

تحلیل علمی آزمون‌ها بر اساس نظریه کلاسیک و نظریه صفت مکنون

محمد علی احمدوند^{*1}

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

* نویسنده مسئول: Ahmadvandma@yahoo.com

تاریخ دریافت مقاله 88/5/15 تاریخ آغاز بررسی مقاله 88/5/20 تاریخ پذیرش مقاله 88/8/11

چکیده

میانگین نمرات و پاسخ نامه های دانشجویان تقریباً 12 الی 14 بود، در حالی که میانگین نمره های داده شده به آنان خیلی بالاتر از میانگین های مذکور بوده است. ضرایب دشواری هر 4 آزمون خارج از سطح بهینه و ضریب تمیز در آزمون انگیزش و هیجان بالاترین و روانشناسی رشد پایینترین بود. در هر 4 آزمون این شاخص ها خارج از سطح بهینه قرار داشت. آزمون ها از روایی و پایایی مناسب برخوردار نبودند. نتیجه پژوهش ضرورت برگزاری کارگاه های تهیه آزمون و ارزیابی را تأیید می کند.

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

تحلیل علمی آزمون ها، متداولترین وسیله برای ارزشیابی یادگیری و پیشرفت تحصیلی و دانشگاه ها است. برای بررسی کیفیت ارزشیابی در دانشگاه تربیت معلم، چهار آزمون روانسنجی، انگیزش و هیجان، احساس و ادراک و روانشناسی رشد از بین مجموعه آزمون ها بطور تصادفی انتخاب، و پاسخنامه های آنان تحلیل گردیده اند شاخص های آماری بکار رفته عبارتند از: درجه دشواری، قدرت تمیز، واریانس و مجذور خی 129 سؤال که روایی و پایایی آزمون ها تعیین گردید و به سؤال های پژوهش پاسخ داده شد و فرضیه های آن مورد آزمون قرار گرفت. از 25 تا 75 درصد سؤال های طرح شده در آزمون های نامناسب شناخته شدند.

ارزشیابی، سؤال های مناسبی را تهیه نمایند(شرفی، 1383، گنجی و ثابت 1383). برای تهیه سؤال های چند گزینه ای، قواعدی تدوین شده است که لازم است طراحان سؤال های امتحانی از آنها سود جویند. آزمون ها باید از روایی⁵ و پایایی⁶ برخوردار باشند. آزمون، زمانی روایی دارد که برای اندازه گیری آنچه مورد نظر است، مناسب باشد. روایی محتوایی نشان می دهد که نمونه سؤال های مورد استفاده در یک آزمون تا چه حد معرف جامعه سؤال هایی است که می توان از محتوای مورد نظر تهیه کرد. از این رو، باید سؤال های آزمون، نمونه کاملی از هدف ها و محتوای درسی باشند(گنجی، 1381، برنی، 1988). سؤال هایی که طبق قواعد درست تهیه شده باشند به روایی آزمون می افزایند. هر سؤال که از هدف اصلی آزمون جدا افتد به سهم خود از روایی کل آزمون می کاهد. آزمون های متشکل از سؤال های بسیار دشوار یا بسیار ساده از سطح روایی خوبی برخوردار نیستند(کوهن و همکاران، 1996). 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 روش های ساده تعیین پایایی، یکی استفاده از روش دو نیمه کردن آزمون می باشد(بری من، 1999، ویک فیلد، 1996).

تحلیل سؤال ها و مراحل تحلیل. هدف از تحلیل سؤال های آزمون واریسی تک تک سؤال ها و تعیین میزان دقت و نارسایی های آنها است. در تحلیل سؤال های آزمون، نقاط قوت و ضعف یک آزمون و کیفیت همه سؤال های آن تعیین می شود. طراح سؤال پی می برد که چگونه دانش جویان به آزمون و همچنین به تک تک سؤال ها واکنش نشان داده اند و می داند که کدام گزینه نقش خاصی را در ارزشیابی نداشته است. متداول ترین

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 می گیرد. اندازه گیری¹ و ارزشیابی² یادگیری بصورت آزمون های چهار گزینه ای و یا تشریحی رخ می دهد. نمرات دانشجویان اعلام می شود و اوراق امتحانی به دایره امتحانات جهت بایگانی سپرده می شود و هیچگونه تحلیلی پیرامون این امتحانات صورت نمی گیرد تا سهولت و دشواری آزمون ها و کیفیت تک تک سؤال ها از نظر شاخص های آماری تعیین گردد و روشن شود که آیا هدف های ارزیابی تحقق یافته اند یا خیر؟ با مساعدت و کمک مسئولین آموزشی و پژوهشی دانشکده، تعدادی از سؤال های آزمون های اجرا شده بر اساس نظریه کلاسیک³ و نظریه مکنون⁴ در قالب طرح پژوهشی مصوب، تحلیل گردید.

