارد، استفاده از نظریه کلاسیک نسبت به راش چند وجهی مبتنی بر اعتبار جزئی، می‌تواند آگاهی‌دهنده‌تر باشد، اما در مورد تحلیل ارزیاب و دشواری بخش‌ها یا سوالات آزمون، تحلیل‌های راش چند وجهی آگاهی‌دهنده‌تر از نتایج کلاسیک خواهد بود. از سوی دیگر با توجه به توزیع توانایی داوطلبان، انتظار می‌رود که بیشترین فراوانی نمرات در سطح متوسط باشد، اما همان‌گونه که در نمودار یک مشخص است، بیشترین فراوانی در نمرات پائین‌تر از میانگین قرار دارد. تفاوت توزیع فراوانی نمرات با توزیع فراوانی توانایی افراد، نشان می‌دهد که تعیین سطح توانایی افراد در نظریه کلاسیک تحت تاثیر وجوه مختلف از جمله دشواری سوالات و همچنین عملکرد ارزیاب قرار گیرد و افرادی که از نظر پارامتر توانایی در یک محدوده قرار دارند، ممکن است نمرات متفاوتی را دریافت کنند. اگر چه ترتیب افراد براساس نتایج کلاسیک و راش چند وجهی یکسان است، اما نمرات داده شده به داوطلبان متناسب با توانایی واقعی آن‌ها نبود که دلیل آن مستقل نبودن نمره افراد از سختی سوالات و عملکرد ارزیاب در نظریه کلاسیک اندازه‌گیری است. با توجه به اینکه نظریه راش چند وجهی می‌تواند توانایی افراد را بگونه‌ای برآورد کند که تحت تاثیر عملکرد ارزیاب قرار نداشته باشد، و از سوی دیگر عدم امکان بررسی و حذف اثر ارزیاب بر روی نمره افراد در نظریه کلاسیک اندازه‌گیری، نظریه راش چند وجهی در بررسی وضعیت داوطلبان، آگاهی‌دهنده‌تر از نظریه کلاسیک اندازه‌گیری خواهد بود.

## References

- Allen, M. J., & Yen, W. M. (۲۰۰۲). *Introduction to Measurement Theory*. Prospect Heights, IL: Waveland Press. Translated in persian, ۲۰۱۳, SAMT
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education (Eds.). (۲۰۱۴). *Standards for educational and psychological testing*. American Educational Research Association. Translated in persian, ۲۰۱۹, Tehran university.
- Andrich D. & Marais I. (۲۰۱۹). *A course in rasch measurement theory: measuring in the educational social and health sciences*. Springer. <https://doi.org/10.1007/978-981-13-7496-8>
- de Ayala, R. J. (۲۰۲۲). *The theory and practice of item response theory* (۲st ed.). Guilford Press.
- Eckes, T. (۲۰۱۵). *Introduction to Many-Facet Rasch Measurement: Analyzing and Evaluating Rater-Mediated Assessments* (۲nd ed.). New York: Peter Lang.
- Embretson, Susan & Reise, S.. (۲۰۰۰). *Item Response Theory For Psychologists*. Translated in persian, Roshd, ۲۰۰۹, Tehran.
- Ezanloo, B., & Hajatpour, S. (۲۰۲۳). An Investigation of the Evaluators' Ratings of the Performance Exams in the Field of Arts Using Multi-Faceted Rasch Model. *Educational Measurement and Evaluation Studies*, ۱۲(۴۲), ۱۰۰-۱۲۳. doi: ۱۰/۲۲۰۳۴/emes.۲۰۲۳/۵۲۸۱۶۱.۲۳۴۴
- Hambleton, Ronald K. & Swaminathan, Hariharan. & Rogers, H. Jane. (۱۹۹۱). *Fundamentals of item response theory*. Newbury Park, Calif: Sage Publications, Translated in persian, ۲۰۱۰, Alameh tabai university.
- Hombo, C. M., Donoghue, J. R., & Thayer, D. T. (۲۰۰۱). A simulation study of the effect of rater designs on ability estimation (ETS Research Report No. RR-۰۱-۰۵). Princeton, NJ: Educational Testing Service.
- Keeves J. P. (۱۹۹۷). Introduction: Advances in Measurement in Education. In Keeves J. P. (Ed.), *Educational research methodology and measurement: an international handbook* (۲nd ed.) (pp. ۷۰۵-۷۱۲). Pergamon.
- Lee, M., & Cha, D. (۲۰۱۶). A comparison of generalizability theory and many facet Rasch measurement in an analysis of mathematics creative problem solving test. *Journal of Curriculum Evaluation*, ۱۹(۲), ۲۵۱-۲۷۹
- Li, G., Pan, Y., & Wang, W. (۲۰۲۱). Using Generalizability Theory and Many-Facet Rasch Model to Evaluate In-Basket Tests for Managerial Positions. *Frontiers in psychology*, ۱۲, ۶۶۰۵۵۳. <https://doi.org/10.3389/fpsyg.2021.660553>
- Polat, M., Sölpük Turhan, N., & Toraman, Çetin . (۲۰۲۲). Comparison of Classical Test Theory vs. F Theory in writing assessment. *Pegem Journal of Education and Instruction*, ۱۲(۲), ۲۱۳-۲۲۵. <https://doi.org/10.4۷۷۵۰/pegegog.۱۲/۰۲.۲۱>
- Robitzsch, A., & Steinfeld, J. (۲۰۱۸). Item response models for human ratings: Overview, estimation methods, and implementation in R. *Psychological Test and Assessment Modeling*, ۶۰(۱), ۱۰۱-۱۳۹.

