

2022

The accuracy of estimating the item difficulty parameter of the one parameter logistic model using the pairwise method in light of sample size and test length change

Mohamed Ahmad Ali Al-Yassin
Yarmouk University, Irbid, Jordan., maha_jawa2001@yahoo.com

Follow this and additional works at: <https://digitalcommons.aaru.edu.jo/jjoas-h>



Part of the [Education Commons](#)

Recommended Citation

Al-Yassin, Mohamed Ahmad Ali (2022) "The accuracy of estimating the item difficulty parameter of the one parameter logistic model using the pairwise method in light of sample size and test length change," *Jordan Journal of Applied Science-Humanities Series*: Vol. 31: Iss. 2, Article 4.
Available at: <https://digitalcommons.aaru.edu.jo/jjoas-h/vol31/iss2/4>

This Article is brought to you for free and open access by Arab Journals Platform. It has been accepted for inclusion in Jordan Journal of Applied Science-Humanities Series by an authorized editor. The journal is hosted on [Digital Commons](#), an Elsevier platform. For more information, please contact rakan@aarj.edu.jo, marah@aarj.edu.jo, u.murad@aarj.edu.jo.

The accuracy of estimating the item difficulty parameter of the one parameter logistic model using the pairwise method in light of sample size and test length change.

دقة تقدير معلمة صعوبة الفقرة للنموذج اللوجستي أحادي المعلمة باستخدام طريقة المزاوجة في ضوء تغير حجم العينة وطول الاختبار

Mohamed Ahmad Ali Al-Yassin^{1*}, Amal Ahmad Mahmoud Al Zoubi.²
Yarmouk University, Irbid, Jordan.¹²

ARTICLE INFO

Article history:

Received 31 Jan 2021

Accepted 15 Mar 2021

Published 01 Apr 2022

<https://doi.org/10.35192/jjoas-h.v31i2.314>

*Corresponding author at Yarmouk University, Jordan.

Mohamed Ahmad Ali Al-Yassin.

Email: maha_jawa2001@yahoo.com.

Keywords:

Pairwise method

Item response theory

Test length

Sample Size

الكلمات المفتاحية:

طريقة المزاوجة

نظرية استجابة الفقرة

طول الاختبار

حجم العينة

ABSTRACT

This study aimed at investigating The accuracy of estimating the item difficulty parameter of the one parameter logistic model using the pairwise method in light of sample size and test length change. To achieve the study objective, the data of binary responses were generated by five different sample size (50, 100, 250, 500, 1000) examinee, and four different length test (20, 30, 40, 70) item using simra function in the pairwise package in R program by determining the level of the ability and difficulty parameter is between (-3, 3) normally distributed with a mean 0 and a standard deviation 1. To answer the questions of the study, the researcher used R program for estimating the difficulty parameter using the pairwise method, and estimating the ability using the weighted likelihood estimation method. The results of the study revealed statistically significant differences at $\alpha = 0.05$ between the standard errors of the estimates of the difficulty parameter due to the sample size for sample size 1000 examinee. The results of the study revealed statistically significant differences at $\alpha = 0.05$ between the standard errors of the estimates of the ability parameter due to the sample size and test length and interaction. The study recommended conduct studies which compare between the likelihood method, the Bayes method, and the pairwise method in the accuracy of estimating the parameters of the item.

هدفت الدراسة إلى الكشف عن دقة تقدير معلمة صعوبة الفقرة للنموذج أحادي المعلمة باستخدام طريقة المزاوجة في ضوء تغير حجم العينة وطول الاختبار. ولتحقيق هدف الدراسة، تم توليد بيانات ثنائية الاستجابة لأربعة مستويات من طول الاختبار (20، 30، 40، 70) فقرة، وخمسة مستويات من حجم العينة (50، 100، 250، 500، 1000) مفحوص باستخدام دالة simra الموجودة في حزمة المزاوجة، وذلك بتحديد مستوى قدرة المفحوصين وصعوبة الفقرات بحيث تكون بين (-3، 3) موزعة توزيعاً طبيعياً بوسط حسابي 0 وانحراف معياري 1. استخدمت برمجية لغة R (Version 3.6.3) لتقدير معلمة الصعوبة للفقرات باستخدام طريقة المزاوجة، وتقدير القدرة باستخدام طريقة الأرجحية الموزونة. وأظهرت نتائج الدراسة وجود فروق ذات دلالة احصائية عند مستوى الدلالة $\alpha = 0.05$ بين الأخطاء المعيارية لتقديرات معلمة الصعوبة تعزى لمتغير حجم العينة ولصالح حجم العينة 1000 مفحوص. وأشارت النتائج إلى وجود فروق ذات دلالة احصائية عند مستوى الدلالة $\alpha = 0.05$ بين الأخطاء المعيارية لتقديرات معلمة القدرة تعزى لمتغير حجم العينة وطول الاختبار والتفاعل بينهما. وأوصت الدراسة بإجراء دراسات تقارن بها بين طريقة الأرجحية العظمى وطريقة بيز وطريقة المزاوجة في دقة تقدير معلم الفقرة.

