

تحلیل علمی آزمون‌ها بر اساس نظریه کلاسیک و نظریه صفت مکنون

محمد علی احمدوند^{۱*}

(۱) استادیار. عضو هیأت علمی دانشگاه آزاد اسلامی واحد تهران جنوب

* نویسنده مسئول: Ahmadvandma@yahoo.com

تاریخ دریافت مقاله ۸۸/۵/۱۵ تاریخ آغاز بررسی مقاله ۸۸/۵/۲۰ تاریخ پذیرش مقاله ۸۸/۸/۱۱

برگزاری کارگاه‌های تهیه آزمون و ارزیابی را تأیید می‌کند.

کلید واژه گان: ارزشیابی، ضریب دشواری، ضریب تمیز، نظریه صفت مکنون.

چکیده

تحلیل علمی آزمون‌ها، متداولترین وسیله برای ارزشیابی یادگیری و پیشرفت تحصیلی و دانشگاه‌ها است. برای بررسی کیفیت ارزشیابی در دانشگاه تربیت معلم، چهار آزمون روانسنجی، انگیزش و هیجان، احساس و ادراک و روانشناسی رشد از بین مجموعه آزمون‌ها بطور تصادفی انتخاب، و پاسخنامه‌های آنان تحلیل گردیده‌اند. شاخص‌های آماری بکار رفته عبارتند از: درجه دشواری، قدرت تمیز، واریانس و مجذور خی ۱۲۹ سؤال که روایی و پایایی آزمون‌ها تعیین گردید و به سؤال‌های پژوهش پاسخ داده شد و فرضیه‌های آن مورد آزمون قرار گرفت. از ۲۵ تا ۷۵ درصد سؤال‌های طرح شده در آزمون‌های نامناسب شناخته شدند. میانگین نمرات و پاسخ‌نامه‌های دانشجویان تقریباً ۱۲ الی ۱۴ بود، در حالی که میانگین نمره‌های داده شده به آنان خیلی بالاتر از میانگین‌های مذکور بوده است. ضرایب دشواری هر ۴ آزمون خارج از سطح بهینه و ضریب تمیز در آزمون انگیزش و هیجان بالاترین و روانشناسی رشد پایینترین بود. در هر ۴ آزمون این شاخص‌ها خارج از سطح بهینه قرار داشت. آزمون‌ها از روایی و پایایی مناسب برخوردار نبودند. نتیجه پژوهش ض رورت

مقدمه

در دانشکده تربیت معلم دانشگاه آزاد اسلامی واحد تهران جنوب، میزان یادگیری و پیشرفت تحصیلی دانشجویان از طریق برگزاری امتحانات متمرکز در پایان هر نیم سال انجام می‌گیرد. اندازه‌گیری^۱ و ارزشیابی^۲ یادگیری بصورت آزمون‌های چهارگزینه‌ای و یا تشریحی رخ می‌دهد. نمرات دانشجویان اعلام می‌شود و اوراق امتحانی به دایره امتحانات جهت بایگانی سپرده می‌شود و هیچگونه تحلیلی پیرامون این امتحانات صورت نمی‌گیرد تا سهولت و دشواری آزمون‌ها و کیفیت تک تک سؤال‌ها از نظر شاخص‌های آماری تعیین گردد و روشن شود که آیا هدف‌های ارزیابی تحقق یافته اند یا خیر؟ با مساعدت و کمک مسئولین آموزشی و پژوهشی دانشکده، تعدادی از سؤال‌های آزمون‌های اجراشده بر اساس نظریه کلاسیک^۳ و نظریه مکنون^۴ در قالب طرح پژوهشی مصوب، تحلیل گردید.

توجیه علمی موضوع. متداولترین ابزار اندازه‌گیری که در تعدادی از گروه‌های آموزشی نظیر روان‌شناسی، مشاوره، آموزش ابتدایی، مدیریت و برنامه‌ریزی و... بکار می‌رود، آزمون‌های چهارگزینه‌ای است. در چند دهه اخیر ارزشیابی آموزشی را سرمایه‌گذاری برای انسانها و پیشرفت آنان تلقی کرده‌اند و از آنان به عنوان فرایندی منظم برای توصیف کردن و هدایت نمودن کیفیت یادگیری و پیشرفت تحصیلی و اطمینان‌یابی از چگونگی فعالیت‌های آموزشی بهره‌برده‌اند (لنرگان، ۱۳۸۱). برای ارزشیابی شایسته لازم است طراحان سؤال‌های آزمون‌ها با هدف‌های آموزشی آشنا باشند و بتوانند به کمک بودجه‌بندی مطالب، هدف‌ها را در قالب‌های مشخص شده، تعیین کنند و در حوزه شناختی با توجه به طبقه‌های دانش، فهمیدن، کاربستن، تحلیل، ترکیب و

