



تدوین و هنجاریابی آزمون پیشرفت تحصیلی زیست‌شناسی

حسن مشتاقیان ابرقویی*^۱، بهرام جوکار^۲

تاریخ دریافت: ۱۴۰۱/۰۱/۱۷ تاریخ پذیرش: ۱۴۰۱/۰۵/۱۵

از صفحه ۲۹ تا ۴۶

چکیده:

هدف پژوهش حاضر، تدوین و هنجاریابی آزمون زیست‌شناسی پایه دهم دوره متوسطه به منظور اندازه‌گیری سطح آموخته‌ها و تعیین پراکندگی پیشرفت تحصیلی دانش‌آموزان در این ماده درسی بود. برای نیل به هدف کلی تحقیق، هر دو مدل کلاسیک آزمون و نظریه سؤال-پاسخ مورد استفاده قرار گرفت. ابزار اولیه پژوهش شامل ۱۵۰ سؤال چهارگزینه‌ای از محتوای زیست‌شناسی پایه دهم رشته تجربی بود که بر روی نمونه‌ای به حجم ۳۰۰ دانش‌آموز دختر و پسر اجرا شد. به دنبال تجزیه و تحلیل سؤال‌های آزمون فرم‌های نهایی آزمون تدوین شد. ابزار نهایی دو فرم موازی ۵۰ سؤالی بود که بر روی نمونه هنجاری به حجم ۹۳۸ دانش‌آموز دختر و پسر شهرستان شیراز اجرا گردید. ضریب پایایی برآورد شده برای ثبات درونی فرم‌های آزمون به ترتیب ۰/۸۹ و ۰/۸۸ به دست آمد. بر اساس تحلیل عاملی، هر دو فرم آزمون از یک عامل کلی اشیاع بودند. نتایج نشان داد که بین میانگین نمرات پسران و دختران تفاوت معناداری وجود ندارد؛ بنابراین هنجارهای استاندارد و هم‌صدک برای همه آزمودنی‌ها محاسبه شد. یافته‌های حاصل از تجزیه و تحلیل نظریه سؤال-پاسخ نشان داد که بیش از ۹۲ درصد از سؤال‌ها به‌طور قابل توجهی با مدل لجستیک سه پارامتری برازش دارند. تابع آگاهی برآورد شده آزمون یک منحنی زنگوله‌ای شکل بود و در محدوده توانایی ۰/۵- تا ۲/۵+ بیشترین اطلاعات را ارائه می‌داد؛ همچنین حداکثر میزان آگاهی دهندگی فرم‌های آزمون در سطح توانایی ۱/۵+ به دست آمد که نشان می‌دهد آزمون تدوین شده تخمین دقیق‌تری از نمره واقعی آزمودنی‌هایی ارائه می‌دهد که سطح توانایی آن‌ها بالاتر از متوسط پیوستار توانایی است.

کلمات کلیدی: آزمون پیشرفت زیست‌شناسی، مدل کلاسیک اندازه‌گیری، نظریه سؤال-پاسخ، هنجاریابی.

۱. دکتری سنجش و اندازه‌گیری؛ دبیر زیست‌شناسی دبیرستان باقرالعلوم، آموزش و پرورش ناحیه ۴ شیراز * kavir311@gmail.com

۲. استاد روانشناسی تربیتی، عضو هیات علمی دانشکده روانشناسی و علوم تربیتی دانشگاه شیراز.

مقدمه و بیان مسئله

از امتحانات غیررسمی معلمان گرفته تا آزمون‌های استاندارد شده که در سطح ملی اجرا می‌شود، سنجش^۱ برای مدت‌های طولانی بخش جدایی‌ناپذیر از فرآیند آموزشی بوده است (پلگرینو، چودوفسکی و گلیرز^۲، ۲۰۰۱: ۱۹). بر طبق استانداردهای آزمون‌های آموزشی و روان‌شناختی^۳، سنجش را می‌توان «هرگونه روش سیستماتیک به دست آوردن اطلاعات از آزمون‌ها و دیگر منابع که به منظور استنباط درباره ویژگی‌های افراد، اشیا یا برنامه‌ها به کار می‌رود، تعریف کرد» (سیمون، ارسیکان و روسو^۴، ۲۰۱۳: ۱۷۲). به‌طور کلی هدف از سنجش کمک به یادگیری، اندازه‌گیری پیشرفت فردی دانش‌آموزان و ارزیابی برنامه‌ها است (پلگرینو و همکاران، ۲۰۰۱: ۳۷). اگرچه ظاهراً چنین تصور می‌شود که سنجش مرحله پایانی فعالیت‌های آموزشی معلم است؛ اما واقعیت کنونی این است که سنجش و اندازه‌گیری، اغلب تعیین‌کننده نوع فعالیت‌های آموزشی معلم و یادگیری دانش‌آموزان است و می‌تواند تأثیر بسزایی در فرآیند یاددهی-یادگیری داشته باشد (سیزک^۵، ۱۹۹۳). کاملاً روشن است که سنجش یک بخش جدا شده از سیستم آموزشی نیست. آنچه اندازه گرفته می‌شود و اینکه اطلاعات به چه شکل استفاده شود تا حد زیادی به برنامه درسی و روش‌های آموزشی مورد استفاده وابسته است. با وجود این، سنجش تأثیر قوی بر هر دو برنامه درسی و شیوه‌های آموزشی دارد (لین^۶، ۲۰۰۰؛ پلگرینو و همکاران، ۲۰۰۱). نتایج حاصل از چندین مطالعه عمده درباره تأثیر آزمون^۷ روی آموزش و برنامه درسی، نشان داده است اطلاعاتی که بر اساس آن آزمون‌ها ساخته می‌شوند، تأثیر معناداری نه تنها بر روی تدریس، بلکه بر روی برنامه آموزشی و انگیزش دانش‌آموزان دارد (سالمون-کاکس^۸، ۱۹۸۱؛ کروکس^۹، ۱۹۸۸؛ فلیپس^{۱۰}، ۲۰۱۲). از نظر بلک و همکاران^{۱۱} (۲۰۰۳) فعالیت‌های سنجش اگر اطلاعاتی به عنوان بازخورد برای معلمان و دانش‌آموزان فراهم کند، می‌تواند به اصلاح فرآیندهای یاددهی-یادگیری بینجامد. آن‌ها به معلمان، دانش‌آموزان و والدین آن‌ها در تعیین سطح و کیفیت یادگیری دانش‌آموزان کمک می‌کنند. آن‌ها به معلمان کمک می‌کنند تا درک کنند که چگونه آموزش را بر اساس شواهدی از یادگیری دانش‌آموزان منطبق سازند. همچنین به مدیران و سرپرستان کمک می‌کند تا سطح پیشرفت فردی دانش‌آموزان را در کلاس، مدرسه یا منطقه تعیین کنند؛ و در آخر آنکه سنجش به سیاست‌گذاران و عموم مردم کمک می‌کند تا به ارزیابی تأثیر سیستم‌های آموزشی بپردازند (پلگرینو و همکاران، ۲۰۰۱: ۱۹).

1. Assessment

2. Pellegrino, Chudowsky & Glaser

3. The Standards for Educational and Psychological Testing

4. Simon, Ericikan & Rousseau

5. Cizek

6. Linn

7. testing

8. Salmon-Cox

9. Crooks

10. Phelps

11. Black et al



علی‌رغم نقش سنجش در جهت‌بخشی به فعالیت‌های یاددهی-یادگیری، امروزه با انتقادهای بسیاری روبرو شده است. برخی مواقع این انتقادهای به ابزارهای سنجش بوده است؛ اما بیشتر انتقادهای به سنجش، نه ناشی از خود آزمون‌ها، بلکه استفاده‌های ثانویه‌ای است که از نتایج این آزمون‌ها صورت می‌گیرد. زمانی که نتایج آزمون‌ها و امتحانات به‌عنوان اندازه‌های پاسخگویی معیار قضاوت، پاداش و تنبیه معلمان، کارایی مدارس و برنامه‌ها به کار می‌رود، موجب پیامدهای ناخواسته بی‌شماری می‌شود. برای نمونه، معلمان ممکن است با هزینه کردن دیگر اهداف آموزشی، وقت زیادی را صرف آماده‌سازی دانش آموزان برای امتحان کنند؛ ضمن آنکه ممکن است تقلب رادر بین افراد رواج دهد (کورتز، ۲۰۰۸). همچنین استفاده از این نوع سنجش، موجب کم‌سواد شدن دانش آموزان و یادگیری سطحی می‌شود (سیزک، ۱۹۹۳).