توجیه علمی موضوع. متداولترین ابزار اندازه گیری که در تعدادی از گروه های آموزشی نظیر روان شناسی، مشاوره، آموزش ابتدایی، مدیریت و برنامه ریزی و... بکار می رود، آزمون های چهار گزینه ای است. در چند دهه اخیر ارزشیابی آموزشی را سرمایه گذاری برای انسانها و پیشرفت آنان تلقی کرده اند و از آنان به عنوان فرایندی منظم برای توصیف کردن و هدایت نمودن کیفیت یادگیری و پیشرفت تحصیلی و اطمینان یابی از چگونگی فعالیت های آموزشی بهره برده اند (بازرگان، 1381). برای ارزشیابی شایسته لازم است طراحان سؤال های آزمون ها با هدف های آموزشی آشنا باشند و بتوانند به کمک بودجه بندی مطالب، هدف ها را در قالب های مشخص شده، تعیین کنند و در حوزه شناختی با توجه به طبقه های دانش، فهمیدن، کاربرستن، تحلیل، ترکیب و

¹. Measurement

². Evaluation

³. Classic theory

⁴. Latent trait theory

⁵. Validity

⁶. Reliability

حداکثر اطلاع را درباره تفاوت بین آزمودنی‌ها به دست می‌دهند (سلیمی زاده، 1377). ضریب تمیز⁷ قدرت سؤال را در تمایز گذاری یا تشخیص بین گروه قوی و گروه ضعیف آزمون شوندگان مشخص می‌کند. هر قدر این ضریب بزرگتر باشد قوه تمیز آن بیشتر است (پیراون لاین نت⁸؛ 2009).

در تحلیل سؤال‌های آزمون، علاوه بر تعیین ضریب‌های دشواری و تمیز برای هر سؤال، بررسی نحوه پراکندگی پاسخ‌های مربوط به گزینه‌های انحرافی هر سؤال نیز ضروری است. گزینه‌های انحرافی سؤال‌ها 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 گزینه‌های غلط را انتخاب کنند. مشکل‌ترین قسمت در ساخت یک سؤال گزینه‌های غلط آن است که با وجود غلط بودن، به پاسخ درست بیشتر شباهت داشته باشند. در صورتی یک سؤال به خوبی عمل می‌کند که افراد ضعیف بیشتر از افراد گروه قوی گزینه‌های انحرافی آن سؤال را انتخاب نمایند (سیف، 1382؛ ردی⁹، 2009).

تحلیل سؤال بر اساس نظریه صفت مکنون. در نظری کلاسیک، ضریب دشواری، قدرت تمیز، واریانس و شاخص‌های دیگر آماری تعیین می‌گردند و بر اساس این شاخص‌ها می‌توان سؤال امتحانی را تحلیل نمود. گرچه به روش کلاسیک، ایراداتی وارد شده است. اما می‌تواند برای طراحان سؤال راهگشا باشد و آنها را به سوی ساخت سؤال‌های مناسب‌تر رهنمون شود. منتقدین روش کلاسیک مدعی‌اند که ویژگی‌های سؤال‌ها به گروه آزمونی وابسته هستند. به طور مثال، اگر سؤال‌ها ساده باشند آزمودنی‌ها قوی‌تر و اگر سؤال‌ها دشوارتر باشند آزمودنی‌ها ضعیف‌تر ارزشیابی می‌شوند. اگر آزمون برای افراد قوی اجرا شود، دشواری کمتری را نشان