ارزشیابی، سؤال‌های مناسبی را تهیه نمایند (شریفی، ۱۳۸۳، گنجی و ثابت ۱۳۸۳). برای تهیه سؤال‌های چند گزینه‌ای، قواعدی تدوین شده است که لازم است طراحان سؤال‌های امتحانی از آنها سود جویند. آزمون‌ها باید از روایی^۵ و پایایی^۶ برخوردار باشند. آزمون، زمانی روایی دارد که بخواهد اندازه‌گیری آنچه مورد نظر است، مناسب باشد. روایی محتوایی نشان می‌دهد که نمونه سؤال‌های مورد استفاده در یک آزمون تا چه حد معرف جامعه سؤال‌هایی است که می‌توان از محتوای مورد نظر تهیه کرد. از این رو، باید سؤال‌های آزمون، نمونه کاملی از هدف‌ها و محتوای درسی باشند (گنجی، ۱۳۸۱، برنی، ۱۹۸۸). سؤال‌هایی که طبق قواعد درست تهیه شده باشند به روایی آزمون می‌افزایند. هر سؤال که از هدف اصلی آزمون جدا افتد به سهم خود از روایی کل آزمون می‌کاهد. آزمون‌های متشکل از سؤال‌های بسیار دشوار یا بسیار ساده از سطح روایی خوبی برخوردار نیستند (کوهن و همکاران، ۱۹۹۶). یک آزمون زمانی پایایی دارد که اگر آن را در فاصله زمانی کوتاه چندین بار به گروه واحدی از دانشجویان بدهیم نتایج حاصل نزدیک به هم باشد. از روش‌های ساده تعیین پایایی، یکی استفاده از روش دو نیمه کردن آزمون می‌باشد (بری من، ۱۹۹۹، ویک فیلد، ۱۹۹۶).

تحلیل سؤال‌ها و مراحل تحلیل. هدف از تحلیل سؤال‌های آزمون واریسی تک تک سؤال‌ها و تعیین میزان دقت و نارسایی‌های آنها است. در تحلیل سؤال‌های آزمون، نقاط قوت و ضعف یک آزمون و کیفیت همه سؤال‌های آن تعیین می‌شود. طراح سؤال‌پی می‌برد که چگونه دانش‌جویان به آزمون و همچنین به تک تک سؤال‌ها واکنش نشان داده‌اند و می‌داند که کدام گزینه نقش خاصی را در ارزشیابی نداشته است. متداول‌ترین مورد استفاده اطلاعات بدست آمده از تحلیل سؤال‌ها،

1. Measurement
2. Evaluation
3. Classic theory
4. Latent trait theory

5. Validity
6. Reliability

ها به دست می دهند (سلیمی زاده، ۱۳۷۷). ضریب تمیز^۷ قدرت سؤال را در تمایز گذاری یا تشخیص بین گروه قوی و گروه ضعیف آزمون شوندگان مشخص می کند. هر قدر این ضریب بزرگتر باشد قوه تمیز آن بیشتر است (پیراون لاین نت^۸، ۲۰۰۹).

در تحلیل سؤال های آزمون، علاوه بر تعیین ضریب های دشواری و تمیز برای هر سؤال، بررسی نحوه پراکندگی پاسخ های مربوط به گزینه های انحرافی هر سؤال نیز ضروری است. گزینه های انحرافی سؤال ها باید طوری تهیه شوند که بتوانند افراد ضعیف را به خود جلب کنند. هدف از طرح این سؤال این است که افراد گروه قوی گزینه درست و افراد گروه ضعیف یکی از گزینه های غلط را انتخاب کنند. مشکل ترین قسمت در ساخت یک سؤال گزینه های غلط آن است که با وجود غلط بودن، به پاسخ درست بیشتر شباهت داشته باشند. در صورتی یک سؤال به خوبی عمل می کند که افراد ضعیف بیشتر از افراد گروه قوی گزینه های انحرافی آن سؤال را انتخاب نمایند (سیف، ۱۳۸۲؛ ردی^۹، ۲۰۰۹).

تحلیل سؤال بر اساس نظریه صفت مکنون. در نظری کلاسیک، ضریب دشواری، قدرت تمیز، واریانس و شاخص های دیگر آماری تعیین می گردند و بر اساس این شاخص ها می توان سؤال امتحانی را تحلیل نمود. گرچه به روش کلاسیک، ایراداتی وارد شده است. اما می تواند برای طراحان سؤال راهگشا باشد و آنها را به سوی ساخت سؤال های مناسب تر رهنمون شود. منتقدین روش کلاسیک مدعی اند که ویژگی های سؤال ها به گروه آزمونی وابسته هستند. به طور مثال، اگر سؤال ها ساده باشند آزمودنی ها قوی تر و اگر سؤال ها دشوارتر باشند آزمودنی ها ضعیف تر ارزشیابی می شوند. اگر آزمون برای افراد قوی اجرا شود، دشواری کمتری را نشان

