

## تلفیق رویکرد ریکیسی و رویکرد برنامه‌نویسی ریاضی در طراحی خزانه‌های سؤال بهینه برای سنجش انطباقی کامپیوتری

مریم مقدسین<sup>۱</sup>

تاریخ دریافت: ۹۵/۰۱/۱۴

تاریخ پذیرش: ۹۵/۱۰/۱۵

### چکیده

سنجش انطباقی کامپیوتری (CAT) شیوه‌ای از سنجش توانایی است که دقت برآورد توانایی را افزایش می‌دهد و بدون از دست دادن دقت اندازه‌گیری آزمون، طول آن را کاهش می‌دهد. با این وجود، سنجش انطباقی در صورتی خوب عمل می‌کند که دارای خزانه سؤالی باشد که در آن تعداد کافی سؤال با کیفیت مناسب وجود داشته باشد. بسیاری از محققان خاطرنشان کردند که برای ساخت خزانه سؤالی برای (CAT)، نه تنها اندازه خزانه سؤال مهم است، بلکه توزیع پارامترهای سؤال‌های خزانه نیز از اهمیت به‌سزایی برخوردار است. با این وجود، تحقیقات اندکی در مورد این که چگونه این ویژگی‌های مطلوب تعیین می‌شود، وجود دارد. هدف اصلی این پژوهش، تلفیق ایده "bin-and-union" برگرفته از رویکرد ریکیسی (۲۰۰۳) که یک روش شبیه‌سازی مونت‌کارلو برای تعیین ویژگی‌های خزانه سؤال است، با رویکرد برنامه‌نویسی ریاضی بوده است. در پژوهش فعلی، این روش برای ساخت یک خزانه سؤال بهینه برای آزمون سنجش انطباقی ریاضی به کار رفته است. خزانه سؤال پژوهش حاضر بر اساس مدل سه پارامتری مدرج شده است و روش سیمپسون-هتر برای کنترل مواجهه سؤال به کار رفته است. این طرح شامل برآوردهایی از اندازه مطلوب خزانه سؤال و توزیع مطلوب پارامترهای سؤال‌ها و ویژگی‌های غیر آماری آن بوده است. فرآیند طراحی این خزانه شامل تعیین مجموعه‌ای از ویژگی‌های مطلوب خزانه سؤال، با در نظر گرفتن چندین عامل مهمی که ممکن بود بر نتایج مورد نظر طراحی یک خزانه سؤال اثر گذارد، بوده است. عملکرد خزانه‌های سؤال شبیه‌سازی شده و عملیاتی با در نظر گرفتن مجموعه‌ای از ملاک‌های

۱. استادیار گروه روانشناسی بالینی، دانشگاه خوارزمی (نویسنده مسئول) [mmoghadasin@yahoo.com](mailto:mmoghadasin@yahoo.com)

ارزیابی، با یکدیگر مورد مقایسه قرار گرفته‌اند. نتایج ارزیابی نشان داد که مکانیزم به کاررفته برای تعیین ویژگی‌های مطلوب خزانه سؤال به خوبی عمل می‌کند و برای تعیین ویژگی‌های مطلوب خزانه سؤال مناسب است.

**واژگان کلیدی:** خزانه سؤال بهینه، سنجش انطباقی کامپیوتری، رویکرد ریکیسی، رویکرد برنامه‌نویسی ریاضی، مدل حداقل انحرافات وزن‌دار (WDM)

### مقدمه

کیفیت خزانه سؤال، به‌عنوان یک عامل مهم برای افزایش کیفیت اندازه‌گیری در سنجش انطباقی کامپیوتری (CAT)، در نظر گرفته شده است (فلاگرا<sup>۱</sup>، ۲۰۰۰؛ جنسما<sup>۲</sup>، ۱۹۷۷؛ مک‌برید و وایس<sup>۳</sup>، ۱۹۷۶؛ ریکیسی<sup>۴</sup>، ۱۹۷۶؛ ۲۰۰۳؛ وندرلیندن، ادلاید و ولدکمپ<sup>۵</sup>، ۲۰۰۶؛ ولدکمپ و وندرلیندن، ۲۰۰۰؛ اکسینگ و همبلتون<sup>۶</sup>، ۲۰۰۴). حتی در همان اوایل دهه ۱۹۷۰ - آغاز پژوهش‌های مرتبط با CAT - محققان به‌طور ضمنی و یا به‌صراحت اذعان داشتند که ویژگی‌های خزانه سؤال نقش مهمی در دستیابی به بهترین نتایج ممکن در سنجش انطباقی، خواهد داشت (مک‌برید و وایس، ۱۹۷۶). بااین‌وجود، دستورالعمل‌های اندکی در مورد چگونگی ساخت یک خزانه سؤال با کیفیت بالا، برای سنجش انطباقی کامپیوتری قابل‌دستیابی بودند (هی<sup>۷</sup>، ریکیسی، ۲۰۱۰؛ هی و ریکیسی، ۲۰۱۱).

مطالعات قبلی در مورد طراحی خزانه سؤال بهینه برای CAT روی دو رویکرد عمده مبتنی هستند. رویکرد اول، توسط ولدکمپ و وندرلیندن (۲۰۰۰) ایجاد شده است، این رویکرد از روش برنامه‌نویسی ریاضی<sup>۸</sup> برای طراحی خزانه سؤال استفاده می‌کند. در این

- 
1. Flaugher
  2. Jensema
  3. McBride & Weiss
  4. Reckase
  5. Van der linden, Adelaide & Weldkamp
  6. Xing & Hamblton
  7. Hi
  8. mathematical programming

رویکرد، CAT از طریق رویکرد تست سایه<sup>۱</sup>، اجرا می‌شود (وندرلیندن، ریس<sup>۲</sup>، ۱۹۸۸) که تست از طریق برنامه‌نویسی عدد صحیح خطی دو ارزشی<sup>۳</sup> یا برنامه‌نویسی ۰-۱<sup>۴</sup> تشکیل می‌شود. در این رویکرد ویژگی‌های تست بر اساس یک تابع هدف<sup>۴</sup> که در ارتباط با مجموعه‌ای از قیود<sup>۵</sup> خاص بیشینه می‌شود، محقق می‌گردد (وندرلیندن، ۲۰۰۵a). با این وجود، یک محدودیت بالقوه این رویکرد آن است که به نرم‌افزارهای جبر خطی از قبیل CPLEX و LINDO برای به دست آوردن راه‌حل بهینه نیاز دارد که کاربرد این روش را کمی دشوار می‌کند و ممکن است، کدها و معادلات آن برای اکثریت کاربران در دسترس نباشد که در این صورت اگر برنامه نیاز به اصلاح و یا تغییر داشته باشد، کنترلی بر آن نداشته باشند و چه بسا این احتمال وجود دارد که همیشه راه‌حل قابل اجرا و عملی<sup>۶</sup> در دسترس نباشد (چانگ<sup>۷</sup>، ۲۰۰۷؛ روبین<sup>۸</sup> و همکارانش، ۲۰۰۵). همچنین در این رویکرد فرض بر این است که سؤال‌ها از قبل در خزانه موجود هستند و از روی آن‌ها یک خزانه کوچک‌تر تشکیل می‌شود (گو و ریکیسی، ۲۰۰۷). در این رویکرد از ویژگی‌های یک خزانه سؤال موجود به عنوان نقطه شروع استفاده می‌شود (ریکیسی، ۲۰۱۰). به دلیل این که اجرای روش تست سایه در مقیاس بزرگ، به مسئله‌ای بسیار دشوار تبدیل می‌شود، بنابراین، حل مسائل سرهم کردن تست را با مشکل روبرو می‌کند. به منظور دوری از مشکلات روش تست سایه، از مدل حداقل انحرافات وزن‌دار<sup>۹</sup> که نوعی روش برنامه‌نویسی خطی است،

#### 1. Shadow a test (SAT)

روش آزمون سایه، روش پشتیبانی مفیدی برای سرهم کردن آزمون‌های انطباقی است. ایده زیربنایی این روش، حل یک مسئله بزرگ همزمان به صورت یک توالی از مسائل همزمان کوچک‌تر است. این رویکرد بر اساس این پیش فرض شکل گرفت که اگر بخواهیم از یک خزانه بزرگ مجموعه‌ای آزمون سرهم کنیم، ابتدا تعداد آزمون‌هایی که باید سرهم شود مشخص می‌شود

2. Reese
3. binary linear integer programming
4. objective function
5. constraints
6. feasible
7. Chang
8. Robin
9. weighted deviations model (WDM)

به‌عنوان جایگزین استفاده می‌شود. این روش برای سرهم کردن تست‌های چندگانه و تست‌های سنجش انطباقی به روش مؤثرتری عمل می‌کند. در این روش، ابتدا پیش‌بینی جستجوی راه‌حل برای تست کامل صورت می‌گیرد و همزمان هم قابل‌حل بودن و هم بهینه بودن تست را در نظر می‌گیرد. این روش جزء روش‌های شهودی<sup>۱</sup> حل مسائل سرهم کردن تست است (وندرلیندن، ۲۰۰۵a، ۲۰۰۵b). این روش در اصل توسط استوکینگ و سوانسون<sup>۲</sup>، ۱۹۹۳ به دلیل علاقه و نگرانی آن‌ها در مورد کیفیت ضعیف خزانه‌های سؤال در سرهم کردن تست‌های متوالی در مقیاس بزرگ ایجاد شد. مدل حداقل انحرافات وزن‌دار به‌صراحت ویژگی‌های آماری و غیر آماری سؤال را با تعادل مطلوبی بین ویژگی‌های اندازه‌گیری و ساختاری در نظر می‌گیرد. این ویژگی‌ها به‌وسیله وزن‌هایی که توسط طراحان تست انتخاب می‌شود، در مدل وارد می‌گردد. این روش برخلاف روش تست سایه، ویژگی‌های محتوایی را به‌عنوان اهداف<sup>۳</sup> (به‌جای قیود) فرمول‌بندی می‌کند. انحراف از اهداف محتوایی وزن داده می‌شود و در تابع هدف به همراه فاصله آگاهی سؤال از مقدار هدف<sup>۴</sup> قرار می‌گیرد (استوکینگ، سوانسون و پیرمن<sup>۵</sup>، ۱۹۹۳). البته این روش ابزاری بوده است که در بسیاری از رویکردهای برنامه‌نویسی ریاضی استفاده شده است و مدل‌های غیرقابل‌حل را قابل‌اجرا می‌کرده است (بروک، کندریک و مروس<sup>۶</sup>، ۱۹۹۸). در CAT، مدل حداقل انحرافات وزن‌دار سؤال‌هایی را انتخاب می‌کند که به‌طور متوالی کوچک‌ترین مجموع انحرافات وزن‌دار را دارد.

رویکرد دوم، رویکرد اکتشافی<sup>۷</sup> ریکیسی است (ریکیسی، ۲۰۰۳)، این رویکرد برخلاف روش برنامه‌نویسی ریاضی، بسیار سراسر است. همچنین، در مطالعات گوناگون در مورد طراحی خزانه‌های سؤال بهینه برای CAT استفاده شده است (ریکیسی، ۲۰۰۳،

- 
1. Heuristics
  2. Stocking & Swanson
  3. goal
  4. target
  5. Pearlman
  6. Brooke, Kendrick & Meeraus
  7. heuristic

ریکیسی و هی، ۲۰۰۴، ۲۰۰۵، ۲۰۰۹a، ۲۰۰۹b؛ گو<sup>۱</sup> و ریکیسی، ۲۰۰۷). در این رویکرد، استفاده از برنامه‌ریزی اعداد صحیح کنار گذاشته شده است و در آن فرض نمی‌شود که سؤال‌ها از قبل در خزانه وجود دارد. در عوض، در این رویکرد سؤال‌ها برحسب پارامترهای IRT شبیه‌سازی می‌شود تا با برآوردهای اخیر توانایی مطابقت داشته باشد و میزان آگاهی به‌اندازه کافی بهینه‌ای را ایجاد کند. در روش ریکیسی ابتدا، خزانه سؤال هدف بر اساس صفات غیر آماری از قبیل محتوا به خزانه‌های کوچک‌تری تقسیم‌بندی می‌شود. سپس فرآیند CAT شبیه‌سازی می‌شود، به طوری که خزانه‌های سؤال کوچک‌تر به‌طور همزمان ساخته می‌شوند. شبیه‌سازی با یک آزمودنی که به‌طور تصادفی از توزیع مورد انتظار جامعه فرضی در دامنه مشخص شده انتخاب می‌شود، آغاز شده و CAT برای او اجرا می‌گردد. هر سؤال به نحوی شبیه‌سازی می‌شود که سؤال بهینه‌ای بر اساس برآورد جدید توانایی آزمودنی باشد. فرایند مشابهی برای آزمودنی بعدی نیز تکرار می‌شود، سپس، به همین ترتیب، برای کل نمونه مورد نظر این فرآیند ادامه می‌یابد و سؤال‌ها برای نمونه بزرگی از آزمودنی‌ها شبیه‌سازی می‌شود و به خزانه سؤال اضافه می‌گردند و بدین ترتیب بر اساس این روش که به آن روش "bin-and-union" نیز گفته می‌شود، خزانه سؤال بهینه ساخته می‌شود (ریکیسی، ۲۰۰۳، ۲۰۰۹؛ ریکیسی و هی، ۲۰۰۴، ۲۰۰۹a). برخلاف مسئله‌ی سرهم کردن<sup>۲</sup> خزانه سؤال در رویکرد اول که در آن یک خزانه سؤال از یک خزانه بزرگ<sup>۳</sup> در دسترس بر طبق ویژگی‌های مطلوب سرهم می‌شود (وندربلیندن، آریل و ولدکمپ، ۲۰۰۶، وندربلیندن، ۲۰۰۰a، ۲۰۰۰b، ۲۰۰۰c، ۲۰۰۵a، ۲۰۰۵b)، در مسئله طراحی خزانه سؤال در رویکرد دوم، فرض بر این است که هیچ سؤال واقعی در دسترس نیست. از این رو، از آنجایی که در عمل نیز، زمانی که یک خزانه سؤال طراحی می‌شود، هیچ سؤال واقعی در دسترس نیست، طبیعتاً طراحی یک خزانه سؤال که بهینه باشد، هدف مطلوبی است (هی، ریکیسی، ۲۰۱۰). در رویکرد اکتشافی، طراحی خزانه

---

1. Gu  
2. assemble  
3. master

سؤال بهینه به این صورت نیست که از قبل خزانه بزرگی در دسترس باشد، بلکه هیچ سؤال واقعی در دسترس نیست و در آن از روش مونت کارلو برای تعیین ویژگی‌های یک خزانه سؤال بهینه استفاده می‌شود (گو و ریکیسی، ۲۰۰۷).

ادبیات پژوهشی متعدد در مورد طراحی خزانه‌های سؤال بهینه، هر یک از این رویکردها را به صورت مجزا به کار برده‌اند، حال اگر روشی ایجاد شود که از مزایای هر دو رویکرد به طور همزمان استفاده کند، می‌تواند هم از طریق کمی سازی ویژگی‌های آماری و غیر آماری به عنوان قیود تست که از مزایای روش برنامه‌نویسی ریاضی است و هم از تمام مزایای رویکرد ریکیسی استفاده کند و خزانه‌ای بهینه ایجاد کند. حال سؤال مهمی که مطرح می‌شود، این است که برای طراحی یک خزانه سؤال بهینه چه تلاشی باید صورت گیرد؟ بدیهی است که در طراحی یک خزانه سؤال باید ویژگی‌های آماری و غیر آماری سؤال در نظر گرفته شود. برای مثال، توزیع پارامترهای سؤال مطلوب باید چگونه باشد؟ سؤال‌ها موجود در خزانه سؤال CAT باید چه صفاتی داشته باشند؟ از طرف دیگر، پرسش‌هایی از قبیل؛ چه چیزی باعث می‌شود که اندازه خزانه سؤال کافی باشد یا به عبارتی به چند سؤال در خزانه نیاز داریم؟ نیز باید در نظر گرفته شود؛ بنابراین، به طور خلاصه، زمانی که یک خزانه سؤال بهینه طراحی می‌شود، باید حداقل سه عنصر اساسی در نظر گرفته شود، یعنی، ویژگی‌های آماری، ویژگی‌های غیر آماری و اندازه خزانه سؤال. ویژگی‌های آماری شامل پارامترهای سؤال می‌باشند، ویژگی‌های غیر آماری شامل ویژگی‌های محتوایی، توزیع کلید و مهارت‌های شناختی و غیره می‌باشند (هی و ریکیسی، ۲۰۱۰). از این رو، در این مطالعه برای طراحی خزانه سؤال بهینه از تلفیق رویکرد ریکیسی و مدل حداقل انحرافات وزن‌دار استفاده می‌شود و راهنمایی‌های مقدماتی برای تهیه یک خزانه سؤال بهینه برای یک آزمون سنجش انطباقی کامپیوتری که تحت مدل سه پارامتری مدرج شده است و روش کنترل مواجهه سیمپسون-هتر برای کنترل مواجهه بیش از حد در ساخت خزانه سؤال بهینه وارد می‌شود، ارائه خواهد شد. همچنین، عملکرد این خزانه‌ها با یک خزانه سؤال عملیاتی مقایسه می‌شود.

## روش‌شناسی

این پژوهش از دو قسمت نسبتاً مرتبط تشکیل شده است. در قسمت اول شبیه‌سازی‌های مونت‌کارلو ریکسی (۲۰۰۳) و روش برنامه‌نویسی ریاضی WDM، برای طراحی خزانه‌های سؤال بهینه آزمون ریاضی استفاده می‌شود. در قسمت دوم، عملکرد خزانه‌های سؤال بهینه با خزانه سؤال عملیاتی<sup>۱</sup> بر اساس مجموعه‌ای از ملاک‌های ارزیابی مورد مقایسه قرار گرفته است. در این پژوهش از طریق سه روش تصادفی<sup>۲</sup> (R)، روش تصادفی آمیخته و پیش‌بینی<sup>۳</sup> (MRP) و روش حداقل آگاهی آزمون<sup>۴</sup> (MTI) پارامترهای سؤال خزانه‌های سؤال بهینه شبیه‌سازی شده‌اند. تلفیق روش ریکسی (۲۰۰۳) با رویکرد برنامه‌نویسی خطی حداقل انحرافات وزن‌دار، برای طراحی خزانه‌های سؤال بهینه مدرج شده با مدل سه پارامتری لُجستیک به کار رفته است. تعادل محتوایی در هر یک از آزمون‌های CAT، از طریق مدل انتخاب سؤال حداقل انحرافات وزن‌دار (WDM) وارد برنامه‌نویسی شده است. روش کنترل مواجهه سیمپسون-هتر نیز به منظور کنترل مواجهه در مدل وارد شده است (سیمپسون-هتر، ۱۹۸۵). در این پژوهش، یک خزانه عملیاتی در سه حوزه‌ی محتوایی حسابان-دیفرانسیل، هندسه و جبر ساخته شده تا حوزه‌های تجربی معنادار توانایی را اندازه بگیرد. تست‌های عملیاتی برای هر آزمودنی، شامل ۶۰ سؤال توانایی ریاضی است و از لحاظ محتوایی و ویژگی‌های آماری به‌عنوان مبنای طراحی خزانه سؤال بهینه در این پژوهش در نظر گرفته شده‌اند. در این قسمت، علاوه بر شرح موارد مربوط به روش‌شناسی تحقیق، شرح مختصری در مورد روش برنامه‌نویسی ریاضی (WDM) ارائه می‌شود.<sup>۵</sup>

1. Operational
2. Random Procedure (R)
3. Mixed Random and Prediction Procedure (MRP)
4. Minimum Test Information Procedure (MTI)

۵. برای اطلاع از روش "bin-and-union" رویکرد ریکسی، مفهوم bin در مدل سه پارامتری و ایجاد سؤال‌های بهینه در مدل سه پارامتری به مقاله‌ای با عنوان "طراحی خزانه‌های سؤال بهینه برای سنجش انطباقی کامپیوتری با در نظر گرفتن امنیت آزمون" که توسط نویسندگان این مقاله در مجله "مطالعات اندازه‌گیری و ارزشیابی آموزشی" دوره پنجم، شماره دهم به چاپ رسیده است، مراجعه کنید.

