

سنجدش انطباقی با رایانه

* دکتر محمد عسگری

چکیده

منطبق کردن دشواری سوال با سطح ویژگی آزمودنی، کارآترین روش اندازه‌گیری یک ویژگی است. برای این منظور به آزمودنی‌ها، سوال‌هایی ارائه می‌شود که اطلاعات دقیقی در باره میزان و سطح ویژگی آنها، بدست آید. این موضوع اندیشه زیربنایی سنجش انطباقی است و از یافته‌ها و نتایج مدل‌های نظریه سوال-پاسخ، در تحلیل سوال‌های آزمون، متاثر است. سنجش انطباقی از نتایج طبیعی تفکر به شیوه بیز در برآورد و بررسی تابع آگاهی سوال آزمون است، و دارای دلالت‌های ضمنی برای برای تفسیر و تهیه آزمون است. این شیوه دارای تقریب‌های ساده و متنوعی است. سنجش انطباقی با رایانه (سار)، یک روش پیچیده‌تر سنجش انطباقی است. اجرای برنامه سار تکرار فرآیندی دو مرحله‌ای است که در گام نخست، سوالی به آزمودنی عرضه می‌شود که دشواری آن با توانایی فعلی برآورد شده برای او برابر باشد. در گام دوم، توانایی آزمودنی براساس اطلاعات حاصل از پاسخ او به سوال مرحله اول، برآورد می‌شود. این گام‌ها تکرار می‌شوند تا جایی که ملاک‌های پایان دادن به آزمون محقق شوند. برنامه‌ریزی واقعی، اجرای عملی و حفاظت برنامه سار ساده نیست. چالش‌های عملی و پیچیدگی‌های برنامه سار شامل: خزانه‌های سوال، حفاظت آزمون، و مسائل مربوط به آزمودنی‌ها.

واژگان کلیدی: سنجش انطباقی، سنجش انطباقی با رایانه، نظریه سوال-پاسخ

مقدمه

ایدۀ اساسی سنجش انطباقی، منطبق کردن سوال‌های آزمون با سطح صفت (توانایی یا مهارت) آزمودنی، در جریان آزمودن است (اگن^۱، 2007). به عبارت دیگر، این روش مستلزم انتخاب سوال‌های آزمون در هنگام اجرای آن است؛ به طوری که سوال‌هایی که برای آزمودنی اجرا می‌شود از لحاظ دشواری مناسب او باشد. اگر انطباق سوال‌های آزمون با سطح صفت (توانایی یا مهارت) آزمودنی با استفاده از رایانه انجام شود، به آن «سنجش انطباقی با رایانه»^۲ (سار) گفته می‌شود (وای^۳، 2005). در سار، صفحه رایانه در هر زمان، یک سوال که از خزانه بزرگ سوال‌های طبقه‌بندی شده بر اساس محتوا و دشواری انتخاب می‌شود، را برای آزمودنی نمایش می‌دهد. اولین سوال در سنجش انطباقی با رایانه همیشه دارای سطح دشواری متوسط است. انتخاب سوال بعدی با توجه به عملکرد آزمودنی در سوال‌های قبل توسط رایانه صورت می‌گیرد. بدین ترتیب فرد به سوال‌هایی پاسخ می‌دهد که منطبق با صفات زیر بنایی عملکرد یا خصیصه مکنون او باشد. در نهایت مجموع سوال‌های پاسخ داده شده، آزمونی را تشکیل می‌دهد که منطبق با سطح توanایی آزمودنی است. سنجش انطباقی با رایانه (سار) در عمل دارای مراحل زیر است:

1. داشتن خزانه بزرگی از سوال‌ها؛ 2. اجرای آزمون یا تعداد زیادی سوال روی آزمودنی‌ها؛
3. استخراج نتایج آزمون بر اساس مدل‌های نظریه سوال - پاسخ^۴ IRT؛ 4. گرینش سوال‌های آزمون منطبق با سطوح توanایی آزمودنی‌ها (اگن، 2007؛ وای، 2005؛ هورنکه^۵، 2000).

تأکید اصلی این مقاله بر جنبه‌های روانسنجی سنجش انطباقی با رایانه (سار) است. همچنین کاربردهای عملی آن نیز مورد توجه قرار می‌گیرد. این مقاله با توصیفی کلی از سنجش انطباقی شروع شده، پس از آن درباره سار، مزیت‌ها، ویژگی‌های روانسنجی، و تلویحات کاربردی آن بحث می‌شود. در نهایت به نتیجه‌گیری از بحث‌ها پرداخته شده است.

1. Eggen

2. Computerized Adaptive Testing (CAT)

3. Way

4. Item Response Theory (IRT)

5. Hornke

سنچش انطباقی

سنچش انطباقی¹، برازش² یک آزمون را در جریان آزمودن سطح یک صفت (توانایی یا مهارت) یک آزمودنی است (اگن، 2007). این سنچش از یافته‌ها و نتایج مدل‌های نظریه سؤال-پاسخ (IRT) در تحلیل سؤال‌های آزمون‌ها متأثر است (ستاری، 1382). در IRT فرض می‌شود که؛ الف) عملکرد آزمودنی در آزمون را می‌توان با مجموعه‌ای از عواملی که صفات³، صفات مکنون⁴، و یا توانایی‌ها⁵ نامیده می‌شوند، برآورد کرد؛ و ب) ارتباط بین عملکرد آزمودنی در سؤال و مجموعه صفات مذکور با یک تابع تجمعی که ویژگی سؤال نامیده می‌شود، قابل برآورد است (سوتاridona⁶، پرنل⁷، و واجو⁸، 2003).

کارآترین روش اندازه‌گیری یک ویژگی یا نمره واقعی، منطبق کردن سطح دشواری سؤال با سطح ویژگی یا نمره واقعی آزمودنی است؛ به این معنی که سؤال‌های دشوار برای آزمودنی‌هایی که دارای نمره واقعی بالا یا سطح بالایی از ویژگی هستند، و سؤال‌های آسان برای آزمودنی‌هایی که دارای نمره واقعی پایین یا سطح پایینی از ویژگی هستند، اجرا شود. هنگامی که امکان حدس زدن وجود ندارد، کارآترین سؤال، سوالی است که احتمال پاسخ صحیح دادن آزمودنی به آن، برابر 0/5 باشد. زمانی که گروهی از آزمودنی‌ها با تنوع زیادی از ویژگی‌ها با همان آزمون اندازه‌گیری می‌شوند، امکان ندارد که به طور همزمان، آن آزمون حداقل کارآیی را برای کلیه آزمودنی‌ها داشته باشد. در این موارد، بهتر است سؤال‌هایی به آزمودنی داده شود که اطلاعات عمده‌تری را در باره میزان و سطح ویژگی او بدست دهد. برای رسیدن به این هدف، باید سؤال‌های متفاوتی برای آزمودنی‌های متفاوت مطرح شود؛ یعنی اندازه‌گیری انفرادی شود و این همان اندیشه زیربنایی و اساسی سنچش انطباقی است (اگن، 2007).