يعتبر حجم عينة المفحوصين وطول الاختبار المطلوبان لتقديرات معالم نماذج نظرية استجابة الفقرة (Item response theory, [IRT]) عاملان مهمان؛ لأنهما يؤثران على دقة وكفاءة تقديرات معالم الفقرة والقدرة، ويعد التقدير الدقيق للمعلم هو أمر بالغ الأهمية؛ بسبب استخدام نتائج التقديرات لاتخاذ قرارات بشأن المفحوصين في العديد من المجالات التي يمكن أن تؤثر بشكل كبير في حياتهم كالتقدم في مراحل المدرسة والالتحاق بالجامعة والتوظيف. وتعد الموضوعية في القياس من الأمور الأساسية التي تهتم بها العلوم الطبيعية والعلوم الإنسانية على حد سواء، وكلما كان القياس موضوعياً زادت الدقة في فهم ظاهرة موضوع القياس والتنبؤ بها. ويقصد بالموضوعية في القياس النفسي والتربوي ثبات قياس الظاهرة السلوكية باختلاف الأداة المستخدمة في القياس، وألا يتأثر باختلاف المفحوصين الذين يتم تقدير السمة لديهم من خلال هذه الأداة، وأن يتم تدريج الأداة بوحدة قياس تتوافق مع تدريج مستويات الظاهرة السلوكية موضوع القياس (علام، ٢٠٠٥).

ولتلافي العيوب التي ظهرت في النظرية التقليدية في القياس كعدم الموضوعية في القياس، وأنها تقوم على افتراض أن تباين أخطاء القياس هو نفسه لمستويات القدرة جميعها، ومن أجل الوصول إلى قياس أكثر موضوعية، حاول علماء القياس الاستفادة من التقدم التكنولوجي للتوصل إلى طرق سيكومترية جديدة تساعد في التغلب على جوانب القصور في النظرية التقليدية، وتتمثل نظرية استجابة الفقرة (IRT) الاتجاه المعاصر في القياس النفسي والتربوي (Hambleton & Swaminathan, 1985). وتستند نظرية IRT على افتراضات أساسية هي: افتراض أحادي البعد، وافتراض الاستقلال الموضوعي، وافتراض منحني خاصية الفقرة، وافتراض التحرر من السرعة، وانبثقت عن نظرية IRT مجموعة من النماذج التي تعرف بنماذج السمات الكامنة، وتعد النماذج أحادية البعد الأكثر شيوعاً في تصميم الاختبارات والمقاييس التربوية والنفسية وبنائها، وهي الأكثر ملاءمة للفقرات ثنائية الاستجابة، ويعد عدد المعالم الذي توصف الفقرة به الفارق الأساسي بين نماذج السمات الكامنة أحادية البعد (Hambleton & Swaminathan, 1985).