کاربرد مدل انتخاب سؤال حداقل انحرافات وزن‌دار (WDM) در تعیین قیود محتوایی. در این پژوهش از روش برنامه‌نویسی خطی حداقل انحرافات وزن‌دار برای تعیین محتواها و ایجاد تعادل محتوایی در خزانه‌های سؤال استفاده شده است. ابتدا، محتوای آزمون CAT، توسط متخصصان موضوعی مشخص شد و پس از توافق کامل میان ۵ متخصص، محتواها به کدهای کمی تبدیل شدند. محتواها به سه مجموعه‌ی اصلی (حسابان-دیفرانسیل، هندسه و جبر) تقسیم‌بندی و به دنبال آن هر یک از مجموعه‌ها به زیرمجموعه‌های معین (به ترتیب، ۱۸، ۱۶ و ۱۱ محتوا) تقسیم‌بندی شدند. سپس از طریق روش برنامه‌نویسی ریاضی کدهای مربوط به هر یک از محتواها، وارد برنامه طراحی خزانه سؤال بهینه شدند. از طریق این روش تست‌های سنجش انطباقی برای ۶۰۰۰ نفر سرهم شدند. در این روش، ابتدا پیش‌بینی جستجوی راه‌حل برای تست کامل صورت می‌گیرد و همزمان هم قابل حل بودن و هم بهینه بودن تست در نظر گرفته می‌شود. این روش جزء روش‌های شهودی حل مسائل سرهم کردن تست است. با کاربرد مدل حداقل انحرافات وزن‌دار به صراحت ویژگی‌های آماری و غیر آماری سؤال‌ها با تعادل مطلوبی بین ویژگی‌های اندازه‌گیری و ساختاری در نظر گرفته می‌شود. این ویژگی‌ها از طریق وزن‌هایی که توسط طراحان تست انتخاب می‌شود، در مدل وارد می‌گردد. به عبارت دیگر، ویژگی‌های محتوایی به عنوان اهداف فرمول‌بندی می‌شوند. انحراف از اهداف محتوایی وزن داده می‌شود و در تابع هدف به همراه فاصله آگاهی سؤال از مقدار هدف قرار داده می‌شود. انتخاب سؤال‌ها در CAT، بر اساس رویکرد WDM طوری تنظیم می‌شود که سؤال‌هایی انتخاب شوند که به‌طور متوالی کوچک‌ترین مجموع انحرافات وزن‌دار را داشته باشند. برای انتخاب یک سؤال از سه گام پیروی می‌شود: (۱). اگر سؤالی که قبلاً در تست نبوده به تست اضافه شود، انحراف برای هر یک از قیود محاسبه شود. (۲). انحرافات وزن‌دار در میان همه قیود جمع شود. (۳). در پایان، سؤالی با کوچک‌ترین مجموع وزن‌دار انحرافات انتخاب شود.

در این روش مدل‌یابی،  $x_i$  متغیر تصمیم‌گیری و  $i = 1, \dots, N$  سؤال‌ها را نشان می‌دهد. اگر سؤال در تست وارد شود،  $x_i = 1$  و اگر سؤال از تست خارج شود  $x_i = 0$ . همچنین، در این مدل  $j = 1, \dots, J$  صفات تست به همراه قیود غیر روان‌سنجی را نشان



می‌دهد و حدود پایین و بالای تعداد سؤال‌هایی که در آزمون دارای چنین ویژگی‌هایی هستند را به ترتیب با  $L_j$  و  $U_j$  مشخص می‌کند، البته ممکن است گاهی با یکدیگر برابر باشد. همچنین، اگر سؤال  $i$  دارای ویژگی  $j$  باشد،  $a_{ij} = 1$ ؛ و اگر سؤال  $i$  دارای ویژگی  $j$  نباشد،  $a_{ij} = 0$ . تعداد سؤال‌ها در خزانه،  $W_j$  وزن اختصاص داده شده به هر قید،  $d_{L_j}$  و  $d_{U_j}$  به ترتیب کسری حد پایین و مازاد حد بالا،  $e_{L_j}$  و  $e_{U_j}$ ، به ترتیب، اضافی حد پایین و کسری حد بالا و  $d_\theta$  انحراف از آگاهی هدف برای یک آزمودنی را نشان می‌دهد. دو جدول ۱ و ۲ به صورت خلاصه اطلاعات مربوط به توابع هدف و قیود مربوط به آن را نشان می‌دهد. قیود تست به عنوان ویژگی‌های غیر آماری یا غیر روان‌سنجی، به همراه ویژگی‌های آماری وارد شبیه‌سازی‌های روش اکتشافی مرحله قبل می‌شود. سپس، انحرافات از این قیدها برای هر یک از ۶۰۰۰ تعداد CAT که از کل خزانه بهینه سرهم می‌شود، محاسبه می‌گردد. به‌طور کلی، این مرحله تلفیقی از دو رویکرد برنامه‌نویسی ریاضی و رویکرد اکتشافی است.

جدول ۱. اطلاعات مربوط به قیود و وزن‌های آزمون‌ها در مورد به حداقل رساندن انحرافات از قیود

$\text{minimize } \sum_{j=1}^J W_j d_{L_j} + \sum_{j=1}^J W_j d_{U_j} + W_\theta d_\theta \rightarrow (\text{objective})$	تابع هدف: به حداقل رساندن میزان انحرافات وزندار
---	--

در ارتباط با قید زیر

$$\sum_{i=1}^{60} a_{ij} x_i + d_{L_j} - e_{L_j} = L_j \Rightarrow j = 1, \dots, J$$

$$\sum_{i=1}^{60} a_{ij} x_i - d_{U_j} + e_{U_j} = U_j \Rightarrow j = 1, \dots, J$$

$$\sum_{i=1}^{60} I(\theta) x_i + d_\theta - e_\theta = \infty$$

جدول ۲. اطلاعات مربوط به قیود و وزن های آزمون های CAT در مورد بیشینه کردن آگاهی تست

تابع هدف: به حداکثر رساندن تابع هدف در ارتباط با قیود زیر		تابع هدف: $\maximize \sum_{i=1}^I I_i(\hat{\theta}^{(g-1)})x_i \rightarrow (objective)$	
حد اکثر	حد اقل	وزن	کد قید
۲۵	۲۵	N1	طول تست
۲۰	۲۰	N2	$\sum_{i=1}^I x_i = 60 \rightarrow test - length$
۱۵	۱۵	N3	
۱۸	۱۸	$z_1 = arithmetic$	
۱۶	۱۶	$z_2 = geometry$	$\sum_{i=1}^{25} z_1 = 18, \sum_{i=1}^{20} z_2 = 16, \sum_{i=1}^{15} z_3 = 11$
۱۱	۱۱	$z_3 = algebra$	
۳	۱	For example : $z_{1-1} = 1$ $z_{1-2} = 2$ $\vdots$ $z_{3-11} = 1$	تعداد سؤال ها در زیرمجموعه های تست $\sum_{i \in I_1} x_i \leq n_1^{(a)} z_1, \sum_{i \in I_1} x_i \geq n_1^{(l)} z_1, \sum_{i \in I_2} x_i \leq n_2^{(a)} z_2, \sum_{i \in I_2} x_i \geq n_2^{(l)} z_2, \sum_{i \in I_3} x_i \leq n_3^{(a)} z_3, \sum_{i \in I_3} x_i \geq n_3^{(l)} z_3$
۷	۱	سه حوزه ی شناختی: $h_1$ به کار بستن $h_2$ تجزیه و تحلیل $h_3$ ترکیب	تعداد سؤال ها در هر سطح شناختی $\sum_{i \in C_h} x_i \geq n^{(l)}_h$ $\sum_{i \in C_h} x_i \leq n^{(a)}_h$

ابزارهای پژوهش حاضر از قرار زیر است:

۱). نحوه طراحی خزانة سؤال عملیاتی. در این پژوهش، خزانة سؤال عملیاتی مربوط به آزمون ریاضی است. این آزمون توانایی حل مسائل حساب، هندسه و جبر را برای دانش آموزانی که در سال آخر دبیرستان تحصیل و خود را برای کنکور سراسری ریاضی فیزیک آماده می کرده اند، طراحی شده است و به عنوان مبنایی برای محک و ارزیابی خزانة های سؤال شبیه سازی شده این مطالعه، به کار رفته است. در این پژوهش، برنامه CAT از طریق زبان PHP نوشته و از پایگاه داده MySQL برای ذخیره سازی سؤال ها در خزانة استفاده شده است. سؤال های طراحی شده برای خزانة سؤال عملیاتی توسط نرم افزار BILOG مدرج شده و سؤال هایی که با مدل سه پارامتری لگجستیک برازش داشتند در خزانة ذخیره شدند (دی آیالا، ۲۰۰۹). همچنین، مفروضات IRT در مورد هریک از

سؤال‌ها بررسی شد. سؤال‌هایی که تک‌بعدی بودن و استقلال موضعی را دارا بودند، در خزانه ذخیره شدند (همبلتون و همکاران، ۱۹۹۱؛ دی‌آیالا، ۲۰۰۹). البته سؤال‌ها طوری طراحی شده بودند که هر کدام مفهوماً کاملاً مستقلی نسبت به سؤال‌های دیگر بسنجند، از این‌رو، اغلب آن‌ها مفروضات IRT را برقرار می‌کردند؛ بنابراین، به هر یک از آن‌ها وزن محتوایی جداگانه‌ای بر اساس نظر متخصصین محتوایی و موضوعی داده شد. روش انتخاب سؤال در فرایند CAT، روش پیشینه آگاهی (MI) بوده است. همچنین، برای صرفه‌جویی در زمان محاسبه، از یک جدول آگاهی نیز استفاده شده است. برای به‌دست آوردن برآورد اخیر توانایی هر آزمودنی، قبل از این که دو پاسخ صحیح و غلط در الگوی پاسخ او مشاهده شود، از روش میانگین پسین (MAP) (اوون، ۱۹۷۵) استفاده گردید، به طوری که پیشین مورد انتظار از توزیع نرمال پیروی می‌کرد. پس از این که دو پاسخ صحیح و غلط در الگوی پاسخ آزمودنی مشاهده شد، برای برآورد توانایی از شیوه پیشینه درست‌نمایی (MLE) استفاده گردید. برآورد  $\theta$  بعد از سؤال آخر، به‌عنوان نمره آزمودنی به حساب آمد. همچنین روش سیمپسون-هتر برای کاهش بیش‌مواجهه سؤال‌هایی که از میزان آگاهی بالایی برخوردار بودند استفاده شد و نرخ مواجهه هدف برابر با  $\frac{1}{3} = \frac{0}{33}$  قرار داده شد. هر آزمون بعد از تعداد ثابتی سؤال، یعنی ۶۰ سؤال، با قیود محتوایی مشخص که در جدول ۱ و ۲ آورده شده است، به اتمام می‌رسد. در این مطالعه تعادل محتوایی به‌عنوان یک عامل مهم در ساخت خزانه سؤال بهینه در نظر گرفته شده است؛ بنابراین، تعامل بین دو عامل کنترل محتوایی سؤال‌ها و کنترل مواجهه سؤال‌های ارائه‌شده به آزمودنی‌ها مورد بررسی قرار گرفته شده است. توانایی اولیه برای هر فرد روی صفر تنظیم شد و برنامه CAT به نحوی برنامه‌ریزی شد که برای همه افراد، سؤال یکسانی که پارامتر دشواری آن صفر ( $b = 0$ ) باشد، اجرا کند. خزانه سؤال عملیاتی شامل ۹۲۱ سؤال (۴۵۵ سؤال حساب دیفرانسیل، ۲۵۸ سؤال هندسه و ۲۰۸ سؤال جبر) بود که توسط ۱۵ طراح ساخته و توسط ۱۰ طراح دیگر ارزیابی موضوعی و تخصصی گردید. سپس

به‌وسیله ۳ روان‌سنج که در درس ریاضی نیز تخصص داشتند بررسی و پس از رفع مشکلات محتوایی و مشکلات احتمالی گزینه‌ها، سؤال‌ها بر روی ۵۰۰ نفر از آزمودنی‌هایی که به همان جامعه تعلق و در اجرای مداد-کاغذی آن نیز شرکت داشته‌اند مدرج‌سازی گردید. این خزانه عملیاتی به شکلی طراحی شد تا ملاک‌های توصیف‌شده توسط وندرلیندن (۲۰۰۰a) را داشته باشد!

۲. روش شبیه‌سازی خزانه سؤال بهینه. برنامه‌های شبیه‌سازی‌شده از طریق نسخه اصلی برنامه (MATLAB (MathWorks, 2014)، به‌منظور شبیه‌سازی الگوی خزانه سؤال و ارزیابی خزانه‌های سؤال شبیه‌سازی‌شده و عملیاتی ایجاد شد. همچنین، برای جستجوی بهترین راه‌حل بهینه برای تعیین قیود محتوایی آزمون، از نرم‌افزار جبر خطی GAMS استفاده شد. شبیه‌سازی الگوی خزانه سؤال بهینه در گام‌های زیر انجام گردید:

**گام اول: مدل یابی کردن شیوه‌های CAT:** از آنجا که هدف این پژوهش ساخت خزانه سؤال بهینه برای سنجش مهارت ریاضی بوده است، در شبیه‌سازی خزانه‌های بهینه نیز تمام ویژگی‌های روان‌سنجی آزمون عملیاتی، به‌دقت وارد گردید. طول آزمون مانند آزمون‌های عملیاتی، ثابت (۶۰ سؤال) قرار داده شد. تعادل محتوایی در برنامه از طریق وارد کردن بهترین وزن‌های محتوایی که توسط حل معادلات جبری به‌دست آمد، در برنامه‌ی شبیه‌ساز وارد شد. روش انتخاب سؤال‌ها روش MI به همراه جدول آگاهی در نظر گرفته شد. به‌منظور برآورد توانایی آزمودنی در طول اجرای آزمون، پیش از اینکه آزمودنی در الگوی پاسخ خود حداقل دو پاسخ صحیح و غلط ایجاد کند، از روش برآورد اوون (۱۹۷۵) برای میانگین پسین<sup>۲</sup> استفاده شد، پس از ایجاد حداقل دو پاسخ صحیح و غلط در الگوی پاسخ، از روش بیشینه‌درست‌نمایی استفاده گردید. همچنین، برای بررسی تعامل

---

۱. این ملاک‌ها عبارت‌اند از: الف) خزانه باید به‌اندازه کافی بزرگ باشد تا این اجازه را به ما بدهد که چندین هزار خرده‌آزمون همپوش از سؤال‌ها به دست آید ب) سؤال‌ها دامنه کاملی از سطوح دشواری نسبت به جامعه موردنظر را پوشش دهند ج) خزانه باید شامل ترکیب مناسبی از سؤال‌ها با ضرایب تشخیص بالا و پایین باشد تا هزینه ساخت سؤال را با در نظر گرفتن دقت آزمون کاهش دهد.

۲. توزیع پیشین توانایی با میانگین صفر و انحراف استاندارد یک در نظر گرفته شد.

عامل کنترل مواجهه (S-H) با تعادل محتوایی آزمون، یک شبیه‌سازی بدون وارد کردن عامل S-H و شبیه‌سازی دیگر با وارد کردن عامل S-H (با وجود تمام شرایط مساوی) انجام گرفت.

**گام دوم: ایجاد وزن‌های محتوایی:** در این مرحله تمام قیود آماری و غیر آماری سؤال‌ها، توسط متخصصان موضوعی به صورت کمی تعیین شد و بهترین راه‌حل بهینه از طریق برنامه‌نویسی خطی با نرم‌افزار GAMS جستجو گردید. پس از این که بهترین وزن‌های محتوایی مشخص شد، در برنامه‌ی شبیه‌سازی CAT با رویکرد ریکرسی تلفیق شد.

**گام سوم: ایجاد جامعه آزمون دهندگان:** از آنجا که، خزانه سؤال عملیاتی برای آزمودنی‌هایی که فرض می‌شود، دارای توزیع توانایی نرمال هستند طراحی گردید، در شبیه‌سازی خزانه سؤال بهینه نیز از توزیع نرمال پیروی شد. دو توزیع حجم نمونه در شبیه‌سازی خزانه سؤال بهینه به کار رفت: به این معنی که خزانه‌های سؤال بهینه با یک نمونه شبیه‌سازی و با نمونه‌ای دیگر ارزیابی شدند. الف) تعداد ۶۰۰۰ توانایی ( $\theta$ ) از توزیع  $N(0,1)$  به طور تصادفی انتخاب شد و به عنوان توانایی واقعی آزمودنی‌ها وارد تحلیل گردید. این نمونه به منظور تعیین ویژگی‌های خزانه سؤال بهینه به کار رفت ب). پیوستار توانایی در دامنه -۴ تا +۴ با فواصل ۰/۱۲۵ به ۶۵ طبقه تقسیم و در هر یک از این سطوح ۵۰۰ آزمودنی قرار گرفت (در کل ۳۲۵۰۰ آزمودنی). این نمونه به منظور ارزیابی عملکرد کلی شبیه‌سازی‌ها و محاسبه آماره‌های مشروط توانایی مورد استفاده قرار گرفت.

**گام چهارم: ایجاد پارامترهای سؤال:** برای هر آزمون CAT، سؤال اول طوری طراحی شد که برای سطح توانایی صفر بهینه باشد. بعد از هر پاسخ، سؤال‌های بهینه‌ای برای برآورد اخیر توانایی تولید شد. البته فرض بر این بود که سؤال‌ها بر پایه مدل سه پارامتری لگجستیک مدرج شده‌اند؛ بنابراین، پارامترهای  $a$ ،  $b$  و  $c$  از طریق سه روش  $MRP$ ،  $R$  و  $MTI$  تولید شدند. در هر سه روش، پارامتر  $c$  بر اساس توزیع بتا<sup>۱</sup> تولید شد. پارامترهای  $a$  بسته به

۱. میانگین و واریانس پارامتر  $c$  برابر با میانگین و واریانس پارامتر  $c$  در خزانه عملیاتی بود. بهترین توزیعی که با این پارامتر برازش داشت توزیع بتا بود.

برآورد اخیر توانایی و روش ایجاد پارامتر (MRP، MTI و R)، ایجاد گردیدند و پارامترهای b طوری تولید شدند که سؤال حداکثر میزان آگاهی در برآورد توانایی جدید را ایجاد کند.