به‌رحال صرف‌نظر از جنبه مثبت یا منفی آزمون‌ها، همان‌طور که جیپس^{۱۲} (۲۰۰۳) اظهار می‌کند؛ این یک حقیقت است که آزمون تأثیر قدرتمندی بر شیوه‌ای که معلمان وظیفه‌شان را به اجرا درمی‌آورند دارد؛ به‌طوری‌که امروزه تأثیر سنجش و اندازه‌گیری بر جریان آموزش معلم، به‌پدید آمدن مفهومی به نام «آموزش جهت داده‌شده به‌وسیله اندازه‌گیری»^{۱۳} انجامیده است. منظور از آموزش جهت داده‌شده به‌وسیله اندازه‌گیری، یک رویکرد کلی نسبت به آموزش است که در آن هدف اصلی معلم بالا بردن سطح عملکرد دانش آموزان در یک روش خاص اندازه‌گیری است (لفرانسوا، ۲۰۰۰؛ نقل از سیف؛ ۱۳۸۴). پوفام^{۱۴} (۱۹۸۷) به‌عنوان مدافع «آموزش جهت داده‌شده به‌وسیله اندازه‌گیری»، آن را به‌عنوان مقرون به‌صرفه‌ترین روش بهبود کیفیت آموزش عمومی قلمداد می‌نماید. از نظر پوفام، آموزش جهت داده‌شده توسط اندازه‌گیری، می‌تواند یک نیروی بسیار مؤثر برای بهبود آموزشی باشد و با مراقبت‌های روشن می‌تواند از اثرات منفی بالقوه آن، از جمله تحریف برنامه‌های آموزشی جلوگیری کرد (جیپس، ۲۰۰۳: ۳۱-۳۳).

از جمله ابزارهای سنجش و اندازه‌گیری آموزشی، می‌توان به آزمون‌های پیشرفت اشاره کرد. آزمون پیشرفت، یک روش سیستماتیک برای تعیین میزان آموخته‌های دانش آموزان و اندازه‌گیری پرونداد یادگیری است. به کمک آزمون‌های پیشرفت می‌توان موفقیت دانش آموزان را در طول زمان پیگیری نمود، موفقیت تحصیلی گروه‌های مختلف (کلاس‌ها/مدارس) را مقایسه نمود، نتایج آن را در تصمیم‌گیری‌های خطیر مورد استفاده قرارداد، نقاط قوت و ضعف دانش آموزان را اندازه‌گیری کرد و اثربخشی برنامه آموزشی را معلوم نمود (ایساکس، زارا، هربرت، کومبس و اسمیت^{۱۵}، ۲۰۱۳). از نظر گرونلاند استفاده از این آزمون‌ها نه تنها موجب بهبود تصمیمات آموزشی می‌شوند بلکه از طریق: ۱) بهبود انگیزش دانش آموز، ۲) افزایش نگهداری و انتقال آموخته‌ها، ۳) افزایش خود فهمی دانش آموز و ۴) فراهم نمودن بازخورد در ارتباط با کارایی آموزشی، می‌تواند به یادگیری کمک کند (گرونلاند^{۱۶}، ۱۹۸۸: ۵). با این وجود، حدود تأثیرگذاری و مشارکت آزمون‌های پیشرفت تحصیلی در بهبود یادگیری و آموزش، به میزان زیادی به اصول زیربنایی حاکم بر ساخت و استفاده از آن‌ها بستگی دارد. به عبارت دیگر، آزمون‌های پیشرفت می‌تواند توجه دانش آموزان را به اهداف

12. Gipps

13. Measurement-Driven Instruction

14. Popham

15. Smith

16. Gronlund

اساسی آموزش هدایت کند یا آنان را از آن دور نماید. موجب شود تا دانش آموزان بر جنبه‌های محدود از محتوای دوره تمرکز کنند یا توجه آن‌ها را به حیطه‌های مهم آن معطوف کند. یادگیری سطحی را پاداش دهند یا درک عمیق را تکلیف کنند؛ و در امر تصمیم‌گیری آموزشی؛ اطلاعاتی قابل اعتماد و یا اطلاعات سودار و مخدوش فراهم کند (همان: ۷).

علی‌رغم کاربرد وسیع آزمون‌های پیشرفت و اهمیت آن‌ها در ارزیابی و راهنمایی دانش آموزان، بسیاری از معلمان در زمینه ساخت آزمون‌های مناسب، آموزش ندیده‌اند یا آموزش اندکی دارند (گرونلاند؛ ۱۹۸۸). نتایج برخی مطالعات نشان می‌دهد که اکثر آزمون‌های معلم ساخته، فاقد اعتبار^{۱۷} و پایایی^{۱۸} است و به نظر می‌رسد بسیاری از معلمان فاقد مهارت‌های ساخت آزمون هستند (آل - ویلیامز^{۱۹}، ۲۰۰۲؛ بیکر^{۲۰}، ۲۰۰۳). در این رابطه پوفام اظهار می‌کند: آزمون‌های معلم ساخته، بیشتر بر سطوح پایین‌تر مهارت‌های شناختی تأکید دارند و عمدتاً به یادآوری مطالب قبلاً آموخته شده می‌پردازند (۱۹۸۷). به نظر می‌رسد که بیشتر سؤال‌های آزمون‌های معلم ساخته تنها نیاز به یادآوری اطلاعات بسیار جزئی و خاص دارد و همچنین آموخته‌های دانش آموزان بیش از آنکه مبتنی بر مهارت‌های پیچیده و بالاتر شناختی باشد شامل محفوظات و یادآوری مطالب جزئی و کم‌اهمیت است (هومن، ۱۳۷۲).

مطالعات تأیید می‌کنند که سبک سنجش بر روی سبک یادگیری دانش آموزان تأثیر خواهد گذاشت. کراکس (۱۹۸۸) نشان داد که دانش‌آموزانی که رویکرد سطحی را برای یادگیری به‌کار می‌برند، سازگاری کمتری با ارزیابی‌هایی داشته‌اند که رویکرد یادگیری عمیق را بکار می‌گیرند؛ از طرف دیگر دانش‌آموزانی که می‌توانند با موفقیت رویکرد عمیق را در یادگیری به‌کار برند به‌سادگی می‌توانند خود را با رویکرد سطحی سازگار نمایند. وی نتیجه می‌گیرد که توسعه استراتژی یادگیری عمیق از طریق تأکید بیشتر بر سؤالات سطوح بالاتر حوزه شناختی در ارزیابی حاصل می‌شود. شکی نیست که می‌توان دانش آموزان را به‌جای یادگیری طوطی‌وار به‌سوی یادگیری فرآیندهای بالاتر شناختی که مستلزم درک فعال مفاهیم و اصول بنیادی و حل مسئله است هدایت کرد. یکی از راه‌های تحقق این امر استفاده از بازخوردهای مناسب حاصل از به‌کارگیری شیوه‌های مطلوب سنجش و ارزیابی است. نتایج حاصل از یک تحلیل چندجانبه که توسط فیلیدو و روسو^{۲۱} انجام شده نشان داده است که «استفاده مکرر از سؤال‌های سطح بالای شناختی ضمن آموزش، بر پیشرفت دانش آموزان تأثیر مثبت دارد» (گیج و برلایندر/ ترجمه خوئی نژاد، ۱۳۷۴: ۷۴۲). فردریکسون^{۲۲}، با توجه به تأثیر آزمون‌ها معتقد است: یک وظیفه مهم مریبان و روان‌شناسان، ساخت و توسعه ابزارهایی است که به‌خوبی منعکس‌کننده طیف وسیعی از هدف‌های آموزشی است و کاربرد شیوه‌هایی برای استفاده از آن ابزارها در بهبود فرآیند آموزشی است (جیپس، ۲۰۰۳: ۳۲).

یکی از انواع آزمون‌های پیشرفت تحصیلی، آزمون‌های دروس ویژه می‌باشند. این نوع آزمون‌ها برای اندازه‌گیری آموخته‌های دانش آموزان در یک درس و یا برنامه آموزشی ویژه تهیه و استاندارد می‌شوند و ابزار مفیدی برای ارزشیابی سطح آموخته‌ها و نقاط قوت و ضعف دانش آموزان کلاس در مقایسه با

17. validity

18. reliability

19. Alele-Williams

20. Baker

21. Filido & Roso

22. Frederiksen



سطح محلی؛ منطقه‌ای یا ملی به شمار می‌آیند. (گلاور و برونینگ / ترجمه خرازی، ۱۳۷۵) تحلیل پرسش‌های این گونه آزمون‌ها به معلمان کمک می‌کند تا هدف‌های تدریس خود را بهتر بشناسند و یا آن دسته هدف‌هایی را که در عمل کمتر مورد توجه آن‌ها بوده، مدنظر قرار دهند؛ فرآیند تدریس خود را اصلاح کنند و انتظارات خود را از دانش‌آموزان روشن نمایند؛ یادگیری دانش‌آموزان را تسهیل و توجه آن‌ها را به برنامه درسی متمرکز نمایند (دفتر همکاری‌های علمی بین‌المللی وزارت آموزش و پرورش، ۱۳۷۹). از جمله برنامه‌های درسی ویژه دوره متوسطه در نظام آموزش و پرورش ایران، زیست‌شناسی است. هدف غایی از این برنامه درسی، رشد همه‌جانبه فرد در ابعاد زیستی و تعامل آن با سایر ابعاد رشد و تسلط دانش‌آموزان بر چارچوب مفهومی و شیوه‌های پژوهش در قلمرو علم زیست‌شناسی است (سازمان پژوهش و برنامه‌ریزی آموزشی، ۱۳۸۹: ۶).