انتخاب سؤال های بهتر و مناسب تر برای تشکیل فرم نهایی آزمون است. تحلیل سؤال ها وضعیت دانشجویان را از لحاظ درک مطالب و مفاهیم درس و نقاط ضعف تدریس به خوبی روشن می سازد و می تواند در بهبود روش تدریس مؤثر باشد (ویرسما^۱ و دیگران ۱۹۹۰-کرای گید^۲، ۲۰۰۴). در روش کلاسیک تعیین می شود که چند نفر از دانشجویان گزینه درست سؤال را انتخاب کرده اند و هر یک از گزینه های انحرافی چند نفر را از گروه قوی و گروه ضعیف به خود جلب کرده اند. دوانی^۳ در سال ۱۹۶۷ روشی برای تجزیه و تحلیل سؤال ها ارائه داد. وی ابتدا اوراق را از بالاترین نمره به پایین ترین نمره مرتب کرد. آنگاه یک سوم اوراق را به ترتیب از بالا انتخاب کرده و آنها را گروه بالا نامید و یک سوم اوراق را از کمترین نمره به نام گروه پایین در نظر گرفت و بقیه برگه های امتحانی را کنار گذاشت. در مرحله بعد برای هر سؤال، تعداد گزینه هایی را که شاگردان هر دو گروه انتخاب کرده اند جداگانه شمارش و نتایج را در کارتی درج نمود و دو شاخص دشواری و قوه تمیز را برای هر سؤال تعیین نمود (سیف، ۱۳۸۲؛ نفیسی و زند پارسا، ۱۳۷۶). درصد کل آزمون شوندگان ی که به یک سؤال جواب درست می دهند ضریب دشواری^۴ آن سؤال را بدست می دهد. هر اندازه ضریب دشواری یک سؤال بزرگتر باشد آن سؤال آسان تر است. سؤال هایی بهتر هستند که ضریب دشواری آنها از ۱ کمتر و از صفر بیشتر و به عدد ۰/۵ نزدیک باشند. آلن و ین^۵ (۱۹۷۹) معتقدند سطح بهینه دشواری برای سؤال های چهار گزینه ای در وسط ۰/۲۵ و ۱ یعنی در حدود ۰/۶ قرار دارد (سالکانید^۶، سالکانید^۶، ۲۰۰۸). به طور کلی ضرایب دشواری بین ۰/۳ تا ۰/۷ حد اکثر اطلاع را درباره تفاوت بین آزمودنی

1. Wiersma
2. Craighead
3. Downie
4. Difficulty index
5. Allen and yen
6. Salkind

7. Discrimination index

8. Www.Pareonline.net

9. Redie

می دهد تا زمانی که همین آزمون برای افراد ضعیف اجرا شود. همچنین نمره گذاری در روش کلاسیک غیر واقع‌ی است و به هر سؤال ارزش یکسانی تعلق می‌گیرد. در حالی که سؤال‌های مختلف از نظر محتوایی که مورد پرسش قرار می‌دهند و از نظر دشواری یکسان نیستند که بتوان برای آنها ارزشی برابر قایل شد. در حال حاضر روش تحلیل صفت مکنون^۱ یا نظریه سؤال-پاسخ بر این فرض متکی است که یک ویژگی زیربنایی وجود دارد که به شخص امکان می‌دهد تا در یک تکلیف شناختی معین، موفقیت کسب کند. افزون بر این، این پندار وجود دارد که هر چه شخص از این صفت بیشتر برخوردار باشد، در آزمون مربوط عملکرد بهتری خواهد داشت. در نظریه سؤال-پاسخ هر چه توانایی فرد از می‌زان دشواری سؤال بیشتر باشد، احتمال پاسخ درست فرد به سؤال نیز بیشتر می‌شود و هر چه جایگاه فرد بر روی محور توانایی‌ها از میزان دشواری سؤال کمتر بشود احتمال پاسخ درست این فرد به سؤال نیز کمتر و کمتر می‌شود.

با استفاده از نظریه سؤال-پاسخ می‌توان برخی ویژگی‌های سؤال را تعیین کرد. این ویژگی سؤال برای بسیاری مقاصد آزمون‌سازی از جمله تحلیل سؤال‌های آزمون مفیدند. ویژگی‌های هر سؤال به صورت نمودار نشان می‌دهند و آن را منحنی ویژگی سؤال^۲ (آی‌سی‌سی) می‌نامند. منحنی ویژگی سؤال، احتمال پاسخ درست دادن به هر سؤال را به توانایی آزمون‌شونده ربط می‌دهد. به سخن دیگر، منحنی یا خم ویژگی سؤال یک بازنمایی نموداری از رابطه بین احتمال پاسخ درست دادن به یک سؤال و موقعیت آزمون‌شونده در صفت مورد اندازه‌گیری توسط آزمون است (گلاورز، خرازی ۱۳۷۸). از روی منحنی ویژگی سؤال می‌توان ضریب دشواری و تمیز سؤال را تعیین کرد. ضریب دشواری عبارت است از نمره معیاری که در آن ۵۰ درصد آزمون‌شوندگان سؤال

