

# مقایسه مدل‌های اندازه‌گیری (کلاسیک و

## سؤال - پاسخ) از لحاظ برآورد پارامترهای

### سؤال و توانایی

نوشته مهدی فراهانی

معرفی مقاله

نظریه‌های اندازه‌گیری از یک بُعد به دو دسته اساسی تقسیم می‌شود: نظریه کلاسیک و نظریه‌های جدید اندازه‌گیری (IRT)<sup>۱</sup>. این پژوهش از طریق داده‌های حاصل از اجرای سه آزمون ریاضیات، فیزیک و درس فنی که در مورد ۴۳۰ نفر از داوطلبان کنکور داخلی وزارت نیرو اجرا شده بود، به مقایسه نظریات کلاسیک و جدید از لحاظ برآورد پارامترهای سؤال و توانایی پرداخته است.

تحلیل‌های کلاسیک داده‌ها از طریق برنامه SPSS و برآورد پارامترهای سؤال و توانایی بر حسب مدل‌های جدید اندازه‌گیری از طریق نرم افزار رایانه‌ای BILOG در کالیفرنای آمریکا انجام گرفته است. مقایسه مدل‌های مختلف اندازه‌گیری از نظر برآوردهای متفاوت یا مشابه پارامترهای سؤال به وسیله تابع آگاهی و آزمون وابسته، به آزمون و ارزیابی گذاشته شد. نتایج نشان داد که مدل کلاسیک و مدل‌های IRT برآوردهای متفاوتی برای پارامترهای سؤال به دست می‌دهند و برآورد پارامترهای سؤال بر پایه مدل‌های

IRT، به ویژه مدل سه پارامتری، دقیق تر از برآورد پارامترهای سؤال بر اساس مدل کلاسیک و مدل های ساده تر IRT است. مقایسه مدل ها از نظر برآورد متفاوت پارامتر توانایی آزمودنی ها به وسیله آزمون کای اسکوئر آزموده شد. نتایج این مرحله نیز تفاوت مدل ها را در برآورد پارامتر توانایی آشکار ساخت و نشان داد که مدل سه پارامتری نسبت به مدل کلاسیک و حتی مدل های یک و دو پارامتری IRT برآوردهای متفاوت تر و دقیق تری ارائه می دهد.

این مقاله خلاصه ای است از پایان نامه تحصیلی دوره کارشناسی ارشد آقای مهدی فراهانی، کارشناس آموزش وزارت نیرو که با راهنمایی آقای دکتر محمد کاظم سلیمی زاده، عضو هیأت علمی دانشگاه علامه طباطبائی تهیه شده و در اختیار فصلنامه قرار گرفته است. بدین وسیله از ایشان سپاس گذاری می شود.

فصلنامه



پژوهشگاه علوم انسانی و مطالعات فرهنگی  
پرتال جامع علوم انسانی

## مقدمه

مدل یا نظریه، رکن اساسی هر رشته علمی محسوب می‌شود. در اندازه‌گیری و روان‌سنجی نیز نظریه‌ها و مدل‌هایی وجود دارد که از لحاظ تاریخی و سیر روند تکاملی، به دو دسته کلاسیک (نظریه ضعیف نمره حقیقی) و سؤال - پاسخ (به اختصار IRT) تقسیم می‌شوند (۱).

مبانی مدل کلاسیک اندازه‌گیری را در اوایل قرن حاضر، اسپیرمن معرفی و پایه‌ریزی کرد (۲). این مدل سپس در دو کتاب گالیکسن و لرد و ناویک به اوج توسعه و بسط خود رسید (۳ و ۴) اما هم‌زمان با بسط این مدل، ضعف‌های جدی آن بیش‌تر آشکار شد و روان‌سنجان و متخصصان آزمون‌سازی را بیش‌از پیش به سمت مدل‌های جدید سوق داد (۵). از نیمه دوم قرن بیستم به تدریج زمینه ارائه نظریات جدید مطرح شد و کسانی مانند لرد (۶)، راش (به نقل از رایت، ۱۹۷۷) (۷)، برن بام (۸)، رایت و همبلتون (۹) در این مسیر گام‌های مؤثری برداشتند.

هم اینک فعالیت‌های آزمون‌سازی و اندازه‌گیری در زمینه‌های مختلف - از پیشرفت تحصیلی گرفته تا سنجش نگرش‌ها و ... - در آمریکا و اروپا و براساس مدل‌های جدید (IRT) بررسی و ساخته می‌شود و مدل کلاسیک اعتبار بیش‌تر برای بحث در مورد تاریخچه نهضت آزمون‌سازی یا برآورد پارامترهای نظریات جدید به عنوان برآورد اولیه و مقدماتی مورد استفاده قرار می‌گیرد. معلمان هر چند در محدوده کلاس درس به ظاهر کم‌تر می‌توانند از مدل‌های اندازه‌گیری بهره ببرند اما اطلاع از مبانی نظری طراحی، اجرا و تجزیه و تحلیل آزمون‌های پیشرفت تحصیلی و یافته‌های جدید، بینش و بصیرت بهتری در مورد یکی از وظایف مهم معلمان - یعنی سنجش محصلان - در اختیار آن‌ها قرار می‌دهد.

## مطالعه منابع مربوط به موضوع پژوهش

در بعضی از مطالعات قبلی، محققان به مقایسه مدل‌ها از نظر دقت برآورد پارامترهای سؤال و توانایی پرداخته و برخی دیگر، تأثیر حجم نمونه آزمودنی و سؤال یا نقص مفروضات را در برآورد پارامترها بررسی کرده‌اند. با توجه به محدودیت مقاله حاضر، به‌طور عمده نتایج دسته اول (مقایسه مدل‌ها در برآورد پارامترها) ارائه می‌شود.