مورد استفاده اطلاعات بدست آمده از تحلیل سؤال‌ها، انتخاب سؤال‌های بهتر و مناسب‌تر برای تشکیل فرم‌نهایی آزمون است. تحلیل سؤال‌ها وضعیت دانشجویان را از لحاظ درک مطالب و مفاهیم درس و نقاط ضعف تدریس به خوبی روشن می‌سازد و می‌تواند در بهبود روش تدریس مؤثر باشد (ویرسما¹ و دیگران 1990-کرای گید²، 2004). در روش کلاسیک تعیین می‌شود که چند نفر از دانشجویان گزینه درست سؤال را انتخاب کرده‌اند و هر یک از گزینه‌های انحرافی چند نفر را از گروه قوی و گروه ضعیف به خود جلب کرده‌اند. دوانی³ در سال 1967 روشی برای تجزیه و تحلیل سؤال‌ها ارائه داد. وی ابتدا اوراق را از بالاترین نمره به پایین‌ترین 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 برگه‌های امتحانی را کنار گذاشت. در مرحله بعد برای هر سؤال، تعداد گزینه‌هایی را که شاگردان هر دو گروه انتخاب کرده‌اند جداگانه شمارش و نتایج را در کارتی درج نمود و دو شاخص دشواری و قوه تمیز را برای هر سؤال تعیین نمود (سیف، 1382؛ نفیسی و زند پارسا، 1376). درصد کل آزمون شوندگانی که به یک سؤال جواب درست می‌دهند ضریب دشواری⁴ آن سؤال را بدست می‌دهد. هر اندازه ضریب دشواری یک سؤال بزرگتر باشد آن سؤال آسان‌تر است. سؤال‌هایی بهتر هستند که ضریب دشواری آنها از 1 کمتر و از صفر بیشتر و به عدد 0/5 نزدیک باشند. آلن و یین⁵ (1979) معتقدند سطح بهینه دشواری برای سؤال‌های چهار گزینه‌ای در وسط 0/25 و 1 یعنی در حدود 0/6 قرار دارد (سالکانید⁶، 2008). به طور کلی ضرایب دشواری بین 0/3 تا 0/7

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

⁷. Discrimination index

⁸. Www.Pareonline.net

⁹. Redie

را درست پاسخ داده اند و ضریب تمیز برابر با شیب منحنی ویژگی سؤال است (راسیا³، 2006).
 برای تهیه منحنی ویژگی سؤال نسبت یا درصد آزمون شوندگانی که آن سؤال را دست جواب داده اند در رابطه با نوعی ملاک، مثلاً نمره کل آزمون آنها رسم می شود. بر روی محور افقی نمره کل آزمون و بر روی محور عمودی نسبت آزمون شوندگانی که به سؤال پاسخ درست داده اند، مشخص می شود. در این نظریه می تواند ضریب دشواری سؤال را مستقیماً به سطح زیر منحنی ویژگی سؤال ربط داد. هر چه زیر منحنی ویژگی سؤال بیشتر باشد، ضریب دشواری سؤال بزرگ تر است. هر چه منحنی ویژگی سؤال حالت پلکانی بیشتری داشته باشد همبستگی بین آن سؤال و کل آزمون بیشتر است (سالکانید، 2008؛ سیف، 1382).

سؤال های پژوهش

1. توزیع گزینه درست در بین 4 گزینه به چه صورت است؟
2. میانگین نمره های دانشجویان در هر آزمون چقدر است؟
3. شاخص های آماری هر سؤال در آزمون ها چه وضعی دارند؟
4. چند درصد سؤال های هر آزمون نا مناسب هستند؟
5. آزمون ها با توجه به سطح بهینه هر شاخص چه وضعی دارند؟
6. منحنی ویژگی سؤال های انتخابی چه شکلی دارند؟

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 می گیرد. در حالی که سؤال های مختلف از نظر محتوایی که مورد پرسش قرار می دهند و از نظر دشواری یکسان نیستند که بتوان برای آنها ارزشی برابر قایل شد. در حال حاضر روش تحلیل صفت مکنون¹ یا نظریه سؤال-پاسخ بر این فرض متکی است که یک ویژگی زیربنایی وجود دارد که به شخص امکان می دهد 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 می شود و هر چه جایگاه فرد بر روی محور توانایی ها از میزان دشواری سؤال کمتر بشود احتمال پاسخ درست این فرد به سؤال نیز کمتر و کمتر می شود.