-
1. Adaptive testing
 2. Tailoring
 3. Traits
 4. Latent traits
 5. Abilities
 6. Sotaridona
 7. Pornel
 8. Vallejo

سنجدش انطباقی، از نتایج طبیعی تفکر به شیوه بیز درباره برآورده و بررسی تابع آگاهی سئوال آزمون است (ثرندایک¹، 1982؛ ترجمه هومن، 1369). هدف اساسی از بکارگیری روش‌های بیز در اندازه‌گیری، افزایش دقت برآورد نمره واقعی آزمودنی است. این مطلب در واقع استفاده مؤثر از روش بیز در سنجدش برآذش یافته، یعنی تطبیق دقیق دشواری تکلیف آزمون با سطح توانایی فرد است. یعنی کوششی برای طرح ریزی روش‌های بخصوصی از سنجدش است که در هر مرحله آن، سئوال بعدی به گونه‌ای انتخاب شود که با توجه به برآورده کنونی جایگاه آزمون شونده در صفت مکنون مورد اندازه‌گیری، بیشترین مقدار اطلاعات ممکن را بدست دهد. تفکر به شیوه بیز دارای دلالت‌های ضمنی هم برای تفسیر آزمون و هم برای تهییه آزمون است. در مورد تفسیر آزمون، رهنمودی برای تهییه بهترین برآورده ممکن از جایگاه آزمون شونده در یک خصیصه مکنون با توجه به اطلاعاتی که درباره تاریخچه گذشته فرد و کارکرد فعلی او در یک یا چند آزمون در دست داریم؛ در اختیار ما قرار می‌دهد. در مورد تهییه آزمون، برای گزینش مواد آزمون پایه‌ای منطقی به دست می‌دهد که بیشترین مقدار اطلاعات را برای هر مرحله سنجدش فراهم می‌آورد، و در نتیجه درجه دقت برآورده بعدی ما از وضع نسبی فرد، در صفت مکنون با بیشترین سرعت ممکن افزایش می‌دهد. این موضوع در قلب سنجدش انطباقی یا سنجدش برآذش یافته قرار دارد که علاقه و توجه به آن روز افزون است (ثرندایک، 1982؛ ترجمه هومن، 1369؛ آلن² و ین³، ترجمه دلاور، 1374؛ افروز و هومن، 1375). پیشرفت‌های نظریه‌های اندازه‌گیری، بخصوص نظریه سئوال-پاسخ، تکنولوژی رایانه و فناوری اطلاعات و استفاده از آنها در سنجدش، ابزارهای لازم برای سنجدش انطباقی با رایانه را در سال‌های اخیر تسهیل نموده است (اگن، 2007).

سنجدش انطباقی با رایانه

به طور سنتی، ارزشیابی دانش، مهارت‌ها، توانایی‌ها، و سایر ویژگی‌های شخصیتی توسط آزمون‌های کتبی (داد-کاغذی) انجام می‌گرفته است. پیشرفت و توسعه

1. Thorndike

2. Allen

3. Yen

فناوری اطلاعات¹ در دو دهه گذشته آزمودن مبتنی بر رایانه² را هم در پژوهش‌های آموزشی و هم در عمل تسهیل نموده است (تاو³، وا⁴، و چنگ⁵، 2008). در سنچش انطباقی با رایانه که مورد خاصی از سنچش مبتنی بر رایانه است، هر آزمودنی یک آزمون منحصر به فرد که برای سطح توانایی او برازش یافته است، دریافت می‌کند. سار، استفاده از رایانه در انتخاب سوال‌ها، ضمن پاسخگویی آزمودنی به هر سوال است. زمانی که آزمودنی به سوال پاسخ درست می‌دهد، براساس فرمول خاصی، سوال مشکل‌تری به او ارائه می‌شود. چنانچه به سوال پاسخ غلط داده شود، سوال آسان‌تری پیشنهاد می‌گردد. انتخاب سطح دشواری سوال‌های متواالی، براساس عملکرد آزمودنی در هر مرحله صورت می‌گیرد. بعد از پاسخ آزمودنی، توانایی او به روز شده، و سوالی از خزانه سوال برای وی انتخاب می‌شود که بیشترین تناسب را با توانایی جدید او داشته باشد و این روند ادامه پیدا می‌کند تا ملاک‌های پایان دادن به آزمون محقق شوند (ترین تافولو⁶، جرجیادو⁷، و إكونومیدس⁸، 2006). سنچش انطباقی با رایانه (سار) به صورت روز افزون و در مقیاسی بزرگ در برنامه‌های سنچش و آزمون وارد شده است. مزیت‌های سار هم برای سازندگان و هم برای اجراءکنندگان آزمون مفید و امیدوارکننده بوده است. زیرا، به آزمون‌کننده اجازه می‌دهد تا سطح ویژگی آزمودنی را به سرعت بسنجد، بدون اینکه او را وادر کند تا به مجموعه‌ای خسته‌کننده از سوال‌های بسیار ساده و یا مجموعه‌ای ناخوشایند و ناراحت‌کننده از پرسش‌های خیلی مشکل پاسخ دهد (الیا⁹، رولتا¹⁰، زیمنز¹¹، و آباد¹²). (2000)

1. Information Technology (IT)

2. Computer- Based Testing (CBT)

3. Tao

4. Wa

5. Chang

6. Triantafillou

7. Georgiadou

8. Economides

9. Olea

10. Revuelta

11. Ximenez

12. Abad

انگیزه اصلی استفاده از سار، در کارآیی آن نهفته است. در مقایسه با روش‌های سنجش مداد- کاغذی سنتی، افزایش کارآیی سنجش در نتیجه رایانه‌ای کردن شیوه آزمودن، با تأکید بر کارآیی اندازه‌گیری است. نتایج مطالعات نشان می‌دهد که در سار و با دقت یکسان، به سوال‌های کمتری در مقایسه با روش‌های غیرانطباقی، نیاز هست. این سنجش به حدود 50 تا 60 درصد سوال‌های در مقایسه با روش‌های غیرانطباقی نیاز دارد. امتیازهای دیگر سار نسبت به روش‌های سنتی عبارت‌اند از:

1. هر فردی یک آزمون منحصر به فردی را پاسخ می‌دهد که هم محتوا و هم طول آزمون می‌تواند از فردی به فرد دیگر متفاوت باشد؛ 2. سار برای هر فردی بهینه است؛ به گونه‌ای که حداقل دو پیامد مطلوب خواهد داشت: (الف) از آنجا که سوال‌های کمتری برای دستیابی به دقتی مشابه، در سار نیاز است، کارآیی اندازه‌گیری افزایش می‌یابد، (ب) هر آزمودنی می‌تواند با سطح توانایی خود در چالش باشد؛ این چالش اثر تحریکی و تجربه شده‌ای به عنوان تقویت بر توانایی آنها دارد؛ 3. چون برای اندازه‌گیری توانایی هر فردی به سوال‌های کمتری نیاز است، بنابراین در هزینه و وقت آزمودنی‌ها، تهیه‌کنندگان آزمون، و مجریان آن صرفه‌جویی می‌شود؛ 4. از آنجا که سوال‌های آزمون در چارچوب نظریه IRT تهیه می‌شود، بنابراین آزمون دارای تعدادی ویژگی شناخته شده است؛ بر همین اساس کیفیت احتمالی اندازه‌گیری بهتر است؛ 5. نمره‌گذاری آن عینی و بدون خطاست و امکان نمره‌گذاری و بازنورد مستقیم وجود دارد؛ 6. دستیابی به انواع سوال‌های جدید امکان‌پذیر است؛ 7. اجرای آزمون کارآیی بیشتری دارد؛ 8. برنامه‌ریزی و مکان‌یابی آزمون انعطاف‌پذیر است؛ 9. اینمنی آزمون بیشتر است؛ 10. آزمودنی‌ها انگیزه بیشتری دارند؛ و 11. امکانات بیشتری برای بکارگیری نظریه‌های جدید آزمون وجود دارد (اگن، 2007).

اجرای برنامه سنجش انطباقی با رایانه، اساساً تکرار یک فرایند دو مرحله‌ای است. در گام اول، سوالاتی به آزمودنی ارائه می‌شود که دشواری آن با توانایی فعلی (ذاتی) برآورده شده برای او برابر باشد. در گام دوم، توانایی آزمودنی بر اساس اطلاعات حاصل از پاسخ او به سوال مرحله اول برآورده می‌شود. سپس این دو گام تکرار می‌شوند تا جایی که یکی از ملاک‌های پایان دادن به آزمون محقق شوند؛ ملاک‌هایی مثل تعداد معینی سوال، یا با سطوحی از دقت اندازه‌گیری مورد انتظار. از طریق این

فرآیند الگوریتمی¹ سنچش انطباقی با رایانه توانایی برآورد شده برای آزمودنی، تثیت² تثیت² می‌شود. اگرچه، سنچش انطباقی ایده نظری نسبتاً ساده‌ای دارد، اما برنامه‌ریزی برنامه‌ریزی واقعی، اجرای عملی و حفاظت یک برنامه سنچش انطباقی با رایانه بسیار پیچیده است. چالش‌های عملی که این پیچیدگی را به وجود می‌آورند و باید در تهیه برنامه‌های سار مورد توجه قرار گیرند عبارت‌اند از: تهیه و نگهداری خزانه‌های سؤال³، حفاظت آزمون⁴، و مسائل مربوط به آزمودنی‌ها⁵ (الیا و همکاران، 2000؛ وای، 2005؛ اگن، 2007).