**گام پنجم: ایجاد داده‌های پاسخ:** پاسخ‌های آزمودنی‌ها به دنبال هر سؤالی که بر طبق مدل سه پارامتری لگستیک ایجاد شد، تولید گردید. در مدل سه پارامتری لگستیک احتمال اینکه آزمودنی j به سؤال i پاسخ صحیح دهد، به صورت زیر محاسبه می‌شود:

$$P_i(\theta_j) \equiv c_i + (1 - c_i)(1 + \exp[-1.7a_i(\theta_j - b_i)])^{-1} \quad (1)$$

احتمال اینکه فرد j ( $j=1, 2, \dots, J$ ) با پارامتر  $\theta$ ، به سؤال i ( $i=1, 2, \dots, I$ ) پاسخ صحیح دهد را نشان می‌دهد؛  $a_i$  پارامتر شیب،  $b_i$  پارامتر دشواری و  $c_i$  پارامتر حدس سؤال i است. از آنجا که توانایی واقعی آزمودنی در شبیه‌سازی معلوم بود، بعد از اجرای هر سؤال برای آزمودنی،  $P_i(\theta_j)$  محاسبه گردید. پس از آن، عدد تصادفی  $m_{ij}$  از توزیع یکنواخت  $U(0, 1)$  استخراج و با  $P_i(\theta_j)$  مقایسه می‌شد. اگر  $m_{ij}$  برابر یا کمتر از  $P_i(\theta_j)$  بود، پاسخ برابر با ۱، در غیر این صورت برابر با صفر در نظر گرفته می‌شد.

**گام ششم: تعدیل پس از شبیه‌سازی:** برای هر ترکیبی از روش‌ها و متغیرهای مستقل ۱۰ تکرار صورت گرفت تا جایی که برآورد نسبتاً ثابتی از خزانه سؤال بهینه به دست آمد. از ۱۰ تکرار الگوها و تعدادهای مواجهه سؤال، قبل انجام تعدیل پس از شبیه‌سازی، میانگین گرفته شد.

**جامعه آزمون CAT عملیاتی:** این آزمون برای تمام دانش‌آموزان مقطع پیش‌دانشگاهی که خود را برای کنکور سراسری ریاضی سال ۱۳۹۳ آماده می‌کردند، قابلیت اجرا داشت. البته با این فرض که توزیع این جامعه نرمال با میانگین صفر و انحراف معیار یک بوده است.

**نمونه آزمون CAT عملیاتی:** این آزمون بر روی ۳۵۰ نفر از دانش‌آموزانی که خود را برای آزمون کنکور سراسری ریاضی سال ۱۳۹۳ آماده می‌کردند و در کنکورهای

آزمایشی مرحله‌ای نیز شرکت داشتند، در فواصل فروردین‌ماه تا خرداد ۱۳۹۳ (قبل از کنکور ۹۳) به صورت آنلاین اجرا شد. این نمونه به صورت تصادفی سیستماتیک (از آنجا که لیست کاملی از اسامی شرکت‌کنندگان کنکورهای آزمایشی مرحله‌ای در دسترس بود و اسامی به صورت تصادفی لیست شده بودند، از این روش استفاده شد) انتخاب شد که از یک توزیع نرمال توانایی، با میانگین  $0/17$  و انحراف استاندارد  $0/95$  پیروی می‌کرد. **متغیرهای مستقل پژوهش**. در همه طرح‌های خزانه سؤال دو متغیر مستقل دست‌کاری شدند. روش طراحی پارامترهای سؤال، روش کنترل مواجهه. نرخ مواجهه هدف بر اساس روش سیمپسون-هتر برابر با  $(0/33)$  قرار داده شد. این نرخ در خزانه سؤال عملیاتی نیز برابر با  $0/33$  قرار گرفت (برای جزئیات بیشتر جدول ۳ را ببینید).

جدول ۳. طرح شبیه‌سازی خزانه‌های سؤال بهینه

طول آزمون	۶۰
توزیع توانایی	$N(0,1)$
کنترل مواجهه	عدم کنترل مواجهه روش سیمپسون-هتر (نرخ مواجهه هدف $0/33$ )
مدل تصادفی (R)	مدل تصادفی آمیخته و پیش‌بینی (MRP) مدل کمینه آگاهی آزمون (MTI)
پهنای bin	b-bin: $0/2$ a-bin: $\Delta\alpha^2 = \Delta^2 I_{\max} = 0/4$
تعداد محتوایی	حساب-دیفرانسیل
آزمون سه محتوایی	هندسه
	جبر

**ارزیابی خزانه سؤال شبیه‌سازی شده**. در این پژوهش عملکرد خزانه‌های سؤال بهینه شبیه‌سازی شده با خزانه‌های سؤال عملیاتی بر اساس مجموعه‌ای از ملاک‌های تجربی ارزیابی و مقایسه شد. توزیع نمونه‌گیری دوم مطالعه (توزیع  $32500$  نفری) در ارزیابی

خزانه سؤال مورد استفاده قرار گرفت. ملاک‌های ارزیابی خزانه سؤال به شرح زیر است (چانگ و یینگ<sup>۱</sup>، ۱۹۹۹؛ ریکیسی و هی، ۲۰۰۵):

۱- آگاهی شرطی آزمون<sup>۲</sup>: آگاهی آزمون در هر یک از سطوح توانایی برابر است با مجموع کل آگاهی هر یک از سؤال‌های آزمون در آن سطح توانایی. چون آزمون‌های این مطالعه دارای طول ثابت ۲۰ (I=۱، ۲، ...، I) بودند، آگاهی آزمون، به‌عنوان شاخص کارایی<sup>۳</sup> آزمون در نقاط مختلف توانایی در نظر گرفته شد. هرچه میزان آگاهی آزمون در یک سطح توانایی بیشتر باشد، کارایی آزمون در آن سطح نسبت به سایر سطوح توانایی نیز بیشتر است. برای سطح توانایی j ام (j=۱، ۲، ...، ۶۵) آگاهی آزمون در زیر ارائه شده، نمادهای استفاده شده در این رابطه قبلاً تعریف شده‌اند:

$$I(\theta_j) = \sum_{i=1}^I a_{ij}^2 \frac{(P_{ij} - c_i)^2 \cdot q_{ij}}{(1 - c_i)^2 \cdot p_{ij}} \quad (2)$$

۲- خطای استاندارد شرطی اندازه‌گیری<sup>۴</sup> (CSEM): این شاخص میزان خطای اندازه‌گیری برآورد توانایی را در هر یک از سطوح توانایی واقعی ( $\theta$ ) محاسبه می‌کند:

$$SEM(\theta_j) = \sqrt{\frac{1}{N_i} \sum_{i=1}^{N_i} (\hat{\theta}_{ij} - \bar{\theta}_{ij})^2} \quad (3)$$

اگر  $\theta_j$ ، و توانایی ام (j=۱، ۲، ...، ۶۵) در پیوستار -۴ تا +۴ (یعنی؛ +۴، +۳/۸۷۵، ... -۳/۸۷۵، -۴) را نشان دهد، i هر یک از آزمودنی‌ها در  $\theta_j$  و  $N_i=500$ ، تعداد کل تکرارهای CAT اجرا شده در  $\theta_j$  است.  $\hat{\theta}_{ij}$  ( $\hat{\theta}_{ij}=1, 2, \dots, 500$ ) برآورد  $\theta_{ij}$  و  $\bar{\theta}_{ij} = \frac{1}{N_i} \sum_{i=1}^{500} \hat{\theta}_{ij}$  میانگین ۵۰۰ برآورد  $\theta_{ij}$  ( $\hat{\theta}_{ij}$ ) در  $\theta_j$  است.

- 
1. Chang&Ying
  2. Conditional test information
  3. Efficacy index
  4. Conditional standard error of measurement



۳- اریب و میانگین مجذور خطا<sup>۱</sup> (MSE): در معادله‌های ۴ و ۵،  $N$  تعداد کل آزمودنی‌ها در تمام نقاط ثابت توانایی (۶۵ نقطه ثابت توانایی از -۴ تا ۴) یعنی، برابر با  $(\sum_{j=1}^{65} \sum_{i=1}^{500} N_{ij} = 32500)$  آزمودنی است و  $\hat{\theta}_i$  برآورد کننده  $\theta_i$  آمین آزمودنی با سطح توانایی واقعی  $\theta_i$  است:

$$Bias = \frac{1}{N} \sum_{i=1}^N (\hat{\theta}_i - \theta_i) \quad (۴)$$

$$MSE = \frac{1}{N} \sum_{i=1}^N (\hat{\theta}_i - \theta_i)^2 \quad (۵)$$

۴- اریب شرطی<sup>۲</sup> و میانگین مجذور خطای شرطی<sup>۳</sup> (CMSE): در اصل در معادله‌های ۶ و ۷ مقدار  $\theta_j$  همان مقادیر ثابت (یعنی؛ +۴، +۳/۸۷۵، ...، -۳/۸۷۵، -۴) است، که اینجا برآوردهای  $\hat{\theta}_{ij}$  در هر سطح توانایی، از مقدار ثابت  $\theta_j$  کم می‌شود:

$$Conditional\ Bias(\theta_j) = \frac{1}{N_i} \sum_{i=1}^{N_i} (\hat{\theta}_{ij} - \theta_j) \quad (۶)$$

$$CMSE(\theta_j) = \frac{1}{N_i} \sum_{i=1}^{N_i} (\hat{\theta}_{ij} - \theta_j)^2 \quad (۷)$$

۵- کجی توزیع نرخ مواجهه سؤال؛ آماره کای دو که توسط چانگ و بینگ (۱۹۹۹) ارائه شده، برای اندازه‌گیری میزان کجی توزیع مواجهه سؤال به کاررفته و برابر است با:

$$X^2 = \sum_{i=1}^n \frac{(r_i - \frac{L}{n})^2}{\frac{L}{n}} \quad (۸)$$

- 
1. Bias and mean square error
  2. Conditional Bias
  3. Conditional mean square error
  4. Skewness of item exposure rate distribution

در این معادله  $T_i$ ، نسبت نرخ مشاهده شده  $I$  امین سؤال،  $L=20$ ، طول آزمون و  $n$  تعداد سؤال‌های خزانه است. معادله ۸ اختلاف بین نرخ مواجهه سؤال مشاهده شده و ایده آل را محاسبه می‌کند و مقدار اثربخشی استفاده از خزانه سؤال را نیز تعیین می‌نماید. مقدار کای دو کوچک نشان می‌دهد که بیشتر سؤال‌ها استفاده شده‌اند.

۶- درصد سؤال‌های بیش مواجهه شده<sup>۱</sup>: نرخ مواجهه یک سؤال را می‌توان به‌عنوان نسب تعداد دفعات اجراهای سؤال به تعداد کل آزمودنی‌ها در نظر گرفت. در مجموع، سطح متوسط نرخ مواجهه سؤال مناسب است. نرخ بالای مواجهه یک سؤال بدین معناست که خطر فاش شدن سؤال برای آزمودنی‌های بعدی افزایش می‌یابد. اگر چنین باشد، هم امنیت و هم روایی<sup>۲</sup> آزمون به دلیل نرخ بالای مواجهه سؤال مورد تهدید قرار می‌گیرد؛ بنابراین، درصد سؤال‌های بیش مواجهه شده، به‌عنوان ملاک مهمی برای ارزیابی موفقیت برنامه CAT در نظر گرفته می‌شود (هو و چانگ<sup>۳</sup>، ۲۰۰۱).

۷- درصد سؤال‌های کم مواجهه شده<sup>۴</sup>: نرخ کم مواجهه شدن یک سؤال بدین معناست که یک سؤال بندرت در برنامه CAT مورد استفاده قرار گیرد. خزانه سؤالی که سؤال‌های بسیار زیادی با نرخ مواجهه خیلی پایینی دارد، دارای فایده کمی است. دو موضوع به‌صرفه بودن طراحی سؤال‌ها و مناسب بودن شیوه انتخاب آن‌ها، به دلیل نرخ مواجهه کم سؤال به چالش کشیده می‌شوند. نرخ مواجهه پایین‌تر از ۰/۰۲ به‌عنوان سؤال کم مواجهه شده در نظر گرفته می‌شود (هو و چانگ<sup>۳</sup>، ۲۰۰۱، ریکیسی<sup>۵</sup>، ۲۰۰۹).

۸- نرخ همپوشی آزمون<sup>۶</sup>: نرخ همپوشی آزمون، عبارت است از تعداد مورد انتظار سؤال‌های مشترکی که به دو آزمودنی که به‌طور تصادفی نمونه‌گیری شدند، ارائه می‌شود، تقسیم بر طول مورد انتظار آزمون<sup>۶</sup>. به‌طور ایده آل، تعداد سؤال‌های مشترک بین دو

- 
1. Percentage of overexposed items
  2. Validity
  3. Hau & Chang
  4. Percentage of underexposed items
  5. Test overlap rate
  6. Expected test length

آزمودنی که به‌طور تصادفی نمونه‌گیری شدند، باید حداقل باشد (چانگ و بینگ، ۱۹۹۹؛ چن، آنکنمان، اسپری<sup>۱</sup>، ۱۹۹۹):

$$\begin{aligned}\bar{T} &= \frac{\sum_{i=1}^n \binom{m_i}{2}}{L \binom{N}{2}} \\ &= \frac{\sum_{i=1}^n m_i(m_i - 1)}{LN(N - 1)}\end{aligned}\quad (9)$$

در این رابط  $N$ ، تعداد CAT‌های (با طول ثابت) اجراشده،  $L$ ، تعداد سؤال‌های در هر یک از CAT‌ها،  $n$  تعداد سؤال‌های خزانه و  $m_i$  تعداد دفعاتی است که سؤال  $i$  برای همه  $N$  تعداد CAT اجرا شده است.

سؤال‌های زیر در این پژوهش بررسی می‌شوند:

- ۱) آیا روش شبیه‌سازی مونت کارلو ریکیسی (۲۰۰۳) می‌تواند با روش برنامه‌نویسی ریاضی (WDM) تلفیق شود؟
- ۲) عملکرد خزانه سؤال بهینه برای CAT زمانی که در الگوریتم انتخاب سؤال، علاوه بر تعادل محتوایی، مواجهه بیش‌ازحد سؤال وارد می‌شود در مقابل زمانی که مواجهه بیش‌ازحد وارد نمی‌شود، چگونه است؟
- ۳) آیا خزانه‌های سؤال بهینه‌ای که طراحی می‌شوند بهتر از خزانه سؤال عملیاتی برحسب ملاک تجربی عمل می‌کنند؟

### یافته‌های پژوهش

در این قسمت نتایج مربوط به شبیه‌سازی خزانه‌های بهینه و مقایسه عملکرد آن‌ها با خزانه سؤال عملیاتی شرح داده می‌شود. همان‌طور که در جدول ۴ نمایش داده شده است، در این پژوهش برای ساخت خزانه‌های سؤال بهینه‌ای که در آن تعادل محتوایی نیز از اهمیت برخوردار است، دو عامل (سه روش ایجاد سؤال بهینه و عدم اعمال شیوه کنترل

مواجهه سیمپسون-هتر) دست‌کاری شده‌اند، بنابراین، تعداد شش الگوی خزانه سؤال بهینه به وجود آمده است. به خزانه‌های سؤال بهینه‌ای که بر اساس روش 'bin-and-union' ایجاد شدند، به اختصار ROP<sup>۲</sup> گفته می‌شود (هی و ریکسی، ۲۰۱۱). همچنین، خزانه سؤال عملیاتی به اختصار با OP<sup>۳</sup> نشان داده می‌شود (گو و ریکسی، ۲۰۰۷). پهنای b-bin ها در تمام خزانه‌ها برابر با ۰/۲ و میزان تغییر آگاهی در پارامتر a برابر با ۰/۴ قرار داده شد.

جدول ۴. ترکیب عامل روش ایجاد سؤال بهینه و شیوه کنترل مواجهه برای

ساخت خزانه‌های سؤال بهینه با تعادل محتوایی

شیوه کنترل مواجهه سیمپسون-هتر Exposure control- ) Simpson-Hetter procedure	روش ایجاد سؤال بهینه			
	item generation methods without- Exposure control Exposure control	R ROP_1 ROP_4	MRP ROP_2 ROP_5	MTI ROP_3 ROP_6

**خزانه‌های سؤال بهینه بدون کنترل مواجهه بیش از حد سؤال.** در این مرحله خزانه‌های بهینه با توجه به ایجاد تعادل محتوایی بر اساس روش برنامه‌نویسی WDM و بدون عامل کنترل مواجهه سیمپسون-هتر (S-H) شبیه‌سازی شده‌اند. روش ایجاد سؤال بهینه از طریق سه روش R، MRP، MTI با پهنای  $b\text{-bin} = 0.2$  و میزان آگاهی یا پهنای پارامتر a برابر با  $a\text{-bin}: \Delta a^2 = 2\Delta I_{Maximum} = 0.4$  بوده است. جدول ۵، اندازه‌ها و خلاصه آماره‌های مربوط به پارامترهای سؤال در خزانه‌های سؤال بهینه و عملیاتی را در سه محتوای مشخص ارائه می‌کند. نتایج نشان می‌دهد که خزانه‌های سؤال بهینه شامل حداقل تعداد سؤال می‌باشند. البته یکی از دلایل آن می‌تواند این قضیه باشد که در ساخت آن‌ها هیچ نوع کنترل مواجهه‌ای صورت نگرفته است. توزیع پارامترهای سؤال در ماتریس‌های مربوط به هر یک از خزانه‌ها نشان می‌دهد<sup>۴</sup> که همه خزانه‌های بهینه دارای سؤال‌هایی با دامنه

۱. برای یادآوری به صفحه ۶ در مقدمه مراجعه کنید.

2. Range-optimal item pool

3. Operational item pool

۴. به دلیل محدودیت صفحات از آوردن ماتریس پارامترها خودداری می‌شود و تنها به ذکر نتایج استخراج شده از این ماتریس‌ها اکتفا می‌شود.