با استفاده از نظریه سؤال-پاسخ می توان برخی ویژگی های سؤال را تعیین کرد. این ویژگی سؤال برای بسیاری مقاصد آزمون سازی از جمله تحلیل سؤال های آزمون مفیدند. ویژگی های هر سؤال به صورت نمودار نشان می دهند و آن را منحنی ویژگی سؤال² (آی سی سی) می نامند. منحنی ویژگی سؤال، احتمال پاسخ درست دادن به هر سؤال را به توانایی آزمون شونده ربط می دهد. به سخن دیگر، منحنی یا خم ویژگی سؤال یک بازنمایی نموداری از رابطه بین احتمال پاسخ درست دادن به یک سؤال و موقعیت آزمون شونده در صفت مورد اندازه گیری توسط آزمون است (گلاورز، خرازی 1378). از روی منحنی ویژگی سؤال می توان ضریب دشواری و تمیز سؤال را تعیین کرد. ضریب دشواری عبارت است از نمره معیاری که در آن 50 درصد آزمون شوندگان سؤال

¹. Latent trait theory

². ICC (Item Characteristic Curve)

³. Rasiah

فرضیه های تحقیق

1. تفاوت دشواری یا سهولت آزمون‌ها قابل توجه است.
2. شاخص های آماری آزمون‌ها در سطح بهینه قرار ندارند.
3. درصد قابل توجهی از سؤال های هر آزمون نامناسب هستند.
4. پایایی و روایی آزمون روانسنجی بهتر از آزمون های دیگر است.
5. طراحان با طرح سؤال های مناسب و علمی آشنایی کافی ندارند.

هدف پژوهش، پاسخگویی به سؤال‌ها و آزمون فرضیه‌ها است. نتایج پژوهش در اختیار طراحان آزمون‌ها گذاشته می‌شود تا در بهبود بخشیدن به کیفیت سؤال‌های خود و نگهداری سؤال‌های مناسب در بانک سؤال و حذف سؤال‌های نامناسب اقدام کنند.

روش تحقیق

از بین گروه‌های آموزشی دانشکده تربیت معلم، گروه روان‌شناسی با داشتن بالاترین آزمون‌های چهارگزینه‌ای انتخاب شد و از بین 20 آزمون اجرا شده 4 آزمون احساس و ادراک، انگیزش و هیجان، روان‌سنجی، روان‌شناسی رشد به صورت نمونه و به قید قرعه برگزیده

شدند. پاسخنامه‌های امتحانی که توسط 297 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 سؤال‌ها از نظر ساخت (تنه سؤال، گزینه انحرافی، گزینه کلید) و رعایت نکات استاندارد مورد توجه قرار گرفت و از روش‌های آماری تحلیل سؤال بهره‌برده شد و شاخص‌های آماری نظیر ضریب دشواری، ضریب تمیز، واریانس و مجذورخی برای هر سؤال محاسبه گردید و سؤال‌ها از حیث گزینه‌های انحرافی و شاخص‌ها در سه دسته مناسب، قابل اصلاح و نامناسب تفکیک شدند. میانگین آزمون‌ها در گروه‌های قوی، میانه و ضعیف برای آزمون‌ها محاسبه گردیدند. توزیع نمرات و نمای آنها به دست آمد. پایایی از راه دو نیمه کردن آزمون‌ها تعیین شد و رابطه بین ضریب دشواری و قدرت تمیز با استفاده از نظریه صفت مکنون برای آزمون‌ها با شکل، نشان داده شد و منحنی ویژگی سؤال نیز در چند مورد رسم گردید.

یافته‌ها

برای تحلیل سؤال‌های آزمون‌ها از کارت تحلیل سؤال استفاده شد که در اینجا اولین سؤال روانشناسی رشد به عنوان نمونه درج می‌گردد:

جدول شماره 1. کارت تحلیل سؤال

عنوان آزمون: روان شناسی رشد		شماره سؤال: 1					
متن سؤال							
مرحله چهارم رشد اخلاقی کلبِرگ 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		18		2		20	
گروه پایین	1	14	4	1		20	
شاخص های محاسبه شده				$N_1=14$	$N_U=18$	$N=40$	
	$X^2=5$	$pq=0/16$	$q=0/20$	$P=0/80$	$P_1=70\%$	$P_U=90\%$	$d=0/20$

تحلیل: 18 نفر از گروه قوی و 14 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.