ولقد زاد الاهتمام باستخدام نماذج نظرية IRT لتقدير معالم الفقرات وقدرات المفحوصين، لامتلاكها العديد من الميزات كإعطاء تقديرات دقيقة لمعلم الفقرة والقدرة، ولكنها تحتوي على بعض العيوب؛ إذ تتطلب استخدام حجوم عينات كبيرة للحصول على دقة لتقديرات المعالم (Lord, 1968). حيث أن حجوم العينات الكبيرة غالباً غير متوفرة في البيئات التعليمية، لذلك فإن المهتمين بالقياس بحثوا عن طرق جديدة وبرمجيات أكثر تطوراً لاستخدامها في تقدير المعالم وفق نظرية IRT بوجود عينات صغيرة، ومنها طريقة المزاوجة Pairwise Method من خلال برمجية لغة R (Heine & Tarnai, 2015). وتعد النماذج ثنائية الاستجابة من أشهر النماذج أحادية البعد استخداماً في بناء الاختبارات والمقاييس، كالنموذج اللوجستي أحادي أو ثنائي أو ثلاثي المعلمة (de Grijter & Van der Kamp, ٢٠٠٥)، ويفترض النموذج أحادي المعلمة أن لكل فقرة معلمة صعوبة خاصة بها، بينما لجميع الفقرات القدرة التمييزية نفسها بين المفحوصين، كما يفترض عدم لجوء المفحوصين للتخمين عند الإجابة على الفقرات (Hambleton, 1989).

ويمكن وصف هذا النموذج بالمعادلة الرياضية الآتية:

$$P_i(\theta) = \frac{1}{1 + e^{-Da(\theta - b_i)}} ; i = 1, 2, 3, \dots, n \quad (1)$$

حيث $P_i(\theta)$: احتمال الاستجابة الصحيحة للمفحوص ذو القدرة θ على الفقرة (أ)، b_i : معلم الصعوبة للفقرة (أ)، a : معلمة القدرة للمفحوص، d : معلم تمييز الفقرة وهو ثابت لجميع الفقرات في هذا النموذج، D : عامل التدريج وهو ثابت قيمته ١,٧.

ويعتمد احتمال إجابة المفحوص على الفقرة إجابة صحيحة في النموذج أحادي المعلمة على قدرة المفحوص ومعلمة الصعوبة، ولأن قيمة احتمال إجابة المفحوص على الفقرة غير معلومة، وإجابات المفحوصين على الفقرات معلومة، فإن أساليب تقدير المعالم تهدف إلى تحديد قيمة القدرة لكل مفحوص وقيمة معلمة الصعوبة للفقرة من خلال الاعتماد على إجابات المفحوصين على الفقرات. وهناك عدة عوامل تؤثر في دقة تقدير معالم الفقرات تأثيراً مباشراً مثل: طرق التقدير، وحجم العينة للمفحوصين، وطول الاختبار، وتوزيع السمة الكامنة، وجودة البيانات، والتوزيعات لمعلم الفقرة (Barnes & Wise, 1991). وأشار هامبلتون وسواميناثان (Hambleton & Swaminathan, 1985) أن هناك ثلاثة طرق رئيسة لتقدير معالم الفقرة والقدرة وفق نماذج نظرية IRT تستخدم طريقة الأرجحية العظمى وهي: الأرجحية العظمى المشتركة Joint Maximum Likelihood (JML) والأرجحية العظمى الهامشية Marginal Maximum Likelihood (MML) والأرجحية العظمى المشروطة Conditional Maximum Likelihood (CML)، وأشار لورد (Lord, 1986) إلى أن طريقة بيز (Bayes) تعد من الأساليب الإحصائية المستخدمة لتقدير معالم الفقرات والقدرة.

ويؤدي التقدير المتزامن لكل من معلم الفقرة و قدرة المفحوص في نماذج نظرية IRT إلى تعقيدات إحصائية في التقدير، مما جعل إجراء التقدير محط تركيز أساسي لأبحاث القياس النفسي (Birnbaum, 1969). وأشار وورم (Warm, ١٩٨٩) إلى طريقة تقدير الأرجحية الموزونة Weighted Likelihood Estimation (WLE) التي تعتمد على توزين دالة الأرجحية لتقدير معلمة قدرة المفحوص للتقليل من التحيز في تقدير القدرة بطريقة الأرجحية العظمى، ولا تعتبر طريقة الأرجحية الموزونة من طرق تقدير بيز، لأنه لم يتم وضع افتراضات حول توزيع معلمة القدرة.