محدودی از سطوح دشواری هستند. دلیل این امر این است که زمانی که قيود محتوایی در تعامل با ویژگی‌های روان‌سنجی قرار می‌گیرند، خزانه‌های سؤال بهینه شبیه‌سازی دارای ویژگی‌های روان‌سنجی دقیق‌تری می‌شوند، به طوری که دامنه دشواری سؤال‌ها محدودتر می‌شود. خزانه‌های سؤال بهینه دارای میانگین دشواری بالاتری نسبت به خزانه‌های عملیاتی هستند. خزانه بهینه MTI در هر سه محتوا دارای حداقل تعداد سؤال است، ولی تفاوت زیادی با خزانه‌های MRP ندارد. با این وجود، خزانه‌های بهینه MRP دارای بیشترین مقدار پارامتر  $a$  هستند. ولی خزانه‌های بهینه  $R$  و MTI دارای میانگین پارامتر  $a$  مشابهی هستند. خزانه‌های MTI به دلیل ماهیت ایجاد سؤال، دارای حداقل میزان پراکندگی در پارامتر  $a$  هستند. توزیع خزانه‌های سؤال  $R$  نسبت به دو خزانه دیگر دارای یک توزیع یکنواخت‌تری در سراسر ماتریس پارامترها است، این نتیجه به دلیل ماهیت روشی است که پارامترهای سؤال را ایجاد می‌کند. در این روش، پارامترها در سراسر ماتریس پراکنده می‌شوند. توزیع پارامتر دشواری و تشخیص سؤال‌ها در این روش بسیار مشابه خزانه عملیاتی است و در تمام محتواها دارای مقادیر پارامتر متنوع‌تری است. اما سؤال‌های دشوار در خزانه‌های بهینه MRP دارای پارامتر ضریب تشخیص بالاتری هستند و برعکس سؤال‌ها آسان دارای پارامترهای ضریب تشخیص متوسط یا پایین‌تری هستند. این نتایج باعث می‌شود که تعداد آزمون‌هایی که در خزانه‌های  $R$  از قيود محتوایی تخطی می‌کنند، در سرتاسر پارامتر توانایی یکنواخت باشد. خزانه‌های MRP در پارامترهای توانایی بالاتر از متوسط، دارای تخطی از قيود کمتری هستند. بررسی نتایج عملکرد این خزانه‌ها در جدول ۶ آورده شده است. برآورد توانایی در هر سه خزانه بهینه و عملیاتی، دارای سطح معینی از اریب مثبت است، ولی مقدار این اریب‌ها بسیار کوچک است. میانگین مجذور خطا (MSE) در خزانه‌های سؤال بهینه کوچک‌تر از خزانه سؤال عملیاتی است. در این میان خزانه‌های سؤال بهینه MTI عملکرد بهتری در شاخص (MSE) نشان می‌دهند، زیرا تعامل میزان حداقل آگاهی در هر سطح با ویژگی‌های محتوایی مورد توجه قرار گرفته است و بنابراین، توانایی را با دقت بیشتری برآورد می‌کنند. خزانه‌های سؤال بهینه دارای نرخ همپوشی پایین‌تری هستند، با وجود این که دارای سؤال‌ها کمتری می‌باشند. نمودار ۱ نشان می‌دهد که نرخ

همپوشی تست در سطوح توانایی زیر ۲- در خزانه‌های عملیاتی و بهینه تقریباً مشابه است. ولی در بقیه سطوح نرخ همپوشی خزانه‌های بهینه کمتر از عملیاتی است. خزانه‌های بهینه در مقدار همپوشی تست‌ها شباهت زیادی به یکدیگر دارند. دلیل این شباهت به خاطر ویژگی مشترک همه آن‌ها در ایجاد پارامتر  $b$  است، زیرا پارامتر  $b$  ارتباط مستقیمی با انتخاب و سرهم شدن تست‌ها دارد. در این پژوهش سؤال عملیاتی همانند آزمون‌های مشابه دیگر در نوع خود، در دو انتهای سطوح توانایی دارای نرخ مواجهه کم نیست، زیرا این برنامه طوری طرح‌ریزی شده که نقاط برش در سطوح بالای توانایی باشند پس طبیعتاً به سؤال‌های دشوارتری نیاز دارد. از این رو، تنها در سطوح پایین توانایی نرخ همپوشی تست کمتر است. همچنین، خزانه‌های بهینه درصد خیلی کوچک‌تری از کم مواجهه شدن سؤال‌ها را نسبت به خزانه عملیاتی دارند. نرخ سؤال‌های بیش مواجهه شده در هر سه خزانه تقریباً مشابه و بیشتر از خزانه عملیاتی است. نمودارهای ۲ تا ۵ درصد‌های مواجهه سؤال در هر یک از سطوح توانایی را برای خزانه‌های عملیاتی و بهینه نشان می‌دهد. در همه خزانه‌های سؤال، سؤال‌ها خیلی آسان و خیلی دشوار که به ترتیب در سطوح پایین و بالای توانایی ارائه می‌شوند، دارای نرخ‌های مواجهه کوچک‌تری هستند. نتایج نشان می‌دهد که در همه خزانه‌ها، سؤال‌های با سطوح دشواری متوسط بیشترین قابلیت استفاده را دارند. در این نمودارها نرخ‌های مواجهه سؤال، به تفکیک محتواها در هر چهار خزانه گزارش شده است. شکل نمودارها الگوی جالبی از مواجهه را در هر یک از خزانه‌ها نشان می‌دهد. در خزانه‌های عملیاتی، سؤال‌ها کم مواجهه شده در سرتاسر طول پیوستار توانایی پراکنده شدند. خزانه‌های بهینه  $R$  دارای الگویی با توزیع نرمال در میزان مواجهه سؤال‌ها هستند. خزانه‌های  $MRP$  در سطوح بالای توانایی دارای نرخ مواجهه کمتری هستند. در پایان، خزانه‌های بهینه  $MTI$  تنها در یک بازه‌ی کوچک در متوسط پارامتر توانایی دارای مواجهه بیش‌ازحد می‌باشند و در سطوح دیگر میزان مواجهه کم می‌شود. جدول ۷ نرخ‌های مواجهه سؤال‌ها را به تفکیک محتواها در هر چهار خزانه سؤال نشان می‌دهد. نتایج جدول ۷ گویای این مطلب است که در محتوای اول، دوم و سوم خزانه‌های  $MRP$  دارای بیشترین نرخ مواجهه هستند، زیرا ویژگی‌های پارامترهای سؤال‌ها در خزانه‌های عملیاتی (مانند میزان

ضریب همبستگی بین پارامترهای  $a$  و  $b$  در تعامل با پارامتر توانایی) با ویژگی‌های این روش ایجاد سؤال تعامل برقرار کرده و سؤال‌هایی با ضریب تشخیص بالا در سطوح بالای توانایی ایجاد می‌کنند که با توجه به ماهیت آزمون (میانگین دشواری آزمون بالاتر از متوسط است)، این نتیجه به دست می‌آید. همچنان که در نمودار ۶ ملاحظه می‌شود، میانگین آگاهی خزانه‌های سؤال در خزانه‌های  $MTI$ ،  $R$  و خزانه سؤال عملیاتی در تمام سطوح کاملاً مشابه است. خزانه  $MRP$  در سطوح توانایی زیر ۲ کاملاً مشابه خزانه‌های دیگر است ولی در سطوح توانایی بالای ۲ متوسط میزان آگاهی بسیار بیشتر می‌شود که به دلیل ماهیت روش ایجاد پارامترهای سؤال است. نتایج مربوط به آگاهی تست‌ها نمایانگر این است که اگر در ساخت خزانه‌های بهینه سؤال به ویژگی‌های کاربردی آزمون‌های عملیاتی توجه شود، نتایج خزانه‌های بهینه و عملیاتی بسیار به یکدیگر شبیه می‌شود، با این تفاوت که اندازه خزانه‌های سؤال بهینه بسیار کمتر از خزانه سؤال عملیاتی است.

نمودار ۷، ۸ و ۹ به ترتیب، خطای استاندارد شرطی اندازه‌گیری (CSEM)، اریب شرطی و میانگین مجذور خطا (CMSE) را در هر چهار خزانه سؤال نشان می‌دهد. خطاهای استاندارد اندازه‌گیری الگویی مشابه با خزانه‌های عملیاتی دارند. ولی در تمام سطوح توانایی کمتر از خزانه عملیاتی است. نمودار ۸ نشان می‌دهد که در خزانه عملیاتی میزان اریب در اغلب سطوح توانایی بیشتر از خزانه‌های بهینه است. خزانه  $MTI$  نسبت به خزانه‌های دیگر از میزان اریب کمتری برخوردار است. دلیل این نتیجه تعامل بین ویژگی‌های محتوایی و حداقل میزان آگاهی آزمون است. همچنین، نمودار ۹ میانگین مجذور خطا را در سطوح متفاوت توانایی نشان می‌دهد. نتایج این نمودار نشان می‌دهد که  $MSE$  هر سه خزانه بهینه کوچک‌تر از خزانه سؤال عملیاتی است. به خصوص خزانه  $MTI$  از حداقل مقدار  $MSE$  برخوردار است.

جدول ۵. اندازه خزانه سؤال و آماره های پارامتر سؤال، بدون S-H (b-bin=0.2)، با تعادل محتوا

خزانه سؤال	اندازه خزانه	a				b				c			
		میانگین	انحراف استاندارد	حداکثر	حداقل	میانگین	انحراف استاندارد	حداکثر	حداقل	میانگین	انحراف استاندارد	حداکثر	حداقل
Content 1 (arithmetic)													
OP	۴۵۵	۱/۰۸۹	۰/۸۸۴۴	۳/۰۴۵	۰/۱۶۶	-۰/۰۳۹	۳/۹۸۱	-۳/۵۹۶	۰/۱۴۵	۰/۸۰۱	۰/۴۱۷۹	۰/۰۰۵	
ROP_1	۳۳۰	۱/۰۸۹	۰/۲۷۷	۲/۰۸۸	۰/۴۶	۰/۰۱۶	۳/۶۵	-۲/۸۲	۰/۱۴۶	۰/۰۷	۰/۴۱۰۲	۰/۰۰۱	
ROP_2	۱۷۸	۲/۰۱۲	۰/۸۸۱	۳/۱۲	۰/۹۵	۰/۰۲۰۱	۳/۸۲۵	-۳/۱۱۸	۰/۱۴۸	۰/۰۷۲	۰/۴۳	۰/۰۲۴	
ROP_3	۱۷۴	۱/۵۶۲	۰/۲۵۳	۲/۳۲۴	۰/۹۸۷	۰/۰۱۴۶	۳/۸۴۱	-۳/۵۹۲	۰/۱۴۲	۰/۰۶۴	۰/۴۶	۰/۰۰۱	
Content 2 (geometry)													
OP	۲۵۸	۱/۲۰۶	۰/۲۲۴۵	۲/۹۳	۰/۲۴۵	-۰/۰۴۸۲	۳/۸۴	-۳/۴۵۸	۰/۱۸۴	۰/۰۹۷	۰/۴۷۶	۰/۰۹۱	
ROP_1	۱۹۹	۱/۸۳۲	۰/۲۶۳	۲/۶۵۴	۰/۵۲۶	-۰/۰۸۵۵	۳/۵۵	-۳/۵۵۴	۰/۱۷۱	۰/۰۷۴	۰/۴۲	۰/۰۰۱	
ROP_2	۱۶۸	۱/۸۸۲	۰/۲۷۸	۲/۹۳۵	۰/۹۲۱	۰/۰۱۲۸	۳/۷۴	-۳/۴۹	۰/۱۸۳	۰/۰۶۸	۰/۴۶	۰/۰۰۱۵	
ROP_3	۱۶۷	۱/۸۵۹	۰/۲۷۲	۲/۵۱	۱/۱۵	-۰/۰۱۳۵	۳/۱۰۴	-۳/۱۰۴	۰/۱۷۸	۰/۰۸۱	۰/۴۸	۰/۰۰۲	
Content 3 (algebra)													
OP	۲۰۸	۱/۳۵۶	۰/۲۴۷	۲/۸۸۹	۰/۵۳۸	۰/۰۳۴۴	۳/۶۸۹	-۳/۴۰۹	۰/۱۷۴	۰/۰۷۵۴	۰/۴۸۹	۰/۰۰۴	
ROP_1	۱۸۳	۱/۷۸۳	۰/۲۶۸	۲/۷۸۹	۰/۳۳۵	۰/۰۲۳۸	۳/۵۳۲	-۳/۴۵۶	۰/۱۸۱	۰/۰۶۲	۰/۴۹۷	۰/۰۰۳	
ROP_2	۱۵۷	۲/۱۴۸	۰/۲۹۱	۳/۰۴۶	۰/۹۵	۰/۰۱۸۴	۳/۵۸۹	-۳/۱۴	۰/۱۸۹	۰/۰۷۵	۰/۴۶۵	۰/۰۰۲	
ROP_3	۱۵۵	۱/۷۹۴	۰/۲۷۹	۲/۴۸۹	۰/۹۰۲	-۰/۰۱۰۹	۳/۵۹۸	-۳/۰۸۹	۰/۱۸۴	۰/۰۷۴۸	۰/۴۵	۰/۰۰۱	



جدول ۶: خلاصه‌ی آماره‌های عملکرد خزانه سؤال، بدون S-H ( $b\text{-bin}=0.2$ )، با تعادل محتوا

آماره‌ها	OP	R	MRP	MTI
Bias	۰/۰۰۵۲	۰/۰۰۲۱	۰/۰۰۰۷۳	۰/۰۰۱۹۹
MSE	۰/۰۱۷۴۵	۰/۰۱۲۶۷	۰/۰۱۱۳۸	۰/۰۰۰۷۶
کجی نرخ مواجهه	۹۸/۸۵۶	۵۱/۴۰۲	۴۸/۴۳۶۵	۴۳/۴۹۴
نرخ همپوشی سؤال	۰/۵۱۷۳	۰/۳۷۶۵	۰/۳۸۴۴	۰/۳۹۴۱
درصد سؤال‌های با نرخ مواجهه بزرگ‌تر از $\frac{1}{3}$	۸/۱۶۹٪	۹/۶۳۴٪	۱۱/۵۲۴٪	۱۰/۰۰۳٪
درصد سؤال‌های با نرخ مواجهه کوچک‌تر از 0.02	۵۱/۵۲۶٪	۲۵/۵۳۱٪	۲۱/۱۵٪	۱۹/۱۸۹٪
درصد تست‌هایی که از قیود تست تخطی دارند	۵۴/۸٪	۳/۰۰۲٪	۲/۰۰۱۳٪	۲/۰۰۱۵٪
اندازه خزانه سؤال	۹۲۱	۶۱۲	۵۰۳	۴۹۶

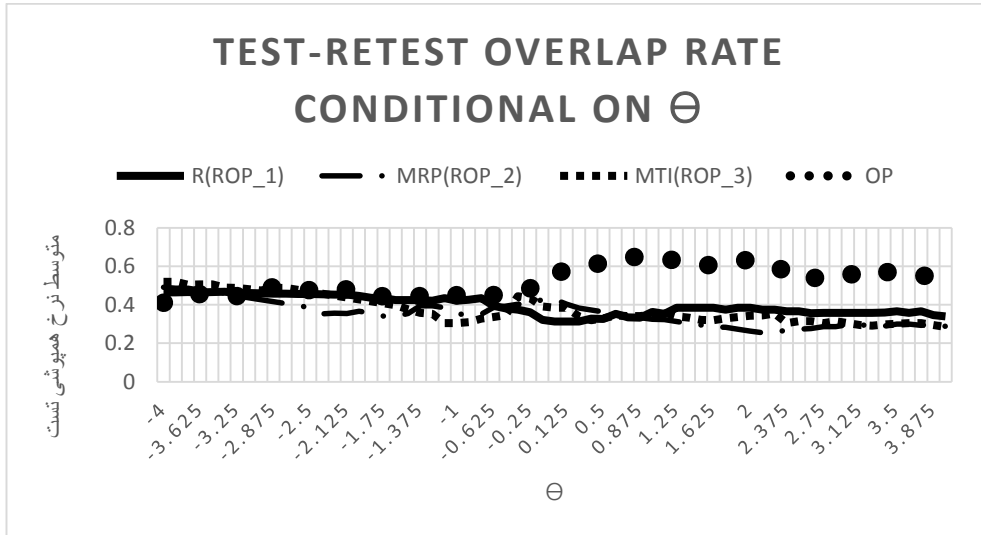
جدول ۸ و ۹ نتایج مربوط به میزان تخطی از قیود را در تست‌ها برحسب پارامترهای توانایی نشان می‌دهد. در CAT عملیاتی، ۳۵۰ تست سرهم شدند و در CAT های شبیه‌سازی شده، ۶۰۰۰ تست سرهم شدند، برای قابل‌مقایسه کردن میزان تخطی‌ها، درصد فراوانی نسبی آن‌ها محاسبه شد. مطابق با نتایج این دو جدول، خزانه عملیاتی در تمام سطوح توانایی میزان تخطی از قیود بیشتری نسبت به خزانه‌های بهینه دارد. در خزانه‌های R میزان تخطی‌ها در دو دامنه توانایی بیشتر است و تقریباً الگویی مشابه با خزانه عملیاتی دارد. در خزانه MRP و MTI در سطوح پایین توانایی میزان تخطی از قیود بیشتر است. در خزانه MRP در سطوح بالای توانایی به دلیل وجود سؤال‌ها بیشتر، تخطی از قیود در تست‌هایی که سرهم می‌شود، به حداقل خود می‌رسد.

جدول ۷. درصد سؤال‌ها بیش مواجهه شده و کم مواجهه شده

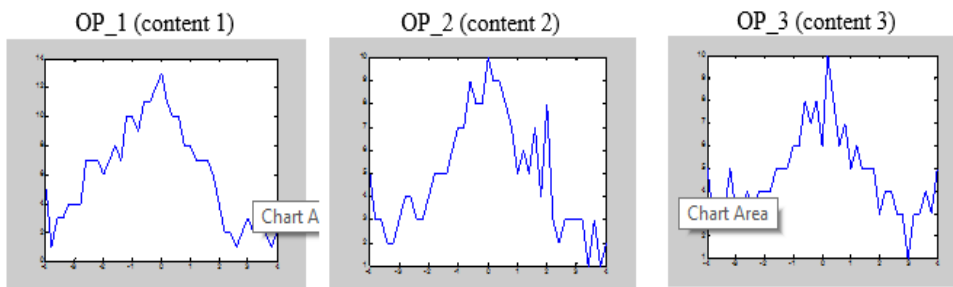
در هر یک از محتواها برای خزانه‌های سؤال بدون S-H

آماره	محتوای اول				محتوای دوم				محتوای سوم			
خزانه‌های سؤال	MTI	MRP	R	OP	MTI	MRP	R	OP	MTI	MRP	R	OP
درصد سؤال‌های با نرخ مواجهه بزرگ‌تر از $\frac{1}{3}$	۱۱/۷۷٪	۱۳/۴۶٪	۹/۹۲٪	۹/۸۱٪	۸/۵۲٪	۹/۹۷٪	۹/۵۸٪	۷/۵۸٪	۸/۳۱٪	۹/۹۷٪	۷/۸۳٪	۷/۰۸٪
اندازه خزانه	۱۷۴	۱۷۸	۲۳۰	۴۵۵	۱۶۷	۱۶۸	۱۹۹	۲۵۸	۱۵۵	۱۵۷	۱۸۳	۲۰۸

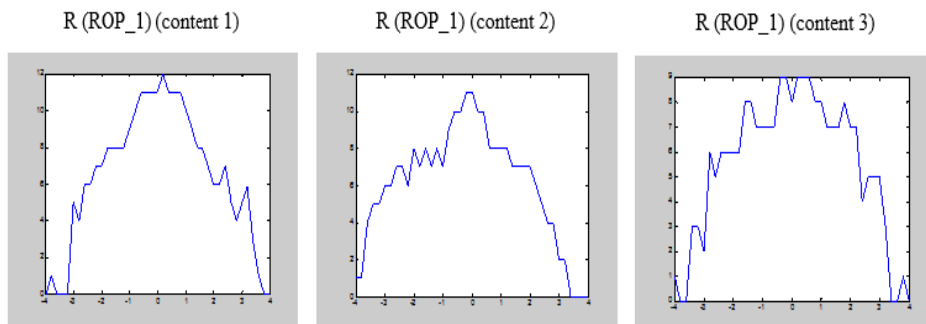




نمودار ۱. نرخ همپوشی تست مشروط به  $\Theta$  بدون **S-H** ( $b\text{-bin: } 0.2$ ) با تعادل محتوایی