129 سؤال 4 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 3 آزمون دیگر از کیفیت بهتری برخوردار بود. 75 درصد سؤال های این آزمون مناسب ارزیابی شد. آزمون روانسنجی رتبه دوم و آزمون حساس و ادراک رتبه سوم و آزمون روان شناسی رشد رتبه چهارم را به دست آورد.

جدول شماره 2. 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
75 درصد	29	4	3	22	انگیزش و هیجان
57 درصد	30	6	7	17	روانسنجی
55 درصد	40	10	8	22	احساس و ادراک
43 درصد	30	12	5	13	روان شناسی رشد

آزمون روانشناسی رشد از لحاظ انتخاب سؤال، ضعیف است و 57 درصد سؤال قابل اصلاح و نامناسب دارد.

جدول شماره 3. مقایسه میانگین های نمرات دانشجویان در 4 آزمون

نما	میانگین گروه ضعیف	میانگین گروه میانه	میانگین گروه قوی	میانگین کل	آزمون ها
14	10/50	14/52	17/17	14/18	روان شناسی
14	10/35	14/35	17/64	14/11	انگیزش و هیجان
14	10/80	14/10	16/90	13/93	احساس و ادراک
14	9/87	12/12	14/83	12/27	روانشناسی رشد

آزمون روانشناسی رشد دارای کمترین میانگین و آزمون روانسنجی دارای بالاترین میانگین است. با توجه به میانگین کل نمرات در 4 آزمون، روانسنجی بالاترین میانگین و روانشناسی رشد پایین‌ترین میانگین را به خود

جدول شماره 4. مقایسه محل گزینه درست در بین چهار گزینه آزمون‌ها

آزمون‌ها	الف	ب	ج	د
احساس و ادراک	6	10	16	8
انگیزش و هیجان	3	6	11	9
روان‌سنجی	8	10	7	5
روانشناسی رشد	6	11	6	7

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

گزینه درست باید به نسبت 25 درصد بین 4 گزینه توزیع شود. در آزمون احساس و ادراک، گزینه کلید بیشتر در قسمت ج قرار گرفته است (40 درصد). در آزمون انگیزش و هیجان نیز، گزینه کلید بیشتر در قسمت ج قرار گرفته است (38 درصد). در آزمون روان‌سنجی و روان‌شناسی رشد، گزینه کلید بیشتر در قسمت ب قرار گرفته است (به ترتیب 33 درصد و 37 درصد).

جدول شماره 5. مقایسه شاخص‌های ضریب دشواری و ضریب تمیز و واریانس در 4 آزمون

آزمون‌ها	ضریب دشواری	ضریب تمیز	واریانس
احساس و ادراک	0/72	0/31	0/16
انگیزش و هیجان	0/68	0/36	0/18
روان‌سنجی	0/68	0/33	0/18
روانشناسی رشد	0/65	0/26	0/18

از 30 سؤال، تعداد 13 سؤال مناسب، 5 سؤال قابل اصلاح و 12 سؤال نامناسب، تشخیص داده شد. با توجه به منحنی ویژگی سؤال برای سؤال شماره 1، نسبت پاسخدهی در هر دو گروه بالا و پایین زیاد است. به منظور اجتناب از طولانی شدن مقاله به درج منحنی ویژگی سؤال شماره 1 از آزمون روان‌شناسی رشد اکتفا می‌شود.

از بین 3 شاخص فقط واریانس آزمون‌ها در حد بهینه قرار داشت.

برای 30 سؤال روان‌شناسی رشد، میانگین کل نمرات دانشجویان 12/27 و نمای نمرات آزمون 14 شد. طرح سؤال، تمایل بیشتری در نهادن گزینه کلید در قسمت ب داشته است (36/6 درصد).

تحلیل: دانشجویان دارای نمرات 16 و 14 در حد بالایی به سؤال پاسخ درست نداده اند. دانشجویان دارای نمره 18 به سؤال جواب نداده اند و از نمره 20 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 30 - 0 است).
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.