هرچند آموزش زیست‌شناسی می‌تواند نقش مهمی در رشد تفکر علمی و روحیه حقیقت‌جوئی در دانش‌آموزان ایفا کند؛ به نظر می‌رسد ضعف دانش‌آموزان در این ماده درسی جدی است. بر طبق یافته‌های بین‌المللی نتایج مطالعه تیمز^{۲۳} ۲۰۱۱، دانش‌آموزان ایرانی پایه هشتم در درس زیست‌شناسی در بین ۴۲ کشور مورد مطالعه، رتبه ۲۲ را به خود اختصاص داده‌اند. متوسط نمره مقیاس زیست‌شناسی پایه هشتم این دانش‌آموزان ۴۶۶ بوده است که از میانگین کلی (۴۷۴) به‌طور معنی‌داری کمتر است. همچنین از نظر نمره مقیاس محتوای آزمون، دانش‌آموزان ایرانی پایه هشتم بالاترین میانگین را در حوزه شناختی دانش (۴۷۹) و کمترین میانگین را در حوزه شناختی کاربرد (۴۷۰) به دست آورده‌اند که هر دو به‌طور معنی‌داری متفاوت از میانگین سه حوزه (دانش، کاربرد، استدلال) (۴۷۴) است. شواهد داخلی نیز نشانگر چنین کاستی‌هایی است. برای نمونه، میانگین نمرات امتحانات هماهنگ کشوری در سال ۱۳۹۴ برای دانش‌آموزان رشته تجربی دختر ۱۲/۷۱ و پسر ۱۱/۴۸ بوده است؛ و نزدیک به ۱۵ درصد دانش‌آموزان رشته تجربی، شکست تحصیلی را تجربه نموده و مردود شده‌اند (مرکز سنجش وزارت آموزش و پرورش، ۱۳۹۴).

برای تغییر این پیامدهای خطیر، حداقل سه گزینه وجود دارد: تغییر سیستم آموزشی، تغییر آزمون‌ها و شیوه‌های سنجش و یا تغییر هردو. در اینجا، تأکید ما بر راه دوم است. البته این به معنای فقدان شایستگی گزینه‌های اول و سوم نیست؛ و باید توجه داشت که آزمون‌ها خودشان به تنهایی موجب بهبود آموزش و یادگیری نمی‌شوند؛ و بسیاری از عوامل دیگر از جمله کیفیت برنامه درسی، مهارت و تجربه معلمان و حمایت دانش‌آموزان در خارج از کلاس، بر آموزش و یادگیری تأثیر می‌گذارند (پلگرینو و همکاران، ۲۰۰۱: ۳۱) اما همان‌طور که اسکات^{۲۴} استدلال می‌کند، آزمون‌ها را می‌توان، به‌عنوان ابزاری مترقی، منصفانه، منطقی و معقول برای اصلاحات آموزشی به کار برد (۲۰۱۱). در همین رابطه، وولفولک (۲۰۰۴) بیان می‌دارد، اگر آزمون‌ها مشخص کنند که معلمان واقعاً چه چیزهایی را آموزش می‌دهند و دانش‌آموزان چه چیزهایی را می‌آموزند - که واقعاً چنین است - پس راه بهبود آموزش، راهی مستقیم اما رو به بالا است: ارزیابی توانایی‌ها و عادات مهم و اساسی؛ بنابراین با ساخت و به‌کار بردن آزمون‌هایی که به‌جای مطالب سطحی، فرآیندهای بالاتر

23. Third International Mathematics and Science Study (TIMSS)

24. Scott

شناختی را مورد تأکید و سنجش قرار دهد و استفاده از آن‌ها جهت ارزشیابی دانش آموزان می‌تواند در فرآیند یادگیری آن‌ها تغییر بنیادی و مطلوب ایجاد کرد. بی‌شک «کوشش در جهت اصلاح برنامه درسی مدارس تنها زمانی می‌تواند مؤثر باشد که امتحانات و دیگر روش‌های ارزشیابی تقویت گردد، لذا طراحی امتحانات جدید به‌عنوان بخش اساسی از فرآیند اصلاح آموزش برنامه درسی مدارس تصدیق می‌گردد» (دایره المعارف تعلیم و تربیت، جلد ۸: ۱۱۱).

با توجه به آنچه مطرح شد، ضرورت در دسترس داشتن آزمون‌های استاندارد شده پیشرفت تحصیلی احساس می‌شود؛ لذا پژوهش حاضر در پی آن بود که به‌منظور اندازه‌گیری پیشرفت و مهارت‌های اساسی دانش آموزان در برنامه درسی زیست‌شناسی، آزمونی بر اساس اهداف و محتوای آموزشی این ماده درسی برای پایه دهم دوره متوسطه نظری تدوین نماید؛ تا بدین وسیله ابزار مناسبی جهت اندازه‌گیری و سنجش آموخته‌های دانش آموزان فراهم نماید.

مبانی نظری تدوین و استانداردسازی آزمون‌های پیشرفت بر اصول و روش‌های روان‌سنجی نهاده شده است. امروزه دوروش، مدل کلاسیک آزمون (گالیکسون^{۲۵}، ۱۹۵۰) و نظریه سؤال-پاسخ^{۲۶} (IRT)، برای ساختن آزمون‌ها و تفسیر نمرات، به متخصصان اندازه‌گیری کمک کرده است (همبلتون^{۲۷}، ۱۹۸۹). در سال‌های اخیر به دلیل محدودیت‌های نظریه کلاسیک، کاربرد آن کاهش یافته و از طرفی به دلیل ظهور کامپیوتر و نرم‌افزارها، استفاده از IRT رواج یافته است. همبلتون، راجرز و سوامینتان^{۲۸} (۱۹۹۱) سه محدودیت عمده نظریه کلاسیک آزمون تغییرپذیری شاخص‌های سؤال، وابستگی توانایی آزمودنی به آزمون، عدم امکان پیش‌بینی عملکرد آزمودنی در یک سؤال آزمون-راش‌سناسایی کرده‌اند. علاوه بر این، مفاهیم و تعاریف اساسی تئوری کلاسیک غیرقابل آزمون هستند و به سادگی، درست فرض می‌شوند و هیچ راهی برای تعیین تجربی ارتباط بین مفروضات نظریه کلاسیک آزمون با واقعیت وجود ندارد (همبلتون و کوک، ۱۹۷۷؛ همبلتون و همکاران، ۱۹۹۱؛ دی مارس، ۲۰۱۰). این در حالی است که مطلوب آن است که (الف) آماره‌های سؤال وابسته به گروه نباشند، (ب) نمرات آزمون به دشواری آزمون وابسته نباشد، (ج) مدل‌های آزمون، مبنایی برای تطبیق سؤال‌های آزمون با سطح توانایی فراهم کنند، (د) مدل‌های آزمون مبتنی بر فرضیات غیرقابل قبول نباشد. در حال حاضر شواهد قابل توجهی وجود دارد که نشان می‌دهد این ویژگی‌های مطلوب و موارد دیگر را می‌توان در چارچوب نظریه سؤال-پاسخ، به دست آورد (همبلتون و سوامینتان، ۱۹۸۹). در تئوری سؤال-پاسخ فرض می‌شود که (الف) زیربنای عملکرد آزمودنی در آزمون، یک توانایی یا ویژگی واحد است و (ب) رابطه بین عملکرد آزمون‌شونده در هر سؤال و توانایی اندازه‌گیری شده توسط آزمون را می‌توان با یک تابع افزایشی تکنوا^{۲۹} توصیف کرد. این تابع یک تابع مشخصه سؤال^{۳۰} نامیده می‌شود و احتمال پاسخگویی صحیح آزمودنی‌ها را در سطوح مختلف توانایی ارائه می‌دهد. آزمون‌شوندگانی که توانایی بیشتری دارند، نسبت به آزمون‌شوندگان با توانایی کمتر، احتمال بیشتری برای پاسخ صحیح به سؤال دارند (دی مارس^{۳۱}، ۲۰۱۰). توابع مشخصه سؤال، همان‌طور که معمولاً در مدل‌های تک‌بعدی^{۳۲} آزمون مرسوم است، معمولاً با یک، دو یا سه پارامتر توصیف می‌شوند (دی آیلالا^{۳۳}، ۲۰۰۹). مدل سه

25. Gulliksen

26. Item Response Theory

27. Hambleton

28. Hambleton, Rogers & Swaminathan

29. monotonic

30. Item Characteristic Function

31. DeMars

32. unidimensional



پارامتری پیچیده‌ترین مدل IRT است. در این مدل سه پارامتر دشواری، تمیز و حدس وجود دارد که باید برآورد شود. یوری^{۳۴} (۱۹۷۷) ضمن مقایسه مدل‌های یک، دو و سه پارامتری نشان داد که مدل سه پارامتری به‌بترین وجه قادر است آزمون‌های چندگزینه‌ای را با استفاده از داده‌های واقعی آزمون توصیف کند.

مبتنی بر اهداف این مطالعه و به‌منظور سنجش پیشرفت تحصیلی و سطح توانایی دانش‌آموزان در زیست‌شناسی، با توجه به اهداف آموزشی و برنامه‌های درسی این موضوع، آزمونی تدوین و هنجاریابی شد. از آنجاکه جهت تولید آزمون، باید از ویژگی‌های ساختی و فنی آن اطمینان حاصل کرد؛ سؤال‌های پژوهشی زیر مدنظر قرار گرفت.