را درست پاسخ داده‌اند و ضریب تمیز برابر با شیب منحنی ویژگی سؤال است (راسیا، ۲۰۰۶). برای تهیه منحنی ویژگی سؤال نسبت یا درصد آزمون‌شوندگانی که آن سؤال را دست‌جواب داده‌اند در رابطه با نوعی ملاک، مثلاً نمره کل آزمون آنها رسم می‌شود. بر روی محور افقی نمره کل آزمون و بر روی محور عمودی نسبت آزمون‌شوندگانی که به سؤال پاسخ درست داده‌اند، مشخص می‌شود. در این نظریه می‌تواند ضریب دشواری سؤال را مستقیماً به سطح زیر منحنی ویژگی سؤال ربط داد. هر چه زیر منحنی ویژگی سؤال بیشتر باشد، ضریب دشواری سؤال بزرگ‌تر است. هر چه منحنی ویژگی سؤال حالت پلکانی بیشتری داشته باشد همبستگی بین آن سؤال و کل آزمون بیشتر است (سالکانید، ۲۰۰۸؛ سیف، ۱۳۸۲).

سؤال‌های پژوهش

۱. توزیع گزینه درست در بین ۴ گزینه به چه صورت است؟
۲. میانگین نمره‌های دانشجویان در هر آزمون چقدر است؟
۳. شاخص‌های آماری هر سؤال در آزمون‌ها چه وضعی دارند؟
۴. چند درصد سؤال‌های هر آزمون نامناسب هستند؟
۵. آزمون‌ها با توجه به سطح بهینه هر شاخص چه وضعی دارند؟
۶. منحنی ویژگی سؤال‌های انتخابی چه شکلی دارند؟

فرضیه‌های تحقیق

^۱ . Latent trait theory

^۲ . ICC (Item Characteristic Curve)

^۳ . Rasiah

شدند. پاسخنامه های امتحانی که توسط ۲۹۷ نفر از دانشجویان تکمیل شده بود از اداره امتحانات دریافت گردید و با تدوین جدول مشخصات برای هر سؤال و نحوه پاسخدهی دانشجویان، کارت تحلیل سؤال تهیه گردید.

همه سؤال ها از نظر ساخت (تنه سؤال، گزینه انحرافی، گزینه کلید) و رعایت نکات استاندارد مورد توجه قرار گرفت و از روش های آماری تحلیل سؤال بهره برده شد و شاخص های آماری نظیر ضریب دشواری، ضریب تمیز، واریانس و مجذورخی برای هر سؤال محاسبه گردید و سؤال ها از حیث گزینه های انحرافی و شاخص ها در سه دسته مناسب، قابل اصلاح و نامناسب تفکیک شدند. میانگین آزمون ها در گروه های قوی، میانه و ضعیف برای آزمون ها محاسبه گردیدند. توزیع نمرات و نمای آنها به دست آمد. پایایی از راه دو نیمه کردن آزمون ها تعیین شد و رابطه بین ضریب دشواری و قدرت تمیز با استفاده از نظریه صفت مکنون برای آزمون ها با شکل، نشان داده شد و منحنی ویژگی سؤال نیز در چند مورد رسم گردید.

یافته ها

برای تحلیل سؤال های آزمون ها از کارت تحلیل سؤال استفاده شد که در اینجا اولین سؤال روانشناسی رشد به عنوان نمونه درج می گردد:

۱. تفاوت دشواری یا سهولت آزمون ها قابل توجه است.

۲. شاخص های آماری آزمون ها در سطح بهینه قرار ندارند.

۳. درصد قابل توجهی از سؤال های هر آزمون نامناسب هستند.

۴. پایایی و روایی آزمون روانسنجی بهتر از آزمون های دیگر است.

۵. طراحان با طرح سؤال های مناسب و علمی آشنایی کافی ندارند.

هدف پژوهش، پاسخگویی به سؤال ها و آزمون فرضیه ها است. نتایج پژوهش در اختیار طراحان آزمون ها گذاشته می شود تا در بهبود بخشیدن به کیفیت سؤال های خود و نگهداری سؤال های مناسب در بانک سؤال و حذف سؤال های نامناسب اقدام کنند.

روش تحقیق

از بین گروه های آموزشی دانشکده تربیت معلم، گروه روان شناسی با داشتن بالاترین آزمون های چهار گزینه ای انتخاب شد و از بین ۲۰ آزمون اجرا شده ۴ آزمون احساس و ادراک، انگیزش و هیجان، روان سنجی. روان شناسی رشد به صورت نمونه و به قید قرعه برگزیده

جدول شماره ۱. کارت تحلیل سؤال

عنوان آزمون: روان شناسی رشد		شماره سؤال: ۱					
متن سؤال							
مرحله چهارم رشد اخلاقی کلبرگ کدام است؟							
الف. جهت گیری هدف- وسیله ای							
ب. جهت گیری حفظ کردن نظم اجتماعی*							
ج. جهت گیری اصول اخلاقی همگانی							
د. جهت گیری قرارداد اجتماعی							
گروهها	الف	ب	ج	د	سفید	جمع	ملاحظات
گروه بالا	۱۸			۲		۲۰	
گروه پایین	۱	۱۴	۴	۱		۲۰	

N=40		N _U =18		N ₁ =14		شاخص های محاسبه شده	
P _U =90%	P ₁ =70%	P=0/80	q=0/20	pq=0/16	X ² =5	d=0/20	

تحلیل: ۱۸ نفر از گروه قوی و ۱۴ نفر از گروه ضعیف گزینه درست را انتخاب کرده اند.
سؤال آسان است اما با توجه به پراکندگی خوب و معنا دار بودن مجذور خی بعنوان سؤال آسان، می توان در اول آزمون از آن استفاده کرد.