الف) تهیه و نگهداری خزانه سؤال در سنچش انطباقی با رایانه: در این باره، تعدادی موضوع کاربردی مؤثر بر اندازه‌گیری، باید شناسایی شوند. این شناخت به تهیه خزانه سؤال کمک می‌کند و هنگام اجرا، به اجرای نتیجه آزمون کمک می‌کند تا نتایج حاصل از آزمون را به صفت بنیادی اندازه‌گیری شده ربط دهد. این موضوعات عبارت‌اند از: 1. اندازه خزانه سؤال و کنترل آن، 2. تعیین ابعاد خزانه سؤال، 3. مدل‌های پاسخ، 4. جابجایی و اصلاح سؤال‌ها 5. اضافه کردن سؤال به خزانه سؤال، 6. حفظ هماهنگی مقیاس‌ها.

1. اندازه خزانه سؤال برای سنچش انطباقی: برای سنچش انطباقی با رایانه، وجود خزانه سؤال با تعداد سؤال‌های کافی و برآذش یافته با یک مدل IRT ضروری است. همه سؤال‌ها باید به خوبی آماده اندازه‌گیری صفتی باشند که، می‌بایست اندازه‌گیری شود. یک خزانه سؤال خوب، تمام جوانب مرتبط با صفتی را که قرار است با استفاده از سار اندازه‌گیری شود، به خوبی پوشش می‌دهد. در اصل، معیارهای سؤال‌های خوب در سار خیلی متفاوت از سؤال‌های خوب در آزمون‌های کتبی (مداد- کاغذی) نیست (اگن، 2007). قبل پیشنهاد شده بود که حداقل به خزانه سؤالی با 100 سؤال، برای سنچش انطباقی با رایانه ضروری است (یوری، 1977، به نقل الیا و همکاران، 2000). براین اساس، همه آزمون‌سازان تلاش می‌کردند آزمون‌هایی تهیه نمایند که حداقل 100 سؤال داشته باشد. اما سه عامل باعث شد که تعداد سؤال‌های خزانه

1. Algorithm

2. Converge

3. Item Pools

4. Test Security

5. Examinee Issues

سؤال افزایش چشمگیری یابد: نخست این که، در چند دهه اخیر آزمون‌های سنتی بهبود چشمگیری داشته‌اند، به گونه‌ای که سطح عملکرد آزمون‌های آموزشی تقریباً دقیقی به اندازه سنجش انطباقی دو مرحله‌ای دارند. دوم، الزام‌ها و اجرهایی که بر شیوه‌های انتخاب سوال‌های سنجش انطباقی تحمیل شده‌اند (مثل الزام‌های محتوایی و منطقی آزمون) ضرورت ایجاد خزانه سوال‌های بزرگتری را نسبت به نبود چنین الزام‌هایی، به منظور دستیابی به اطلاعات مشابه، به وجود آورده است. سوم این که، اگر قرار باشد آزمونی که تهیه شده است برای مدتی طولانی مورد استفاده قرار گیرد، باید دارای خزانه سوال زیادی باشد، تا این‌مانی سوال‌ها به خطر نیافتد. امروزه بسیاری از آزمون‌های انطباقی خزانه سوال‌هایی با بیشتر از 1000 سوال در اختیار دارند و با استفاده از چرخش دوره‌ای سوال‌ها، تعداد سوال‌های خزانه را به 2000 سوال می‌رسانند. بدیهی است که، محدود کردن دایره شمول یک آزمون، نخستین عامل در تعیین این مطلب است که چند سوال برای آزمون لازم است. در تعیین تعداد سوال‌های خزانه تمایز بین باید و واقعیت کارآیی سنجش انطباقی ضروری است (الیا و همکاران، 2000؛ وای، 2005).

2. تعیین ابعاد پاسخ‌های سوال‌های خزانه‌های سوال یک بعدی نیستند و بنابراین پاسخ‌های سوال‌ها نیز چنین هستند. بنابراین در بررسی ابعاد پاسخ‌ها، تعامل بین آزمودنی‌ها و سوال‌های آزمون بررسی می‌شود. لذا ضروری است که در تعریف جامعه مورد نظر روش و صریح اقدام شود. اگر تحلیل ابعاد با داوطلبان غیرعلقه مند و یا دانشجویان قبل از آموزش انجام شود، ابعادی که به دست می‌آید شبیه همان ابعادی نخواهد بود که از اجرای آزمون برای گروهی دانشجویان واقعی و با انگیزش فردی بالا پس از آموزش به دست می‌آید. دو نکته ضروری در بررسی ابعاد پاسخ‌ها عبارت‌اند از: نکته اول در زمان تهیه خزانه سوال باید مد نظر قرار گیرد و آن این ضرورت است که تعیین شود که آیا برای توصیف پاسخ‌های آزمودنی‌ها از مدل یک بعدی و یا چند بعدی سوال پاسخ باید استفاده نمود. این کار را می‌توان با تحلیل ابعاد سوال‌ها انجام داد. نکته دوم آن است که ضرورت دارد، ابعاد فضایی پاسخ‌ها زمانی که محتوای آزمون و یا آزمون شوندگان به صورت منظم تغییر می‌کنند، مورد بررسی قرار گیرند. این موضوع غالب توسط آزمون‌سازان برای تعیین مقیاس اندازه‌گیری مورد غفلت قرار می‌گیرد. روش‌های زیادی برای تعیین ابعاد فضایی پاسخ‌ها وجود دارد که عبارت‌اند از: تحلیل عاملی اکتشافی؛ تحلیل عاملی تأییدی؛

تحلیل عاملی کامل؛ مقیاس پردازی چند بعدی؛ و تکنیک‌های متفاوت وابسته به نظریه‌های سئوال-پاسخ (اگن، 2007؛ الیا و همکاران، 2000؛ وای، 2005).

اگر تحلیل‌ها نشان داد که رویکرد یک بعدی کافی است، استفاده از مدل سئوال پاسخ یک بعدی مفروض به صرفه و عاقلانه است. اما اگر نتایج تحلیل نشان داد که ساختار چند بعدی برای تبیین پاسخ‌های آزمودنی‌ها ضروری است، منطقی است که با دقت ساختار آزمون را مورد توجه قرار داده و برای ساده‌سازی آن تلاش شود به گونه‌ای که مدل‌های سئوال پاسخ چند بعدی ساده‌تری، پاسخ‌ها را تبیین نماید (الیا و همکاران، 2000).

3. مدل‌های سئوال پاسخ: در سار داشتن یک خزانه سئوال با کیفیت بالا ضروری است، اما کافی نیست. سئوال‌ها باید با یک مدل مناسب IRT برآذش یافته و مدرج‌سازی¹ شده باشند. بر اساس داده‌های تجربی حاصل از اجرای عملی سئوال‌ها،

یک مدل IRT برآذش یافته و پارامترهای سئوال و پارامترهای یک یا بیشتر از یک توزیع توانایی (θ) در جامعه برآورده می‌شوند. کار یک مدل IRT، با مدرج‌سازی سئوال‌های آغاز می‌شود. مدرج‌سازی سئوال‌های برآذش یک مدل IRT و برآورده پارامترهای سئوال را در بر می‌گیرد. اگر پارامترهای سئوال‌های برآورده شدند، آنگاه آنها را می‌توان در خزانه سئوال نگهداری و برای سار مورد استفاده قرار داد. در یک سار احتمال عملکرد تတی یک آزمون نقش اساسی در استنباط در باره آزمودنی دارد.

خصوصیات یک مدل IRT خوب، برای استفاده عملی و مناسب در سار عبارت‌اند از:
1. امکان برآورده توانایی (θ تတی) یک آزمودنی با مقیاسی مشترک، و با هر زیر مجموعه‌ای از سئوال‌های مدرج‌سازی شده در خزانه سئوال وجود دارد. بنابراین، ضرورتی ندارد که مجری آزمون، سئوال‌های یکسانی را به آزمودنی بدهد تا برآورده قابل مقایسه‌ای از θ آنها بدست آورد.