۱- آیا آزمون زیست‌شناسی تدوین شده از پایایی و اعتبار کافی برخوردار است تا بتوان از آن به‌عنوان ابزاری پایا و معتبر استفاده کرد؟

۲- آیا محتوای آزمون زیست‌شناسی تدوین شده مبتنی بر تحلیل عاملی از یک عامل کلی اشباع است؟

۳- آیا تفاوتی بین عملکرد دانش‌آموز پسر و دختر در آزمون زیست‌شناسی تدوین شده وجود دارد؟

۴- آزمون زیست‌شناسی تدوین شده تا چه اندازه با مدل لجستیک سه پارامتری در تئوری سؤال- پاسخ برآزش دارد؟

۵- اندازه‌های برآورد پارامترهای سؤال (دشواری، تشخیص، حدس) در آزمون زیست‌شناسی تدوین شده چقدر است؟

۶- میزان آگاهی‌های برآورد شده سؤالات آزمون تدوین شده در چه سطحی قرار دارند؟

۷- هنجارهای آزمون زیست‌شناسی تدوین شده برای دانش‌آموزان مدارس دولتی چگونه است؟

روش پژوهش

طرح پژوهشی حاضر از نوع «تحقیق و توسعه^{۳۵}» بوده است که در زمره پژوهش‌های کاربردی قرار می‌گیرد. برخی از مؤلفه‌های اصلی طرح پژوهش عبارت بود از ۱- تعریف حیطه مورداندازه‌گیری، ۲- بیان اهداف و موارد استفاده ابزار اندازه‌گیری، ۳- تعریف رفتاری دانش و مهارت‌های مورداندازه‌گیری، ۴- تعیین الگوی سؤال‌های آزمون، ۵- اجرای مقدماتی و تجربی آزمون به‌منظور اصلاح و یا حذف سؤالات نامناسب، ۶- هنجاریابی آزمون و ارزیابی پایایی و اعتبار آزمون، ۷- تهیه دستورالعمل‌ها و راهنمای اجرا، نمره‌گذاری و تفسیر آزمون.

شرکت‌کنندگان

جامعه آماری این پژوهش را کلیه دانش‌آموزان پایه دهم دبیرستان‌های دولتی شهر شیراز در سال تحصیلی ۱۳۹۴-۱۳۹۳ تشکیل می‌داد. حجم کل جامعه آماری ۳۹۰۳۹ دانش‌آموز بود. از این جامعه ۳۰۰ دانش‌آموز دختر و پسر به‌صورت تصادفی برای نمونه‌گیری و ۹۳۸ دانش‌آموز (۴۸۲)

33. de Ayala

34. Urry

35. research & development

دختر و ۴۵۶ پسر) برای مرحله نهایی هنجاریابی به روش نمونه‌گیری چندمرحله‌ای^{۴۶} از دانش آموزان مدارس دولتی چهار ناحیه آموزشی شهر شیراز انتخاب شدند. معیار تعیین حجم نمونه در مرحله تجربی، پارامترهای دشواری و تمیز سؤال بود (ثرندایک، ۱۹۸۲ / ترجمه هومن، ۱۳۶۶) و برای نمونه هنجاریابی معیارهای مهم تعیین حجم نمونه عبارت بودند از:

- حجم نمونه برای تحلیل عاملی کافی باشد.
- حجم نمونه برای تجزیه و تحلیل داده‌ها با مدل لجستیک سه پارامتری در نظریه سؤال - پاسخ کافی باشد.

ابزار اندازه‌گیری

برای تدوین سؤالات آزمون ضمن تعیین محتوای کتاب درسی زیست‌شناسی پایه دهم (چاپ ۱۳۹۳) و تعریف حیطه مورد اندازه‌گیری (منطبق با حیطه شناختی بلوم)، به کمک محقق و تیمی از دبیران زیست‌شناسی ۱۵۰ سؤال چهارگزینه‌ای بر اساس جدول مشخصات آزمون و منطبق بر اصول روان‌سنجی طراحی گردید. سؤالات تولیدشده پس از اصلاح و ویراستاری بر روی یک نمونه در دسترس اجرا شد تا از نظر قابل فهم بودن محتوا و مدت زمان لازم برای پاسخگویی، اطمینان حاصل شود. سپس سؤال‌ها در سه فرم موازی - از نظر محتوا و هدف - توزیع شد. جهت اجرای تجربی، آزمون اولیه تدوین شده بر روی ۳۰۰ دانش‌آموز (۱۰۰ دانش‌آموز برای هر فرم) دبیرستان‌های متوسطه دولتی شیراز که به صورت تصادفی ساده انتخاب شده بودند، اجرا شد. پس از جمع‌آوری داده‌ها و ضمن بررسی شواهد اعتباریابی، سؤالات آزمون بر اساس معیارهای فنی زیر مورد تجزیه و تحلیل قرار گرفت:

۱- دشواری هر سؤال حداقل ۰/۲۵ و حداکثر ۰/۷۵ باشد.

۲- شاخص‌های تمیز هر سؤال کمتر از ۰/۳۰ نباشد.

۳- ضریب همبستگی دورشته‌ای هر سؤال حداقل ۰/۲۵ باشد.

۴- گزینه‌های انحرافی هر سؤال دارای قدرت جذب کافی باشند.

بر اساس معیارهای ذکر شده، ۱۰۰ سؤال برای فرم‌های نهایی آزمون انتخاب شدند. سؤال‌های نهایی در دو شکل موازی الف و ب هر یک با ۵۰ سؤال چهارگزینه‌ای توزیع و از ساده به دشوار مرتب شدند. ضمن تعیین شواهد اعتباریابی، فرم‌های نهایی آزمون تدوین شده بر روی نمونه هنجاری به شکلی اجرا شد که فرم‌های آزمون در هر کلاس به صورت تصادفی به آزمودنی‌ها اختصاص یابد. زمان پاسخگویی به مجموعه سؤالات آزمون «۶۰» دقیقه برای هر فرم بود. به هر پاسخ صحیح نمره «۱» و به هر پاسخ اشتباه نمره صفر تعلق گرفت. نمرات خام کل برای هر فرم آزمون از مجموع تعداد پاسخ‌های صحیح به دست آمد. در هر فرم آزمون «۸» درصد از کل سؤال‌ها شامل مفاهیم روش علمی، «۱۷» درصد مفاهیم زیست‌شناسی سلولی و مولکولی، «۱۵» درصد مفاهیم تغذیه و گوارش، «۱۹» درصد مفاهیم تنفس و دفع، «۲۹» درصد مفاهیم گردش مواد و ۱۲ درصد مفاهیم



زیست گیاهی را در برمی گرفت. سؤال‌های آزمون سطح دانش و فرآیندهای شناختی بالاتر (فهم، کاربرد و تحلیل) را مورد ارزیابی قرار می‌داد (جدول ۱).

جدول ۱- ترکیب و توزیع سؤال‌های آزمون زیست‌شناسی بر حسب محتوا و عملکرد

مجموع	تحلیل	کاربرد	فهم	دانش			عملکرد محتوا
				مفاهیم و اصول	روش‌ها و ملاک‌ها	اصلاحات و واقعیت‌ها	
۸			۵		۳		نگرش علمی
۱۱	۱	۲				۸	ساختار سلولی
۶			۲	۱	۱	۲	ساختار شیمیایی
۱۵		۲	۴	۱	۳	۵	تغذیه و گوارش
۸	۲	۱	۴			۱	دفع مواد
۱۲		۳	۱	۱		۷	زیست گیاهی
۲۹	۷	۶	۵	۱	۲	۸	گردش مواد
۱۱		۱	۴	۱	۳	۲	تنفس
۱۰۰	۹	۱۴	۲۷	۵	۱۲	۲۳	جمع

یافته‌های پژوهش

مشخصه‌های توصیفی سؤالات آزمون:

جدول ۲ خلاصه مشخصه‌های دشواری و تمیز سؤالات آزمون مبتنی بر مدل کلاسیک اندازه‌گیری را نشان می‌دهد. بر مبنای این جدول دامنه دشواری سؤالات آزمون بین ۰/۲ تا ۰/۸ گسترده شده بود و همه سؤالات از قدرت تمیز بالاتر از ۰/۳ برخوردار بودند. متوسط دشواری سؤالات آزمون ۰/۴۴ و متوسط قدرت تشخیص آن‌ها ۰/۵۵ به دست آمد؛ بنابراین مبتنی بر معیارهای روان‌سنجی (کاپلان و ساکوزو^{۳۷}، ۲۰۱۸) سؤالات آزمون کمی دشوار اما از قدرت تشخیص خوبی برخوردار بوده است.

جدول ۲- خلاصه دشواری و قدرت تشخیص سؤال‌های آزمون زیست‌شناسی						
قدرت تشخیص سؤال			دشواری سؤال			فرم آزمون
متوسط	حداکثر	حداقل	متوسط	حداکثر	حداقل	آزمون
۰/۵۳	۰/۶۲	۰/۳۱	۰/۴۳	۰/۸۳	۰/۱۹	الف
۰/۵۷	۰/۶۷	۰/۲۹	۰/۴۵	۰/۷۹	۰/۲۳	ب

پایایی آزمون:

از آنجایی که در این مطالعه هر سؤال به صورت دو ارزشی (با مقادیر ۱ یا ۰) امتیازدهی شده بود، فرمول ۲۰ کودر-ریچاردسون^{۳۸} (۱۹۳۷) برای برآورد پایایی محاسبه شد. ضرایب پایایی برای فرم الف و ب به ترتیب «۰/۸۹» (خطای استاندارد اندازه گیری ۴/۳۵) و «۰/۸۸» (خطای استاندارد اندازه گیری ۴/۳۳) برآورد شد (جدول ۳). با توجه به معیار هلمشتاتر^{۳۹} (۱۹۶۴) برای ضرایب پایایی آزمون‌های پیشرفت، ضرایب برآورد شده برای هر دو فرم آزمون قابل قبول بود ضمن آنکه نشان می‌داد که سؤال‌های آزمون زیست‌شناسی از ثبات داخلی رضایت بخشی برخوردار بوده است.