واقترح شوبن (Choppin, ١٩٨٥) تقنيتين يمكن استخدامها لتقدير معالم الفقرات وفق النموذج أحادي المعلمة، فالتقنية الأولى هي الأرجحية العظمى التي تعتمد على خوارزمية تكرارية مثل طريقة نيوتن رافسون، أما التقنية الثانية هي طريقة المزاجية التي لا تعتمد على خوارزمية تكرارية. ويعتبر نموذج ثيرستون للمقارنات الزوجية أساساً لطريقة المزاجية، والذي أشار إلى هذه الطريقة هو العالم راش (Rasch, 1960)، وطورها العالم شوبن (Choppin, 1985) لتصبح طريقة عملية لمعايرة الفقرات داخل بنوك الأسئلة، واستخدامها لتقدير معلمة صعوبة الفقرات. إن المبدأ الأساسي في تطبيق المقارنات الزوجية هو الحكم على العديد من الأشياء (مثل الفقرات) على شكل أزواج وتحديد أي من الزوجين مفضلين من قبل المفحوص، ثم الحصول على مقياس متري لمجموعة من الفقرات، ومن مميزات طريقة المزاجية: تجنب العمليات الحسابية الطويلة في تحليل البيانات الكبيرة، ومقارنة فقرتين فقط في كل مرة، والتعامل مع مصفوفات البيانات غير المكتملة دون أي متاعب بناءً على أسلوبها الرقمي، وتقدير معالم النموذج المستخدم بوجود بيانات مفقودة، وتعتبر دقتها في تقدير معلمة الصعوبة أفضل من دقة التقدير باستخدام طريقة الأرجحية العظمى عند استخدامها مع العينات الصغيرة (Heine & Tarnai, 2015).

ويمكن اشتقاق طريقة المزاجية من خلال معادلة النموذج اللوجستي الذي صاغها راش (Rasch, 1960) كما في معادلة ٢:

$$P(x_{vi}) = \frac{e^{x_{iv}(\theta_v - b_i)}}{1 + e^{(\theta_v - b_i)}} \quad (2)$$

حيث x_{vi} : استجابة المفحوص v على الفقرة i (٠، ١)، θ_v : قدرة المفحوص v ، b_i : صعوبة الفقرة i .

وعلى فرض أن عينة من المفحوصين ($v=1 \dots n$)، يجب أن تجيب على فقرتين (J)، وبتحقيق افتراض الاستقلالية العشوائية المشروطة للفقرتين المقدمتين للمفحوص وفق

نموذج راش، فإن الإجابات الأربعة المتوقعة على الفقرتين ستكون إما (٠، ١، ٢)، وتعطي احتمالات الإجابة لها كالتالي:

$$P(x_{vi} = 0, x_{vj} = 0) = \frac{1}{1 + e^{(\theta_v - b_i)}} \times \frac{1}{1 + e^{(\theta_v - b_j)}} \quad (3)$$

$$P(x_{vi} = 1, x_{vj} = 0) = \frac{e^{(\theta_v - b_i)}}{1 + e^{(\theta_v - b_i)}} \times \frac{1}{1 + e^{(\theta_v - b_j)}} \quad (4)$$

$$P(x_{vi} = 0, x_{vj} = 1) = \frac{1}{1 + e^{(\theta_v - b_i)}} \times \frac{e^{(\theta_v - b_j)}}{1 + e^{(\theta_v - b_j)}} \quad (5)$$

$$P(x_{vi} = 1, x_{vj} = 1) = \frac{e^{(\theta_v - b_i)}}{1 + e^{(\theta_v - b_i)}} \times \frac{e^{(\theta_v - b_j)}}{1 + e^{(\theta_v - b_j)}} \quad (6)$$

حيث تشير معادلة ٣ إلى أن المفحوص أجاب إجابة خاطئة عن كلا الفقرتين، ومعادلة ٤ تشير إلى أن المفحوص أجاب إجابة صحيحة عن الفقرة i وإجابة خاطئة عن الفقرة j ، كما أن معادلة ٥ تشير إلى أن المفحوص أجاب إجابة خاطئة عن الفقرة i وإجابة صحيحة عن الفقرة j ، وتشير معادلة ٦ إلى أن المفحوص أجاب إجابة صحيحة عن كلا الفقرتين. ولإيجاد الاحتمال المشترك لفقرتين الذي يحقق الدرجة ١ من خلال جمع المعادلتين (٤، ٥) التي تعطى بمعادلة ٧:

$$P(x_{vi} + x_{vj} = 1) = \frac{e^{(\theta_v - b_i)}}{\left((1 + e^{(\theta_v - b_i)}) \times (1 + e^{(\theta_v - b_j)}) \right)} + \frac{e^{(\theta_v - b_j)}}{\left((1 + e^{(\theta_v - b_i)}) \times (1 + e^{(\theta_v - b_j)}) \right)} \quad (7)$$