با در نظر گرفتن میانگین ضریب دشواری ($P=0/65$) می توان ادعا کرد ضریب تمیز آزمون بین $0 - 62 +$ تا $0/62 -$ خواهد بود.
برای 40 سؤال آزمون احساس و ادراک، میانگین کل نمرات دانشجویان برابر $13/93$ و نمای آزمون نمره 14 شد. طراح سؤال، تمایل بیشتری در جای دادن گزینه کلید در قسمت (ج) داشته است. از 40 سؤال مناسب تشخیص داده شد. ضریب همبستگی سؤال های آزمون $0/16$ محاسبه شد که در سطح $5/05$ 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 شاخص های دیگر، آزمون احساس و ادراک در حد مناسب ارزیابی شد. برای 29 سؤال آزمون انگیزش و هیجان، میانگین کل $41/11$ و نمای بدست آمد. طراح سؤال، گزینه کلید را با درصد بالاتری در قسمت ج جای داده است (37 درصد). از 29 سؤال آزمون، 4 سؤال نامناسب، 3 سؤال قابل اصلاح و 22 سؤال مناسب شناسایی گردید. ضریب همبستگی سؤال های 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 1 آزمون، نسبتاً ساده و قابل استفاده و سؤال شماره 13 نسبتاً 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 شاخص های دیگر، آزمون انگیزش و هیجان در سطح مناسب توصیف گردید.
برای سؤال 30 سؤال روانسنجی میانگین کل نمرات $14/18$ و نمای نمرات آزمون 4 شد. طراح سؤال، تمایل بیشتری در نهادن گزینه کلید، در قسمت ب داشته است (33 درصد) از 30 سؤال، تعداد 17 سؤال مناسب، 7 سؤال قابل اصلاح و 6 سؤال مناسب، شناخته شد. با توجه به منحنی ویژگی سؤال، سؤال شماره 1، نسبت پاسخدهی در گروه پایین کم و در گروه بالا زیاد است و این نسبت از نمره 12 به بالا افزایش می یابد. سؤال از حیث نظریه صفت مکنون مناسب بود. برای سؤال شماره 30، از نمره 13 به بالا پاسخدهی افزایش یافت. از حیث نظریه صفت مکنون، سؤال مطلوب و قابل نگهداری در بانک سؤال است.

بود. با توجه به رابطه بین ضرایب و همچنین شاخص‌های دیگر، آزمون روانشناسی در حد نسبتاً آسان و قدرت تمیز نسبتاً پایین ارزیابی گردید.

سادگی سوق یافته‌اند. فرضیه پژوهش مبتنی بر ضریب دشواری مناسب آزمون‌ها مورد تأیید قرار گرفت.

با توجه به درصد بالای سؤال‌های نامناسب، در آزمون‌ها، روایی و پایایی آنها برابر فرضیه پژوهش مورد تأیید قرار نگرفته و بر خلاف انتظار آزمون روان‌سنجی روایی و پایایی بالاتری از آزمون‌های دیگر نداشت. در حالی که طراح آزمون که سنجش و اندازه‌گیری را تدریس می‌کند آزمون‌ی برگزار کرده که از کیفیت لازم برخوردار نیست. آزمون‌های چهارگانه از لحاظ ضریب تمیز نیز خارج از سطح بهینه هستند و نمی‌توانند به خوبی دو گروه دانشجویان قوی و ضعیف را از هم جدا سازند. فرضیه دیگر پژوهش نیز مبتنی بر ضریب تمیز پایین آزمون‌ها تأیید گردید. آزمون‌ها از لحاظ واریانس در سطح قابل قبولی قرار دارند و فرضیه پژوهش در این زمینه تأیید نگردید. چنانچه سؤال‌های قابل اصلاح آزمون‌ها، مورد بررسی قرار گیرند و سؤال‌های بهتری جایگزین سؤال‌های نامناسب شوند می‌توان آزمون‌ها را به سوی کیفیت بهتر برای ارزشیابی علمی توان آموزشی دانشجویان هدایت کرد.

تحلیل سؤال‌های آزمون بر اساس صفت مکنون تا حدودی با تحلیل کلاسیک سؤال‌ها همخوانی دارد و دانشجویان دارای نمرات بالاتر، بهتر توانسته‌اند به سؤال‌ها پاسخ درست بدهند.

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

ضریب دشواری آزمون روان‌سنجی ($P=0/68$) کمی بیشتر از سطح بهینه و ضریب تمیز ($d=0/33$) نیز پایین‌تر از سطح بهینه و واریانس آزمون ($V=0/18$) قابل قبول

بحث و نتیجه‌گیری

برای تهیه و تدوین سؤال چهارگزینه‌ای به آموزش و مهارت نیاز هست. در طرح هر سؤال باید قواعد آزمون‌سازی رعایت گردد و طراح با توجه به هدفهای آموزشی و طبقات چندگانه شناختی، سؤال خود را تدوین نماید و دقت لازم را در تهیه گزینه‌های انحرافی بکاربرد. آزمون‌های معلم‌ساخته، پس از اجرا نیازمند تحلیل و بازنگری هستند تا طراح از میان روایی و پایایی و شاخص‌های آماری آزمون خود آگاه گردد.