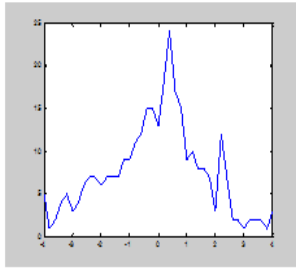


نمودار ۲. درصد سؤال‌ها بیش مواجهه شده در خزانه عملیاتی برحسب محتواهای سه‌گانه

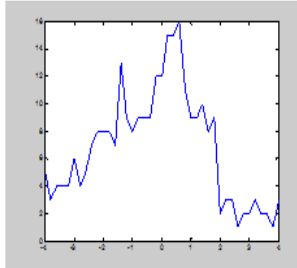


نمودار ۳. نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده در سطوح توانایی بر اساس روش R

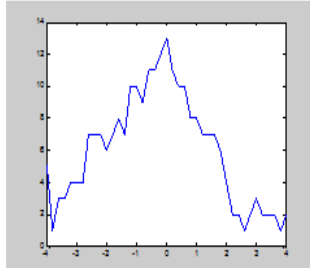
MRP (ROP\_2) (content 1)



MRP (ROP\_2) (content 2)

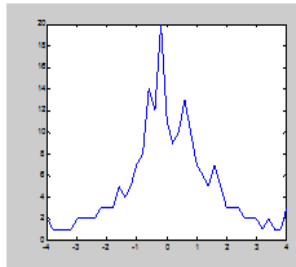


MRP (ROP\_2) (content 3)

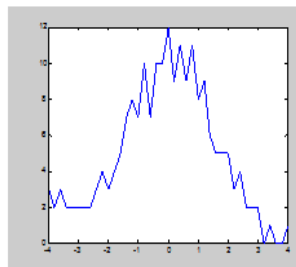


نمودار ۴. نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده در سطوح توانایی بر اساس روش MRP

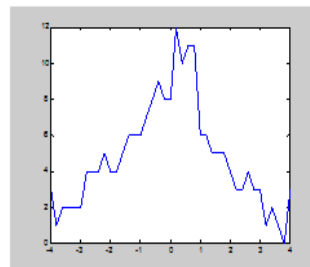
MTI (ROP\_3) (content 1)



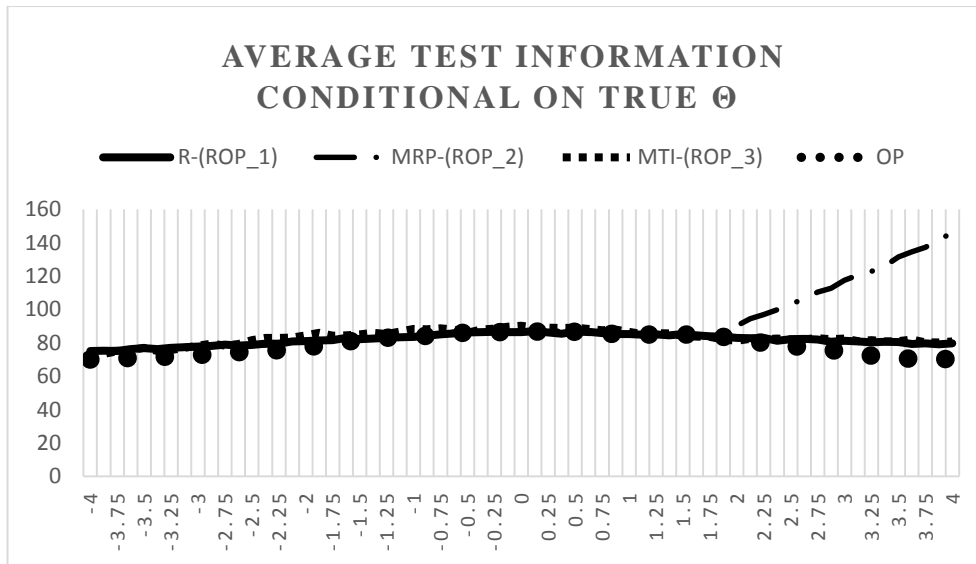
MTI (ROP\_3) (content 2)



MTI (ROP\_3) (content 3)

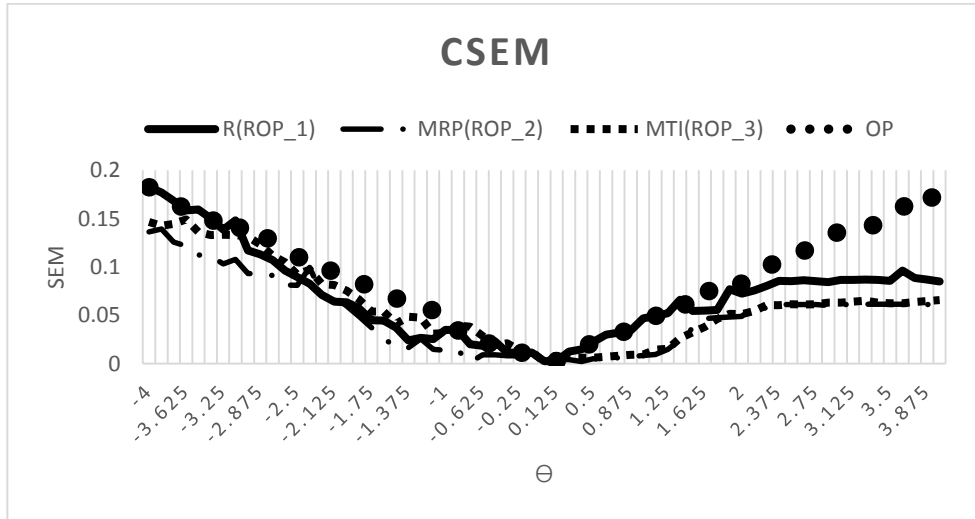


نمودار ۵. نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده بر اساس روش MTI

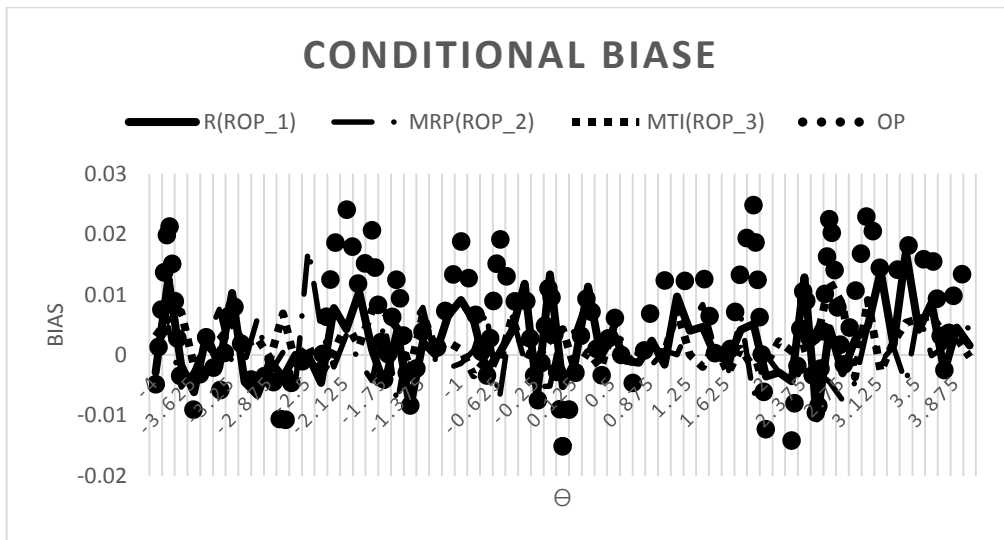


نمودار ۶. متوسط آگاهی تست مشروط به توانایی واقعی در خزانه‌های سؤال

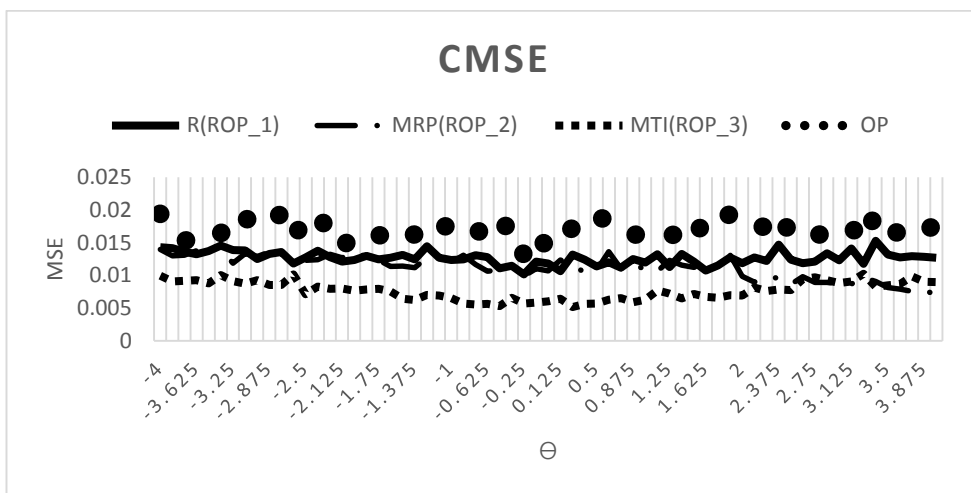
بدون S-H ( $b$ -bin: 0.2)، با تعادل محتوا



نمودار ۷. خطای استاندارد اندازه‌گیری (CSEM) در خزانه‌های سؤال بدون S-H ( $b\text{-bin: } 0.2$ ), با تعادل محتوایی



نمودار ۸. اربیب شرطی (conditional Bias) در خزانه‌های سؤال بدون S-H ( $b\text{-bin: } 0.2$ ), با تعادل محتوایی



نمودار ۹. میانگین مجذور خطا (CMSE) در خزانه‌های سؤال بدون S-H ( $b$ -bin: 0.2)، با تعادل محتوایی

**خزانه‌های سؤال بهینه با کنترل مواجهه بیش از حد سؤال.** در این مرحله، خزانه‌های سؤال با روش کنترل مواجهه S-H، شبیه‌سازی شدند. جدول ۱۰ اندازه‌ها و خلاصه آماره‌های مربوط به پارامترهای سؤال در خزانه‌های سؤال بهینه و عملیاتی را در هر سه محتوای ارائه می‌کند. نتایج نشان می‌دهد که با وجود کنترل مواجهه سؤال در ساخت خزانه‌های سؤال بهینه، هنوز این خزانه‌ها شامل حداقل تعداد سؤال می‌باشند. در این مرحله هم همانند خزانه‌هایی که بدون کنترل مواجهه ایجاد شدند، هر سه خزانه بهینه دارای سؤال‌های با دامنه محدودی از سطوح دشواری هستند و میانگین پارامتر دشواری بالاتری دارند. در این مرحله نیز خزانه‌های بهینه MTI در هر سه محتوا دارای حداقل تعداد سؤال است. خزانه‌های MRP (ROP\_5) نسبت به MRP (ROP\_2) دارای تعداد سؤال‌های بسیار بیشتری هستند. در صورتی که مقایسه خزانه‌های R (ROP\_4) نسبت به R (ROP\_1) و MTI (ROP\_6) نسبت به MTI (ROP\_3) به این اندازه تفاوت نشان نمی‌دهد. در این مرحله نیز خزانه‌های MTI دارای حداقل میانگین پارامتر  $a$  می‌باشند و خزانه‌های بهینه MRP دارای بیشترین مقدار پارامتر  $a$  هستند. نوع توزیع پارامترهای سؤال در هر سه خزانه سؤال بهینه مشابه خزانه‌های سؤال در مرحله قبل است، با این تفاوت که میانگین ضریب

تشخیص در هر سه خزانه بالاتر است. بررسی نتایج عملکرد این خزانه‌ها در جدول ۱۱ ارائه شده است. برآورد توانایی در هر سه خزانه بهینه و عملیاتی، دارای سطح معینی از اریب مثبت است، ولی مقدار این اریب‌ها بسیار ناچیز است. میانگین مجذور خطا (MSE) در خزانه‌های سؤال بهینه کوچک‌تر از خزانه سؤال عملیاتی است. در میان خزانه‌های سؤال بهینه، MRP عملکرد بهتری در این شاخص نشان می‌دهد، زیرا دارای سؤال‌هایی با ضریب تشخیص بالاتری هستند؛ بنابراین، برآورد توانایی را با دقت بیشتری برآورد می‌کند. همچنین، خزانه‌های سؤال بهینه دارای نرخ همپوشی تست پایین‌تری هستند، با وجود این که دارای سؤال‌ها کمتری می‌باشند. نمودار ۱۰ نشان می‌دهد که نرخ همپوشی تست در تمام سطوح توانایی کمتر از خزانه‌های عملیاتی است. این نرخ در سرتاسر پارامتر توانایی به صورت یکنواخت توزیع شده است. در این مرحله نیز سه خزانه بهینه شباهت زیادی در همپوشی تست‌ها دارند. همچنین، خزانه‌های بهینه درصد خیلی کوچک‌تری از بیش مواجهه شدن سؤال‌ها را نسبت به خزانه عملیاتی دارند. نرخ سؤال‌ها بیش مواجهه شده در خزانه‌های MTI کمتر از خزانه‌های دیگر است. نمودارهای ۱۱ تا ۱۳ درصد‌های مواجهه سؤال در هر یک از سطوح توانایی را نشان می‌دهد. الگوی این نمودارها نشان می‌دهد که توزیع مواجهه سؤال‌ها در این مرحله یکنواخت‌تر از مرحله‌ی قبل است. جدول ۱۲ نرخ‌های مواجهه سؤال‌ها را به تفکیک محتواها در هر چهار خزانه سؤال نشان می‌دهد. نتایج این جدول گویای این مطلب است که در محتوای اول، دوم و سوم خزانه‌های R دارای بیشترین نرخ مواجهه هستند و با این که دارای بیشترین تعداد سؤال است، دارای نرخ بالاتری از مواجهه است. در این مرحله نیز خزانه‌های MTI دارای کمترین نرخ مواجهه هستند. البته نتایج خزانه‌های MTI و MRP در نرخ مواجهه سؤال، بسیار مشابه است. این نتایج گویای این قضیه است که با وجود این که خزانه‌های MTI دارای تعداد سؤال کمتری هستند، نرخ‌های مواجهه کمتری دارند. نرخ سؤال‌ها کم مواجهه شده در این نوع خزانه‌ها نیز کمتر از بقیه است، همه این نتایج گویای بهینه بودن این نوع خزانه‌ها از نظر استفاده بهینه از سؤال‌ها است. همچنان که در نمودار ۱۴ ملاحظه می‌شود، میانگین آگاهی خزانه‌های سؤال در خزانه‌های بهینه و خزانه سؤال عملیاتی در تمام سطوح کاملاً مشابه

است. مطابق نتایج جداول ۱۳ و ۱۴ خزانه عملیاتی در تمام سطوح درصد فراوانی نسبی تخطی از قیود بیشتری نسبت به خزانه‌های بهینه دارد. در خزانه‌های R میزان تخطی‌ها در دو دامنه‌ی توانایی بیشتر است و تقریباً الگویی مشابه با خزانه عملیاتی دارد. در خزانه MRP و MTI در سطوح پایین توانایی میزان تخطی از قیود بیشتر است. در خزانه MRP در سطوح بالای توانایی به دلیل وجود سؤال‌ها بیشتر، تخطی از قیود در تست‌هایی که سرهم می‌شود، به حداقل خود می‌رسد. در هر سه خزانه سؤال بهینه، میزان تخطی از قیود محتوایی تست‌ها بیشتر از زمانی است که کنترل مواجهه سؤال وارد نشده است، این نتیجه دلیلی بر این امر است که وارد کردن کنترل مواجهه S-H بر انتخاب سؤال‌ها تأثیر می‌گذارد و این امکان وجود دارد که برنامه CAT، سؤالی را برای اجرا انتخاب کند که کاملاً با قیود محتوایی هماهنگ نباشد و در عمل میزان این تخطی‌ها را بیشتر کند. نمودارهای ۱۵، ۱۶ و ۱۷ به ترتیب خطای استاندارد شرطی اندازه‌گیری (CSEM)، اریب شرطی و میانگین مجذور خطا (CMSE) را در هر چهار خزانه سؤال نشان می‌دهند. خطاهای استاندارد اندازه‌گیری در سطوح توانایی زیر صفر الگویی مشابه با خزانه‌های عملیاتی دارند، با این وجود مقادیر آن کمتر از خزانه عملیاتی است. ولی در سطوح توانایی بالاتر از صفر، به‌خصوص در خزانه‌های MRP و MTI این خطا به صفر نزدیک می‌شود و دارای الگوی متفاوتی می‌شود. نمودار ۱۶ نشان می‌دهد که در خزانه عملیاتی میزان اریب در اغلب سطوح توانایی بیشتر از خزانه‌های بهینه است. خزانه MRP نسبت به خزانه‌های دیگر از میزان اریب کمتری برخوردار است. همچنین، نمودار ۱۷ میانگین مجذور خطا را در سطوح متفاوت توانایی نشان می‌دهد. نتایج این نمودار نشان می‌دهد که MSE هر سه خزانه بهینه کوچک‌تر از خزانه سؤال عملیاتی است. به‌خصوص خزانه MRP از حداقل مقدار MSE برخوردار است.