وللحصول على الاحتمالية المشروطة بالدرجة ١ على الفقرة i ، من خلال قسمة معادلة ٤ على المعادلة ٧ وتعطى بمعادلة ٨:

$$P(x_{vi} = 1 | x_{vi} + x_{vj} = 1) = \frac{\frac{e^{(\theta_v - b_i)}}{1 + e^{(\theta_v - b_i)}} \times \frac{1}{1 + e^{(\theta_v - b_j)}}}{\frac{e^{(\theta_v - b_i)}}{\left((1 + e^{(\theta_v - b_i)}) \times (1 + e^{(\theta_v - b_j)}) \right)} + \frac{e^{(\theta_v - b_j)}}{\left((1 + e^{(\theta_v - b_i)}) \times (1 + e^{(\theta_v - b_j)}) \right)}} \quad (8)$$

ويمكن تبسيط معادلة ٨ رياضياً، وإلغاء معلمة القدرة كما في معادلة ٩:

$$P(x_{vi} = 1 | x_{vi} + x_{vj} = 1) = \frac{e^{b_i}}{e^{b_i} + e^{b_j}} \quad (9)$$

وبالمثل يمكن صياغة الاحتمال المشروط بالدرجة ١ على الفقرة ج، ويمكن تقدير الاحتمال المشروط بإجابة المفحوص إجابة صحيحة على الفقرة أ وإجابة خاطئة على الفقرة ج في معادلة ٩ من خلال بيانات حقيقية، ويمكن التعبير عنها من خلال معادلة ١٠:

$$f_{i,j} = \frac{e^{b_i}}{e^{b_i} + e^{b_j}} \quad (10)$$

حيث أن: عدد المفحوصين الذين أجابوا إجابة صحيحة على الفقرة أ وإجابة خاطئة على الفقرة ج، $f_{j,i}$: عدد المفحوصين الذين أجابوا إجابة صحيحة على الفقرة ج وإجابة خاطئة على الفقرة أ. $n_{i,j} = f_{i,j} + f_{j,i}$: عدد المفحوصين الذين حققوا مجموع الدرجة ١ على الفقرتين ويمكن اعتبارها حجم عينة المفحوصين n ، وبالمثل يمكن إعادة صياغة الاحتمال المشروط بالدرجة ١ على الفقرة ج. وبالتالي يمكن تقدير صعوبة الفقرة من خلال الاحتمالية المشروطة، إذ يمكن إعادة كتابة معادلة ٩ كما في معادلة ١١:

$$\frac{f_{i,j}}{f_{j,i}} = \frac{e^{b_i}}{e^{b_j}} \quad (11)$$

وبأخذ اللوغاريتم الطبيعي للمعادلة السابقة لينتج نسب لوغاريتم مرجح النجاح Odds كما في معادلة ١٢:

$$\ln\left(\frac{f_{i,j}}{f_{j,i}}\right) = b_i - b_j \quad (12)$$

وبالمثل يمكن كتابة المعادلة للفقرة ج كما يأتي (Heine & Tarnai, 2015):