۱۲۹ سؤال ۴ آزمون بر اساس محاسبه شاخص ها
ارزیابی شد و مناسب بودن، قابل اصلاح و ویژگی نا
مناسب بودن را دارد.
طبق بررسی های انجام شده، آزمون انگیزش و هیجان
از لحاظ درصد سؤال های مناسب نسبت به ۳ آزمون

جدول شماره ۲. وضعیت آزمون ها بر حسب سؤال های مناسب، قابل اصلاح و نا مناسب

درصد سؤال های مناسب	تعداد سؤال	سؤال های نا مناسب	سؤال های قابل اصلاح	سؤال های مناسب	آزمون
۷۵ درصد	۲۹	۴	۳	۲۲	انگیزش و هیجان
۵۷ درصد	۳۰	۶	۷	۱۷	روانشنجی
۵۵ درصد	۴۰	۱۰	۸	۲۲	احساس و ادراک
۴۳ درصد	۳۰	۱۲	۵	۱۳	روان شناسی رشد

آزمون روانشناسی رشد از لحاظ انتخاب سؤال،
ضعیف است و ۵۷ درصد سؤال قابل اصلاح و نا مناسب
دارد.

جدول شماره ۳. مقایسه میانگین های نمرات دانشجویان در ۴ آزمون

نما	میانگین گروه ضعیف	میانگین گروه میانه	میانگین گروه قوی	میانگین کل آزمون ها	روان شناسی
۱۴	۱۰/۵۰	۱۴/۵۲	۱۷/۱۷	۱۴/۱۸	روان شناسی
۱۴	۱۰/۳۵	۱۴/۳۵	۱۷/۶۴	۱۴/۱۱	انگیزش و هیجان
۱۴	۱۰/۸۰	۱۴/۱۰	۱۶/۹۰	۱۳/۹۳	احساس و ادراک
۱۴	۹/۸۷	۱۲/۱۲	۱۴/۸۳	۱۲/۲۷	روانشناسی رشد

آزمون روانشناسی رشد دارای کمترین میانگین و
آزمون روانسنجی دارای بالاترین میانگین است. با توجه
به میانگین کل نمرات در ۴ آزمون، روانسنجی بالاترین
میانگین و روانشناسی رشد پایین ترین میانگین را به خود

اختصاص داد. نمای نمرات در هر ۴ آزمون نمره ۱۴ بود.
طراحان سؤال با آگاهانه و یا نا آگاهانه گزینه درست را
در بین ۴ گزینه (الف، ب، ج، د) به نسبت تقریباً متوازی
جای نداده اند و این خود یک نقص تلقی می شود.

جدول شماره ۴. مقایسه محل گزینه درست در بین چهار گزینه آزمون ها

آزمون ها	الف	ب	ج	د
احساس و ادراک	۶	۱۰	۱۶	۸
انگیزش و هیجان	۳	۶	۱۱	۹
روان سنجی	۸	۱۰	۷	۵
روانشناسی رشد	۶	۱۱	۶	۷

گزینه درست باید به نسبت ۲۵ درصد بین ۴ گزینه توزیع شود. در آزمون احساس و ادراک، گزینه کلید بیشتر در قسمت ج قرار گرفته است (۴۰ درصد). در آزمون انگیزش و هیجان نیز، گزینه کلید بیشتر در قسمت ج قرار گرفته است (۳۸ درصد). در آزمون روان سنجی و روان شناسی رشد، گزینه کلید بیشتر در قسمت ب قرار گرفته است (به ترتیب ۳۳ درصد و ۳۷ درصد).

ضرایب دشواری آزمون ها متفاوت هستند. برخی از آنها دارای سادگی زیاد و برخی تا حدی دشواری دارند. آزمون احساس و ادراک در مقایسه با آزمون های دیگر ساده تر و آزمون رشد دشوارتر ارزیابی شد. قدرت تمیز آزمون انگیزش و هیجان بیشتر از بقیه بود.