یکی از نتایج مطالعات مربوط به مقایسه دقت برآورد پارامترها در مدل‌های اندازه‌گیری حاکی از آن است که مدل منطقی سه پارامتری نسبت به مدل‌های یک و دو پارامتری با آزمون‌های ۲۰ سؤالی در برآورد جایگاه افراد در صفت مکنون و رتبه‌بندی آزمودنی‌ها براساس صفت مورد سنجش دارای قدرت بیش‌تری است. نکته دیگر آن‌که

پژوهشگران با تقسیم آزمودنی‌ها به دو گروه بالا و پایین ( $\theta_L = -2/5$  تا  $0$  و  $2/5$  تا  $\theta_U = 0$ ) و مقایسه دقت برآورد جایگاه حقیقی افراد بر حسب مدل‌های IRT، نشان دادند که افزایش دقت و قابلیت پیش‌بینی توانایی افراد توسط مدل سه پارامتری نسبت به سایر مدل‌ها، در گروه‌های با توانش پایین بیش‌تر از گروه‌های با توانش بالاست؛ زیرا امکان استفاده از حدس و شانس کاذب در گروه‌های با توانایی کم بیش‌تر است و تنها مدل سه پارامتری در برآورد پارامترها به عامل  $C_i$  (مجاناب پایین ICC) توجه می‌کند. مقایسه مدل‌های یک و دو پارامتری (بررسی تأثیر پارامتر قدرت تشخیص سؤال در برآورد توانایی افراد) تمایز و تفاوت جدی نشان نداد (۱۰).

راید برای تحلیل داده‌های چهار خرده آزمون مورد استفاده - شامل محاسبات عددی، تجسم فضای سه بعدی، خزانه لغات و استدلال ریاضی - از نرم‌افزارهای BILOG و TESTFACT بهره‌گرفت و بر اساس یک نمونه ۴۰۶ نفره از آزمودنی‌ها و با حذف سؤال‌هایی که کم‌تر از ۷۵ درصد آزمودنی‌ها به آن‌ها پاسخ داده بودند، نتیجه گرفت که هر چهار خرده آزمون از لحاظ سطح دشواری سؤال‌ها و قدرت تشخیص افراد در سطوح مختلف توانایی، تفاوت داشته است. بنابراین، مدل‌های یک و دو پارامتری از لحاظ برآورد پارامترها تفاوت معنی‌دار نشان داده‌اند (۱۱).

هومن در تحقیقی با استفاده از آزمون تهران - استنفرد - بینه (TSB) (۱۲) به بررسی و مقایسه برآورد پارامترهای دشواری و توانایی پرداخت و در واقع، توانمندی مدل راش را در برآورد پارامترها در شرایط نقض مفروضات بررسی کرد. نتیجه کلی تحقیق مذکور آن است که مدل راش برای برآورد پارامتر دشواری سؤال‌ها وقتی به سؤال‌ها تفاوت دارد، مناسب نیست اما برای برآورد توانایی افراد مناسب و خوب است (۱۳).

در مورد مدل دو پارامتری و تفاوت برآورد توانایی بر پایه آن و مدل کلاسیک، انصارین به پژوهش دست زد. او با استفاده از داده‌های حاصل از اجرای آزمون هوش تهران - استنفرد - بینه (TSB) به برآورد منحنی ویژه سؤال‌ها، پارامترهای دشواری و قدرت تشخیص سؤال‌ها و توانایی آزمودنی‌ها اقدام کرد. نتیجه آن‌که نمرات خام یکسان دارای برآورد یکسانی از توانایی و موقعیت آزمودنی بر روی پیوستار مکنون نبودند (۱۴) (۱۵).

دیوجی طی یک بررسی، کاربرد مدل یک پارامتری راش را برای سؤال‌های چند گزینه‌ای مورد بررسی قرار داده است اما به رغم استفاده از مدل راش برای برآورد پارامترها در سؤال‌های چندگزینه‌ای، به عقیده دیوجی به علت نبود پارامتر حدس و قدرت تشخیص یکسان در مدل راش، این کاربردها درست نیست (۱۶). البته علاوه بر دیوجی،

محققان دیگری هم به نامناسب بودن استفاده از مدل راش برای سؤال‌های چندگزینه‌ای اشاره کرده‌اند؛ از جمله، آندرسن عدم برازش مدل را به نابرابری قدرت تشخیص نسبت داده است. همبلتون و تراب نشان داده‌اند که مدل دو پارامتری، توزیع نمرات را بهتر از مدل راش پیش‌بینی می‌کند (۹).

همبلتون و موری از طریق نمودار باقی‌مانده برای آزمون‌های ریاضی، برازش مدل سه پارامتری و عدم برازش مدل تک پارامتری را برای داده‌ها مطرح کرده‌اند (۱۶).

یکی از موضوع‌ها و نکات مهم در مورد مدل‌های IRT، حجم نمونه آزمودنی و سؤال است. این دو عامل، به‌ویژه در مدل سه پارامتری، می‌تواند بر برآورد پارامترهای سؤال و توانایی تأثیرات جدی داشته باشد. همبلتون و کوک در یک مطالعه با انتخاب آزمون‌هایی با سه طول ۱۰، ۲۰ و ۸۰ سؤالی و نمونه‌هایی با حجم ۵۰، ۲۰۰ و ۱۰۰۰ آزمودنی به بررسی اثرات حجم نمونه آزمودنی و ویژگی‌های خزانه سؤال و تعداد سؤال بر خطای استاندارد برآورد توانایی اقدام کردند. طول آزمون ۱۰ سؤالی حداقل طول ممکن برای یک آزمون و آزمون ۸۰ سؤالی نیز از طول‌های متداول تست است. در مورد حجم نمونه آزمودنی ۵۰ و ۱۰۰۰ نفر نیز به همین‌گونه استدلال کرده‌اند. سؤال‌ها از دو خزانه سؤال استخراج شدند. در عمل، سؤال‌های خزانه ۱ دارای دامنهٔ عریض‌تری برای پارامترهای دشواری و قدرت تشخیص سؤال بودند. پارامتر حدس هر دو خزانه سؤال ۰/۲۵ در نظر گرفته شد. جمع‌بندی نتایج این مطالعه به صورت زیر است:

۱. حجم نمونه پاسخ‌گویان و طول آزمون، دو عامل بسیار مهم در دقت منحنی‌های  $SE(\theta)$  است. موارد نقض و استثنای این امر به نوسانات نمونه‌گیری مربوط می‌شود.
۲. در کرانه‌های پیوستار توانایی، دقت منحنی‌های  $SE(\theta)$  حتی با وجود نمونه‌های بزرگ آزمودنی، بسیار پایین است.
۳. در اکثر موارد با نمونه‌های ۲۰۰ آزمودنی و ۲۰ سؤال، دقت برآورد خطای استاندارد توانایی قابل قبول خواهد بود. البته این نکته بیش‌تر در دامنهٔ وسط توانایی [۱-، +۱] صادق است.
۴. افزایش طول تست از ۱۰ به ۲۰ سؤال بیش از افزایش آن از ۲۰ به ۸۰ سؤال، دقت  $SE$  را بهبود می‌بخشد.
۵. در مورد حجم نمونه نیز افزایش افراد از ۵۰ به ۲۰۰، بیش از ۲۰۰ به ۱۰۰۰ نفر دقت برآورد  $SE$  را ارتقا می‌دهد (۱۷).

لرد طی یک مطالعه و ضمن مقایسه مدل‌های یک و دو پارامتری IRT در برآورد نمرهٔ حقیقی آزمودنی‌ها، تلاش کرده است تأثیر حجم نمونه را بررسی کند. داده‌های مطالعه

شامل پاسخ ۳۰۰۰ دانش آموز کلاس ششم به آزمون خزانه لغات متروپولیتن با برنامه LOGIST تجزیه و تحلیل شده است. نتایج مطالعه نشان داد وقتی حجم نمونه کوچک باشد، پارامتر قدرت تشخیص (z) سؤال‌ها و پارامتر مجانب یا حدس سؤال‌ها (C<sub>i</sub>) را نمی‌توان به دقت تعیین کرد. از این رو، در بعضی موقعیت‌های معین و محدود و با حجم نمونه کوچک‌تر از ۱۰۰ یا ۲۰۰ آزمودنی، برآورد کننده نمره حقیقی X آزمودنی در مدل راش (یک پارامتری) می‌تواند اندکی بهتر از برآورد کننده نمره حقیقی بر پایه مدل دو پارامتری باشد (۶).

### بیان اهداف و فرضیه‌های پژوهش

هدف اصلی پژوهش حاضر مقایسه مدل‌های اندازه‌گیری (کلاسیک و سؤال - پاسخ) از لحاظ برآوردهای متفاوت یا مشابهی است که برای پارامترهای سؤال‌های آزمون و توانایی آزمودنی‌ها به دست می‌دهند.

فرضیه‌های این پژوهش که آزمون آن‌ها مورد توجه است، عبارت اند از:

۱. کاربرد مدل‌های سؤال - پاسخ (با تعداد پارامتر مناسب) برای برآورد دقیق مشخصات سؤال بر مدل کلاسیک برتری دارد.
۲. برآورد توانایی آزمودنی‌ها با استفاده از مدل‌های IRT (با تعداد پارامتر مناسب) از برآورد توانایی افراد بر پایه مدل کلاسیک دقیق‌تر است.
۳. در صورت وجود عامل حدس در پاسخ دادن به سؤال‌ها، افزودن پارامتر حدس (C<sub>i</sub>) به مدل IRT، مدل برازنده‌تری برای داده‌ها ایجاد می‌کند.

### روش اجرای پژوهش

الف. آزمودنی‌ها

جامعه این پژوهش را همه داوطلبان آزمون ورودی دوره‌های داخلی وزارت نیرو تشکیل می‌دهد و نمونه تحقیق تعداد ۵۵۳ نفر از داوطلبانی هستند که در آزمون داخلی گزینش دانشجو (مورخ ۱۳۷۳/۳/۶) در رشته قدرت (مقطع کاردانی) شرکت کرده‌اند. برای نمونه برداری همه شرکت‌کنندگان انتخاب شدند. در واقع، شرکت‌کنندگان در آزمون مذکور نمونه‌ای از همه داوطلبان فرض شده‌اند. حجم نمونه اولیه ۵۵۳ نفر بود اما تعدادی از پاسخ‌نامه‌ها به علت مخدوش و غیرقابل استفاده بودن از نمونه حذف شد و حجم نمونه نهایی به ۴۳۰ نفر کاهش یافت.

## ب. روش‌های آماری

برای تجزیه و تحلیل داده‌ها ابتدا در ایران مشخصه‌های کلاسیک سؤال‌ها و آزمون‌ها محاسبه شد. سپس داده‌ها که شامل پاسخ‌های ۴۳۰ آزمودنی به سؤال‌های چهارگزینه‌ای سه خرده آزمون ۲۰ سؤالی بود، از طریق نظام شبکه جهانی اطلاعات (Internet) به یکی از دانشگاه‌های آمریکا (دانشگاه UCLA) ارسال گردید. پس از تحلیل داده‌ها با نرم‌افزار بای‌لوگ، خروجی کامپیوتر شامل برآورد پارامترهای سؤال و توانایی افراد، به ایران فرستاده شد.<sup>۱۳</sup> خلاصه اطلاعات در مورد خرده آزمون ریاضیات به همراه نمودارهای دشواری و قدرت تشخیص سؤال‌ها در پایان مقاله ارائه شده است (نمودار شماره ۲۰۱ و جدول شماره ۴).