2. دشواری سئوال با مقیاس مشترکی مثل توانایی θ آزمودنی‌ها بیان می‌شود. بنابراین، امکان برآذش یک آزمون با هر سطحی از توانایی θ یک آزمودنی وجود دارد.
3. آگاهی یک سئوال عملکردی از توانایی θ است، در نتیجه تابع آگاهی می‌تواند به عنوان اساسی برای برآذش سئوال انتخابی به کار گرفته شود. با این وجود، ویژگی‌های مطلوب فوق، فقط زمانی درست‌اند که برآذش سئوال‌های انتخاب شده برای مدل IRT لحاظ شده باشند (اگن، 2007).

بسیاری از پژوهشگران انتخاب مدل ویژه سوال پاسخ را به مجادله تبدیل نموده‌اند. اگرچه، هیچ دلیل نظری و یا عملی وجود ندارد که باید فقط به یک مدل پاسخ محدود شد. برای نمونه، اگر فردی بخواهد یک مقیاس اندازه‌گیری با تعداد محدودی از افراد بسازد، مدل یک پارامتری^۱ مطلوب‌ترین مدل احتمالی برای استفاده اوست (Lord², 1983)، حتی اگر انتظار برآزنده بودن سوال‌های خزانه با مدل یک پارامتری را نداشته باشیم. زمانی که حجم نمونه بزرگتری در دسترس باشد، و فرد قصد افروden تعدادی سوال به سوال‌های خزانه داشته باشد، بدون این که ویژگی‌های^۳ مقیاس اندازه‌گیری آسیب بیند، هیچ چیزی مانع از به کارگیری مدل سه پارامتری^۳ برای سوال‌های جدید نیست. به همان نسبتی که به سوال‌ها افروده می‌شود، مقیاس نیز مطلوب‌تر می‌شود (Alia & Hamkaran, 2000).

برای دستیابی به کیفیتی بالا در مدرج‌سازی سوال‌های، انتخاب مدل IRT انتخاب روش برآوردهای آماره‌های برازش، ماهیت و حجم نمونه آزمودنی‌ها نقش مرتبط به هم، و مهمی دارد (Fischer⁴ و Molenaar⁵, 1995؛ ون در لیندن⁶ و Hambleton⁷, 1996). در انتخاب مدل IRT، هر چه مدل انتخابی ساده‌تر باشد، به حجم نمونه کمتر نیاز است، و روش‌های آماری بهتری برای برآوردهای برازش، آزمودن مدل وجود دارد. برای برآوردهای پارامترها با دقیقی معقول، و برای آزمودن مدل با توانی بالا در مدل‌های یک پارامتری، دو پارامتری، و سه پارامتری؛ به ترتیب حداقل به 500, 200، و 1000 پاسخ آزمودنی در هر سوال نیاز است. به عبارت دیگر، دستیابی به برازش خوب برای یک مدل ساده‌تر خیلی مشکل است. این بدین معنی است که برخی سوال‌های باید از خزانه حذف شوند، که این موضوع روایی خزانه سوال را تهدید می‌کند (Agn, 2007؛ ون در لیندن و Glas⁸, 2000).

دو روش احتمالی متداول برای برآوردهای پارامترهای سوال وجود دارد. نخست روش بیشینه احتمال حاشیه‌ای⁹ (MML) است. در این روش فرض می‌شود که

1. One-parameter logistic model (1PL)

2. Lord

3. Three-parameter logistic model (3PL)

4. Fischer

5. Molenaar

6. Van der Linden

7. Hambleton

8. Glas

9. Marginal maximum likelihood (MML)

نمونه‌های تصادفی از جامعه‌ای با توزیع شناخته شده‌ای از توانایی انتخاب شده است. (به صورت معمول فرض می‌شود توزیع توانایی در جامعه نرمال است). با روش بیشینه احتمال حاشیه‌ای (MML)، پارامترهای سؤال و پارامترهای توزیع توانایی همزمان با هم برآورد می‌شوند. روش دوم، روش بیشینه احتمال موقعیتی^۱ (CML) است. اگر از این روش استفاده شود، امکان برآورد همزمان پارامترهای سؤال و پارامترهای توانایی وجود ندارد، در روش بیشینه احتمال موقعیتی، به مفروضه‌هایی در باره توزیع توانایی آزمودنی‌ها نیاز نیست؛ فقط به نمونه‌هایی از جامعه نیاز است. این موضوع در عمل بسیار مطلوب است. زیرا، در آموزش و پرورش دستیابی به نمونه‌های تصادفی واقعی به آسانی امکان پذیر نیست. انتخاب نمونه‌هایی با روش خوش‌های روایی برآورد پارامترهای سؤال را با استفاده از روش بیشینه احتمال موقعیتی تهدید نمی‌کند. با این وجود، روش برآورد (CML)، برای هر مدلی قابل استفاده نیست. از این روش در مدل یک پارامتری و گاهی اوقات در مدل دو پارامتری می‌توان استفاده کرد (اگن، 1990؛ ورهلست^۲، گلاس^۳، و ورسترالن^۴ 1995). اگر از (CML) برای برآورد کردن پارامترهای سؤال و برآش مدل آزمون استفاده شود، مدرج‌سازی با برآورد پارامترهای یک توزیع توانایی به صورت جداگانه تکمیل می‌شود (اگن، 2007).

4. جابه‌جایی، اصلاح و آزمودن مجدد: در تهیه یک خزانه سؤال و مقیاس اندازه‌گیری برای استفاده در یک آزمون انطباقی، ضروری است شیوه‌هایی را تعیین و به کار گرفته شود که به حذف و کنار گذاشتن سؤال‌هایی که عملکرد خوبی ندارند، کمک نماید. این موضوع ارتباطی به مدل‌های سؤال پاسخ انتخاب شده ندارد. به محض این که سؤال‌های ضعیف شناسایی شدند، باید آنها را از خزانه سؤال کنار گذاشت؛ زیرا سؤال‌های ضعیف علت اصلی خطای برآورد توانایی و تصمیم‌گیری درباره آزمودنی‌ها هستند. یک روش برای بررسی سؤال‌هایی که برآزنده نیستند، رسم منحنی سؤال-پاسخ^۵ تجربی و مقایسه آن با منحنی نظری^۶ است. مقایسه منحنی

1. Conditional maximum likelihood (CML)

2. Verhelst

3. Glas

4. Verstralen

5. Item response curve

6. Theoretical curve

نظری با منحنی تجربی در اکثر موارد رویکردی بسیار ساده برای تشخیص بصری سوالهای ضعیف است. فقط از طریق عدم پذیرش سوالهای ضعیف و شیوه اصلاحی است که می‌توان خزانه سوالی نگهداری نمود که به هماهنگی اندازه‌گیری کمک نماید. در آزمون انطباقی این مطلب مهمی است که هرچه آزمون کوتاه‌تر باشد، عملکرد غیر عادی سوالهای ضعیف، برآن تأثیر بیشتری خواهد داشت (الیا و همکاران، 2000).

5. افزایش سوالهای خزانه سوال: زمانی که خزانه سوالی از سوالهای مدرج شده¹ برای یک مقیاس اندازه‌گیری بخصوص، در اختیار است، می‌توان از طریق هر طرح لنگر²، سوالهایی را، به خزانه سوال اضافه نمود. معمولاً گروهی از آزمودنی‌ها به مجموعه‌ای از سوالهای مدرج شده قبلی و تعدادی سوال جدید مدرج شده، پاسخ می‌دهند. پس از آن یک روش لنگر برای مدرج‌سازی سوالهای جدید به کار می‌رود. بکارگیری این روش در یک آزمون انطباقی بسیار منطقی است. زیرا؛ طی آن فرد می‌تواند سوالهای جدید را در میان آزمون جدید وارد نماید (الیا و همکاران، 2000).