$$\ln\left(\frac{f_{j,i}}{f_{i,j}}\right) = b_j - b_i \quad (13)$$

وبعد عملية تقدير معالم الفقرات وقدرات المفحوصين، فإنه يجب التحقق من مدى دقة التقدير، فقد أشار لورد (Lord, 1980) إلى وجود معايير مختلفة للكشف عن دقة تقدير المعالم ومن أبرزها: الخطأ المعياري في التقدير، ودالة المعلومات للفقرات والاختبار، والفاعلية النسبية للاختبار. يعد الخطأ المعياري في التقدير مؤشراً إحصائياً جيداً على دقة تقدير معلمة القدرة ومعالم الفقرة، ويرتبط الخطأ المعياري في التقدير عكسياً مع الجذر التربيعي لدالة معلومات الاختبار، ويكون الخطأ المعياري في التقدير أقل ما يمكن عند مستويات القدرة التي تناظر أقصى معلومات للاختبار (Hambleton et al., 1991). ويشير وورم (Warm, 1978) إلى أن الخطأ المعياري في التقدير هو القيمة المتوقعة للانحراف المعياري لأخطاء التقدير، وكلما قلت قيمته كان مؤشراً على دقة التقدير. وكلما قل مقدار الخطأ المعياري في التقدير، فإن ذلك يؤدي إلى زيادة الثبات وبالتالي زيادة دقة الاختبار (Reeve & Fayers, 2004). طور العالم إيفرون (Efron, 1979) طريقة لاختيار عينات بالإرجاع سميت بطريقة البوتستراب (Bootstrap) وهي اختبار مجموعة من عينات عشوائية التي تشتمل على (n) من العناصر المسحوبة بالإرجاع بشكل عشوائي من (N) من البيانات الأصلية، وتقوم هذه الطريقة في تقدير التباينات والأخطاء المعيارية وفترات الثقة والقيمة الاحتمالية. ويمكن تقدير الأخطاء المعيارية لمعلم الفقرات باستخدام طريقة البوتستراب عند استخدام طريقة المزوجة في تقدير معالم الفقرات، إذ أنها لا تستند إلى دالة الأرجحية العظمى في تقدير معالم الفقرات (Finch & French, 2019). إن عدد عينات البوتستراب التي يجب سحبها تختلف باختلاف الغاية التي استخدمت من أجلها تقنية البوتستراب، مثلاً قد تبلغ عدد عينات البوتستراب ١٠٠٠ عينة. ودرس كلاً من إيمون وماسوم (Imon & Masoom, 2005) مسألة تحديد عدد عينات البوتستراب المطلوب سحبها من المجتمع، وأشاروا إلى عدم وجود اتفاق بين الباحثين الإحصائيين حول عدد التكرارات الكافية التي تحقق الغاية من عملية تقدير المعالم، ويمكن تقدير الخطأ المعياري في تقدير معلمة الصعوبة كما في المعادلة الآتية:

$$S(\theta_b) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\theta_b - \theta_B)^2} \quad (14)$$

حيث $S(\theta_b)$: الخطأ المعياري لتقدير صعوبة الفقرة، θ_b : معلمة صعوبة الفقرة الحقيقية، θ_B : معلمة صعوبة الفقرة المقدر، B: عدد عينات البوتستراب (Heine & Tarnai, 2015).

يتناول هذا الجزء مراجعة للدراسات ذات الصلة بمشكلة البحث، وقد تم عرض هذه الدراسات وفق التسلسل الزمني، فمن الدراسات الأجنبية ذات الصلة بمشكلة الدراسة، دراسة هيونغ وآخرون (Huang et al., 2001) التي هدفت إلى مقارنة فعالية بعض البرمجيات الجاهزة (Bilog- MG₃, PIC) في معايرة فقرات الاختبارات المخصصة لاختبار المفحوصين لمهن معينة والتي تتضمن فقراتها أجزاء لحقل الاختصاص وأخرى مشتركة، وتركز هذه الدراسة في جانب منها على فحص أثر حجم عينة المفحوصين على دقة تقدير معالم الفقرات، ولتحقيق هذا الغرض تم استخدام النموذج اللوجستي ثلاثي المعلمة لمعايرة ٣٦٠ فقرة باستخدام طريقة المحاكاة حيث تم افتراض أن أول ١٢٠ فقرة تشكل الجزء المشترك من الاختبار بينما باقي الفقرات شكلت أربعة اختبارات يتضمن كل منها ٦٠ فقرة، وقد تم اختيار حجم العينة ما بين ٢٥٠ مفحوص و ١٥٠٠ مفحوص، واستخدم معيار جذر معدل مربعات الفروق (RMSD) لمقارنة دقة تقديرات المعالم. أظهرت النتائج أن أخطاء التقدير لمعلمة الصعوبة ومعلمة التمييز تكون أكبر عندما يقل حجم عينة المفحوصين وذلك عند استخدام برمجية Bilog-MG₃.