جدول ۳- ضرایب پایایی و خطاهای استاندارد آزمون زیست‌شناسی

خطای استاندارد	ضریب پایایی	انحراف استاندارد	میانگین	تعداد سؤال	حجم گروه نمونه			فرم آزمون
					کل	دختر	پسر	
۴/۳۵	۰/۸۸۶	۹/۴۰	۲۲/۸۵	۵۰	۴۶۸	۲۲۲	۲۴۶	الف
۴/۳۲	۰/۸۸۴	۹/۲۴	۲۲/۷۰	۵۰	۴۷۰	۲۲۶	۲۴۴	ب

اعتبار آزمون:

اگرچه چندین روش خاص برای اعتبار یابی آزمون‌های آموزشی شرح داده شده است، اما اعتبار محتوا^{۴۰} در سنجش پیشرفت تحصیلی حائز اهمیت بسیار است؛ بنابراین در این پژوهش علاوه بر اعتبار سازه (از طریق تحلیل عاملی)، از اعتبار محتوایی نیز استفاده شد. اعتبار محتوایی آزمون از ابتدای تدوین سؤالات، ضمن تعریف دقیق حیطه محتوایی و تهیه جدول مشخصات آزمون (جدول ۱) و تدوین سؤالات آزمون بر اساس اهداف و محتوای آموزشی در آن تنیده شد. همچنین به منظور تعیین تجربی اعتبار محتوایی آزمون و برآورد میزان تطابق محتوای سؤالات با حیطه مورد سنجش از ضریب توافق (۰/۶۲) بین نظرات کارشناسان گروه‌های آموزشی استفاده شد که نشانگر متناسب بودن سؤالات با محتوای تدریس است. نتایج تحلیل عاملی بر اساس روش مؤلفه اصلی نیز نشان داد که «۱۵/۷» درصد از کل واریانس در فرم الف و «۱۵/۶» درصد از کل واریانس در فرم ب توسط یک عامل کلی که کاملاً متمایز از بقیه عوامل است توضیح داده می‌شود. همچنین ۶۸/۷ درصد از واریانس مشترک بین سؤالات در فرم الف و ۶۹/۳ درصد از واریانس مشترک بین سؤالات در فرم ب با این عامل کلی تبیین می‌شود (جدول ۴).

38. Kuder – Richardson

39. Helmstadter

40. content-related validity evidence



جدول ۴- ارزش ویژه، درصد واریانس و درصد واریانس تراکمی سؤال‌های آزمون زیست‌شناسی

عامل	ارزش ویژه		درصد واریانس		درصد واریانس تراکمی	
	فرم الف	فرم ب	فرم الف	فرم ب	فرم الف	فرم ب
۱	۷/۸۴۰	۷/۷۸۷	۱۵/۶۸	۱۵/۵۷	۱۵/۶۸	۱۵/۵۷
۲	۱/۸۳۶	۱/۸۹۰	۳/۶۷	۳/۷۸	۱۹/۳۵	۱۹/۳۵
۳	۱/۷۳۸	۱/۵۴۶	۳/۴۸	۳/۰۹	۲۲/۸۳	۲۲/۴۵

تفاوت بین گروه‌ها

برای بررسی اینکه آیا عامل جنسیت در نمره کل آزمون زیست‌شناسی دخیل است یا خیر، از آزمون t برای گروه‌های مستقل استفاده شد. بر اساس یافته‌ها احتمال اینکه اختلاف میانگین بین دو گروه ناشی از شانس باشد بسیار زیاد بود و شواهد کافی برای رد فرضیه صفر در سطح معنی‌داری $0/05$ وجود نداشت؛ بنابراین بین عملکرد پسران و دختران در آزمون زیست‌شناسی تفاوت معنی‌داری وجود نداشت (جدول ۵). نبود تفاوت معنی‌دار بین میانگین نمرات پسران و دختران مؤید این است که سازه اندازه‌گیری شده توسط آزمون زیست‌شناسی، مستقل از جنسیت آزمودنی‌ها است.

جدول ۵- آزمون تفاوت بین عملکرد دختران و پسران در فرم‌های آزمون زیست‌شناسی

فرم آزمون	پسران		دختران		df	t
	M	SD	M	SD		
الف	۲۲/۷۳	۱۰/۶۳	۲۲/۹۷	۸/۰۷	۴۲۵/۳	*۰/۷۸
ب	۲۲/۳۴	۱۰/۳۹	۲۳/۰۳	۸/۰۲	۴۲۴/۷	*۰/۴۳

* $P < 0/05$

هندجار آزمون

در پژوهش حاضر از یک شاخص نسبی (نمرات مشتق شده) جهت ارائه مجموعه هندجارها استفاده شد. برای این منظور، نمرات خام آزمون با استفاده از توزیع فراوانی به نمرات صدک و نمرات استاندارد نرمال شده (نمرات Z و نمرات T با میانگین «۵۰» و انحراف استاندارد «۱۰») تبدیل شدند. به‌عنوان نمونه، جدول ۶ خلاصه هندجار آزمون زیست‌شناسی برای فرم الف را نشان داده است. بر اساس این جدول، دانش‌آموزی با نمره خام ۴۶ بالاتر از ۹۶ درصد افراد جامعه قرار می‌گیرد و نمره استاندارد این فرد $1/74$ واحد انحراف استاندارد بالاتر از میانگین جامعه قرار می‌گیرد. همچنین این فرد دارای نمره $67/4$ در مقیاس T است.

جدول ۶- خلاصه هنجار آزمون پیشرفت علوم زیستی (فرم الف آزمون)

نمره خام	درصد فراوانی	رتبه درصدی	Z نمره	T نمره
۷	۰/۹	۰/۵۳	-۲/۶۰۳	۲۳/۹۷
۸	۰/۴	۱/۱۸	-۲/۲۸۹	۲۷/۱۱
۹	۱/۱	۱/۹۲	-۲/۰۸۶	۲۹/۱۴
۱۰	۱/۱	۲/۹۹	-۱/۸۹۳	۳۱/۰۷
۱۱	۱/۷	۴/۳۸	-۱/۷۱۶	۳۲/۸۴
:	:	:	:	:
۴۶	۰/۹	۹۶/۰۵	۱/۷۴۰	۶۷/۴۰
۴۷	۱/۵	۹۷/۲۲	۱/۸۹۳	۶۸/۹۳
۴۸	۰/۲	۹۸/۰۸	۲/۰۴۰	۷۰/۴۱
۴۹	۰/۹	۹۸/۶۱	۲/۱۶۲	۷۱/۶۲
۵۰	۱/۱	۹۹/۵۷	۲/۵۳۱	۷۵/۳۱

همترازسازی^{۴۱}

پس از اطمینان از اینکه هر دو فرم آزمون دارای شاخص‌های آماری (تعداد سؤال‌ها، میانگین، انحراف معیار، ضریب پایایی و خطای استاندارد) یکسان هستند، برای اینکه نمره آزمون شونده تحت تأثیر فرمی که در آن شرکت کرده قرار نگیرد از روش‌های همترازسازی خطی و غیرخطی (انتقال بر اساس سطح) پیشنهاد شده توسط ثرندایک (۱۹۸۲) استفاده شد. برای همترازسازی خطی، پس از محاسبه میانگین و انحراف معیار هر یک از گروه‌های شرکت‌کننده در فرم‌های آزمون، فرم الف (A) به‌عنوان آزمون لنگر انتخاب شد و با استفاده از معادله «(۱)» نمره معادل ب (B) محاسبه شد (جدول ۷).

$$\text{معادله ۱} \quad x_B = \bar{x}_B + s_B / s_A (x_A - \bar{x}_B)$$

جایی که «XB» و «XA» میانگین هستند و «SB» و «SA» انحرافات استاندارد هستند.

بر اساس معادله ۱ جدول همترازسازی دو فرم آزمون تهیه شد که خلاصه‌ای از آن در جدول ۷ منعکس شده است. مطابق این جدول: نمره خام ۹ در فرم ب معادل با نمره خام ۸/۹۲ در فرم الف است که نشانگر ۲/۲۰ واحد انحراف استاندارد پایین‌تر از میانگین جامعه (نمره Z) و برابر با نمره ۲۸/۰۸ در مقیاس نمره T (مقیاسی با میانگین ۵۰ و انحراف معیار ۱۰) است.

41. equating



جدول ۷- خلاصه همترازسازی خطی آزمون پیشرفت علوم زیستی (تبدیل نمره خام فرم ب به فرم الف)

نمره خام ب	نمره خام الف	Z نمره	T نمره
۵	۴/۸۵	-۲/۹۸۵	۲۰/۱۵
۶	۵/۸۷	-۲/۶۹۳	۲۳/۰۷
۷	۶/۸۹	-۲/۵۳۲	۲۴/۶۸
۸	۷/۹۰	-۲/۳۷۱	۲۶/۲۹
۹	۸/۹۲	-۲/۱۹۲	۲۸/۰۸
:	:	:	:
۴۶	۴۶/۵۶	۱/۷۶۷	۶۷/۶۷
۴۷	۴۷/۵۷	۱/۸۹۵	۶۸/۹۵
۴۸	۴۸/۵۹	۲/۰۴۲	۷۰/۴۲
۴۹	۴۹/۶۱	۲/۱۶۴	۷۱/۶۴
.	.	.	.