6. حفظ هماهنگی مقیاس: اگر یک مقیاس چند سال مورد استفاده قرار گیرد، میانگین‌ها، درصدها، و سایر ویژگی‌های نمونه در طول این سال‌ها تغییر می‌کنند. طرح لنگر با پارامترهای ثابت می‌تواند این تغییرات محیطی را به گونه‌ای ثبیت کند که مقیاس اندازه‌گیری در پارامترهای برآورده شده، انحراف خیلی زیادی نداشته باشد. اگرچه، منطقی است مطالعه‌ای انجام شود و طی آن برخی از سوالهایی که قبل از مدرج شده‌اند، مجدداً مدرج شوند تا مشخص شود که انحراف ناچیز است و در محدوده بازگشت‌هایی است که ما از خطای نمونه‌گیری انتظار داریم. به علاوه برای حفظ هماهنگی مقیاس می‌توان شرایط اجرایی آزمون را محدود کرد. موضوعاتی مثل تغییرات اجرایی، محدودیت‌های زمانی، تناسب نمونه و سایر شرایط محیطی که علت ناپایداری در مقیاس اندازه‌گیری می‌شوند، را تعریف و مشخص نمود. اگر می‌خواهیم مقیاسی داشته باشیم که در دوره‌ای طولانی از زمان ثابت و پایدار باقی بماند نیازمند

1. Calibrated

2. Linking design

کنترل یا محاسبه این عوامل هستیم (الیا و همکاران، 2000). از آنجا که تمام مدل‌های IRT فرض می‌کنند که سوال‌های آزمون از نوع سوال‌هایی هستند که توانایی آزمودنی را می‌سنجند، لذا مهم‌ترین عامل ممکن است سرعت باشد. اگر در آزمون محدودیت زمان پاسخگویی اعمال شود، سوال‌های آخر آزمون دشواری بالا و سوال‌های نخست دشواری پایینی خواهد داشت و این فرآیند دوری در تکرار مجدد آزمون‌ها باعث فاصله گرفتن مقیاس اندازه‌گیری آزمون جدید، از آزمون قبلی خواهد شد. ساده‌ترین روش برای حل این مشکل آن است که الف) تغییر مکان هر سوال در حوزه بخصوصی از آزمون امکان‌پذیر شود. ب) فقط از آزمون‌هایی استفاده شود که محدودیت زمانی برای پاسخگویی به آنها اعمال نشود (سوتاریدونا و همکاران، 2003).

ب) حفاظت آزمون: حفاظت آزمون موضوعی جالب در تمام برنامه‌های آزمودن است. اگر ایمنی آزمون به خطر بیافتد، مهم نیست که چگونه ویژگی‌های روان‌سنجدی آزمون تغییر می‌کند، بلکه مهم آن است که روایی تفسیر نمرات آن درست نیست. موقفيت سنچش انطباقی با رایانه (سار) به سلامت خزانه سوال‌های آن بستگی دارد. به میزانی که سوالی و یا سوال‌هایی به وسیله آزمودنی‌های قبلی شناخته شوند، پارامترهای برآورده شده برای آن سوال (برآورده شده از طریق یک مدرج سازی)، برای مدتی طولانی قابل استفاده نخواهد بود. به نسبتی که آزمودنی‌های بیشتری محتوای آزمون را شناسایی کنند، پارامتر دشواری سوال‌های آن آزمون ساده‌تر شده و به قسمت پایین مقیاس توانایی انتقال می‌یابند، پارامتر تمیز سوال‌های آن آزمون به سمت صفر متماطل شده و پارامتر حدس زدن پاسخ درست سوال‌های به گونه‌ای غیرمنطقی افزایش می‌یابد (کالتون¹، 1998؛ الیا و همکاران، 2000). به همین دلیل ضروری است که انتخاب و نمره‌گذاری سار محترمانه باشد. دو موضوع جالب و اساسی درباره محترمانه بودن سوال‌ها در سار وجود دارد: افشا شدن سوال‌ها، سرقت سوال‌ها.

۱) افشا شدن سوال‌ها: در آزمون‌های مداد- کاغذی (کتبی) به طور معمول تمام آزمودنی‌ها را می‌توان به طور همزمان آزمون کرد. دفترچه‌های آزمون به اندازه کافی برای آزمون هر یک از آزمودنی‌ها با قیمت مناسب تهیه، و در اختیار آنها قرار داد.

برعکس، سنجش مبتنی بر رایانه پرهزینه است. این احتمال وجود دارد که تعداد رایانه‌ها کمتر از آزمودنی‌هایی باشد که قرار است آزمون شوند. این مطلب بر این دلالت دارد که برخی از آزمودنی‌ها قبل از دیگران آزمون می‌شوند. بعلاوه یکی از مزیت‌های سنجش مبتنی بر رایانه که قابلیت آن را برای ضرورت آزمودن مشخص می‌نماید، آن است که هر یک از آزمودنی‌ها در زمان‌های متفاوتی آزمون می‌شوند و به تعداد زیادی مجری نیاز نیست. این مطلب گویای آن است که یک تمرین عمومی برای آزمودنی‌ها در طول اجرای آزمون وجود دارد که با خودشان درباره سوال‌های آزمون نجوا کنند، بخصوص زمانی که پیامدهای آن برای عملکرد در آزمون بالا باشند. دانش‌آموزانی که سوال‌ها را مطالعه می‌نمایند و از دوستان خود درباره آنها مطالبی شنیده‌اند و پس از آن به سوال‌های آزمون پاسخ می‌دهند، به صورت بالقوه نسبت به سایرین برتری خواهند داشت و توانایی آنها با سوگیری مثبت برآورد خواهد شد. یک راه حل این مشکل، استفاده از خزانه سوال‌های بزرگ است، به گونه‌ای که کم کردن هر کدام از این مجموعه از خزانه، چنین افشاگری‌هایی را تحت تأثیر قرار ندهد. مشکل مرتبط دیگری که در مدارس رخ می‌دهد، این است که معلمان در باره سوال‌های آزمون بخصوصی که از دانش‌آموزان پیش‌بینی شده‌اند پرس‌وجو می‌کنند و بعد به سرعت محتوای آزمون را یاد می‌دهند، در این صورت عملکرد دانش‌آموزان افزایش می‌باید (الیا و همکاران، 2000؛ سوتاریدونا و همکاران، 2003؛ وای، 2005).

(2) سرقت سوال‌ها: تمام برنامه‌های آزمون باید در باره سرقت سوال‌های آزمون حساس باشند. برخی از طرفداران سار معتقدند که چنین آزمون‌هایی ذاتاً از اینمنی بیشتری برخوردارند؛ زیرا از فرم‌های آزمون هیچ کپی وجود ندارد که بتوان آن را سرقت کرد و یا فتوکپی از آن برداشت. با این حال سار در مقابل سرقت سوال‌هایش کاملاً آسیب‌پذیر است (کالتون، 1998؛ الیا و همکاران، 2000).

ج) مسائل مربوط به آزمودنی‌ها: اجرای نندگان سار باید راه حل‌های مناسبی برای تعدادی از مسائل فنی و عملی پیدا کنند. در یک سار که برای ارزیابی توانایی افراد مورد استفاده قرار می‌گیرد، آزمون‌کنندگان باید نسبت به مسائل بالقوه‌ای که ممکن است در حین اجرای یک سار برای آزمودنی‌ها اتفاق بیافتد، هشیار بوده و در این باره آینده‌نگری نمایند.