جدول شماره ۵. مقایسه شاخص های ضریب دشواری و ضریب تمیز و واریانس در ۴ آزمون

واریانس	ضریب تمیز	ضریب دشواری	آزمون ها
۰/۱۶	۰/۳۱	۰/۷۲	احساس و ادراک
۰/۱۸	۰/۳۶	۰/۶۸	انگیزش و هیجان
۰/۱۸	۰/۳۳	۰/۶۸	روان سنجی
۰/۱۸	۰/۲۶	۰/۶۵	روانشناسی رشد

از بین ۳ شاخص فقط واریانس آزمون ها در حد بهینه قرار داشت. برای ۳۰ سؤال روان شناسی رشد، میانگین کل نمرات دانشجویان ۱۲/۲۷ و نمای نمرات آزمون ۱۴ شد. طرح سؤال، تمایل بیشتری در نهادن گزینه کلید در قسمت ب داشته است (۳۶/۶ درصد).

از ۳۰ سؤال، تعداد ۱۳ سؤال مناسب، ۵ سؤال قابل اصلاح و ۱۲ سؤال نامناسب، تشخیص داده شد. با توجه به منحنی ویژگی سؤال برای سؤال شماره ۱، نسبت پاسخدهی در هر دو گروه بالا و پایین زیاد است. به منظور اجتناب از طولانی شدن مقاله به درج منحنی ویژگی سؤال شماره ۱ از آزمون روان شناسی رشد اکتفا می شود.

نمودار شماره ۱. منحنی ویژگی سؤال شماره ۱ روان شناسی رشد

ضریب دشواری آزمون روان شناسی رشد ($P=0/65$) کمی بیشتر از سطح بهینه بود. ضریب تمیز آزمون ($d=0/26$) پایین تر از سطح بهینه و واریانس آزمون ($V=0/18$) قابل قبول تلقی شد. با توجه به رابطه بین ضریب دشواری و ضریب تمیز آزمون و همچنین شاخص های دیگر، آزمون های روانشناسی رشد در حد نسبتاً دشواری با قدرت تمیز پایین ارزیابی شد.

تحلیل: دانشجویان دارای نمرات ۱۶ و ۱۴ در حد بالایی به سؤال پاسخ درست نداده اند. دانشجویان دارای نمره ۱۸ به سؤال جواب نداده اند و از نمره ۲۰ به بعد نسبت پاسخدهی به سؤال افزایش یافته است. با توجه به سطح زیر منحنی و پلکانی نبودن نمودار، سؤال آسان شناخته می شود (ملاک نمره ۳۰ - ۰ است).

نمودار شماره ۲. رابطه بین ضریب دشواری و قدرت تمیز آزمون رشد

ساده ارزیابی شد. ضریب دشواری آزمون انگیزش و هیجان ($P=0/68$) خارج از سطح بهینه و سطح تمیز ($d=0/36$) خارج از بهینه و واریانس ($V=0/18$) قابل قبول تعیین گردید. با توجه به رابطه بین ضریب دشواری و ضریب تمیز آزمون و همچنین شاخص‌های دیگر، آزمون انگیزش و هیجان در سطح مناسب توصیف گردید. برای سؤال ۳۰ سؤال روانسنجی میانگین کل نمرات $14/18$ و نمای نمرات آزمون ۴ شد. طراح سؤال، تمایل بیشتری در نهادن گزینه کلید در قسمت ب داشته است (۳۳ درصد) از ۳۰ سؤال، تعداد ۱۷ سؤال مناسب، ۷ سؤال قابل اصلاح و ۶ سؤال مناسب، شناخته شد. با توجه به منحنی ویژگی سؤال، سؤال شماره ۱، نسبت پاسخدهی در گروه پایین کم و در گروه بالا زیاد است و این نسبت از نمره ۱۲ به بالا افزایش می‌یابد. سؤال از حیث نظریه صفت مکنون مناسب بود. برای سؤال شماره ۳۰، از نمره ۱۳ به بالا پاسخدهی افزایش یافت. از حیث نظریه صفت مکنون، سؤال مطلوب و قابل نگهداری در بانک سؤال است.

ضریب دشواری آزمون روان سنجی ($P=0/68$) کمی بیشتر از سطح بهینه و ضریب تمیز ($d=0/33$) نیز پایین تر از سطح بهینه و واریانس آزمون ($V=0/18$) قابل قبول بود. با توجه به رابطه بین ضرایب و همچنین شاخص‌های دیگر، آزمون روانشناسی در حد نسبتاً آسان و قدرت تمیز نسبتاً پایین ارزیابی گردید.

با در نظر گرفتن میانگین ضریب دشواری ($P=0/65$) می‌توان ادعا کرد ضریب تمیز آزمون بین $62 - 0 +$ تا $0/62 -$ خواهد بود.

برای ۴۰ سؤال آزمون احساس و ادراک، میانگین کل نمرات دانشجویان برابر $13/93$ و نمای آزمون نمره ۱۴ شد. طراح سؤال، تمایل بیشتری در جای دادن گزینه کلید در قسمت (ج) داشته است. از ۴۰ سؤال مناسب تشخیص داده شد. ضریب همبستگی سؤال‌های آزمون $0/16$ محاسبه شد که در سطح $0/05$ معنی دار بود و ضریب همبستگی بین دو گروه هم‌تا $0/97$ محاسبه شد.