- Taylor, C. S. (۲۰۱۳). *Validity and validation*. Oxford University Press, Translated in persian, ۲۰۱۰, Alameh tabai university.
- Wind, S., & Hua, C. (۲۰۲۲). *Rasch Measurement Theory Analysis in R* (۱st ed.). Chapman and Hall/CRC. <https://doi.org/10.1201/9781003174666>.
- Wolf. R.M. (۱۹۹۷). Rating Scales. In Keeves J. P. (Ed.), *Educational research methodology and measurement: an international handbook* (۲nd ed.) (pp. ۹۵۸-۹۶۵). Pergamon.

## پیوست ۱.

جدول ۱. فروانی رتبه‌های داده شده توسط هر ارزیاب

نمره	خواندن ملودی در گام‌های ماژور				خواندن ملودی مینور				وزن خوانی ترکیبی و ساده			
	ارزیاب ۱	ارزیاب ۲	ارزیاب ۳	ارزیاب ۴	ارزیاب ۱	ارزیاب ۲	ارزیاب ۳	ارزیاب ۴	ارزیاب ۱	ارزیاب ۲	ارزیاب ۳	ارزیاب ۴
۰	۱۰	۸	۸	۹	۱۱	۸	۱۰	۹	۳	۲	۲	۱
۰/۲۵	-	۶	۳	-	-	۴	۴	-	۱	۱	-	-
۰/۵۰	۱۰	۵	۱۰	۱۲	۴	۳	۷	۴	۳	۲	۳	۳
۰/۷۵	-	-	-	-	-	-	-	-	-	-	-	-
۱	۶	۳	۵	۳	۵	۳	۲	۴	۴	۲	۱	۱
۱/۲۵	-	-	-	-	-	-	-	-	-	-	-	-
۱/۵۰	-	-	-	۱	-	-	۲	-	۲	۵	۲	۵
۱/۷۵	-	-	-	-	-	-	-	-	-	-	-	-
۲	۲	۳	۳	۱	۱	۲	۱	۱	۲	۱	۲	۹
۲/۲۵	-	-	-	-	-	-	-	-	-	-	-	-
۲/۵۰	۲	۴	۴	۴	۴	۴	۳	۴	۲	۴	۴	۵
۲/۷۵	-	-	-	-	-	-	-	-	-	-	-	-
۳	۵	۴	۴	۳	۱	۲	۱	۳	۴	۱۳	۳	۳
۳/۲۵	-	-	-	-	-	-	-	-	-	-	-	-
۳/۵۰	۱	۴	۳	۳	۳	۱	۴	۴	۳	۵	۴	۴
۳/۷۵	-	-	-	-	-	-	-	-	-	-	-	-
۴	۱	-	-	۲	۴	۴	۲	۴	-	۴	۳	۵
۴/۲۵	-	-	-	-	-	-	-	-	-	-	-	-
۴/۵۰	-	۱	۱	۱	۱	۱	-	-	۱	۳	۵	۶
۴/۷۵	-	-	-	-	-	۱	-	-	-	۲	-	-
۵	۴	۳	۳	۳	۷	۷	۸	۳	۳	۷	۸	۳

جدول ۲. طبقه‌بندی نمرات آزمون سُرایش

طبقه	دامنه	طبقه	دامنه
۱	۰	۷	$(۲/۵, ۳]$
۲	$(۰, ۰/۵]$	۸	$(۳, ۳/۵]$
۳	$(۰/۵, ۱]$	۹	$(۳/۵, ۴]$
۴	$(۱, ۱/۵]$	۱۰	$(۴, ۴/۵]$
۵	$(۱/۵, ۲]$	۱۱	$(۴/۵, ۵]$
۶	$(۲, ۲/۵]$		