وأجرى الدرايين (٢٠٠١) دراسة هدفت إلى البحث في فعالية النموذج اللوجستي أحادي المعلمة في دقة تقدير قدرة المفحوص، وصعوبة الفقرة باختلاف حجم عينة المفحوصين والاختبار. تكونت عينة المفحوصين في الدراسة من (٥٠، ١٠٠، ٥٠٠) مفحوص، وعدد فقرات الاختبار (٢٥، ٥٠، ٣٠٠) فقرة. أشارت نتائج الدراسة إلى وجود فروق ذات دلالة إحصائية لتفاعل كل من حجم عينة المفحوصين وطول الاختبار على دقة تقدير قدرة المفحوص. في حين لم تظهر النتائج وجود فروق إحصائية في دقة تقدير معلمة القدرة تعزى لحجم عينة المفحوصين، ووجود فروق إحصائية في دقة تقدير معلمة القدرة تعزى لطول الاختبار لصالح طول الاختبار الكبير، أما فيما يتعلق بدقة تقدير معلمة الصعوبة، فقد بينت النتائج بأن هناك فروقا إحصائية تعزى لتفاعل كل من حجم العينة وطول الاختبار، وفروق إحصائية في دقة تقدير معلمة الصعوبة تعزى لحجم عينة المفحوصين وطول الاختبار لصالح العينات الكبيرة للمفحوصين والاختبار. وفي دراسة عبانة (٢٠٠٤) التي هدفت إلى دراسة أثر حجم العينة وطريقة انتقائها وعدد الفقرات وطريقة انتقائها على دقة تقدير معالم الفقرة والقدرة وفق النموذج اللوجستي ثلاثي المعلمة. اختار اختبار القدرة العقلية المؤلف من أربع اختبارات فرعية هي: اختبار المفردات، واختبار المتشابهات، واختبار المتضادات، واختبار الحساب، وطبقت ثلاث اختبارات بأطوال مختلفة (٢٠، ٣٠، ٧١) فقرة على ثلاثة عينات من المفحوصين (٢٥٠، ٥٠٠، ١٠٠٠) مفحوص، وقد استخدمت برمجية Bilog لتحليل البيانات وتقدير معالم فقرات الاختبار وقدرة المفحوصين. أشارت النتائج إلى أن دقة تقديرات معالم الفقرة تزداد بزيادة حجم عينة المفحوصين، وأشارت النتائج إلى أن دقة تقديرات معلمة القدرة تزداد بزيادة طول الاختبار.

وأجرى فيتزباترك (Fitzpatrick, 2008) دراسة هدفت إلى معرفة أثر تخفيض طول اختبارات التهيئة على معدل إتقان الطلبة. تكونت الاختبارات في الدراسة من حجوم مختلفة (٥، ١٠، ١٥) فقرة طبقت على عينات مختلفة من المفحوصين تتراوح ما بين (٣٣٠٠) مفحوص إلى (٢٦٠٠٠) مفحوص، وتم تحليل البيانات وفق النموذج اللوجستي أحادي المعلمة. أظهرت نتائج الدراسة زيادة في تباينات الاختلافات في معدلات إتقان المفحوصين عند استخدام أطوال اختبارات قصيرة، كما أن هناك اختلافات وعدم استقرار معلمة الصعوبة عند استخدام اختبار طوله أقل من ١٥ فقرة، وأوصت الدراسة باستخدام اختبارات طولها أكبر من ١٥ فقرة لزيادة الاستقرار في تقدير معالم الفقرات. وأجرى دي لوتوري وهونغ (De la Torre & Hong, 2010) دراسة هدفت إلى التعرف على أثر حجم العينة على دقة تقدير معالم الفقرة والقدرة في اختبارات مطورة حسب نماذج نظرية IRT. واستخدمت الدراسة نموذج (IRT - HO) higher order item response theory في توليد مجموعة من البيانات ضمن ظروف اختبار مختلفة. حيث تم توليد البيانات لحجوم عينات مكونة من ٥٠٠ مفحوص و ١٠٠٠ مفحوص، وطول اختبارات مكونة من ١٠ فقرات و ٢٠ فقرة باستخدام طريقة مونتي كارلو من أجل التعرف على أثر حجم العينة في دقة تقدير معالم الفقرة والقدرة في الاختبار. وأشارت النتائج إلى أن حجم العينة وطول الاختبار يؤثران في تقدير معالم الفقرة ولصالح العينة ١٠٠٠ مفحوص والاختبار ٢٠ فقرة، حيث كان الوسط الحسابي للخطأ المعياري في تقدير معلمة الصعوبة يتراوح بين ٠,٠٦ و ٠,٠٢. وأشارت النتائج أيضاً إلى عدم تأثير دقة تقدير معلمة القدرة بحجم العينة وتأثيرها بطول الاختبار، حيث كانت الأخطاء المعيارية أقل ما يمكن عند استخدام طول اختبار ٢٠ فقرة مقارنة مع طول الاختبار ١٠ فقرات، وكان الوسط الحسابي للخطأ المعياري في تقدير معلمة القدرة يتراوح بين ٠,٥٨ و ٠,٤٨.