جدول ۱۰. اندازه‌های خزانه سؤال و آماره‌های پارامتر سؤال، با S-H (b-bin=0.2). با تعادل محتوا

نوع سؤال	اندازه خزانه	a				b				c					
		میانگین	استاندارد	انحراف	حداکثر	حداقل	میانگین	استاندارد	انحراف	حداکثر	حداقل	میانگین	استاندارد	انحراف	حداکثر
Content 1 (arithmetic)															
OP	۴۵۵	۱/۰۸۹	۰/۲۸۴۴	۳/۰۴۵	۰/۱۶۶	۳/۰۴۵	۰/۲۸۴۴	۳/۰۴۵	۰/۱۶۶	۳/۰۴۵	۰/۲۸۴۴	۳/۰۴۵	۰/۱۶۶	۳/۰۴۵	۰/۲۸۴۴
ROP_4	۲۶۱	۱/۰۸۹	۰/۳۰۱	۳/۰۰۴	۰/۹۰۲	۳/۰۰۴	۰/۳۰۱	۰/۹۰۲	۰/۳۴۶	۰/۹۰۲	۰/۳۰۱	۰/۳۴۶	۰/۹۰۲	۰/۳۴۶	۰/۳۰۱
ROP_5	۲۴۹	۲/۲۸۹	۰/۳۰۱۲	۳/۱۸۶	۰/۹۳۸	۳/۱۸۶	۰/۳۰۱۲	۰/۹۳۸	۰/۵۴۹	۰/۹۳۸	۰/۳۰۱۲	۰/۵۴۹	۰/۹۳۸	۰/۳۰۱۲	۰/۵۴۹
ROP_6	۲۱۴	۱/۶۱۳	۰/۱۶۶۷	۲/۵۲۸	۱/۱۸۹	۲/۵۲۸	۰/۱۶۶۷	۱/۱۸۹	۰/۲۱۸۰	۱/۱۸۹	۰/۱۶۶۷	۰/۲۱۸۰	۱/۱۸۹	۰/۱۶۶۷	۰/۲۱۸۰
Content 2 (geometry)															
OP	۲۵۸	۱/۲۰۶	۰/۲۲۴۵	۲/۹۳	۰/۲۴۵	۲/۹۳	۰/۲۲۴۵	۰/۲۴۵	۰/۴۸۲	۰/۲۴۵	۰/۴۸۲	۰/۴۸۲	۰/۲۴۵	۰/۴۸۲	۰/۲۴۵
ROP_4	۳۳۹	۱/۸۹۳	۰/۲۹۴	۲/۸۵۴	۰/۹۰۵	۲/۸۵۴	۰/۲۹۴	۰/۹۰۵	۰/۹۶۸	۰/۹۰۵	۰/۹۶۸	۰/۹۰۵	۰/۹۶۸	۰/۹۰۵	۰/۹۶۸
ROP_5	۳۳۶	۲/۰۹۳	۰/۳۰۴	۳/۱۸۳	۱/۰۸۶	۳/۱۸۳	۰/۳۰۴	۱/۰۸۶	۰/۲۱۵	۱/۰۸۶	۰/۳۰۴	۰/۲۱۵	۱/۰۸۶	۰/۳۰۴	۰/۲۱۵
ROP_6	۲۰۶	۱/۵۴۲	۰/۲۸۶	۲/۵۰۲	۱/۰۸۷	۲/۵۰۲	۰/۲۸۶	۱/۰۸۷	۰/۳۶۸	۱/۰۸۷	۰/۳۶۸	۰/۳۶۸	۱/۰۸۷	۰/۳۶۸	۰/۳۶۸
Content 3 (algebra)															
OP	۲۰۸	۱/۳۵۶	۰/۲۴۷	۲/۸۸۹	۰/۵۳۸	۲/۸۸۹	۰/۲۴۷	۰/۵۳۸	۰/۳۲۴	۰/۵۳۸	۰/۳۲۴	۰/۳۲۴	۰/۵۳۸	۰/۳۲۴	۰/۳۲۴
ROP_4	۱۹۹	۱/۷۶۹	۰/۲۹۸	۳/۰۰۶	۰/۹۱۸	۳/۰۰۶	۰/۲۹۸	۰/۹۱۸	۰/۴۲۸	۰/۹۱۸	۰/۴۲۸	۰/۴۲۸	۰/۹۱۸	۰/۴۲۸	۰/۴۲۸
ROP_5	۱۹۲	۲/۰۲۳	۰/۲۹۹	۲/۷۸۶	۰/۹۴۸	۲/۷۸۶	۰/۲۹۹	۰/۹۴۸	۰/۴۶۱	۰/۹۴۸	۰/۴۶۱	۰/۴۶۱	۰/۹۴۸	۰/۴۶۱	۰/۴۶۱
ROP_6	۱۸۸	۱/۷۲۹	۰/۲۷۴	۲/۵۲۳	۱/۱۰۴	۲/۵۲۳	۰/۲۷۴	۱/۱۰۴	۰/۴۹۶	۱/۱۰۴	۰/۴۹۶	۰/۴۹۶	۱/۱۰۴	۰/۴۹۶	۰/۴۹۶

جدول ۱۱. خلاصه‌ی آماره‌های عملکرد خزانه سؤال با S-H (b-bin=0.2)، با تعادل محتوا

آماره‌ها	OP	R	MRP	MTI
Bias	۰/۰۰۵۲	۰/۰۰۱۲	۰/۰۰۰۱۱	۰/۰۰۰۷۲
MSE	۰/۰۱۷۴۵	۰/۰۱۲۶۵	۰/۰۰۷۳	۰/۰۱۲۲۴
کجی نرخ مواجهه	۹۸/۸۵۶	۳۰/۲۵۲	۲۶/۷۴۱	۲۶/۹۹۴
نرخ همپوشی سؤال	۰/۵۱۷۳	۰/۲۲۴۳	۰/۲۳۶۷	۰/۲۵۰۱
درصد سؤال‌های با نرخ مواجهه بزرگ‌تر از $\frac{1}{3}$	۸/۱۶۹٪	۵/۱۵۳٪	۴/۶۳۵٪	۳/۸۶۲٪
درصد سؤال‌های با نرخ مواجهه کوچک‌تر از 0.02	۵۱/۵۲۶٪	۲۳/۰۸۶٪	۲۴/۴۲٪	۱۷/۳۲٪
درصد تست‌هایی که از قیود تست تخطی دارند	۵۴/۸٪	۳/۴٪	۲/۵٪	۲/۶٪
اندازه خزانه سؤال	۹۲۱	۶۹۹	۶۷۷	۶۰۸

جدول ۱۲. درصد سؤال‌ها بیش مواجهه شده و کم مواجهه شده در هر محتوا

برای خزانه‌های سؤال با S-H

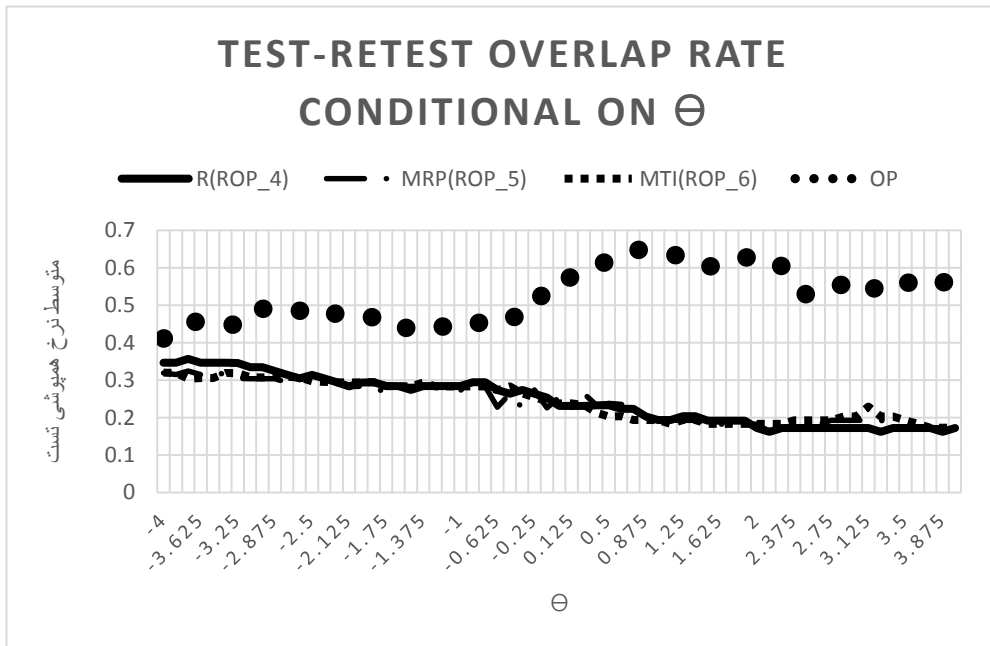
آماره	محتوای ۱				محتوای ۲				محتوای ۳			
خزانه‌های سؤال	MTI	MRP	R	OP	MTI	MRP	R	OP	MTI	MRP	R	OP
درصد سؤال‌های با نرخ مواجهه بزرگ‌تر از $\frac{1}{3}$	۴/۱۸۹٪	۴/۸۷٪	۵/۶۴٪	۹/۸۱٪	۳/۷۵٪	۴/۹۲٪	۵/۵۲٪	۷/۵۸٪	۳/۸۲٪	۴/۸۶٪	۵/۲۶٪	۷/۸۱٪
اندازه خزانه	۲۱۴	۲۴۹	۲۶۱	۴۵۵	۲۵۸	۳۳۹	۳۳۹	۲۵۸	۲۰۶	۲۰۶	۲۰۸	۱۸۸

جدول ۱۳. تعداد تست‌ها بر اساس تخطی از قیود محتوایی در CAT‌ها، بدون عامل کنترل مواجهه S-H

خزانه سؤال	OP	R (ROP_4)	MRP (ROP_5)	MTI (ROP_6)	کل
CAT total	۳۵۰	۶۰۰	۶۰۰	۶۰۰	۱۹۳
سطوح توانایی	OP	R	MRP	MTI	کل
۱	۱	۳	۵	۶	-۰/۵
۲	۲	۲	۴	۱۲	-۱
۳	۳	۳	۴	۱۲	-۱/۵
۴	۴	۴	۱۱	۱۲	-۲
۵	۵	۵	۱۶	۱۵	-۲/۵
۶	۶	۶	۲۴	۱۹	-۳
۷	۷	۷	۱۶	۲۴	-۳/۵
۸	۸	۸	۲۲	۲۷	-۴
۹	۹	۹	۲۹	۲۷	-۴
۱۰	۱۰	۱۰	۲۹	۲۷	-۴
۱۱	۱۱	۱۱	۲۹	۲۷	-۴
۱۲	۱۲	۱۲	۲۹	۲۷	-۴
۱۳	۱۳	۱۳	۲۹	۲۷	-۴
۱۴	۱۴	۱۴	۲۹	۲۷	-۴
۱۵	۱۵	۱۵	۲۹	۲۷	-۴
۱۶	۱۶	۱۶	۲۹	۲۷	-۴
۱۷	۱۷	۱۷	۲۹	۲۷	-۴
۱۸	۱۸	۱۸	۲۹	۲۷	-۴
۱۹	۱۹	۱۹	۲۹	۲۷	-۴
۲۰	۲۰	۲۰	۲۹	۲۷	-۴
۲۱	۲۱	۲۱	۲۹	۲۷	-۴
۲۲	۲۲	۲۲	۲۹	۲۷	-۴
۲۳	۲۳	۲۳	۲۹	۲۷	-۴
۲۴	۲۴	۲۴	۲۹	۲۷	-۴
۲۵	۲۵	۲۵	۲۹	۲۷	-۴
۲۶	۲۶	۲۶	۲۹	۲۷	-۴
۲۷	۲۷	۲۷	۲۹	۲۷	-۴
۲۸	۲۸	۲۸	۲۹	۲۷	-۴
۲۹	۲۹	۲۹	۲۹	۲۷	-۴
۳۰	۳۰	۳۰	۲۹	۲۷	-۴
۳۱	۳۱	۳۱	۲۹	۲۷	-۴
۳۲	۳۲	۳۲	۲۹	۲۷	-۴
۳۳	۳۳	۳۳	۲۹	۲۷	-۴
۳۴	۳۴	۳۴	۲۹	۲۷	-۴
۳۵	۳۵	۳۵	۲۹	۲۷	-۴
۳۶	۳۶	۳۶	۲۹	۲۷	-۴
۳۷	۳۷	۳۷	۲۹	۲۷	-۴
۳۸	۳۸	۳۸	۲۹	۲۷	-۴
۳۹	۳۹	۳۹	۲۹	۲۷	-۴
۴۰	۴۰	۴۰	۲۹	۲۷	-۴
۴۱	۴۱	۴۱	۲۹	۲۷	-۴
۴۲	۴۲	۴۲	۲۹	۲۷	-۴
۴۳	۴۳	۴۳	۲۹	۲۷	-۴
۴۴	۴۴	۴۴	۲۹	۲۷	-۴
۴۵	۴۵	۴۵	۲۹	۲۷	-۴
۴۶	۴۶	۴۶	۲۹	۲۷	-۴
۴۷	۴۷	۴۷	۲۹	۲۷	-۴
۴۸	۴۸	۴۸	۲۹	۲۷	-۴
۴۹	۴۹	۴۹	۲۹	۲۷	-۴
۵۰	۵۰	۵۰	۲۹	۲۷	-۴
۵۱	۵۱	۵۱	۲۹	۲۷	-۴
۵۲	۵۲	۵۲	۲۹	۲۷	-۴
۵۳	۵۳	۵۳	۲۹	۲۷	-۴
۵۴	۵۴	۵۴	۲۹	۲۷	-۴
۵۵	۵۵	۵۵	۲۹	۲۷	-۴
۵۶	۵۶	۵۶	۲۹	۲۷	-۴
۵۷	۵۷	۵۷	۲۹	۲۷	-۴
۵۸	۵۸	۵۸	۲۹	۲۷	-۴
۵۹	۵۹	۵۹	۲۹	۲۷	-۴
۶۰	۶۰	۶۰	۲۹	۲۷	-۴
۶۱	۶۱	۶۱	۲۹	۲۷	-۴
۶۲	۶۲	۶۲	۲۹	۲۷	-۴
۶۳	۶۳	۶۳	۲۹	۲۷	-۴
۶۴	۶۴	۶۴	۲۹	۲۷	-۴
۶۵	۶۵	۶۵	۲۹	۲۷	-۴
۶۶	۶۶	۶۶	۲۹	۲۷	-۴
۶۷	۶۷	۶۷	۲۹	۲۷	-۴
۶۸	۶۸	۶۸	۲۹	۲۷	-۴
۶۹	۶۹	۶۹	۲۹	۲۷	-۴
۷۰	۷۰	۷۰	۲۹	۲۷	-۴
۷۱	۷۱	۷۱	۲۹	۲۷	-۴
۷۲	۷۲	۷۲	۲۹	۲۷	-۴
۷۳	۷۳	۷۳	۲۹	۲۷	-۴
۷۴	۷۴	۷۴	۲۹	۲۷	-۴
۷۵	۷۵	۷۵	۲۹	۲۷	-۴
۷۶	۷۶	۷۶	۲۹	۲۷	-۴
۷۷	۷۷	۷۷	۲۹	۲۷	-۴
۷۸	۷۸	۷۸	۲۹	۲۷	-۴
۷۹	۷۹	۷۹	۲۹	۲۷	-۴
۸۰	۸۰	۸۰	۲۹	۲۷	-۴
۸۱	۸۱	۸۱	۲۹	۲۷	-۴
۸۲	۸۲	۸۲	۲۹	۲۷	-۴
۸۳	۸۳	۸۳	۲۹	۲۷	-۴
۸۴	۸۴	۸۴	۲۹	۲۷	-۴
۸۵	۸۵	۸۵	۲۹	۲۷	-۴
۸۶	۸۶	۸۶	۲۹	۲۷	-۴
۸۷	۸۷	۸۷	۲۹	۲۷	-۴
۸۸	۸۸	۸۸	۲۹	۲۷	-۴
۸۹	۸۹	۸۹	۲۹	۲۷	-۴
۹۰	۹۰	۹۰	۲۹	۲۷	-۴
۹۱	۹۱	۹۱	۲۹	۲۷	-۴
۹۲	۹۲	۹۲	۲۹	۲۷	-۴
۹۳	۹۳	۹۳	۲۹	۲۷	-۴
۹۴	۹۴	۹۴	۲۹	۲۷	-۴
۹۵	۹۵	۹۵	۲۹	۲۷	-۴
۹۶	۹۶	۹۶	۲۹	۲۷	-۴
۹۷	۹۷	۹۷	۲۹	۲۷	-۴
۹۸	۹۸	۹۸	۲۹	۲۷	-۴
۹۹	۹۹	۹۹	۲۹	۲۷	-۴
۱۰۰	۱۰۰	۱۰۰	۲۹	۲۷	-۴

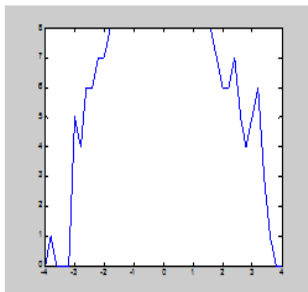
جدول ۱۴. درصد تست‌ها بر اساس تخطی از قیود محتوایی در CAT ها، بدون عامل کنترل مواجهه S-H

کل	۴	۳/۵	۳	۲/۵	۲	۱/۵	۱	۰/۵	۰	-۰/۵	-۱	-۱/۵	-۲	-۲/۵	-۳	-۳/۵	-۴	CAT total	سطوح توانایی جزانه سوال
۱۵۴۸	۰/۰۸	۰/۰۷	۰/۰۷	۰/۰۳	۱/۰۰۴	۰/۰۲	۰/۰۱	۱/۰۰۴	۱/۰۰۱	۱/۰۰۱	۱/۰۰۸	۰/۰۴	۰/۰۴	۰/۰۴	۰/۰۵	۰/۰۴	۰/۰۳	۲۵۰	OP
۰/۰۲۴	۰/۰۰۵	۰/۰۰۳	۰/۰۰۲	۰/۰۰۱۵	۰/۰۰۱۵	۰/۰۰۱	۱/۰۰۲۳	۰/۰۰۰۳	۰/۰۰۰۵	۰/۰۰۱	۰/۰۰۰۵	۰/۰۰۰۵	۰/۰۰۱۳	۰/۰۰۲	۰/۰۰۲۸	۰/۰۰۳۶	۰/۰۰۵	۶۰۰۰	R (ROP_4)
۰/۰۲۵	۰	۰/۰۰۰۳	۰/۰۰۰۳	۰	۰/۰۰۰۵	۰/۰۰۰۶	۰/۰۰۱۵	۰	۰	۰/۰۰۰۸	۰/۰۰۰۶	۰/۰۰۱۸	۰/۰۰۲۶	۰/۰۰۴	۰/۰۰۲۶	۰/۰۰۴۱	۰/۰۰۵	۶۰۰۰	MRP (ROP_5)
۰/۰۲۶	۰/۰۰۱۵	۰/۰۰۱۱	۰/۰۰۰۸	۰	۰/۰۰۰۵	۰/۰۰۰۶	۰/۰۰۰۳	۰/۰۰۰۸	۰	۰/۰۰۰۱	۰/۰۰۰۲	۰/۰۰۰۱	۰/۰۰۰۲	۰/۰۰۰۲۵	۰/۰۰۰۳۱	۰/۰۰۰۴	۰/۰۰۰۴۵	۶۰۰۰	MTI (ROP_6)

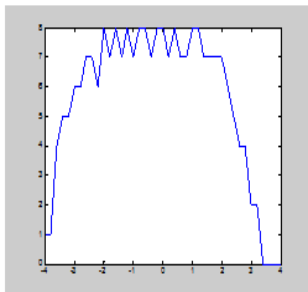


نمودار ۱۰. نرخ همپوشانی تست مشروط به  $\square$  با کنترل S-H ( $b$ -bin: 0.2)، با تعادل محتوا

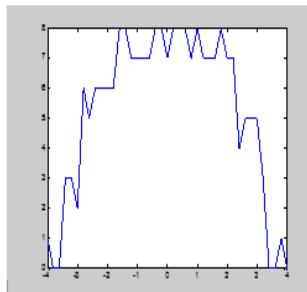
R (ROP\_4) (content 1)



R (ROP\_4) (content 2)

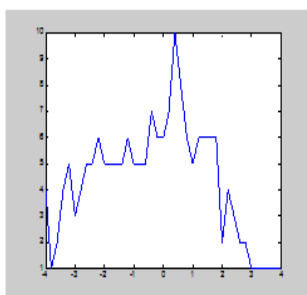


R (ROP\_4) (content 3)

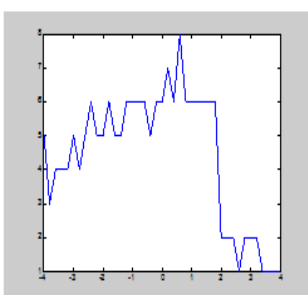


نمودار ۱۱. نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده بر اساس روش R