در نخستین دید اجمالی، تجربه شرکت در یک سار خیلی متفاوت از یک آزمون ستی به نظر نمی‌رسد. یک سئوال بر روی صفحه رایانه ظاهر می‌شود، آزمودنی پاسخ سئوال خود را وارد رایانه می‌کند، و پس از آن سئوال بعدی ظاهر می‌شود و این فرآیند تا آزمون تمام شود، ادامه می‌یابد. هر چند، چند جنبه منحصر به فرد وجود دارد که در تجربه سار عملکرد آزمودنی یا آزمودنی‌ها را تحت تأثیر قرار می‌دهد. نخست اینکه، سنچش مبتنی بر رایانه بسیاری از آزمودنی‌ها، ناآشنناست. سئوال‌هایی که روی صفحه رایانه ظاهر می‌شوند ممکن است برای آزمودنی‌ها از لحاظ خواندن دشوارتر و یا آسان‌تر باشند. لازم است سئوال‌های طولانی‌تر، یعنی سئوال‌هایی که در کل صفحه رایانه گسترده شده‌اند، توسط آزمودنی‌ها بر روی صفحه رایانه بالا یا پایین شوند. ورود پاسخ‌های آزمودنی‌ها با استفاده از صفحه کلید و یا موشواره از وارد کردن پاسخ در یک پاسخنامه دفترچه آزمون، یا سیاه کردن آن در پاسخنامه متفاوت است. دوم اینکه، در یک آزمون ستی آزمودنی‌ها معمولاً به تمام سئوال‌های یک آزمون پاسخ می‌دهند. این موضوع برای آزمودنی‌ها آزادی عمل زیادی در انتخاب سئوال‌ها می‌دهد؛ یعنی می‌توانند از بعضی از سئوال‌ها رد شوند تا آنها را بعداً پاسخ دهند. پاسخ سئوال‌ها را مرور کنند و احتمالاً تغییر دهند. بر عکس، آزمودنی‌های سار کنترل بسیار کمی بر این موضوع دارند؛ زیرا سئوال‌ها به صورت نوعی یک بار اجرا می‌شوند و در آن آزمودنی امکان مرور سئوال‌ها و پاسخ‌های خود را ندارند (سوتاریدونا و همکاران، 2003). در این قسمت ضمن بحث در باره دورنمای آزمودنی‌ها در حین اجرای سار، سه موضوع مربوط به آنها معرفی می‌شود.

(1) مرور سئوال‌ها: یکی از سئوال‌های بحث‌انگیز مرتبط با اجرای سار که باید در باره آن تصمیم‌گیری شود، آن است که آیا به آزمودنی‌ها، اجازه مرور و تغییر پاسخ‌هایی را که قبلاً به سئوال‌ها داده‌اند، داده شود (وای، 2005). پژوهشگران قدیمی حوزه سنچش انطباقی معتقد بودند که اجازه دادن به آزمودنی برای برگشت به پاسخ‌های خود و تغییر آنها اثرات منفی بر کارآیی اندازه‌گیری در سار دارد؛ اگر پاسخی تغییر یافت، سئوال‌های بعدی که در نتیجه این تغییر توسط رایانه انتخاب می‌شوند، طولانی‌تر و بهتر از دیگری نخواهد بود. همچنین برخی از پژوهشگران نگرانند که آزمودنی‌ها برای بهبود نمرات خود در سار از تقلب استفاده کنند. آزمودنی‌ها این کار می‌توانند در ابتدا به سئوال‌ها پاسخ غلط بدهند؛ در نتیجه رایانه به آنها ساده‌ترین سئوال‌های ممکن را عرضه می‌کند. پس از آن در فرصت مرور

سئوال‌ها، پاسخ‌های خود را از غلط به درست تغییر داده و نمرات بالاتری کسب کنند (واینر¹، 2000؛ 1983). به همین دلایل، برخی از برنامه‌های سار، امکان برگشت و مرور پاسخ‌ها وجود ندارد. مرور سئوال‌ها در سار به عنوان روشی برای افزایش کارآیی سنجش انطباقی مورد بررسی قرار گرفته است. مرور سئوال‌ها به زمان اضافی نیاز دارد، پاسخ‌های تغییر داده شده ممکن است خطای استاندارد برآورد توانایی آزمودنی را افزایش دهد. ممکن است آزمودنی‌ها به صورت راهبردی از مرور کردن سئوال‌ها برای افزایش نمرات خودشان استفاده نمایند. اگر چه، واکنش‌های آزمودنی‌ها به سار، مثبت بوده است، اما آنها نارضایتی شدیدی نسبت به فقدان مرور کردن سئوال‌ها ابراز نموده‌اند (باقی²، فرارا³، گابریز⁴، 1994؛ لگ⁵ و بوهر⁶، 1992؛ ویسپول⁷، رکلین⁸، و وانگ⁹، 1994).

اگرچه پژوهش‌های جدید نشان داده‌اند که اگر آزمودنی‌ها اجازه مرور پاسخ‌های خود و تجدید نظر در آنها را داشته باشند، اثر آشکاری بر کارآیی اندازه‌گیری ندارد، و راهبردهای دستیابی به نمرات بالاتر با تسلی به تقلب، صرفاً سودمندی ناچیزی برای آزمودنی‌ها در برخواهد داشت (استون¹⁰ و لونز¹¹، 1994؛ ویسپول¹²، هندریکسون¹³، و بلایر¹⁴، 2000). بر اساس نتایج پژوهش‌ها، به نظر می‌رسد دلیل روانسنجی مناسبی وجود ندارد که در سار به آزمودنی‌ها اجازه مرور و تجدید نظر در پاسخ‌هایشان داده نشود. با این وجود، اگر آزمودنی‌ها اجازه مرور و تجدید نظر در پاسخ‌هایشان را داشته باشند، بر پیچیدگی‌های برنامه سار افزوده می‌شود (وای، 2005).

(2) محدودیت‌های زمانی: در نظر گرفتن محدودیت زمانی معقول و منطقی برای آزمون‌های استاندارد شده قدیمی چالش برانگیز است. اگر زمان آزمون بسیار طولانی شود، نتیجه آن از دست رفتن کارآیی آزمون است. اگر زمان آزمون بسیار کوتاه باشد،

1. Wainer

2. Baghi

3. Ferrara

4. Gabrys

5. Legg

6. Buhar

7. Vispocklin

8. Roklin

9. Wang

10. Stone

11. Lunz

12. Vispoel

13. Hendrickson

14. Bleirler

برخی از آزمودنی‌ها نمی‌توانند تمام سوال‌های یک آزمون را پاسخ گویند و در نتیجه عملکرد آنها در آزمون، پایین‌تر از سطح عملکرد استاندارد، برآورده می‌شود. اگرچه، درنظر گرفتن محدودیت زمانی در سار پیچیده‌تر است. یک دلیل آن است که در سار از نمره دقیق به عنوان ملاکی برای پایان دادن به آزمون استفاده می‌شود، و به همان میزان که یک فرد نمی‌داند چه تعدادی سوال را پاسخ خواهد داد، نمی‌داند که چقدر زمان برای پاسخ دادن به سوال‌ها دارد. تصمیم‌گیری راجع به محدودیت فرصت‌های زمانی لازم در یک سار موضوعی بسیار مهم است. دلیل این موضوع آن است که تحمیل هر نوع محدودیت زمانی متضاد با برنامهٔ سنچش است که به دانش‌آموزان کمک می‌کند تا سطوح بهینه‌ای از عملکرد خود نشان دهند. محدودیت زمان همچنین می‌تواند دشواری‌هایی در مدرج‌سازی سوال‌های آزمون به وجود آورد (اگن، 2007؛ وای، 1998؛ الیا و همکاران، 2000).

تحقیقات مرتبط

پژوهش‌هایی که در بارهٔ اثر مرور کردن سوال‌ها بر عملکرد آزمودنی‌ها در سنچش انطباقی صورت گرفته‌اند، نشانگر آن است که: (الف) زمانی که آزمودنی‌ها حق تعییر دادن و مرور کردن سوال‌ها و پاسخ‌های خود را دارند، عملکرد آنها افزایش می‌یابد؛ (ب) پاداش نمره‌ای که در نتیجهٔ تعییر پاسخ‌ها به دست می‌آید دانش‌آموزان را تشویق می‌کند تا در بارهٔ سوال‌ها دوباره به تفکر پرداخته و یا دوباره آنها را بخوانند. بنابراین، از این موضوع می‌توان نتیجهٔ گیری کرد که ممانعت از مرور سوال‌ها، از بین بردن فرصت تعییر پاسخ‌هایی است که عملکرد آزمون را بهبود می‌بخشد. (ج) نبود فرصت برای مرور پاسخ‌ها اضطراب امتحان را افزایش می‌دهد. بدین معنا که نبود فرصت برای مرور سوال‌ها به معنای فقدان کنترل آزمودنی بر محیط و آزمون است و این مطلب منجر به افزایش اضطراب امتحان می‌شود. همچنین مشخص شده که آزمودنی‌ها زمانی استرس‌های محیطی را بهتر تحمل می‌کنند که، احساس کنند بر محیط کنترل بیشتری دارند. لذا افزایش کنترل بر محیط منجر به کاهش اضطراب و بهبود عملکرد در تکلیف می‌شود (اگن، 2007؛ وای، 1998).