برای تعیین برتری مدل‌های سؤال - پاسخ بر مدل کلاسیک در برآورد دقیق مشخصات سؤال (فرضیه ۱ پژوهش) از یک ابزار نیرومند نظریات جدید اندازه‌گیری به نام تابع آگاهی<sup>۱۴</sup> استفاده شد. تابع آگاهی نظامی است که ورودی آن پارامترهای سؤال و خروجی آن میزان آگاهی‌دهندگی آزمون می‌باشد. متفاوت بودن تابع آگاهی یک آزمون بر حسب مدل‌های اندازه‌گیری، بیانگر تفاوت مدل‌ها از برآوردهایی است که برای پارامترهای سؤال‌های آزمون مذکور محاسبه شده است. برای بررسی معناداری تفاوت تابع آگاهی آزمون‌ها به علت وابسته بودن داده‌ها، از آزمون  $\chi^2$  برای داده‌های وابسته استفاده شد و تابع آگاهی هر آزمون که براساس مدل‌های یک و دو پارامتری و کلاسیک برآورد شده بود، با تابع آگاهی همان آزمون در مدل سه پارامتری مقایسه گردید. معناداری تفاوت این توابع آگاهی با آزمون شد.

به منظور آزمون فرضیه دوم - تفاوت یا عدم تفاوت مدل‌ها در برآورد توانایی و سطح صفت مکنون آزمودنی‌ها - از آزمون مجذور کای از نوع نیکویی برازش استفاده شد تا مدل‌های یک و دو پارامتری و کلاسیک را در مقایسه با مدل سه پارامتری آزمون کند. برای این منظور، ابتدا نمرات به مقیاس استاندارد یا  $Z$  برده شدند و سپس به ۱۲ طبقه از  $[-۲/۵, -۳]$  تا  $[۲/۵, ۳]$  تقسیم گردیدند. آنگاه فراوانی طبقات براساس آزمون نیکویی برازش با درجات آزادی  $k - 1 = 11$  مقایسه و در دو سطح معناداری ۵٪ و ۱٪ بررسی شد.

برای بررسی و تعیین برازش یا عدم برازش هر یک از مدل‌های IRT با داده‌ها، یعنی به منظور آزمون فرضیه ۳ پژوهش مبنی بر برابری بیشتر تر مدل سه پارامتری نسبت به سایر مدل‌های IRT، نرم‌افزار BILOG از نوعی آزمون مجذور کای استفاده می‌کند. فرض صفر آزمون نیکویی برازش مدل - داده‌ها بر تناسب و برازش مدل با داده‌ها تأکید می‌ورزد و در

واقع، تفاوت مدل با داده‌ها را انکار می‌کند و فرض خلاف عدم برازش را بیان می‌دارد. در عمل، سطح احتمال مجذور کای محاسبه شدهٔ هر سؤال با  $0/01$  و  $0/05$  مقایسه می‌شود. در صورت بزرگ‌تر بودن سطح احتمال هر سؤال از  $0/01$  یا  $0/05$ ، نتیجه گرفته می‌شود که سؤال با مدل دارای برازش نسبی ( $0/05 < \alpha < 0/01$ ) یا برازش کامل ( $\alpha > 0/05$ ) است.

#### پ. ابزار گردآوری داده‌ها

ابزار مورد استفاده برای گردآوری داده‌های پژوهش شامل سه آزمون بیست سؤالی ریاضیات، فیزیک و درس فنی به صورت چهارگزینه‌ای بوده است. درس‌های مذکور جزء مواد امتحانی اصلی و مشترک آزمون‌های گزینش دانشجو در وزارت نیرو به شمار می‌روند. سایر مواد امتحانی تعداد آزمودنی کمی دارد. این آزمون‌ها به وسیلهٔ متخصصان حیطه‌های مزبور که معمولاً جزء کارشناسان وزارت نیرو هستند، طراحی و تهیه می‌شود. در مورد نحوهٔ اجرای ابزار تا حد قابل قبولی می‌توان شرایط تقریباً استاندارد و اصولی را برای آن‌ها در نظر گرفت. از جمله در ابتدای دفترچهٔ آزمون دربارهٔ شیوهٔ پاسخ‌گویی به سؤال‌ها، مدت زمان اجرا، وجود نمرهٔ منفی در آزمون، تعداد سؤال‌های هر خرده آزمون و ... توضیحاتی ذکر شده است.

#### یافته‌های پژوهش

برای آزمون فرضیهٔ ۱، مقایسهٔ مدل سه پارامتری و مدل کلاسیک از نظر برآوردهایی که برای پارامترهای سؤال‌ها به دست می‌دهند، از تابع آگاهی خرده آزمون‌ها بهره گرفته شد. تفاوت توابع آگاهی هر آزمون که حاصل مدل‌های مختلف بود، محاسبه و با آزمون  $t$  برای داده‌های وابسته آزمون شد. نتایج نشان داد که نه تنها بین تابع آگاهی هر آزمون براساس مدل کلاسیک تفاوت قابل توجهی دیده می‌شود و این تفاوت حتی در سطح  $0/01 < \alpha$  و با بیش از ۹۹ درصد اطمینان معنادار است بلکه توابع آگاهی مدل‌های دو پارامتری و یک پارامتری نیز برای آزمون از تابع آگاهی هر خرده آزمون براساس مدل کلاسیک در همین سطح تفاوت معناداری دارد. در واقع، با درجهٔ آزادی ۳۲ و مقدار  $t$  مبین  $2/45$  تنها توابع آگاهی آزمون‌ها بین مدل‌های یک و دو پارامتری معنادار نبود و سایر مقادیر همه تفاوت معناداری داشتند. بنابراین، می‌توان گفت علاوه بر مدل سه پارامتری IRT، حتی مدل‌های دو و یک پارامتری نیز پارامترهای سؤال‌ها را دقیق‌تر و مناسب‌تر از مدل کلاسیک برآورد می‌کنند و فرضیهٔ صفر رد و فرضیهٔ اول پژوهش با بیش از ۹۹ درصد اطمینان تأیید می‌شود (جدول شمارهٔ ۱).