ضریب دشواری آزمون احساس و ادراک ($P=0/72$) و خارج از سطح بهینه بود. ضریب تمیز ($d=0/31$) خارج از سطح بهینه بود. واریانس آزمون قابل قبول تلقی گردید. با توجه به رابطه بین ضریب دشواری و ضریب تمیز آزمون و همچنین شاخص‌های دیگر، آزمون احساس و ادراک در حد مناسب ارزیابی شد. برای ۲۹ سؤال آزمون انگیزش و هیجان، میانگین کل $11/14$ و نمای ۱۴ بدست آمد. طراح سؤال، گزینه کلید را با درصد بالاتری در قسمت ج جای داده است (۳۷ درصد). از ۲۹ سؤال آزمون، ۴ سؤال نامناسب، ۳ سؤال قابل اصلاح و ۲۲ سؤال مناسب شناسایی گردید. ضریب همبستگی سؤال‌های آزمون با نمره کل آزمون در سطح $0/01$ معنی دار بود ضریب همبستگی بین دو گروه هم‌تا $0/94$ محاسبه گردید. با توجه به منحنی ویژگی سؤال، سؤال شماره ۱ آزمون، نسبتاً ساده و قابل استفاده و سؤال شماره ۱۳ نسبتاً

نمودار شماره ۳. محل آزمون‌ها روی سطوح بهینه و خارج از بهینه ضریب دشواری

بحث و نتیجه گیری

برای تهیه و تدوین سؤال چهار گزینه ای به آموزش و مهارت نیاز هست. در طرح هر سؤال باید قواعد آزمون سازی رعایت گردد و طراح با توجه به هدفهای آموزشی و طبقات چندگانه شناختی، سؤال خود را تدوین نماید و دقت لازم را در تهیه گزینه های انحرافی بکار ببرد. آزمون های معلم ساخته، پس از اجرا نیازمند تحلیل و بازنگری هستند تا طراح از میان روایی و پایایی و شاخص های آماری آزمون خود آگاه گردد.

درس سنجش و اندازه گیری در گرایش های رشته علوم تربیتی ارایه می شود و دانشجویان کارشناسی ارشد نیز در این گرایش ها، درس ارزشیابی آموزشی را می گذرانند و با شیوه های تهیه، اجرا و استاندارد کردن آزمون ها آشنا می شوند و بقیه دانشجویان و همچنین استادان غیر علوم تربیتی از این شیوه ها نا آگاه هستند. از این رو رشته های دیگر به گونه ای تجربی به تهیه آزمون چهار گزینه ای می پردازند و پس از امتحان نمره را اعلام می کنند و گاهی اوقات نیز، همان سؤال ها را برای نیم سالهای بعد مورد استفاده قرار می دهند.

در این میان تلاشها و زحمات دانشجویان با تعدادی سؤال غیر استاندارد ارزشیابی می شود و بین دانش آموزان تمیز و تشخیص درست بر اساس میزان یادگیری و پیشرفت درسی آنان به عمل نمی آید و از برگذاری امتحانات متمرکز در یک دوره حدود ۱۰ روزه، نتیجه درستی عاید نمی گردد.

برای نتایج این پژوهش آزمون های احساس و ادراک، انگیزش و هیجان، روان سنجی و همچنین روان شناسی رشد از سطح بهینه دشواری دور هستند و به سمت سادگی سوق یافته اند. فرضیه پژوهش مبتنی بر ضریب دشواری مناسب آزمون ها مورد تأیید قرار گرفت. با توجه به درصد بالای سؤال های نا مناسب، در آزمون ها، روایی و پایایی آنها برابر فرضیه پژوهش مورد تأیید قرار نگرفته و بر خلاف انتظار آزمون روان سنجی روایی و پایایی بالاتری از آزمون های دیگر نداشت. در حالی که طراح آزمون که سنجش و اندازه گیری را تدریس می کند آزمونی برگزار کرده که از کیفیت لازم برخوردار نیست. آزمون های چهارگانه از لحاظ ضریب تمیز نیز خارج از سطح بهینه هستند و نمی توانند به خوبی دو گروه دانشجویان قوی و ضعیف را از هم جدا سازند. فرضیه دیگر پژوهش نیز مبتنی بر ضریب تمیز پایین آزمون ها تأیید گردید. آزمون ها از لحاظ واریانس در سطح قابل قبولی قرار دارند و فرضیه پژوهش در این زمینه تأیید نگردید. چنانچه سؤال های قابل اصلاح آزمون ها، مورد بررسی قرار گیرند و سؤال های بهتری جایگزین سؤال های نا مناسب شوند می توان آزمون ها را به سوی کیفیت بهتر برای ارزشیابی علمی توان آموزشی دانشجویان هدایت کرد.

تحلیل سؤال های آزمون بر اساس صفت مکنون تا حدودی با تحلیل کلاسیک سؤال ها همخوانی دارد و دانشجویان دارای نمرات بالاتر، بهتر توانسته اند به سؤال ها پاسخ درست بدهند.