برای همترازسازی غیرخطی نمره خام مربوط به صدک های انتخابی (۲، ۵، ۱۰، ۲۵، ۵۰، ۷۵، ۹۰، ۹۵ و ۹۸) آزمودنی‌ها در هر دو فرم آزمون برابر در نظر گرفته شد (جدول ۸). مطابق این جدول فردی که در هر دو فرم آزمون نمره ۲۱ کسب نماید در صدک ۵۰ قرار می‌گیرد به این معنا که نمره وی از نمره نیمی از افراد جامعه بالاتر است.

جدول ۸- همترازسازی غیرخطی فرم‌های آزمون زیست‌شناسی

صدک برگزیده	ام ۲	ام ۵	ام ۱۰	ام ۲۵	ام ۵۰	ام ۷۵	ام ۹۰	ام ۹۵	ام ۹۸
فرم الف	۹	۱۱	۱۳	۱۶	۲۱	۲۷	۳۷	۴۵	۴۸/۶۲
فرم ب	۱۰	۱۱/۵۵	۱۳	۱۶	۲۱	۲۸	۳۵	۴۴/۴۵	۴۸/۵۸

ارزیابی خوبی برازش

یکی از مهم‌ترین موضوعات در نظریه سؤال- پاسخ، مناسب بودن و برازش^{۴۲} یک مدل تابع ریاضی با داده‌های جمع‌آوری شده از یک اجرای آزمایشی است. در این تحقیق برای ارزیابی برازش مدل

42. fitness

لجستیک سه پارامتری با داده‌های آزمون زیست‌شناسی از نرم‌افزار ASCAL^{۴۳} استفاده شد. ASCAL از یک رویکرد بیشینه درستی‌مایی مشترک^{۴۴} با توزیع پیشین^{۴۵} و آزمون کای اسکور جهت برازش سؤال‌ها با مدل استفاده می‌کند. بر اساس نتایج حاصل: از ۱۰۰ سؤال هر دو فرم آزمون، تنها ۸ سؤال با مدل برازش نداشت و ۹۲٪ سؤال‌ها با مدل ریاضی برازش خوبی نشان داد. عدم برازش برخی سؤالات ممکن است به دلیل وجود تعدادی سؤال ساده (اثر کف) و تعدادی سؤال دشوار (اثر سقف) بوده باشد که طبیعتاً نسبت به دیگر سؤالات از قدرت تشخیص کمتری برخوردارند (امبرتسون و رایس، ۲۰۰۰/ ترجمه شریفی و همکاران، ۱۳۸۸).

برآورد پارامترهای سؤال

مبتنی بر مدل ۳ پارامتری نظریه سؤال-پاسخ، برای هر سؤال آزمون سه پارامتر دشواری سؤال، شیب یا تمیز سؤال و حدس کاذب برآورد می‌شود که در ادامه برآورد پارامترهای مذکور توصیف شده است. ۱- پارامتر دشواری سؤال: (bi) نشان‌دهنده نقطه‌ای از مقیاس توانایی است که در آن آزمون‌شونده احتمال $(1+ci)/2$ برای پاسخ صحیح به یک ماده آزمون دارد (همبلتون، ۱۹۸۹). محدوده دشواری سؤال‌های آزمون برای فرم الف، بین $0/989$ - (سؤال ۱) تا $1/954$ + (سؤال ۴۵) و برای فرم ب بین $0/962$ - (سؤال ۲) تا $2/162$ (سؤال ۴۹) برآورد شد. متوسط پارامتر دشواری سؤال‌ها برای فرم الف $0/74$ + و برای فرم ب $0/73$ + به دست آمد؛ که نشان می‌دهد آزمون نسبتاً دشوار بوده است. ۲- پارامتر تمیز سؤال: (ai) نشان‌دهنده قدرت تشخیص سؤال جهت تفکیک افراد قوی و ضعیف است. دامنه پارامتر تمیز (شیب) سؤال برای فرم الف بین $0/404$ + (سؤال ۲) تا $1/993$ (سؤال ۳۵) و برای فرم ب بین $0/408$ + (سؤال ۴) تا $2/182$ (سؤال ۴۸) برآورد شد. متوسط تمیز سؤال برای فرم الف، $1/13$ و برای فرم ب، $1/16$ به دست آمد که نشان می‌دهد سؤالات آزمون از قدرت تمیز کافی برخوردار بوده است.

۳- پارامتر حدس کاذب (ci): احتمال دادن پاسخ صحیح از روی حدس به سؤال توسط آزمون‌شوندگانی با کمترین توانایی را نشان می‌دهد (فالکنر-باند و ولز^{۴۶}، ۲۰۱۶). دامنه پارامتر حدس از حداقل $0/10$ + (سؤال ۲) فرم ب تا حداکثر $0/27$ + (سؤال ۲۷) فرم الف برآورد شد. متوسط پارامتر حدس برای فرم الف، $0/23$ + و برای فرم ب، $0/21$ + به دست آمد.

منحنی مشخصه سؤال و آزمون

منحنی مشخصه سؤال (ICC) یک تابع ریاضی است که احتمال موفقیت در یک سؤال را به توانایی اندازه‌گیری شده توسط مجموعه سؤال‌ها یا آزمونی که شامل سؤال‌ها است مرتبط می‌کند (همبلتون و سوامینتان، ۱۹۸۵) در این مطالعه ICC‌ها به شکل توابع توزیع لجستیک سه پارامتری بر اساس معادله ۳ برآورد شدند.

43. A Microcomputer Program for Estimating Logistic IRT Item Parameters

44. joint maximum likelihood

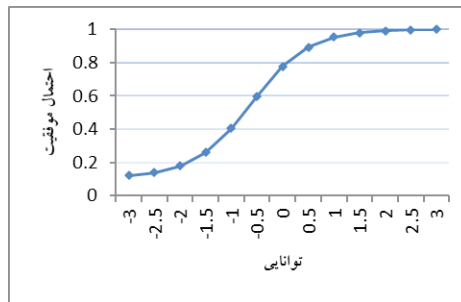
45. prior distributions

46. Faulkner-Bond and Wells



$$p_i(\theta) = c_i + (1 - c_i) \frac{e^{Da_i(\theta - b_i)}}{1 + e^{Da_i(\theta - b_i)}} \quad (i = 1, 2, \dots, n) \quad \text{معادله ۳}$$

به طوری که $p_i(\theta)$ احتمال این است که یک آزمودنی با توانایی θ به سؤال i به درستی پاسخ دهد، b_i پارامتر دشواری سؤال، a_i پارامتر تمیز سؤال و c_i پارامتر حدس کاذب سؤال است و D مقدار ثابت برابر با $1/7$ است. برای نمونه، شکل ۱ منحنی ویژه سؤال ۴ فرم الف آزمون را نشان داده است. بر اساس این نمودار، احتمال اینکه یک آزمودنی با توانایی متوسط به این سؤال پاسخ صحیح دهد برابر $0/8$ است. این سؤال دارای قدرت تشخیص مناسب ($a_i = 1/0.5$) است، احتمال دادن پاسخ صحیح به سؤال، صرفاً از روی حدس اندک است ($c_i = 0/11$) و دشواری آن کمی پایین‌تر از متوسط ($b = -0/71$) برآورد می‌شود.



شکل ۱- منحنی مشخصه سؤال آزمون زیست‌شناسی (سؤال ۴ فرم الف)

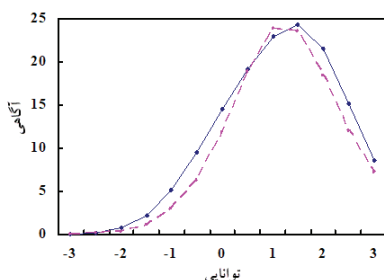
منحنی مشخصه آزمون، نمره خام مورد انتظار در یک آزمون را به خصیصه مکنون (توانایی) مورد اندازه‌گیری مرتبط می‌سازد به عبارت دیگر منحنی ویژه آزمون منعکس‌کننده رابطه تکنوآیین نمره واقعی و نمرات توانایی برای مجموعه خاصی از سؤال‌های آزمون است (همبلتون، ۱۹۸۹). برای محاسبه تابع منحنی ویژه آزمون، مجموع احتمال موفقیت برای سطوح مشخصی از توانایی برای سؤال‌های مختلف محاسبه شد؛ و سپس منحنی ویژه آزمون ترسیم گردید که بر اساس آن می‌توان نمره واقعی فرد در آزمون را برآورد کرد.