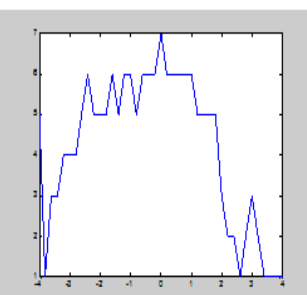
MRP (ROP\_5) (content 1)



MRP (ROP\_5) (content 2)

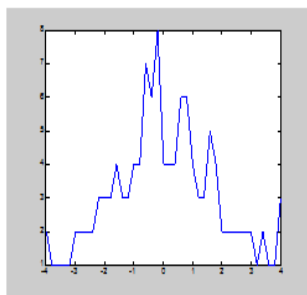


MRP (ROP\_5) (content 3)

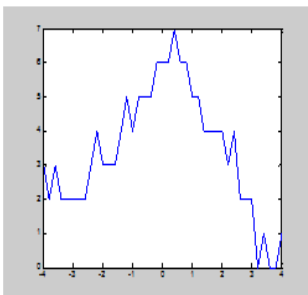


نمودار ۱۲: نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده بر اساس روش MRP

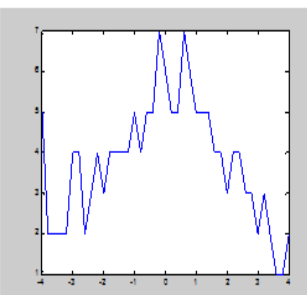
MRP (ROP\_6) (content 1)



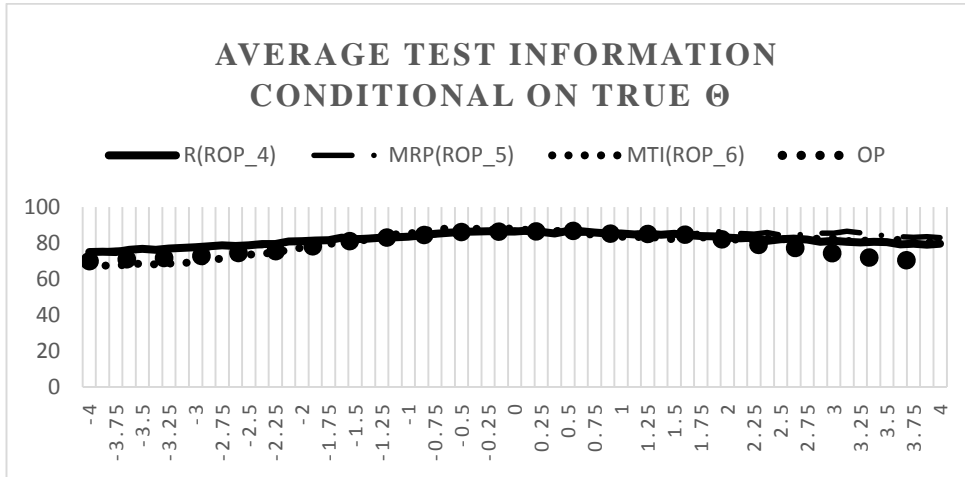
MRP (ROP\_6) (content 2)



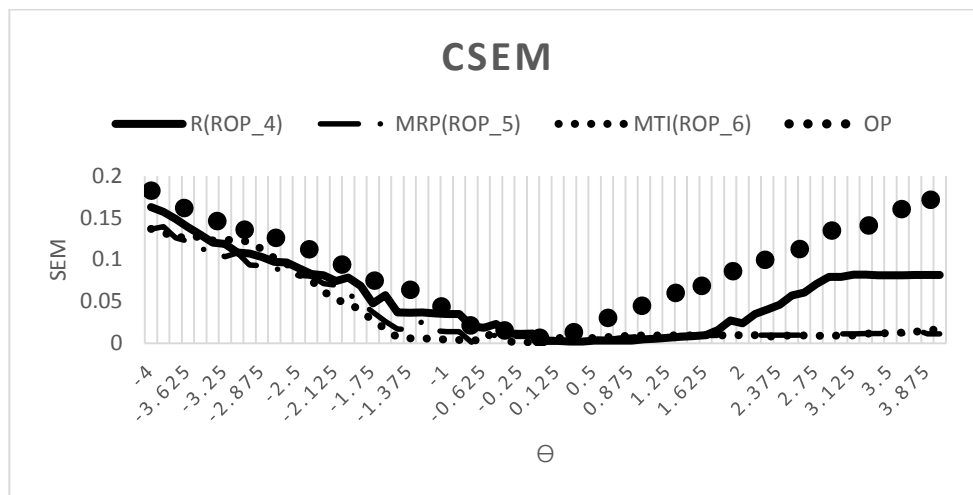
MRP (ROP\_6) (content 3)



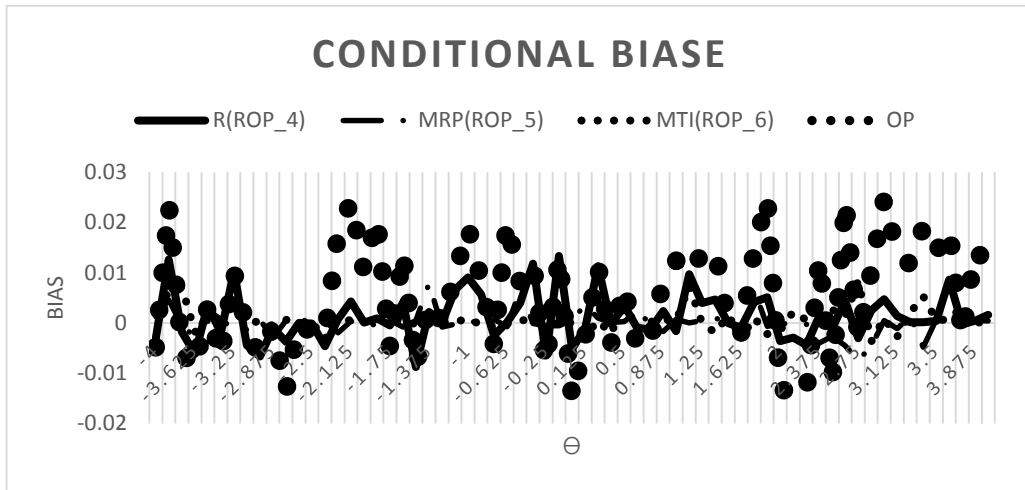
نمودار ۱۳. نمودارهای مربوط به درصد سؤال‌ها بیش مواجهه شده بر اساس روش MTI



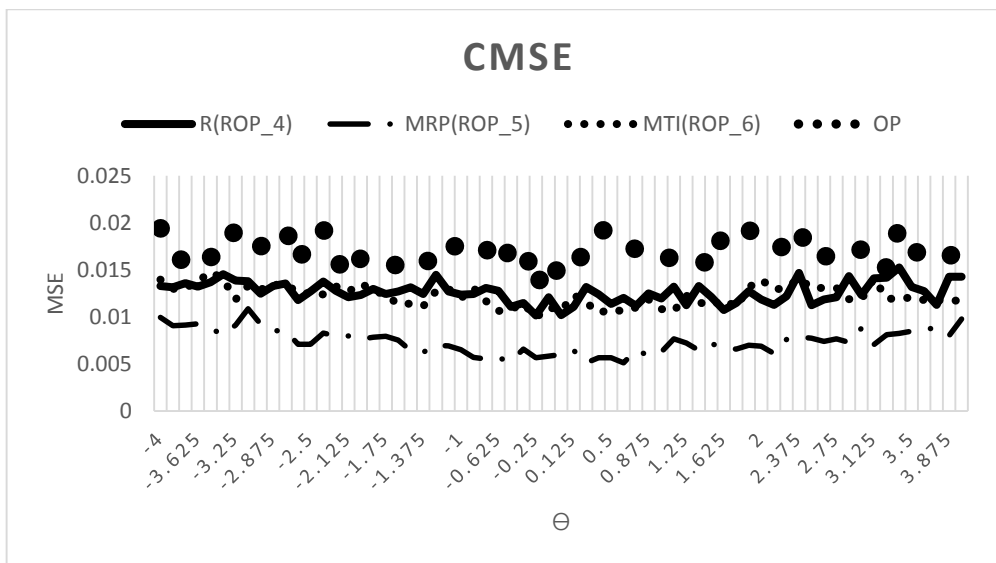
نمودار ۱۴. متوسط آگاهی تست مشروط به توانایی واقعی در خزانه‌های سؤال  
با S-H ( $b$ -bin: 0.2)، با تعادل محتوا



نمودار ۱۵. خطای استاندارد اندازه‌گیری (CSEM) در خزانه‌های سؤال  
با S-H ( $b$ -bin: 0.2)، با تعادل محتوا



نمودار ۱۶: اربیب شرطی (conditional Bias) در خزانه‌های سؤال  
با S-H ( $b\text{-bin}: 0.2$ )، با تعادل محتوا



نمودار ۱۷: میانگین مجذور خطا (CMSE) در خزانه‌های سؤال  
با S-H ( $b\text{-bin}: 0.2$ )، با تعادل محتوا

## نتیجه‌گیری

هدف این پژوهش، ساخت خزانه‌های سؤال بهینه بر اساس رویکرد اکتشافی ریکیسی، با در نظر گرفتن تعادل محتوایی و نرخ مواجهه سؤال‌ها در آزمون بوده است. خزانه‌هایی که با توجه به مدل سه پارامتری لجستیک مدرج شده‌اند. یکی از پیشنهاد‌های پژوهشی ریکیسی در تحقیقات خود این بود که؛ تعادل محتوایی متنوع یکی از مهم‌ترین مؤلفه‌ها در سنجش انطباقی است که می‌توان آن را در الگوریتم‌های انتخاب سؤال CAT گنجانده (هی و ریکیسی، ۲۰۱۰؛ ۲۰۱۱؛ گو و ریکیسی، ۲۰۰۷)، بنابراین در این پژوهش این عامل وارد برنامه‌ی شبیه‌سازی شد، البته لازم به ذکر است که در تحقیقات آن‌ها این عامل در نظر گرفته می‌شد، ولی از طریق پارتیشن‌بندی کردن خزانه بهینه به خزانه‌های کوچک‌تر این کار صورت می‌گرفت و هیچ‌کدام از این تحقیقات این عامل را به صورت کمی وارد تحلیل نکردند. در این پژوهش تعادل محتوایی سؤال‌ها از طریق مدل حداقل انحرافات وزن‌دار ایجاد شد. به طوری که از پیش، کدهای محتوایی سؤال‌ها توسط متخصصین موضوعی، مشخص و به همراه ویژگی‌های آماری وارد برنامه شد. در این پژوهش با دست‌کاری دو عامل: روش ایجاد سؤال بهینه (R, MRP, MTI)، کنترل یا عدم کنترل مواجهه بیش‌ازحد سؤال با روش سیمپسون-هتر (S-H)، ۶ مدل طراحی خزانه سؤال بهینه (ROP\_1, ROP\_2, ROP\_3, ..., ROP\_6) ایجاد شد. همه خزانه‌های سؤال بهینه‌ای که در این پژوهش طراحی شد، صرف‌نظر از عواملی چون کنترل مواجهه و روش ایجاد سؤال بهینه، عملکرد بهتری نسبت به خزانه‌های عملیاتی داشتند. دلیل این امر این است که در مجموع، الگوهای خزانه سؤال بهینه در جستجوی مطلوب‌ترین و مناسب‌ترین ترکیب سؤال‌ها برای تشکیل یک خزانه سؤالی هستند که از طریق آن بتوان تعداد زیادی از تست‌های انطباقی را سرهم کرد. با این وجود، در دنیای واقعی خزانه سؤالی وجود ندارد که به طور مطلق بهینه باشد، زیرا به تعداد عوامل و ترکیب‌های متفاوتی از سؤال‌های موجود در خزانه محدود می‌شود. این دلایل باعث می‌شود که هر یک از این نوع خزانه‌ها دارای صحت و دقت اندازه‌گیری متفاوتی باشند و هر یک از لحاظ بهینه بودن کاملاً از یکدیگر

متفاوت باشند. اما، در کل، هدف کلی برای الگوهای خزانه سؤال بهینه این است که دارای سه ملاک مهم باشند که توسط وندرلیندن (۲۰۰۰a) ارائه شده است. اولین ملاک این است که به اندازه کافی برای سرهم کردن چندین هزار خرده آزمون همپوش، بزرگ باشد. استوکینگ (۱۹۹۴)، مباحث متنوعی در مورد اندازه خزانه سؤال در مورد آزمون‌های ورودی سرنوشت‌ساز که به شکل CAT اجرا می‌شود، مطرح کرد. این قاعده بیان می‌کرد که یک خزانه سؤال CAT برای آزمون‌های سرنوشت‌ساز، ۶ تا ۱۲ برابر طول آزمون CAT باشد. در این پژوهش، اندازه خزانه‌های سؤال بهینه از ۴۹۶ تا ۶۹۹ سؤال در نوسان بود، هیچ‌کدام از خزانه‌ها بیشتر از ۷۲۰ ( $12 * 60 = 720$ ) سؤال نداشتند؛ بنابراین، به عبارت دیگر، روش bin-and-union به اندازه زیادی با توصیه استوکینگ برای ساخت یک خزانه سؤال با اندازه کافی برای یک برنامه CAT سازگار است. این نتایج نشان می‌دهد که ملاک اولی که وندرلیندن در مورد خزانه‌های سؤال بهینه مطرح کرد، در مورد خزانه‌های سؤال بهینه در پژوهش حاضر برقرار است. در مورد ملاک دوم که در مورد حدود دامنه سطوح دشواری سؤال‌های خزانه است، باید ذکر کرد که همه خزانه‌های سؤال بهینه طراحی شده در این پژوهش در دامنه وسیعی از سطوح دشواری سؤال پراکنده شده‌اند. به طوری که توزیع پارامترهای b در تمام ROP ها به صورت یکنواختی در سراسر مقیاس توانایی توزیع شده است، این توزیع برخلاف توزیع نرمال زیربنایی توانایی آزمودنی‌هایی است که آزمون برای آن‌ها شبیه‌سازی شد. به عبارت دیگر، پارامترهای b به طور یکنواختی در طول مقیاس مهارت توزیع شدند، به طوری که با توزیع توانایی استفاده شده در شبیه‌سازی مطابقت ندارند این ویژگی نیز در تحقیقات هی و ریکیسی (۲۰۱۰) مورد تأیید قرار گرفت، زیرا ماهیت CAT ایجاب می‌کند که توزیع پارامتر b در طول پیوستار توانایی یکنواخت باشد. در کل، ویژگی‌های پارامتر b سؤال پیشنهاد شده در این مطالعه با نتایج یوری (۱۹۷۷) و جنسما (۱۹۷۷) نیز سازگار است. ملاک سوم، هزینه ساخت و طراحی سؤال‌ها است که در این مطالعه به طور مستقیمی برآورد نشده است. با این وجود، یک مکانیسم مکنون برای برآورد غیرمستقیم هزینه نوشتن سؤال‌ها و تلاش برای به حداقل رساندن آن، در روش طراحی خزانه سؤال ریکیسی مورد استفاده قرار گرفته



است. ایجاد سؤال‌ها بهینه با روش‌های  $R$ ،  $MRP$  و  $MTI$  به رابطه بین پارامترهای سؤال و هزینه طراحی سؤال، به‌طور همزمان توجه می‌کند. البته تفاوت این سه شیوه ایجاد پارامترهای سؤال در نگاه اول در مفروضات آن‌ها است. روش  $MTI$  بر این فرض استوار است که سؤال‌هایی با ضرایب تشخیص بالا (با مقدار  $a$  بالا) هزینه طراحی بالایی دارند؛ بنابراین، روش  $MTI$  تلاش می‌کند تا تعداد سؤال‌هایی با ضرایب تشخیص بالا را از طریق شبیه‌سازی کردن سؤال‌هایی که شرط حداقل آگاهی تست را داشته باشند، محدود کند (گو و ریکسی، ۲۰۰۷). در دو روش دیگر، یعنی روش  $R$  و قسمت  $R$  روش  $MRP$ ، بر تصادفی‌سازی رابطه بین پارامترها توجه دارند؛ بنابراین از لحاظ هزینه طراحی به‌صرفه نیستند، چون به تعداد بالایی سؤال نیاز خواهند داشت و سؤال‌هایی با توزیع یکنواخت پارامتر دشواری و تشخیص ایجاد می‌کند. اما روش  $P$  یا پیش‌بینی (در روش  $MRP$ )، فرض می‌کند که طراحی سؤال‌هایی با مشخصات معین از قبل تعیین شده و مورد نیاز، کم‌هزینه‌تر از طراحی سؤال‌ها بدون در نظر گرفتن این ویژگی‌ها است. این روش هزینه ایجاد سؤال را از طریق مدل‌یابی کردن مشخصات سؤال (مانند، رابطه بین پارامترهای  $IRT$  سؤال‌ها) به حداقل می‌رساند و سؤال‌هایی شبیه‌سازی می‌کند که مشابه با سؤال‌های موجود در خزانه عملیاتی است. همچنین، بهتر است که خزانه‌های سؤال برای آزمون‌های  $CAT$ ، دقت اندازه‌گیری مشابهی در سراسر دامنه توانایی ایجاد کنند. بررسی دقیق‌تر به خزانه عملیاتی نشان می‌دهد که سؤال‌هایی با ضرایب تشخیص بالاتر در دامنه پارامتر  $b$  بین  $1/72$  تا  $4$  قرار می‌گیرند. در عمل، زمانی که از خزانه عملیاتی به‌طور فراوانی استفاده می‌شود، سؤال‌هایی با ضرایب تشخیص بالا که به احتمال زیادی توسط الگوریتم‌های  $CAT$  انتخاب می‌شوند، کنار گذاشته می‌شوند، طراحی سؤال‌های با چنین ضرایب تشخیص بالایی برای جایگزینی با سؤال‌ها قبلی بسیار دشوار است؛ بنابراین، این نتیجه ممکن است این شک را ایجاد کند که آیا در این مرحله عملکرد یکسانی بر روی سطوح توانایی یکسان با مراحل قبل، می‌تواند به‌آسانی تکرار شود. با این وجود، خزانه‌های بهینه‌ای که از طریق روش ریکسی طراحی می‌شوند، دارای سؤال‌ها بیشتری با ویژگی‌های یکنواخت در طول پیوستار توانایی هستند، در نتیجه این نوع خزانه‌ها عملکرد بهتری از نظر دقت و صحت طبقه‌بندی، در

حداکثر سطوح توانایی مکنون ایجاد می‌کنند. در مجموع، خزانه‌های سؤالی که به‌طور بهینه طراحی می‌شوند، عملکرد بهتری نسبت به خزانه‌های سؤال عملیاتی در شاخص‌های ارزیابی خزانه‌های سؤال ایجاد می‌کنند. این نتایج، مشابه با تحقیقات یوری (۱۹۷۷)، گو و ریکیسی (۲۰۰۷)، هی و ریکیسی (۲۰۱۰) بوده است، در این پژوهش‌ها نیز این ویژگی‌ها به‌عنوان ویژگی‌های برتر خزانه‌های سؤال CAT توصیه شدند. پارامترهای a، در همه ROP ها در حدود ۱/۱ با حداقل مقداری بیشتر از ۸/۸. به اوج خود می‌رسند. نتایج حاصل از خروجی این شبیه‌سازی‌ها حاکی از آن است که تلفیق روش مونت کارلو ریکیسی و روش برنامه‌نویسی حداقل انحرافات وزن‌دار (WDM) امکان‌پذیر است. همچنین، با توجه به این که ملاک‌های وندرلیندن (۲۰۰۰a) در مورد خزانه‌های بهینه برای CAT در این پژوهش وجود داشته است، این شبیه‌سازی با موفقیت انجام گرفته است.