در یک پژوهش دو نسخه از آزمون خزانه لغات انگلیسی برای اسپانیایی زبان‌ها (یک نسخهٔ انطباقی و نسخهٔ دیگر غیرانطباقی) در نمونه‌ای از دانشجویان اسپانیایی سال اول دانشگاه مورد استفاده قرار گرفت. اثرات نوع آزمون (آزمون انطباقی با رایانه

در برابر آزمون غیرانطباقی با رایانه) و شرایط مرور کردن (کسانی که اجازه مرور کردن سؤال و پاسخ را داشتند در مقابل کسانی که این اجازه را نداشتند) بر روی چند متغیر روان‌شناختی آزمایش شد. متغیرهای بین آزمودنی‌ها قبل و پس از مرور کردن به منظور مطالعه اثرات مرور کردن بر روی متغیرهای روان‌شناختی و روان‌سنجدی برای شرایط مرور کردن در هر دو آزمون اندازه‌گیری شدند. دو دسته اصلی نتایج پس از مرور کردن بدست آمد:

- الف) تعداد پاسخ‌های درست و توانایی برآورد شده افزایش معنی‌داری داشت؛
- ب) میزان اضطراب موقعیتی کاهش یافت؛ ج) تفاوت خطاهای اندازه‌گیری معنی‌دار نبود. تعامل اثرات (نوع آزمون با آزادی عمل در مرور کردن) معنی‌دار نبود. پیشنهاد ضمنی این نتایج، آن است که باید در تهیه آزمون‌های انطباقی با رایانه شرایط مرور سؤال‌ها و پاسخ‌ها را برای آزمون شوندگان فراهم نمود (ولیا و همکاران، 2000).
- پژوهش‌ها نشان داده‌اند که اگر آزمودنی‌ها اجازه مرور سؤال‌ها و پاسخ داشته باشند، الف) حدود 60 درصد آزمودنی‌ها حداقل یکی از پاسخ‌های خود را تغییر می‌دهند؛ ب) فقط درصد کوچکی (بین 2 تا 5 درصد) از پاسخ‌ها تغییر داده می‌شوند؛ ج) از بین پاسخ‌های تغییر داده شد، حدود 50 درصد از آنها از غلط به درست تغییر می‌یابند؛ د) از بین آزمودنی‌هایی که پاسخ خود را تغییر می‌دهند، بین 42 تا 52 درصد از آنها سطح توانایی خود را در آزمون بهبود می‌دهند. فقط درصد کوچکی از این آزمودنی‌ها (بین 10 تا 15 درصد) سطح توانایی خود را کاهش می‌دهند؛ ه) میزان ناچیزی از دقت از بین می‌رود؛ و) همبستگی بین توانایی برآورد شده قبل و پس از مرور بیشتر از 0/98 است و میانگین تفاوت‌ها بین 0/07 تا 0/20 است؛ ز) همبستگی منفی و معنی‌دار بین سطح اضطراب و توانایی وجود دارد. تعامل بین اضطراب و مرور معنی‌دار نبود؛ ح) اثر معنی‌داری در زمان آزمون دارد (مرور، زمان آزمون را بین 10 تا 37 درصد افزایش می‌دهد)؛ ط) زمانی که توانایی به عنوان متغیر مستقل در نظر گرفته شود، مشخص شده است کسانی که از بالاترین سطح توانایی برخوردارند، حداقل تغییرات را در پاسخ‌های خود اعمال می‌کنند و این تغییرات اکثراً از غلط به درست و تعداد کمی از آنها از درست به غلط صورت می‌گیرد. بنابراین کسانی که از بالاترین سطح توانایی برخوردارند از بیشترین مزیت مرور بهره‌مند می‌شوند (ولیا و همکاران، 2000).

در مطالعه‌ای اولیا و همکاران (2000) نشان دادند که 81/52 درصد آزمودنی‌ها از فرصت برای مرور پاسخ‌ها استفاده کردند. در سنچش انطباقی با رایانه 80 درصد آزمودنی‌ها پاسخ‌هایشان را مرور کردند و از بین آنها 66/7 درصد توانایی برآورده شده خود را افزایش دادند، 25 درصد کاهش دادند، و 8/3 درصد تغییری در توانایی برآورد شده خود به وجود نیاوردند. در کل، 13/5 درصد پاسخ‌ها تغییر کرد. در سار 12/6 درصد پاسخ‌ها تغییر کرد و از بین آنها 42 درصد از غلط به غلط تغییر یافت، 43 درصد از غلط به درست، و 15 درصد از درست به غلط.

در ارتباط با محدودیت زمانی؛ هدف عملی آن است که محدودیت زمانی به گونه‌ای مشخص شود که محدودیت اثر معنی‌دارای بر عملکرد دانش‌آموزان نداشته باشد. این موضوع با یافته‌های برخی گروه‌های اقلیت نژادی که زمان بیشتری برای پاسخگویی قایل شده‌اند، پیچیده‌تر شده است (یاقی و همکاران، 1992؛ لگ¹ و بوهر²، بوهر²، 1992؛ زارا³، 1992). اما برخی پژوهش‌ها نشان داده‌اند که دادن زمان پاسخگویی بیشتر به دانش‌آموزان گروه‌های اقلیت در آزمون‌های سنتی عملکرد مرتبط با اکثریت دانش‌آموزان را در آنها افزایش نمی‌دهد. به نظر می‌رسد ارتباط بین محدودیت زمانی و عملکرد در آزمون با اضطراب امتحان بین آزمون‌هایی با محدودیت زمان و بدون محدودیت زمان نشان داده‌اند که این قضاوت‌ها برای آزمودنی‌های با اضطراب بالا بیشتر است (هیلی⁴، 1984؛ سیمن⁵ و اونزیگ بوزی⁶، 1995؛ به نقل، الیا و همکاران، 2000). این یافته‌ها پیشنهاد می‌کند که طولانی کردن محدودیت زمان در یک سار ممکن است برای برخی آزمودنی‌ها مفیدتر از سایر آزمودنی‌ها باشد. تفاوت‌های فردی پیدا شده در بین آزمودنی‌ها، مشخص می‌کند که احتمالاً مشکل است تنها محدودیت زمانی را به عنوان دفاعی منصفانه تلقی کرد. بنابراین، می‌بایست محدودیت‌های زمانی خیلی آزاد منشانه‌تری در نظر گرفت، یا اصلاً محدودیت زمانی قایل نشویم. این نکته را به خاطر داشته باشید که برنامه سار خیلی کوتاه‌تر از آزمون‌های سنتی همتای آن است؛ می‌بایست در دادن وقت برای بازگشت آزمودنی‌ها دقت نمود. در آن صورت استرس وابسته به امتحان در نتیجه آن،

1. Legg

2. Buher

3. Zara

4. Hilly

5. Seaman

6. Onzuegbuzie

کاهش یافته و روایی آزمون افزایش می‌یابد. با گسترش برنامه‌های آزمون در ایالت متحده آمریکا شواهدی وجود دارد که کودکان گروه‌های اقلیت و فقیر دسترس کمتری به رایانه در منزل و مدرسه دارند (ساتون¹، 1997). بخاطر اینکه دسترسی کمتر، دلالت بر تجربه کمتر دارد، ارتباط بین تجربه رایانه و عملکرد در سار اهمیت بیشتری پیدا می‌کند.