درس سنجش و اندازه‌گیری در گرایش‌های رشته علوم تربیتی ارایه می‌شود و دانشجویان کارشناسی ارشد نیز در این گرایش‌ها، درس ارزشیابی آموزشی را می‌گذرانند و با شیوه‌های تهیه، اجرا و استاندارد کردن آزمون‌ها آشنا می‌شوند و بقیه دانشجویان و همچنین استادان غیر علوم تربیتی از این شیوه‌ها نا آگاه هستند. از این رو رشته‌های دیگر به گونه‌ای تجربی به تهیه آزمون چهارگزینه‌ای می‌پردازند و پس از امتحان نمره را اعلام می‌کنند و گاهی اوقات نیز، همان سؤال‌ها را برای نیم‌سالهای بعد مورد استفاده قرار می‌دهند.

در این میان تلاشها و زحمات دانشجویان با تعدادی سؤال غیر استاندارد ارزشیابی می‌شود و بین دانش‌آموزان تمیز و تشخیص درست بر اساس میزان یادگیری و پیشرفت درسی آنان به عمل نمی‌آید و از برگذاری امتحانات متمرکز در یک دوره حدود 10 روزه، نتیجه درستی عاید نمی‌گردد.

برای نتایج این پژوهش آزمون‌های احساس و ادراک، انگیزش و هیجان، روان‌سنجی و همچنین روان‌شناسی رشد از سطح بهینه دشواری دور هستند و به سمت

- گلاورز، جی. ای، برونینگ، ار. اچ. (1378). *روان شناسی تربیتی*. ترجمه علینقی خرازی. تهران: مرکز نشر دانشگاهی.
- نفیسی، غلامرضا و زند پارس، علی حسن. (1376). *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 skills Examination*. (exhcoba. Vol.2.No.1.)
- Raja Isaiah Rasiah. (2009). *Relationship between Item Difficulty and Discrimination indices in True/False Type Multiple choice Questions of a*
- آزمون های مطلوب با روایی و پایایی مناسب، برای اندازه گیری یادگیری و پیشرفت تحصیلی دانشجویان اجتناب ناپذیراست. چون توانایی اعضای علمی گروههای آموزشی در طرح سؤال مطلوب، متفاوت است، پیشنهاد می شود که از افراد قوی و آموزش دیده، در طرح سؤال های امتحانی استفاده بیشتری شود و برای افراد علاقه مند کارگاههای ارزشیابی و تهیه و اجرا و تحلیل آزمون دایر گردد.
- اجرای پژوهش در سطح وسیع تر و در گروههای دیگر آموزشی با استفاده از نرم افزارهای کامپیوتری و استفاده از نظریه سؤال-پاسخ در تحلیل آزمون های اجرا شده، اطلاعات علمی ارزشمند و معتبری در اختیار طراحان سؤال و استادان قرار دهد.
- منابع**
- بازرگان، عباس. (1381). *ارزشیابی آموزشی*. تهران: سمت.
- سیف، علی اکبر. (1382). *روشهای اندازه گیری و ارزشیابی آموزشی*. تهران: دوران.
- سلیمی زاده، محمد کاظم. (1377). *آشنایی با تجزیه و تحلیل سؤالات و آزمون ها*. تهران: سازمان سنجش.
- شریفی، حسن پاشا. (1383). *اصول روانسنجی و روان آزمایی*. تهران: رشد.
- گنجی، حمزه. (1371). *روانشناسی عمومی*. تهران: دانشگاه پیام نور.
- گنجی، حمزه و ثابت، مهرداد. (1383). *روانشنجی*. تهران: ساوالان.

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.

Archive of SID

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

The 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, 4 fields 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 4 tests. The finding of research indicated that between 25 till 57% of the items related to 4 tests are unsuitable and the difficulty index, discriminative index, validity and reliability of the tests are outside the optimized level.

Keywords: Evaluation, difficulty index, discrimination index, latent trait theory.

Archive