همچنین، نتایج نشان می‌دهد که بدون توجه به عامل کنترل مواجهه S-H، خزانه‌های سؤال بهینه‌ای که با کنترل تعادل محتوایی طراحی می‌شوند، دارای دقت اندازه‌گیری خوبی نسبت به خزانه عملیاتی هستند. زمانی که عامل مواجهه سؤال نیز کنترل می‌شود، این میزان دقت اندازه‌گیری بیشتر می‌شود. تعامل دو عامل تعادل محتوایی و کنترل مواجهه S-H نتایج جالبی در مورد دقت اندازه‌گیری به‌دست آورد: میزان اریب در هر چهار خزانه مثبت و بسیار کوچک است. زمانی که مواجهه سؤال کنترل می‌شود، این میزان اریب کمتر از زمانی است که کنترل نمی‌شود. زمانی که عامل S-H در شبیه‌سازی وارد نشده است، خزانه (ROP\_3) MTI دارای بهترین عملکرد از نظر دقت اندازه‌گیری است. خزانه (ROP\_2) MRP و (ROP\_1) R دارای خطاهای اندازه‌گیری تقریباً مشابهی در زمانی که عامل S-H در شبیه‌سازی نمی‌شود هستند. اما زمانی که عامل S-H وارد می‌شود، خزانه (ROP\_6) MTI دارای دقت اندازه‌گیری کمتری نسبت به خزانه (ROP\_3) MTI است، زیرا سؤال‌هایی که دارای میزان آگاهی متناسب با دامنه توانایی موردنظر و کد محتوایی مناسب هستند، بیشتر از ۰/۳۳ ارائه شده و برای کنترل این قضیه، برنامه مجبور است که از سؤال‌های bin های هم‌جوار استفاده کند که باعث می‌شود دقت اندازه‌گیری کاهش یابد. در مقابل، زمانی که عامل S-H وارد می‌شود خزانه (ROP\_5) MRP دارای بهترین عملکرد از نظر

خطای اندازه‌گیری است. در مجموع، هر سه خزانه بهینه در هر دو مرحله بهتر از خزانه‌های سؤال عملیاتی از نظر اندازه خزانه، دقت اندازه‌گیری و نرخ مواجهه سؤال‌ها عمل می‌کنند. در کل، بررسی دقیق‌تر نمودارهای مربوط به دقت اندازه‌گیری در هر یک از سطوح توانایی نشان می‌دهد که خزانه‌های سؤال که با کنترل مواجهه سؤال طراحی می‌شوند، دارای دقت اندازه‌گیری بیشتری نسبت به خزانه‌هایی که بدون کنترل مواجهه طراحی می‌شوند، هستند. این نتیجه به دلیل این است که سؤال‌های اضافه‌شده به خزانه‌های بهینه با کنترل مواجهه S-H دارای سؤال‌هایی با ضرایب تشخیص بالاتری هستند. در کل، به نظر می‌رسد که بدون در نظر گرفتن عامل S-H، خزانه‌های MTI از سؤال‌های موجود در خزانه استفاده بیشتری می‌کنند و دارای حداقل سؤال‌های کم مواجهه شده است. همچنین، از نرخ همپوشی تست کمی - با وجود اینکه دارای حداقل تعداد سؤال هستند، برخوردار است. در مجموع، بدون توجه به عامل S-H، خزانه‌های بهینه MTI نسبت به خزانه‌های R و MRP دارای سؤال‌ها کمتری هستند و از نرخ مواجهه کمتری نیز برخوردارند و از سؤال‌ها استفاده بیشتری می‌کنند. زمانی که عامل S-H در شبیه‌سازی وارد می‌شود، خزانه (ROP\_5) MRP دارای دقت اندازه‌گیری بالاتری است ولی از نظر اقتصادی به صرفه نیست؛ بنابراین، توصیه می‌شود که زمانی که به صرفه بودن طراحی خزانه‌های سؤال، تعادل محتوایی و امنیت آزمون از لحاظ نرخ مواجهه سؤال، عامل بسیار مهمی می‌باشند، برای کاهش تعداد سؤال‌ها مورد نیاز در خزانه CAT از روش MTI استفاده شود. البته نکته دیگری که باید به آن توجه کرد این است که زمانی که عامل S-H وارد می‌شود، در هر سه خزانه سؤال بهینه، میزان تخطی از قیود محتوایی تست‌ها بیشتر از زمانی است که S-H وارد نشده است، این نتیجه دلیلی بر این امر است که وارد کردن کنترل مواجهه سیمپسون-هتر بر انتخاب سؤال‌ها تأثیر می‌گذارد و این امکان وجود دارد که برنامه‌ی CAT، سؤال‌ها را برای اجرا انتخاب کند که از قیود محتوایی تخطی دارد؛ بنابراین، با توجه به این مباحث، می‌توان در مورد چگونگی نحوه عملکرد خزانه‌های سؤال بهینه در هر دو حالت وارد کردن و همچنین، وارد نکردن عامل S-H در شبیه‌سازی، نتیجه‌گیری کرد. البته تصمیم‌گیری نهایی در مورد آن به موقعیت استفاده از خزانه برمی‌گردد.

در مجموع، عملکرد خزانه‌های سؤال بهینه در شاخص‌های ارزیابی تجربی بهتر از خزانه سؤال عملیاتی بوده است.

نتایج این پژوهش نشان داد که تعمیم رویکرد ریکسی در طراحی خزانه‌های سؤال در موقعیت‌هایی که سؤال‌ها با مدل سه پارامتری مدرج می‌شوند و تلفیق آن با رویکرد برنامه‌نویسی ریاضی به‌خوبی عمل می‌کند. در مقایسه با رویکرد برنامه‌نویسی ریاضی، رویکرد ریکسی شیوه CAT را به‌صورت سراسر تری شبیه‌سازی می‌کند. رویکرد ریکسی با شیوه‌های متفاوت انتخاب سؤال و فرآیند برآورد توانایی مطابقت کامل دارد و دارای انعطاف بیشتری نسبت به رویکرد برنامه‌نویسی ریاضی است. همچنین کاربرد این روش بسیار آسان‌تر است و نیاز به نرم‌افزارهای پیچیده ندارد. در این رویکرد قیود مربوط به صفات غیر آماری از قبیل صفات محتوایی، در مرحله‌ی اول طراحی خزانه وارد برنامه‌نویسی می‌شود و خزانه‌ها را به بخش‌های کوچک‌تری بخش‌بندی می‌کند (هی و ریکسی، ۲۰۱۰). در مقابل رویکرد برنامه‌نویسی ریاضی ساختار بندی ریاضی بیشتری دارد و تمام قیود آماری و غیر آماری را به‌صورت کمی درمی‌آورد و سپس بهترین راه‌حل بهینه را از طریق برنامه‌نویسی خطی جستجو می‌کند (وندربلیندن، ۲۰۰۵a). اما این روش به کاربرد رویکرد انتخاب سؤال "تست سایه" در شبیه‌سازی CAT نیاز دارد. رویکرد ریکسی روی تصادفی سازی کردن پارامترهای سؤال در شبیه‌سازی تأکید دارد، درحالی‌که رویکرد برنامه‌نویسی ریاضی روی بهینه‌سازی سؤال‌های "ساختگی" از قبل تعریف‌شده تأکید دارد. اما دو رویکرد در پایان کار شبیه‌سازی به نتایج یکسانی می‌رسند. رویکرد ریکسی در بعضی جنبه‌ها مشابه رویکرد برنامه‌نویسی ریاضی است. یکی از شباهت‌های مهم بین این دو رویکرد این است که روش شبیه‌سازی خزانه سؤال P (در MRP) و رویکرد برنامه‌نویسی ریاضی در مورد کاهش هزینه‌های طراحی مانند یکدیگر عمل می‌کنند، به‌طوری‌که هر دو فرایند طراحی خزانه سؤال هزینه‌ها را به حداقل می‌رسانند. رویکرد برنامه‌نویسی ریاضی تابع هزینه را که معکوس تعداد سؤال‌ها واقعی با

ترکیب معینی از صفات شامل پارامترهای IRT سؤال‌ها است، تعریف می‌کند. این رویکرد فرض می‌کند که ایجاد سؤال‌ها واقعی با ترکیب مشخص پارامترهای سؤال، هزینه کمتری دارند. این ایده مشابه روش P است که در آن روش در فرایند شبیه‌سازی به‌احتمال بیشتری سؤال‌ها در طول خط رگرسیون پارامترهای b روی پارامترهای a سؤال‌ها واقعی، ایجاد می‌شود (گو و ریکسی، ۲۰۰۷). هر یک از این دو رویکرد بر همین اساس می‌توانند ایده‌های مشابهی از این نوع را از یکدیگر اقتباس کنند تا طراحی خزانه سؤال را بهبود دهند. در این پژوهش، ویژگی‌های برتر برنامه‌نویسی ریاضی شامل کدگذاری دقیق ریاضی ویژگی‌های غیر آماری در ارتباط با ویژگی‌های آماری با رویکرد ریکسی تلفیق شده است. برای دوری از دشواری‌های روش انتخاب سؤال تست سایه نیز از مدل انتخاب سؤال حداقل انحرافات وزن‌دار (WDM) استفاده شده است. نتایج نشان داد که از طریق این روش می‌توان در سرهم کردن تست‌های CAT میزان تخطی از قیود محتوایی را به حداقل رساند؛ بنابراین، نتایج این تحقیق نشان داد که تلفیق رویکرد برنامه‌نویسی WDM به همراه رویکرد اکتشافی ریکسی، نتایج مطلوبی ایجاد می‌کند و با کمی‌سازی ویژگی‌های محتوایی به همراه ویژگی‌های روان‌سنجی سؤال‌ها، می‌توان میزان تخطی تست‌ها را از قیود محتوایی کاهش داد. این روش ویژگی‌های آماری و غیر آماری سؤال را با تعادل مطلوبی بین ویژگی‌های اندازه‌گیری و ساختاری در نظر می‌گیرد. در این پژوهش، این ویژگی‌ها به‌وسیله وزن‌هایی که توسط طراحان اولیه تست انتخاب شد، در مدل وارد شد. این روش برخلاف روش تست سایه، ویژگی‌های محتوایی را به‌عنوان اهداف به‌جای قیود فرمول‌بندی می‌کند. انحراف از اهداف محتوایی وزن داده‌شده و در تابع هدف به همراه فاصله آگاهی سؤال از مقدار هدف قرار می‌گیرد. از این طریق می‌توان مدل‌های غیرقابل‌حل را قابل‌اجرا کرد و تست‌هایی با مقدار تخطی کمتری ایجاد کرد. این نتایج هم‌راستا با تحقیقات بروک، کندریک و مروس، ۱۹۹۸؛ استوکینگ، سوانسون و پیرمن، ۱۹۹۳ است.

## منابع

- Brooke, A., Kendrick, D., & Meeraus, A. (1988). GAMS: A user's guide. Redwood City CA: The Scientific Press.
- Chang, H. (2007). Book review: Linear models for optimal test design. *Psychometrika*, 72, 279-281.
- Chang, H. H., & Ying, Z. (1999). Alpha-stratified multistage computerized adaptive testing. *Applied Psychological Measurement*, 23, 211-222.
- Chang, H. H., & van der Linden, W. J. (2003). Optimal stratification of item pools in a-stratified computerized adaptive testing. *Applied Psychological Measurement*, 27, 262-274.
- Cheng, Y., & Chang, H. (2009). The maximum priority index method for severely con- strained item selection in computerized adaptive testing. *British Journal of Mathematical and Statistical Psychology*, 62, 369-383.
- Chen, S. Y., Ankenmann, R. D., & Spray, J. A. (1999). Exploring the relationship between item exposure rate and test overlap rate in computerized adaptive testing (No. ACT-RR-99-5): American College Testing Program, Iowa City, IA.
- De Ayala, R.J. (2009). The theory and practice of item response theory. New York: Guilford Press.
- Flaugher, R. (2000). Item pools. In H. Wainer (Ed.), Computerized adaptive testing: A primer (pp. 37-59). Mahwah, NJ: Lawrence Erlbaum.
- Gu, L. (2007). Designing optimal item pools for computerized adaptive tests with exposure controls. Unpublished doctoral dissertation. Michigan State University.
- Gu, L. & Reckase, M. D. (2007). Designing optimal item pools for computerized adaptive tests with Sympon-Hetter exposure control. Paper Presented at the 2007 GMAC Conference on Computerized Adaptive Testing, Minneapolis, MN.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). Fundamentals of item response theory. Newbury Park CA: Sage.
- Hau, K. T., & Chang, H. H. (2001). Item selection in computerized adaptive testing: Should more discriminating items be used first. *Journal of Educational Measurement*, 38 (3), 249-266.
- He, W., & Reckase, M. (2010). Optimal item pool design for a highly constrained computerized adaptive test. Unpublished doctoral dissertaion. Michigan State University.
- He, W., & Reckase, M. (2011). Optimal item pool design for a highly constrained computerized adaptive test. Paper presented at the National Council on Measurement in Education, Denver, CO.

- Jensema, C. J. (1972). An application of latent trait mental test theory to the Washington Pre-College Testing Program. Unpublished doctoral dissertation. University of Washington, 1972.
- Jensema, C. J. (1977). Bayesian tailored testing and the influence of item bank characteristics. *Applied Psychological Measurement*, 1, 111-120.
- Lord, F. (1980). Applications of item response theory to practical testing problems. Hillsdale, NJ: Lawrence Erlbaum.
- McBride, J. R., & Weiss, D. J. (1976). Some properties of a Bayesian adaptive ability testing strategy (Research Rep No. 76-1). Minneapolis, MN: Psychometric Methods Program, Department of Psychology.
- Owen, R. J. (1975). A Bayesian sequential procedure for quantal response in the context of adaptive mental testing. *Journal of the American Statistical Association*, 70, 351-356.
- Reckase, M. D. (1974). An application of the Rasch simple logistic model to tailored testing. Paper presented at the Annual Meeting of the American Educational Research Association, Chicago, Illinois.
- Reckase, M. D. (1976). The effect of item pool characteristics on the operation of a tailored testing procedure. Paper presented at the spring meeting of the Psychometric Society, Murray Hill, NJ.
- Reckase, M.D. (1989). Adaptive testing: The evolution of a good idea. *Educational Measurement: Issues and Practice*, 8(3), 11-15.
- Reckase, M. D. (2001, September). Item pool design for computerized adaptive tests. Invited small group session at the 6th Conference of the European Association of Psychological Assessment, Aachen, Germany.
- Reckase, M. D. (2003). Item pool design for computerized adaptive tests. Paper presented at the National Council on Measurement in Education, Chicago, IL.
- Reckase, M. D., & He, W. (2004). The ideal item pool for the NCLEX-RN examination— Report to NCSBN: Michigan State University.
- Reckase, M. D., & He, W. (2005). Ideal item pool design for the NCLEX-RN exam. Michigan State University, East Lansing, MI.
- Reckase, M. D. (2009). Optimal Item Pool Design for the 2009 NCLEX Exam. A Report Submitted to National Council of State Boards of Nursing March 2009.
- Reckase, M. D., & He, W. (2009a). Optimal item pool design for the 2009 NCLEX Exam-report to the National Council of State Boards of Nursing (NCSBN): Michigan State University.
- Reckase, M. D., & He, W. (2009b). The influence of item pool quality on the functioning of computerized adaptive tests. Paper presented at the annual meeting of Psychometric Society, Cambridge, U.K.

- Reckase, M. D. (2010). Designing Item Pools to Optimize the Functioning of Computerized Adaptive Test. *Psychological Test and Assessment Modeling*, 52, 2010 (2), 127-141.
- Robin, F., van der Linden, W. J., Eignor, D. R., Steffen, M., & Stocking, M. L. (2005). A comparison of two procedures for constrained adaptive test construction (ETS Research Rep No. RR-04-39). Princeton, NJ: Educational Testing Service.
- Stocking, M. L., & Swanson, L. (1993). A method for severely constrained item selection in adaptive testing. *Applied Psychological Measurement*, 17, 277-292.
- Stocking, M. L., Swanson, L., & Pearlman, M. (1993). Application of an automated item selection method to real data. *Applied Psychological Measurement*, 17, 167-176.
- Stocking, M. L. (1994). Three practical issues for modern adaptive testing item pools (No. ETS- RR-94-5): Educational Testing Service, Princeton, NJ.
- Sympson, J. B., & Hetter, R. D. (1985, October). Controlling item-exposure rates in computerized adaptive testing. Proceedings of the 27th annual meeting of the Military Testing Association (pp. 973-977). San Diego, CA: Navy Personnel Research and Development Center.
- Urry, V. W. (1977). Tailored testing: A successful application of latent trait theory. *Journal of Educational Measurement*, 14, 181-196.
- Van der Linden, W. J. (1998). Optimal assembly of psychological and educational tests. *Applied Psychological Measurement*, 22, 195-211.
- Van der Linden, W. J., & Reese, L. (1998). A model for optimal constrained adaptive testing. *Applied Psychological Measurement*, 22 (3), 259-270.
- Van der Linden, W. J., & Glas, C. A. W. (2000 a). Capitalization on item calibration error in adaptive testing. *Applied Measurement in Education*, 13(1), 35-53.
- Van der Linden, W. J. (2000 b). Constrained adaptive testing with shadow tests. In W. J. van der Linden, & C. A. W. Glas (Eds.), *Computerized adaptive testing: Theory and practice* (pp. 27-52). Boston: Kluwer Academic Publishers.
- Van der Linden, W. J. (2000 c). Optimal assembly of tests with item sets. *Applied Psychological Measurement*, 24, 225-240.
- Van der Linden, W. J. (2005a). A comparison of item-selection methods for adaptive tests with content constraints. *Journal of Educational Measurement*, 42, 283-302.
- Van der Linden, W. J. (2005b). *Linear models for optimal test design*. New York: Springer-Verlag.



- Van der Linden, W. J., Adelaide, A., & Veldkamp, B. P. (2006). Assembling a computerized adaptive testing item pool as a set of linear tests. *Journal of Educational and Behavioral Statistics*, 31(1), 81-100.
- Veldkamp, B. P., & van der Linden, W. J. (2000). Designing item pools for computerized adaptive testing. In W. J. van der Linden, & C. A. W. Glas (Eds.), *Computerized adaptive testing: Theory and practice* (pp. 149–162). The Netherlands: Kluwer Academic Publishers.
- Wise, S., & Kingsbury, G. G. (2000). Practical issues in developing and maintaining a computerized adaptive testing program. *Psicologica*, 21, 135-155. Retrieved from <http://www.uv.es/psicologica/articulos1y2.00/wise.pdf>.
- Xing, D., & Hambleton, R. K. (2004). Impacts of test design, item quality, and item bank size on the psychometric properties of computer-based credentialing examinations. *Educational and Psychological Measurement*, 64(1), 5-21.