جدول شماره ۱ - محاسبه آزمون گروه‌های وابسته برای تفاوت بین توابع آگاهی آزمون‌ها  
برحسب مدل‌های مختلف

تابع آگاهی	آزمون ریاضیات			آزمون فیزیک			آزمون درس فنی		
	MEAN	STD	T(Ob)	MEAN	STD	T(Ob)	MEAN	STD	T(Ob)
DF(1.3)	۱/۶۶	۳/۶۲	۲/۶۲	۱/۸۱	۳/۶۹	۲/۸۲	۱/۵۵	۲/۹۹	۲/۹۷
DF(۲.۳)	۱/۶۴	۳/۷۴	۲/۵۲	۱/۵۷	۳/۲۴	۲/۷۹	۱/۱۷	۲/۳۷	۲/۸۳
DF(۱.۲)	۰/۰۲	۰/۵۷	۰/۱۵	۰/۲۴	۰/۹۴	۱/۴۸	۰/۳۸	۱/۱۹	۱/۸۳
DF (cl. ۳)	۲/۹۲	۴/۷۱	۳/۵۶	۲/۸۱	۴/۵۹	۳/۵۲	۲/۵۸	۳/۷۹	۳/۹۱
DF (cl. ۲)	۱/۲۸	۱/۶۳	۴/۵۰	۱/۲۴	۱/۸۷	۳/۸۲	۱/۴۱	۱/۹۲	۴/۲۳
DF (cl. ۱)	۱/۲۶	۱/۳۳	۵/۴۴	۱/۰۰	۱/۰۵	۵/۴۵	۱/۰۳	۰/۹۷	۶/۱۴

$$n = ۳۳$$

$$df = ۳۲$$

$$t(cr) \cdot ۰/۰۵ = ۱/۶۹$$

$$t(cr) \cdot ۰/۰۱ = ۲/۴۵$$

فرضیه ۲ مدعی است که مدل سه پارامتری نظریه سؤال - پاسخ توانایی افراد را نسبت به مدل کلاسیک اندازه‌گیری به گونه‌ای متفاوت و دقیق‌تر برآورد می‌کند. به این ترتیب، باید برای آزمون این فرضیه، تفاوت مدل سه پارامتری از مدل کلاسیک در برآورد توانایی بررسی شود. به این منظور، آزمون مجدور کای بین مدل سه پارامتری و مدل کلاسیک (البته مدل‌های دو و یک پارامتری نیز) اجرا شد تا تفاوت فراوانی‌های افراد در طبقات مختلف آزمون شود. نتایج نشان داد که نه تنها توانایی برآورد شده برای افراد براساس مدل سه پارامتری از مدل کلاسیک متفاوت است بلکه حتی مدل‌های دو و یک پارامتری نیز در برآورد پارامتر توانایی آزمودنی‌ها با مدل سه پارامتری تفاوت و تمایز دارد. این امر در سطح  $\alpha < ۰/۰۱$  هم معنی‌دار بود. از سوی دیگر، از آن‌جا که مدل سه پارامتری از لحاظ تعداد پارامتر مناسب‌ترین مدل برای داده‌های حاضر محسوب می‌شود و طبق نتایج آزمون خوبی برازندگی، برازش و مناسبت کاملاً آشکاری با داده‌های پژوهش دارد، توانایی برآورد شده برای افراد را که از مدل سه پارامتری به دست آمده است، می‌توان دقیق‌تر و مناسب‌تر از توانایی حاصل از سایر مدل‌های اندازه‌گیری تلقی کرد. بنابراین، فرضیه صفر رد و فرضیه دوم پژوهش حاضر نیز با بیش از ۹۹ درصد اطمینان تأیید می‌شود (جدول شماره ۲).

جدول شمارهٔ ۲- برازندگی مدل - داده‌ها برحسب مدل‌های IRT (تعداد سؤال‌ها)

مدل‌ها	۱ - PAR			۲ - PAR			۳ - PAR		
	بدون برازش	برازش نسبی	برازش کامل	بدون برازش	برازش نسبی	برازش کامل	بدون برازش	برازش نسبی	برازش کامل
آزمون‌ها	$\alpha < .01$	$.01 < \alpha < .05$	$\alpha > .05$	$\alpha < .01$	$.01 < \alpha < .05$	$\alpha > .05$	$\alpha < .01$	$.01 < \alpha < .05$	$\alpha > .05$
آزمون ریاضیات	۱۶	۳	۱	۱۸	۱	۱	۱	۲	۱۷
آزمون فیزیک	۱۹	۱	۰	۱۷	۱	۲	۲	۲	۱۶
آزمون درس فنی	۱۸	۲	۰	۱۶	۱	۳	۱	۲	۱۷

$$r \leq df \geq 9$$

$$1 \leq df \geq 9$$

$$1 \leq df \geq 8$$

جدول شمارهٔ ۳- مجذور کای محاسبه شده بین پارامتر توانایی مدل سه پارامتری

و مدل‌های دو و یک پارامتری و کلاسیک

مدل‌ها	ریاضی		فیزیک		درس فنی	
	دو پارامتری	یک پارامتری	کلاسیک	دو پارامتری	یک پارامتری	کلاسیک
سه پارامتری	۲۰۹	۵۰۰۶	۲۱۹۰۶	۲۲۰۹	۳۷۸	۸۲
دو پارامتری	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴
یک پارامتری	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴
کلاسیک	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴	۲۹/۵	۱۲/۴

$$df = 11$$

$$X^2(cr) . / . 0.5 = 19/675$$

$$X^2(cr) . / . 0.1 = 24/725$$

برای آزمون فرضیهٔ ۳ پژوهش از نوعی آزمون خوبی برازندگی استفاده شد و برازش سؤال‌های آزمون با مدل‌های یک، دو و سه پارامتری بررسی گردید. براساس نتایج به دست آمده، مدل سه پارامتری برازش چشم‌گیری با داده‌ها داشت؛ در حالی که نه تنها مدل دو پارامتری بلکه مدل یک پارامتری نیز بدون برازش شناخته شد. بنابراین، فرضیهٔ سوم پژوهش نیز با بیش از ۹۹ درصد اطمینان و در سطح  $\alpha < 0.01$  معنادار شناخته شد و فرض صفر رد و فرض پژوهش (خلاف) تأیید گردید. بدین ترتیب مدل سه پارامتری، برای داده‌های آزمون‌های مورد استفاده در این پژوهش نسبت به سایر مدل‌های IRT برازنده‌تر و مناسب‌تر شناخته شد (جدول شمارهٔ ۳).