- گنجی، حمزه. (۱۳۷۱). *روانشناسی عمومی*. تهران: دانشگاه پیام نور.
- گنجی، حمزه و ثابت، مهرداد. (۱۳۸۳). *روانشناسی*. تهران: ساوالان.
- گلاورز، جی. ای، برونینگ، ار. اچ. (۱۳۷۸). *روان‌شناسی تربیتی*. ترجمه علینقی خرازی. تهران: مرکز نشر دانشگاهی.
- نفیسی، غلامرضا و زند پارس، علی حسن. (۱۳۷۶). *سنجش و ارزشیابی*. تهران: دانشگاه آزاد اسلامی واحد جنوب.

- Bryman. A. (1991). *Quantitative data analysis*, Rout ledge. London.
- Burney. D. H. (1998). *Research Methods*. Brooksycole Publishing Company Pacific Grove.
- Craighead. W. E. (2004). *The Concies Corsin Encyclopedia of psychology and behavioral science*. Chien R. J. Sweden. M. E. Phillips.
- S. M (1996). *Psychological testing and assessment*, may feet publishing Company. California.
- <http://Pareonline.Net>. (2009). *Basic item analysis for multiple choice tests*. Practical assessment Research and Evaluation Education. Springer Publishing Company.
- <http://redie.Uabc>. (2009). *The level of difficulty and discrimination power of the basic knowledge and*

با اینکه چندین دهه از پیدایش روش تحلیل کلاسیک گذشته است هنوز اجرای آن در نظام ارزشیابی متداول نشده است. امروزه با استفاده از نظریه صفت مکنون به گونه‌ای بهتر و علمی‌تر می‌توان در مورد سؤال‌ها و آزمون‌ها به داوری پرداخت و نسبت به اصلاح آنها و بهبود کیفیت ارزشیابی همت گماشت. استفاده از سؤال و آزمون‌های مطلوب با روایی و پایایی مناسب، برای اندازه‌گیری یادگیری و پیشرفت تحصیلی دانشجویان آج تراب‌ناپذیر است. چون توانایی اعضای علمی گروه‌های آموزشی در طرح سؤال مطلوب، متفاوت است، پیشنهاد می‌شود که از افراد قوی و آموزش دیده، در طرح سؤال‌های امتحانی استفاده بیشتری شود و برای افراد علاقه‌مند کارگاه‌های ارزشیابی و تهیه و اجرا و تحلیل آزمون دایر گردد.

اجرای پژوهش در سطح وسیع‌تر و در گروه‌های دیگر آموزشی با استفاده از نرم افزارهای کامپیوتری و استفاده از نظریه سؤال-پاسخ در تحلیل آزمون‌های اجرا شده، اطلاعات علمی ارزشمند و معتبری در اختیار طراحان سؤال و استادان قرار دهد.

منابع

- بازرگان، عباس. (۱۳۸۱). *ارزشیابی آموزشی*. تهران: سمت.
- سیف، علی اکبر. (۱۳۸۲). *روشهای اندازه‌گیری و ارزشیابی آموزشی*. تهران: دوران.
- سلیمی زاده، محمد کاظم. (۱۳۷۷). *آشنایی با تجزیه و تحلیل سؤالات و آزمون‌ها*. تهران: سازمان سنجش.
- شریفی، حسن پاشا. (۱۳۸۳). *اصول روانسنجی و روان‌آزمایی*. تهران: رشد.

skills Examination.

(exhcoba.Vol.2.No.1.)

- Raja Isaiiah Rasiah. (2009). Relationship *between Item Difficulty and Discrimination indices in True/False Type Multiple choice Questions of a Para-Clinical*. Multidisciplinary paper.Vol.35No.2.
- Salkind, N. (2008). *Encyclopedia of Educational psychology*.
- Wakefield. J. F. (1996). *Educational psychology*. Houghton M. Elfin Company. Boston.
- William. W et al. (1990). *Educational Measurement and Testing*. Ally and Bacon. Boston.

Quarterly Journal of Educational Psychology
Islamic Azad University Tonekabon Branch
Vol. 1, No. 1, autumn 2009, No 1

The Scientific analysis of tests Based on Classic theory and latent trait theory

Ahmadvand. Mohammad Ali*¹

1, 2) Assistant professor. Islamic Azad University. Tehran South Branch

*Corresponding author: Ahmadvandma@yahoo.com

Abstract

Scientific test analysis is most popular means for the evaluation of student at universities. In order to study the quality of evaluation at teacher training faculty, 4fields of study. Namely psychology was randomly chosen. Students answer sheets were analyzed in those 4 tests. The aim of study was to check the reliability and validity, difficulty index and discrimination index of the tests. The hypothesis of research was that statistic indicates of reliability and validity is outside the optimized level in 4tests. The finding of research indicated that between 25 till 57% of the items related to4tests are unsuitable and the difficulty index, discriminative index, validity and reliability of the tests are outside the optimized level.

Key words: Evaluation, difficulty index, discrimination index, latent trait theory.