تابع آگاهی سؤال و آزمون

توابع آگاهی^{۴۷} نقش برجسته‌ای در IRT دارند و درباره وضعیت نسبی فرد در خصیصه مکنون اطلاعاتی به دست می‌دهند. در این پژوهش با استفاده از رابطه برین بام (۱۹۶۸) توابع آگاهی سؤال برای سطوح معینی از توانایی (θ) به دست آمد (معادله ۴).

$$I_i(\theta) = \frac{2.89a_i^2(1 - c_i)}{[c_i + e^{1.7a_i(\theta - b_i)}][1 + e^{1.7a_i(\theta - b_i)}]^2} \quad (i = 1, 2, 3, \dots, n) \quad \text{معادله ۴}$$

جایی که θ (Ii) آگاهی ارائه شده توسط سؤال i در سطح θ است و ai, bi و ci قبلاً تعریف شده بودند. هرگاه مفروضه استقلال موضعی^{۴۸} صادق باشد، تابع آگاهی آزمون برابر با مجموع توابع آگاهی سؤال‌های تشکیل دهنده یک آزمون است (ثرن‌دایک، ۱۹۸۲ / ترجمه هومن، ۱۳۶۹). شکل ۲ نمودار تابع آگاهی هر دو فرم آزمون زیست‌شناسی را نشان داده است. بر اساس نمودار شکل ۲ هر دو آزمون در محدوده گسترده‌ای از توانایی ۰/۵- تا ۲/۵+ اطلاعاتی پیرامون آزمودنی‌ها ارائه می‌کند اما بیشترین اطلاعات فراهم‌شده توسط هر دو آزمون در محدوده توانایی ۱/۵+ به دست می‌آید؛ بنابراین می‌توان نتیجه گرفت آزمون برای دامنه متوسط به بالای توانایی از پایایی بیشتری برخوردار است.



شکل ۲- تابع آگاهی آزمون زیست‌شناسی (فرم الف- خط ممتد- و فرم ب- خط ناپیوسته)

بحث و نتیجه گیری

هدف از این پژوهش ارائه ابزاری برای سنجش آموخته‌ها و آگاهی از پیشرفت تحصیلی زیست‌شناسی دانش‌آموزان پایه اول متوسطه دوم در رشته تجربی بود. بر این اساس دو فرم موازی از آزمون زیست‌شناسی بر اساس مدل‌های کلاسیک و IRT تدوین و هنجاریابی شد. برآورد ضریب پایایی نشان داد که بیش از ۸۸ درصد واریانس نمرات آزمون مربوط به واریانس واقعی صفت مورد اندازه‌گیری و ۱۲ درصد مربوط به عوامل خطا است. نتیجه حاصله در مقایسه با نتایج گزارش شده برای آزمون‌های مشابه نیز رضایت‌بخش است. برای نمونه ضرایب پایایی آزمون استاندارد پیشرفت استنفورد [علوم] بین ۰/۷۱ تا ۰/۹۶ گزارش شد (مهرنر و لهمان^{۴۹}، ۱۹۸۶، ۴۹۶)؛ بنابراین مجموعه سؤالات آزمون زیست‌شناسی از ثبات درونی بالا و رضایت بخشی برخوردار است. نتایج حاصل از برآورد ضریب دشواری و تمیز سؤال‌ها مبتنی بر مدل کلاسیک نشان داد سؤال‌ها از قدرت دشواری (۰/۲) الی (۰/۸) و قدرت تمیز مناسبی (بیشتر از ۰/۳) برخوردار است؛ لذا می‌تواند اطلاعات لازم در مورد تفاوت‌های فردی دانش‌آموزان را در اختیار آزمایش‌کننده قرار دهد (کاپلان و ساکوزو، ۲۰۱۸). اگرچه روش‌های متعددی در اعتباریابی آزمون‌ها شرح داده شده است، اما اعتبار محتوایی آزمون‌های پیشرفت تحصیلی از اهمیت بیشتری برخوردار است^{۵۰} (AERA, APA, & NCME, ۲۰۱۴).

48. Local independence

49. Mehrens and Lehman

50. American Educational Research Association, American Psychological Association, & National Council of Measurement in Education



در ادبیات روان‌سنجی رایج است که جهت تعیین شواهد مرتبط با اعتبار محتوایی آزمون، عموماً نظرات متخصصان حیطه مورد‌سنجش مورد مشورت قرار گیرد (پیترسون و ترگست^{۵۱}، ۱۹۸۹؛ آبراهام، ویلیامسون و وستبروک^{۵۲}، ۱۹۹۴). در این پژوهش با توجه به تدوین سؤالات امتحانی مبتنی بر تعریف دقیق حوزه محتوایی آزمون و اهداف آموزشی توسط معلمان مجرب و متخصصین آموزش و پرورش و برآورد ضریب تطابق محتوای سؤالات با حیطه مورد‌سنجش (۰/۶۲) مبتنی بر نظرات کارشناسان گروه‌های آموزشی، می‌توان نتیجه گرفت که آزمون زیست‌شناسی از اعتبار محتوای کافی برای اندازه‌گیری دانش و مهارت‌های زیستی برخوردار است. نتایج تحلیل عاملی نیز نشان داد که آزمون پیشرفت زیست‌شناسی از یک عامل کلی اشباع‌شده است، بنابراین می‌توان نتیجه گرفت که عملکرد دانش‌آموزان در این آزمون تا حد زیادی تحت تأثیر یک عامل یا توانایی کلی است که آزمون، اساساً برای اندازه‌گیری آن ساخته شده است.

نتایج آزمون t نشان داد که بین عملکرد پسران و دختران تفاوت معناداری در آزمون مشاهده نمی‌شود. این بدان معناست که آزمون نسبت به گروه‌های جنسی بی‌طرفانه عمل می‌کند و عاری از سوگیری است و سازه اندازه‌گیری شده توسط آزمون، مستقل از جنسیت است. نتیجه حاصل با سایر نتایج آزمون‌های مشابه نیز هماهنگ است. برای نمونه در سومین مطالعه بین‌المللی ریاضی و علوم، تفاوت معنی‌داری بین عملکرد دانش‌آموزان دختر و پسر در مجموعه سؤالات زیست‌شناسی به دست نیامد (کیامنش، ۱۳۷۶).

یافته‌های حاصل از تجزیه و تحلیل سؤالات بر اساس مدل IRT نشان داد که برازش بین داده‌های تجربی و مدل نظری رضایت‌بخش است؛ بنابراین، می‌توان نتیجه گرفت که استفاده از مدل لجستیک سه پارامتری قابل توجیه بوده و نتایج حاصل از آن معتبر است (همبلتون و کوک، ۱۹۷۹؛ همبلتون، ۱۹۸۹). متوسط پارامتر تشخیص برای فرم‌های آزمون ۱/۱۳ (فرم الف) و ۱/۱۶ (فرم ب) به دست آمد که نشان می‌دهد سؤال‌های آزمون از قوه تمیز کافی برای تفکیک دانش‌آموزان قوی و ضعیف به‌ویژه در سطوح میانی و بالای توانایی برخوردار است. همچنین بالا بودن قوه تمیز اکثر سؤال‌ها نشان می‌دهد که پاسخگویی به اغلب سؤال‌های آزمون فقط به خصیصه مکنون توانایی و خیلی کم به عوامل دیگر بستگی دارد (ترندایک، ۱۹۸۲ / ترجمه هومن، ۱۳۶۹).

بر اساس یافته‌ها، متوسط پارامتر دشواری برای فرم‌های آزمون $+۰/۷۳$ بود که با توجه به دامنه مقیاس توانایی (± ۲) نشان می‌دهد آزمون تدوین‌شده کمی دشوار بوده و برای دانش‌آموزان بالاتر از متوسط کارایی بیشتری دارد. میانگین پارامتر حدس کاذب برای فرم‌های آزمون $+۰/۲۱$ به دست آمد که کوچک‌تر از مقدار « $۰/۲۵$ » در سؤال‌های چهارگزینه‌ای است. این نشان می‌دهد که دانش‌آموزان از حدس زدن کورکورانه صرف‌نظر کرده و یا گزینه‌های انحرافی به‌خوبی عمل کرده‌اند و توانسته‌اند آزمون‌های با توانایی اندک را به خود جذب کنند (همبلتون، ۱۹۸۹).

بر اساس منحنی‌های آگاهی آزمون، دامنه‌ای از توانایی که در آن، آزمون اطلاعات مناسبی را برای برآورد نمره واقعی دانش‌آموزان فراهم می‌کند در محدوده $(۰/۵ - تا +۲/۵)$ گسترده شده

51. Peterson & Treagust

52. Abraham, Williamson & Westbrook

است و بیشینه آگاهی در سطح توانایی $1/5+$ به دست می‌آید. از آنجایی که تابع آگاهی به‌عنوان عکس مجذور فاصله اطمینان حول برآورد توانایی تعریف می‌شود (برین بام^{۵۳}، ۱۹۶۸) و با خطای استاندارد اندازه‌گیری رابطه معکوس دارد (همبلتون و کوک^{۵۴}، ۱۹۷۷) می‌توان نتیجه گرفت که آزمون پیشرفت زیست‌شناسی تخمین دقیق‌تری از نمره واقعی آزمودنی‌هایی ارائه می‌دهد که سطح توانایی آن‌ها بالاتر از متوسط پیوستار توانایی است. به عبارت دیگر خطای استاندارد اندازه‌گیری برای این سطح از توانایی کوچک‌تر و پایایی نمرات حاصل بیشتر است. همچنین از آنجایی که تابع آگاهی هر دو فرم یکسان است می‌توان نتیجه گرفت که دو فرم آزمون معادل (موازی) یکدیگرند (همبلتون، ۱۹۸۹).

به‌طور کلی بر اساس یافته‌های حاصل از این مطالعه می‌توان نتیجه گرفت که آزمون استاندارد تدوین شده زیست‌شناسی، به‌طور نسبی ابزاری معتبر و قابل اعتماد است که می‌تواند در زمینه آموزش و ارزیابی مورد استفاده قرار گیرد. این آزمون پیشرفت که جهت اندازه‌گیری سطح آموخته‌ها و توانایی دانش‌آموزان توسعه یافته است، ساختاری دارد که می‌تواند الگویی برای آزمون‌های پایا و معتبری باشد که در این زمینه تهیه خواهد شد.