## نتیجه‌گیری

مطابق نتایج تحقیقات و مطالعات قبلی در زمینه نظریه‌های جدید اندازه‌گیری، IRT و مقایسه آن با نظریه کلاسیک آزمون، این مطالعه نیز نشان داد که در سطح بالایی از اطمینان می‌توان گفت:

۱. مدل‌های IRT نسبت به مدل کلاسیک اندازه‌گیری در برآورد پارامترهای سؤال‌ها و توانایی آزمودنی‌ها دارای مناسبت، دقت و برازندگی بیش‌تری است. این نتیجه در مورد آزمون‌های پیشرفت تحصیلی چهارگزینه‌ای که مبنای مطالعه حاضر بوده است، با بیش از ۹۹ درصد اطمینان صدق می‌کند.

۲. در بین مدل‌های نظریه سؤال - پاسخ (IRT) برای برآورد پارامترهای سؤال و توانایی افراد، مدل سه پارامتری در مورد آزمون‌های پیشرفت تحصیلی چندگزینه‌ای این پژوهش، نسبت به مدل‌های یک و دو پارامتری برتری خاصی نشان داد.

۳. در بین مدل‌های یک و دو پارامتری IRT از لحاظ برآورد پارامترهای سؤال براساس داده‌های خرده آزمون‌های حاضر، تفاوت چشم‌گیری ملاحظه نمی‌شود. از این رو می‌توان تأثیرگذاری پارامتر قدرت تشخیص سؤال‌ها را بر برآورد پارامتر دشواری چندان شدید و جدی ندانست. به زبان دیگر، سؤال‌های مورد استفاده در خرده آزمون‌های این مطالعه دارای قدرت تشخیص خیلی متفاوت و مؤثری نیستند اما در مورد برآورد پارامتر توانایی آزمودنی‌ها، می‌توان نتیجه گرفت که براساس داده‌های حاضر، بین مدل‌های یک و دو پارامتری در مقام مقایسه با مدل سه پارامتری، برآورد پارامتر توانایی تا حدی متفاوت است و در مورد آزمون‌های چندگزینه‌ای پارامتر شیب بر روی برآورد پارامتر توانایی تأثیر داشته است.



## پی‌نوشت‌ها

۱. نظریات جدید اندازه‌گیری و روان‌سنجی نخست با اصطلاح صفت‌مکون یا خصیصه‌مکون به شدت پیوند خورد اما نه تدریج با عنوان نظریه سؤال - پاسخ یا Item Response Theory (با علامت اختصاری IRT) یا نظریه منحنی ویژه سؤال رواج یافت که برای مقاصد آزمون‌سازی و تحلیل آماری داده‌ها مناسب‌تر به نظر می‌رسد (همبلتون، ۱۹۹۳).

## 2. Information Function

۳. برای تحلیل داده‌ها در آمریکا به وسیله نرم‌افزار BILOG، علاوه بر مساعی فراوان استاد محترم راهنما آقای دکتر سلیمی‌زاده، نگارنده بر خود لازم می‌داند از عنایت استاد ارجمند دکتر جمال عابدی، عضو هیأت علمی دانشگاه مذکور که تحلیل داده‌ها را مسیر کردند، تشکر و سپاس‌گزاری نماید.



## منابع

- Allen, J.M., & Yen, M.W. (1979). *Introduction to measurement theory*. California: wadsworth.
- ثرندایک، رابرت؛ روان‌سنجی کاربردی، مترجم: حیدرعلی هومن، تهران، دانشگاه تهران، ۱۳۶۹.
- Gulliksen, H. (1950). *Theory of mental tests*. New York: John Wiley & sons.
- Lord, F.M., & Novick, M.R. (1968). *Statistical Theories of mental test scores*. Reading, MA: Addison-wesley.
- Hambleton, R.K., & van der Linden, W.J. (1982). Advances in item response theory and applications: An introduction. *Applied Psychological Measurement*, 6, 373-378.
- Lord, F.M. (1983). Small N justifies Rasch model. In D.J. Weiss (ed.). *New Horizons in testing*. New York: Academic press Inc.
- Wright, B.D. (1977a). Misunderstanding of the Rasch model. *Journal of Educational Measurement*, 14, 219-226.
- Birnbaum, A. (1968). Some Latent trait models and their use in inferring an examinee's