نتیجه‌گیری

۱. برآورد نمره واقعی آزمودنی‌ها از طریق سنجش انطباقی، دقیق‌تر است. بنابراین، استفاده از این روش هم منطقی است و هم منصفانه. موضوعی که باید در آزمون‌سازی مؤسسات آموزشی کشور مثل دانشگاه‌ها و سازمان سنجش آموزش کشور بیشتر مورد توجه قرار گرفته و درباره آن پژوهش‌های لازم و ضروری انجام گیرد.

۲. سنجش انطباقی به آزمون‌کننده اجازه می‌دهد که سطح ویژگی آزمون‌ها را به سرعت بسنجد، بدون اینکه او را وادار کند تا به مجموعه‌ای خسته‌کننده از سوال‌های بسیار ساده یا دسته‌ای ناخوشاپایند و ناراحت‌کننده از پرسش‌های خیلی مشکل پاسخ دهد. بنابراین چنین روشی هم برای آزمودنی‌ها و هم برای آزمون‌کننده و آزمون‌ساز مطلوب است و باید مورد توجه مؤسسات آموزشی، بخصوص آموزش عالی کشور قرار گیرد.

۳. سنجش انطباقی با رایانه که تکرار فرآیند دو مرحله‌ای در سنجش است یکی از روش‌های پیچیده‌تر سنجش است که در آن، در گام نخست سوالی به آزمودنی داده می‌شود که دشواری آن با توانایی فعلی برآورد شده برای او برابر باشد. در گام دوم توانایی آزمودنی براساس اطلاعات حاصل از پاسخ او به سوال مرحله اول برآورد شده؛ این گام‌ها تکرار می‌شوند تا جایی که ملاک‌های پایان دادن به آزمون محقق شود. با لحاظ این مطلب که در دانشگاه‌های کشور امکانات سخت‌افزاری لازم برای اجرای سنجش انطباقی با رایانه ایجاد شده است، این موضوع می‌تواند مورد توجه متخصصان آزمون‌سازی قرار گرفته و از این شیوه که هم کارآیی بیشتری دارد و هم بسیاری از سرخورده‌گی‌های دانشجویان در ارتباط با مشکلات آزمون‌های سنتی می‌کاهد، در دانشگاه‌ها استفاده شود.

4. از آنجا که برنامه‌ریزی واقعی، اجرای عملی و حفاظت از برنامه سنچش انطباقی با رایانه بسیار پیچیده است و چالش‌های عملی خاص خود را دارد، و از طرف دیگر ویژگی‌های فرهنگی، اجتماعی، و اقتصادی جامعه ایرانی در این پیچیدگی‌ها و چالش‌ها مؤثرند. باید پژوهش‌های گستردۀ‌ای در ارتباط با پیچیدگی‌های کاربردی آن صورت گرفته و با لحاظ نتایج و یافته‌های پژوهش‌های مذکور، از این روش در عمل استفاده شود.

5. مطالعاتی به منظور امکان‌سنجدی استفاده از سار در سازمان سنچش آموزش کشور و مؤسسات آموزش عالی انجام شود و براساس یافته‌های حاصل از آنها در بارهٔ استفاده از سار و سایر جوانب آن تصمیم‌گیری شود.

منابع

- افروز، غلامعلی و هومن، حیدرعلی (1375). روش تهیه آزمون هوش. تهران: دانشگاه تهران.
- آلن، مری جی وین، وندی ام (1979). مقدمه‌ای بر نظریه‌های اندازه‌گیری (روانسنجی). ترجمه علی دلاور، 1374، تهران: سمت.
- ژندایک، رابرتس ال (1982). روانسنجی کاربردی. ترجمه حیدر علی هومن، 1369، تهران: دانشگاه تهران.
- ستاری، بهزاد (1382). روانسنجی پیشرفته کاربردی. مشهد: به نشر.

- Baghi, H.; Ferrara, S. F., & Gabrys, R. (1992). *Student attitudes toward computer-adaptive test administrations*. Paper presented at the annual meeting of the American Educational Research Assoc, San Francisco, CA.
- Colton, G. D. (1998). Exam security and high-tech cheating, *The Bar Examiner*, 67(3), 13-35.
- Eggen, T. J. H. M. (1990). Innovative procedures in the calibration of measurements scales. In W.H. Schreiber & K. Ingenkamp (eds). (pp.199- 212). *International developments in large scale assessment*. Windsor, Berkshire: NFER-NELSON
- Eggen, T. J. H. M. (2007). Choices in CAT models in the context of educational testing. In D. J. Weiss (Ed.). *Proceedings of the 2007 GMAC Conference on Computerized Adaptive Testing*. Retrieved [date] from www.psych.umn.edu/psylabs/CATCentral/
- Hornke L. F. (2000). Item response times in computerized adaptive testing. *Psicologica*, 21, 157-173.
- Legg, S. M. & Buhr, D. C. (1992). Computerized adaptive testing with different groups. *Educational Measurement: Issues and practice*, 11(2), 23-27.
- Lord, F.M. (1983). Some test theory for tailored testing. In W.H. Holtzman (Ed.), *Computer-assisted instruction, testing and guidance*. New York: Harper & Row.
- Olea J., J., Revuelta J., Ximenez M.C., & Abad F.H. (2000). Psychometric and psychological effects of review on computerized fixed and adaptive tests. *Psicologica*, 21, 175-189.
- Sotaridonia L. S., Pornel J. B., & Vallejo A. (2003). Some applications of item response theory to testing. *The Philippine Statistician*, 52, 81-92.
- Stone, G. E., & Lunz, M. E. (1994). The effect of review on the psychometric characteristics of computerized adaptive tests. *Applied Measurement in Education*, 7, 211-222.

- Stone, G. E., & Lunz, M. E. (1994). The effect of review on the psychometric characteristics of computerized adaptive tests. *Applied Measurement in Education*, 7, 211-222.
- Sutton, R. E. (1997). Equity and high stakes testing: Implications for computerized testing. *Equity and Excellence in Education*, 30(1), 5-15.
- Tao, Y.-H., Wu, Y.-L., & Chang, H.-Y. (2008). A Practical Computer Adaptive Testing Model for Small-Scale Scenarios. *Educational Technology & Society*, 11(3), 259-274.
- Triantafillou, E., Georgiadou, E., & Anastasias A. (2006). CAT-MD: Computer Adaptive Test on Mobile Devices. Economides University of Macedonia, Egnatias 156, Thessaloniki 54006, GREECE
- Van der Linden, W. J. & Glas, C. A. W. (Eds.) (2000). *Computerized adaptive testing: Theory and practice*. Dordrecht: Kluwer Academic Publishers.
- Van der Linden, W. J. & Hambleton, R. K. (Eds.) (1996). *Handbook of modern item response theory*. New-York: Springer-Verlag.
- Verhelst, N. D., & Glas, C. A. W. (1995). The one-parameter logistic model. In. G. H. Fischer & I. W. Molenaar (Eds.). *Rasch models: Foundations, recent developments, and applications* (pp.215-237). New York: Springer-Verlag.
- Vispoel, W. P., Hendrickson, A. B., & Bleiler, T. (2000). Limiting answer review and change on computer adaptive vocabulary tests: Psychometric and attitudinal results. *Journal of Educational Measurement*, 37, 21-38.
- Vispoel, W. P., Rocklin, T. R., & Wang, T. (1994). Individual differences and test administration procedures: A comparison of fixed-item, computerized adaptive, and self-adapted testing. *Applied Measurement in Education*, 7, 53-59.
- Wainer, H. (1983). Some practical considerations when converting a linearly administered test to an adaptive format. *Educational Measurement: Issues and Practice*, 12, 15-20.
- Wainer, H., (Ed.) (2000). *Computerized adaptive testing: A primer* (2nd Edition). Hillsdale, NJ: Erlbaum.
- Way, W. D. (1998). Protecting the integrity of computerized testing item pools. *Educational Measurement: Issues and Practices*, 17 (4), 17-27.
- Way, W. D. (2005). *Practical Questions in Introducing Computerized Adaptive Testing for K-12 Assessments*. Research Report 05-03.
- Zara, A. R. (1992). *An investigation of computerized adaptive testing for demographically-diverse candidates on the national registered licensure examination*. Paper presented at the annual meeting of the National Council on Measurement in Education, San Francisco, CA.