منابع

- امبرتسون، سوزان ای. و استیون، پی. رایس (۲۰۰۰). نظریه‌های جدید روانسنجی برای روان‌شناسان. ترجمه حسن پاشاشریفی، ولی‌الله فرزاد و همکاران (۱۳۸۸). تهران: رشد.
- ثرندایک، رابرت، ال. (۱۹۸۲). روان‌سنجی کاربردی. ترجمه حیدر علی هومن (۱۳۶۹). تهران: موسسه انتشارات و چاپ دانشگاه تهران.
- سازمان پژوهش و برنامه‌ریزی آموزشی، دفتر برنامه‌ریزی و تألیف کتب درسی. (۱۳۸۹). راهنمای

53. Bimbaum
54. Hambleton & Cook



برنامه‌درس زیست‌شناسی، تهران: مؤلف
سیف، علی اکبر. (۱۳۸۴). سنجش فرایند و فراورده یادگیری. تهران: دوران.
دفتر همکاری‌های علمی بین‌المللی وزارت آموزش و پرورش. (۱۳۷۹). مجموعه گفتارهای
ارزشیابی در آموزش. تهران: مؤلف.
کیامنش، علیرضا. (۱۳۷۶). گزارش سنجش عملکرد در سومین مطالعه بین‌المللی ریاضی و علوم
سال چهارم ابتدایی و سوم راهنمایی. تهران: وزارت آموزش و پرورش.
گلور، جان ای و برونینگ، راجر اچ. (۱۳۷۵). روان‌شناسی تربیتی، اصول و کاربرد آن. ترجمه
علینقی خرازی (۱۳۷۵). تهران: مرکز نشر دانشگاهی
گیج، نیت، ل؛ و برلانیر، دیوید، سی. (۱۳۷۴). روان‌شناسی تربیتی. ترجمه غلامرضا خویی نژاد.
مشهد: حکیم فردوسی.
هومن، حیدر علی. (۱۳۷۲). اندازه‌گیری‌های روانی و تربیتی. تهران: دیبا.

- Abraham, M. R., Williamson, V. M. & Westbrook, S. L. (۱۹۹۴). A cross-age study of the understanding five concepts, *Journal of Research in Science Teaching*, ۱۶۵-۱۴۷: (۲) ۳۱.
- Alele-Williams, G. (۲۰۰۲). Measurement and evaluation in mathematics: The way forward. *Journal of basic sciences*, ۷-۱: (۱) ۱
- American Educational Research Association, American Psychological Association, & National Council of Measurement in Education. (۲۰۱۴). Standards for educational and psychological testing. Washington, DC: Authors
- Baker, J.O. (۲۰۰۳). Testing in modern classroom. London: George. Allen and Unwin
- Birnbaum, A. (۱۹۶۸). Some latent trait models and their use in inferring an examinee's ability. In F. M. Lord & M. R. Novick, *Statistical theories of mental test scores* (pp. ۴۷۲-۳۹۷). Reading, MA: Addison-Wesley.
- Black, Paul; Harrison, Christine; Lee, Clara; Marshall, Bethan and William, Dylan (۲۰۰۳). *Assessment for Learning- putting it into practice*. Maidenhead, U.K.: Open university Press.
- Cizek, Gregory J. (۱۹۹۳). Testing for Learning: A Remonstrance. *Educational Measurement: Issues and Practice*, Volume ۴۲-۴۰: (۴) ۱۲.
- Crooks, T.J. (۱۹۸۸) 'The impact of classroom evaluation practices on students', *Review of Educational Research*, ۴: ۵۸.
- de Ayala, R. J. (۲۰۰۹). *The Theory and Practice of Item Response Theory*, New York: Guilford Publications, Inc.
- DeMars, Christine (۲۰۱۰). *Item Response Theory*. New York: Oxford

University Press, Inc.

Faulkner-Bond Molly and Wells Craig S. (۲۰۱۶). A Brief History of and Introduction to Item Response Theory. In Ronald K. Hambleton and Stephen G. Sireci (Ed.). Educational measurement. New York: The Guilford Press.

Frederiksen, N. (۱۹۸۴) 'The real test bias: Influences of testing on teaching and learning', American Psychologist, ۲۰۲-۱۹۳ : (۳)۳۹.

Gipps, C.V. (۲۰۰۳). Beyond Testing: Towards a Theory of Educational Assessment, London: Washington, D.C. Taylor & Francis e-Library.

Gronlund, N. E. (۱۹۸۸). How to Construct Achievement Tests (۴th ed.). Englewood Cliffs, NJ: Prentice-Hall, Inc.

Gulliksen, H. (۱۹۵۰). Theory of mental tests. John Wiley & Sons Inc.

Hambleton, R. K., Rogers, H. J., & Swaminathan, H. (۱۹۹۱) Fundamentals of item response theory, Newbury Park, Cliff: Sage Publications

Hambleton, R.K. (۱۹۸۹). Principles and selected applications of item-response theory. In R. Linn (Ed.) Educational measurement, (۳th ed.). New York: American Council on Education

Hambleton, R.K. & Cook, L.L. (۱۹۷۷). Latent trait models and their use in the analysis of educational test data, Journal of Educational Measurement, ۹۶-۷۵ : ۱۴.

Helmstadter, G. C. (۱۹۶۴). Principles of Psychological Measurement. New York: Appleton Century Crofts.

International Student Achievement in the TIMSS (۲۰۱۱). Science Content and Cognitive Domains. chapter ۳, TIMSS & PIRLS

International Study Center. Lynch school of education, Boston College.

Isaacs, T., Zara, C., Herbert, G., Coombs, S. and Smith, C. (۲۰۱۳) Key Concepts in Educational Assessment. SAGE Publications Ltd.

Kaplan, Robert M. and Saccuzzo, Dennis P. (۲۰۱۸). Psychological testing: Principles, applications, and issues, (۹th ed.). Boston: Cengage Learning.

Linn, R. (۲۰۰۰). Assessment and accountability. Educational Researcher, ۱۶-۴ : (۲)۲۹.

Mehrens, W. A. and Lehman, H. J. (۱۹۸۶). Using Standardized Test in Education. (۴th ed). New York : Longman.

Pellegrino, J. W., Chudowsky, N. and Glaser, R. (Eds.). (۲۰۰۱) Knowing what Students Know :The Science and Design of Educational Assessment . Washington, DC. National Academy Press.

Peterson, R. F. & Treagust, D. F. (۱۹۸۹). Grade ۱۲- students'



- misconceptions of covalent bonding and structure, *Journal of Chemical Education*, ۴۶۰-۴۵۹ : (۶) ۶۶.
- Phelps, Richard P. (۲۰۱۲). The Effect of Testing on Student Achievement, ۲۰۱۰-۱۹۱۰. *International Journal of Testing*, ۴۳-۲۱ : ۱۲.
- Popham, J. (۱۹۸۷) 'The merits of measurement-driven instruction', *Phi Delta Kappa*, ۸۲-۶۷۹ : ۵.
- Salmon-Cox, L. (۱۹۸۱). Teachers and standardized achievement tests: what's really happening? *Phi Delta Kappan*, ۷۳۶-۷۳۰ : (۱۰) ۶۲.
- Simon, M., Ercikan, K., Rousseau, M., (۲۰۱۳). *Improving large-scale assessment in education :theory, issues and practice* . New York & London: Routledge Flamer.
- Urry, V. W. (۱۹۷۷) Tailored testing: a successful application of latent trait theory. *Journal of Educational Measurement*, ۱۹۶-۱۸۱ : (۲) ۱۴
- Vale, C. D. and Gialluca K. (۱۹۸۵) *ASCAL: A Microcomputer Program for Estimating Logistic IRT Item Parameters*. *Computer Science*

Development and Standardization of Biology Achievement Test

Hassan Moshtaghian Abarghoie¹ *, Bahram Jowkar²

Abstract

The general purpose of this study was to develop and standardize an Achievement test to measure student learning in the biology program at the secondary school. Both classical and IRT models were used to address the research objectives of the study. The preliminary instrument consisted of 150 multiple-choice items that on a sample size of 300 male and female students were performed. The final instrument was two parallel forms of 50 items that were performed on a normative sample of 938 male and female students in Shiraz. The Estimated reliability coefficient for internal consistency with test forms was 0.89, and 0.88, respectively. Based on factor analysis; both forms of the test were an overall factor saturated. Results showed there is no significant difference between the mean scores of boys and girls. So standardized and percentile Norms for all subjects were calculated. Findings from the IRT analysis showed that more than 92% of the items are significantly fitted to Three-Parameter Logistic Model. Test information function was a bell-shaped curve and over a wide ability range from -0.5 to +2.5 provides more information. Also, the maximum information was provided at +1.5 from the ability continuum. Based on this, it can be concluded that the Biology Achievement Test provides a more accurate estimate of the True score of subjects whose ability level is higher than the average ability continuum.

Keyword: Achievement test, Classical model, Item-Response Theory, Standardization

1 -* PhD in Measurement and Assessment; teacher of Baqer Al-Uloom High School, Education District 4 Shiraz -, IRAN
kavir311@gmail.com.

2 - Professor in Educational Psychology; Faculty member of Psychology and Educational Sciences Department, Shiraz University, IRAN