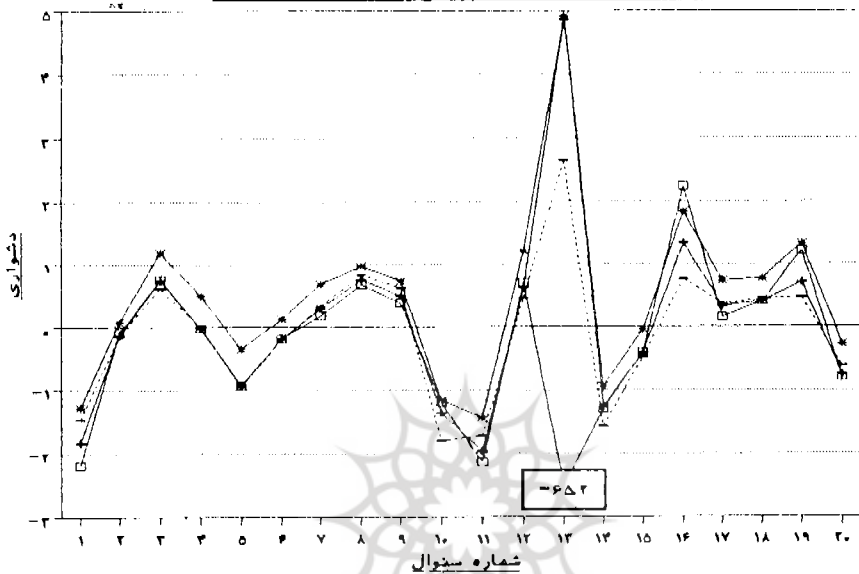
- ability. In F.M. Lord & M.R. Novick. **Statistical theories of mental test scores**, MA: Addison - wesley.
9. Hambleton, R.K. (1993). Principles and selected applications of item response theory. In R.L. Linn (Ed.), **Educational Measurement**. U.S.A.: oryx press.
10. Hambleton, R.K., & Cook, L.L. (1977). Latent trait models and their use in the analysis of educational test data. **Journal of Educational Measurement**, 14, 75-96.
11. Reid, C.A. (1993). **Latent trait modeling of the general test battery used with a rehabilitation client population: An investigation of Model - Data Fit**. Source: DAI-B 54/12, p. 6497. Jun 1994.
۱۲. هومن، حیدرعلی؛ مقایسه مدل تک‌پارامتری راش و مدل دوپارامتری. پایان‌نامه منتشر نشده دانشگاه آزاد اسلامی، ۱۳۷۳.
۱۳. هومن، حیدرعلی؛ هوش آزمای انفرادی تهران - استنفرد - بینه. فصلنامه علوم تربیتی دانشگاه تهران، ویژه‌نامه روان‌سنجی، دوره جدید، سال یکم، شماره ۱-۴.
۱۴. انصارین، علیرضا؛ برآورد خم ویژه سؤال و توانایی آزمودنی‌ها در مقیاس تهران - استنفرد - بینه بر پایه مدل دوپارامتری صفت مکنون. پایان‌نامه منتشر نشده دانشگاه آزاد اسلامی، ۱۳۷۱.
۱۵. هومن، حیدرعلی؛ روش تهیه آزمون هوش، تهران، دانشگاه تهران، ۱۳۷۵.
16. Divgi, D.R. (1986). Does the Rasch model really work for multiple choice items? Not if you look closely. **Journal of Educational Measurement**, 23,283-298.
17. Hambleton, R.K., & Cook, L.L. (1983). The robustness of item response models and effects of test length and sample size on the precision of ability estimates. In D.J. Weiss (Ed.), **New horizons in testing** (pp. 31-49). New York: Academic Press.



پیوست‌ها

شماره ۱

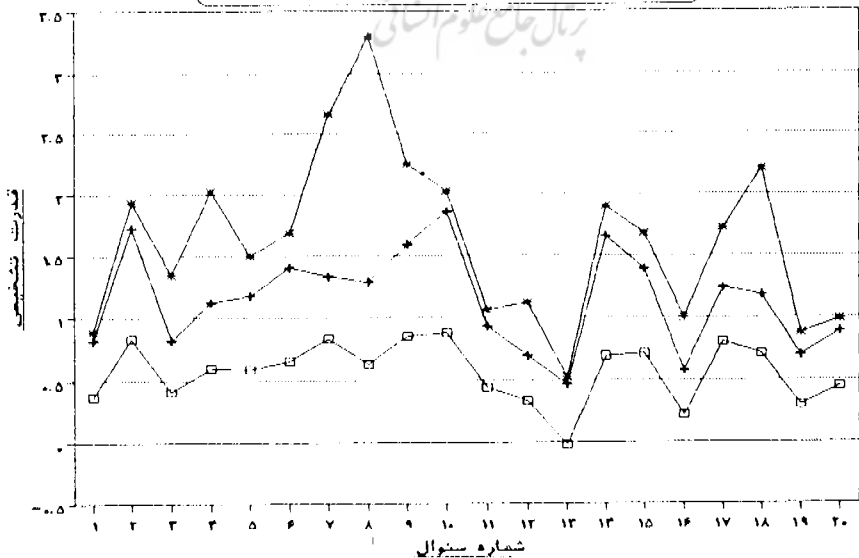
مقایسه پارامتر دشواری سئوالات آزمون ریاضی  
برحسب مدل‌های کلاسیک و سنوال - پاسخ



کلاسیک □ سه پارامتری \* دو پارامتری + یک پارامتری —

شماره ۲

مقایسه پارامتر شیب سئوالات آزمون ریاضی  
برحسب مدل‌های کلاسیک و سنوال - پاسخ



شماره سئوال

جدول شماره ۴ - آمارها و پارامترهای سؤال‌های خرده آزمون ریاضی برحسب مدل‌های کلاسیک و سؤال - پاسخ (IRT)

آزمون ریاضیات	آماره‌های سؤال‌ها (بخش‌ها)						پارامتر دشواری Bi			پارامتر قدرت تشخیص Ai			پارامتر حدس Ci	
	تعداد کوشش	تعداد صحیح	درصد صحیح	دورهنگی	همبستگی نقطه‌ای	دورهنگی نقطه‌ای	کلاسیک	۱ PAR	۲ PAR	۳ PAR	کلاسیک	۲ PAR		۳ PAR
۱	۴۰۲	۳۲۲	۸۰/۱۰	۰/۳۴۷	۰/۴۴۳	۰/۲۰۴	-۱/۴۴۴	-۱/۸۳۲	-۱/۲۲۸		۰/۴۷	۰/۸۲۴	۰/۸۹۵	۰/۲۵۱
۲	۳۵۹	۱۹۲	۵۳/۴۸	۰/۹۴	۰/۵۱	-۰/۹۹	-۰/۱۵۱	-۰/۱۴۵	۰/۶۸	۰/۸۳۳	۰/۸۳۳	۱/۹۳۳	۱/۹۳۳	۰/۱۲
۳	۲۷۹	۱۰۴	۳۷/۲۸	۰/۳۷۹	۰/۴۹۷	۰/۴۵۵	۰/۵۹۴	۰/۷۱۷	۱/۱۷۹	۰/۴۱	۰/۸۲۳	۱/۴۴۲	۱/۴۴۲	۰/۱۸۳
۴	۲۷۸	۱۴۰	۵۰/۳۶	۰/۵۱۲	۰/۴۰۸	-۰/۱۶	-۰/۱۱	-۰/۲۶	۰/۴۹۷	۰/۵۹۱	۱/۱۳	۲/۰۲۷	۲/۰۲۷	۰/۲۲۶
۵	۳۲۱	۲۲۱	۷۱/۹۶	۰/۵۰۹	۰/۴۸۲	-۰/۹۴	-۰/۹۵۹	-۰/۸۳۶	-۰/۴۶	۰/۵۹۱	۱/۱۷۴	۱/۴۹۷	۱/۴۹۷	۰/۲۷
۶	۳۴۲	۱۸۹	۵۵/۲۶	۰/۵۴۶	۰/۴۴۴	-۰/۹۲	-۰/۲۰۱	-۰/۱۹۱	۰/۱۳	۰/۶۵۲	۱/۴۹۸	۱/۶۸۳	۱/۶۸۳	۰/۱۶
۷	۱۹۸	۸۷	۴۳/۹۴	۰/۶۴۱	۰/۱۷۳	۰/۱۷۳	۰/۳۲۳	۰/۷۸	۰/۷۷۳	۰/۸۳۵	۱/۳۳۱	۲/۱۶۵	۲/۱۶۵	۰/۱۸۱
۸	۳۶۰	۱۱۹	۳۳/۰۶	۰/۵۴۴	۰/۴۱۱	۰/۳۵۵	۰/۸۲۱	۰/۶۸۵	۰/۹۱۷	۰/۳۲۳	۱/۲۸۵	۳/۴۰۳	۳/۴۰۳	۰/۱۵۵
۹	۴۰۷	۱۱۴	۲۷/۸۱	۰/۶۵۱	۰/۵۰۹	۰/۳۷۱	۰/۶۱۳	۰/۴۸۵	۰/۷۱۵	۰/۸۵۸	۱/۵۸۳	۲/۴۴۶	۲/۴۴۶	۰/۱۱۸
۱۰	۳۱۲	۴۱۰	۸۵/۱۴	۰/۶۵۹	۰/۴۶۵	-۰/۱۹۷	-۰/۱۸۱۴	-۰/۳۷	-۰/۱۷۳	۰/۸۷۶	۱/۸۶۱	۱/۰۶	۱/۰۶	۰/۲۰۳
۱۱	۳۹۶	۳۳۰	۸۳/۲۳	۰/۴۰۳	۰/۲۷	۰/۲۱۴۸	-۰/۷۳۷	-۰/۸۸۳	-۰/۴۴۴	۰/۴۴	۰/۹۳۱	۱/۰۶	۱/۰۶	۰/۲۵۵
۱۲	۲۵۷	۱۰۳	۴۰/۰۸	۰/۴۱۵	۰/۲۴۸	۰/۷۱۱	۰/۴۵۱	۰/۶۱۹	۱/۲۰۷	۰/۳۲۲	۰/۸۸۹	۱/۱۱۸	۱/۱۱۸	۰/۲۰۵
۱۳	۲۱۲	۱۴	۵/۴۴	-۰/۲۶	-۰/۱۳	-۰/۵۲۱	۲/۶۴۷	۴/۸۵۸	۴/۸۹۴	-۰/۲۶	۰/۴۶۵	۰/۵۳۳	۰/۵۳۳	۰/۲۷
۱۴	۳۶۹	۳۰۴	۸۲/۳۸	۰/۵۷	۰/۳۸۷	-۰/۳۰۸	-۰/۵۸۲	-۰/۳۶۴	-۰/۹۶۶	۰/۹۹۴	۱/۶۵۲	۱/۸۹۵	۱/۸۹۵	۰/۲۲۷
۱۵	۳۴۷	۲۱۷	۶۲/۵۴	۰/۵۸۳	۰/۴۵۱	-۰/۴۱۸	-۰/۴۸۸	-۰/۴۳۷	-۰/۵۲	۰/۷۱۸	۱/۴۹۳	۱/۸	۱/۸	۰/۲۰۱
۱۶	۲۹۷	۹۱	۳۰/۳۴	۰/۳۱	۰/۱۶	۲/۲۴	۰/۷۵۴	۱/۳۲۱	۱/۸۲۷	۰/۲۱۵	۰/۵۷۸	۱	۱	۰/۱۸۳
۱۷	۲۰۷	۹۳	۴۴/۹۳	۰/۶۲۸	۰/۵	۰/۱۴۹	۰/۳۵۴	۰/۳۲۵	۰/۷۱۱	۰/۸۰۷	۱/۴۴	۱/۷۳۴	۱/۷۳۴	۰/۱۶۳
۱۸	۲۴۶	۹۴	۳۸/۳۱	۰/۵۸۳	۰/۴۵۷	۰/۴۹۴	۰/۴۳۷	۰/۴۰۲	۰/۷۱۱	۰/۷۱۸	۱/۱۷۸	۲/۴۰۳	۲/۴۰۳	۰/۱۷۷
۱۹	۱۷۸	۶۲	۳۴/۳۳	۰/۲۹۳	۰/۲۲۷	۱/۰۷	۰/۴۵۹	۰/۷۰۸	۱/۳۰۹	۰/۳۰۶	۰/۷۰۳	۰/۸۸۱	۰/۸۸۱	۰/۱۶۹
۲۰	۴۲۸	۲۱۴	۵۰/۲۴	۰/۴۱۳	۰/۳۲۱	-۰/۸۱۵	-۰/۳۴۴	-۰/۷۶	-۰/۲۸۳	۰/۴۵۳	۰/۸۹۵	۰/۹۸۴	۰/۹۸۴	۰/۲۰۱

MEAN = ۳/۴۰  
STD = ۱۴/۵۹

۸۷



پروشکاه علوم انسانی ومطالعات فرہنگی  
پرتال جامع علوم انسانی