وفي دراسة قام بها جين وآخرون (Chen et al., ٢٠١٤) هدفت إلى تقييم نتائج تحليل نموذج راش باستخدام حجوم عينات صغيرة، ولغرض الدراسة تم استخدام حجوم عينات مختلفة مكونة من (٣٠، ٥٠، ١٠٠، ٢٥٠) مفحوص واختبار مكون من ١٠ فقرات. واستخدم برنامج Mplus لتحليل البيانات. أشارت نتائج التحليل أنه عند استخدام عينات صغيرة (٣٠، ٥٠) مفحوص يؤدي إلى أخطاء معيارية في تقدير المعالم أكبر من نتائج التحليل عند استخدام عينات (١٠٠، ٢٥٠) مفحوص. وأجرى جيانغ وآخرون (Jiang et al., ٢٠١٦) دراسة هدفت إلى معرفة الحجم المناسب لتقدير معالم الفقرة وفق نموذج الاستجابة المتدرجة متعددة الأبعاد Multidimensional Graded Response Model. وتم استخدام طريقة الأرجحية العظمى الهامشية لتقدير المعالم باستخدام برمجية flexMIRT، ولأغراض الدراسة تم توليد بيانات بحجوم عينات مختلفة مكونة من (٥٠٠، ١٠٠٠، ١٥٠٠، ٢٠٠٠) مفحوص، وأطوال اختبارات مختلفة مكونة من (٣٠، ٩٠، ٢٤٠) فقرة. أشارت النتائج إلى أن أقل حجم عينة يقدم تقديرات دقيقة لمعلمة الفقرة هو ٥٠٠ مفحوص عند استخدام اختبارات مكونة من ٣٠ فقرة و ٩٠ فقرة، وعند استخدام اختبار طوله ٢٤٠ فقرة فمن الضروري استخدام عينة مكونة من ١٠٠٠ مفحوص على الأقل، كما وأشارت النتائج إلى أن زيادة حجم العينة أكبر من ١٠٠٠ مفحوص لا يزيد من دقة تقدير المعالم.

وفي دراسة أجراها ساهين وانيل (Sahin & Anil, 2017) هدفت إلى معرفة أثر حجم العينة وطول الاختبار على دقة تقدير معالم نماذج نظرية استجابة الفقرة، وتم استخدام طريقة الأرجحية العظمى الهامشية لتقدير معالم الفقرة. واستخدم برمجية Xcalibre 4.1 لتقدير المعالم. ولتحقيق هدف الدراسة تم بناء ثلاث اختبارات لغوية ذات أطوال مختلفة مكونة من (١٠، ٢٠، ٣٠) فقرة وطبقت على تسعة حجوم عينات مختلفة مكونة من (١٥٠، ٢٥٠، ٣٥٠، ٥٠٠، ٧٥٠، ١٠٠٠، ٢٠٠٠، ٣٠٠٠، ٥٠٠٠) مفحوص. أشارت النتائج إلى أنه يمكن استخدام حجم عينة مكونة من ١٥٠ مفحوص على الأقل مع اختبارات مكونة من (١٠، ٢٠، ٣٠) فقرة لتقدير معلمة الصعوبة بدقة وفق النموذج أحادي المعلمة، وأشارت النتائج إلى أن أوساط الأخطاء المعيارية في تقدير معلمة الصعوبة تتناقص بزيادة حجم العينة حيث تراوحت بين ٠,٣٣ و ٠,٠١.

المعالجة الإحصائية

تم استخدام العديد من المعالجات الإحصائية للإجابة عن أسئلة الدراسة، وذلك وفق الخطوات التالية:

التحقق من افتراضات نموذج نظرية الاستجابة للفقرة أحادي المعلمة

أحادية البعد (Unidimensionality):

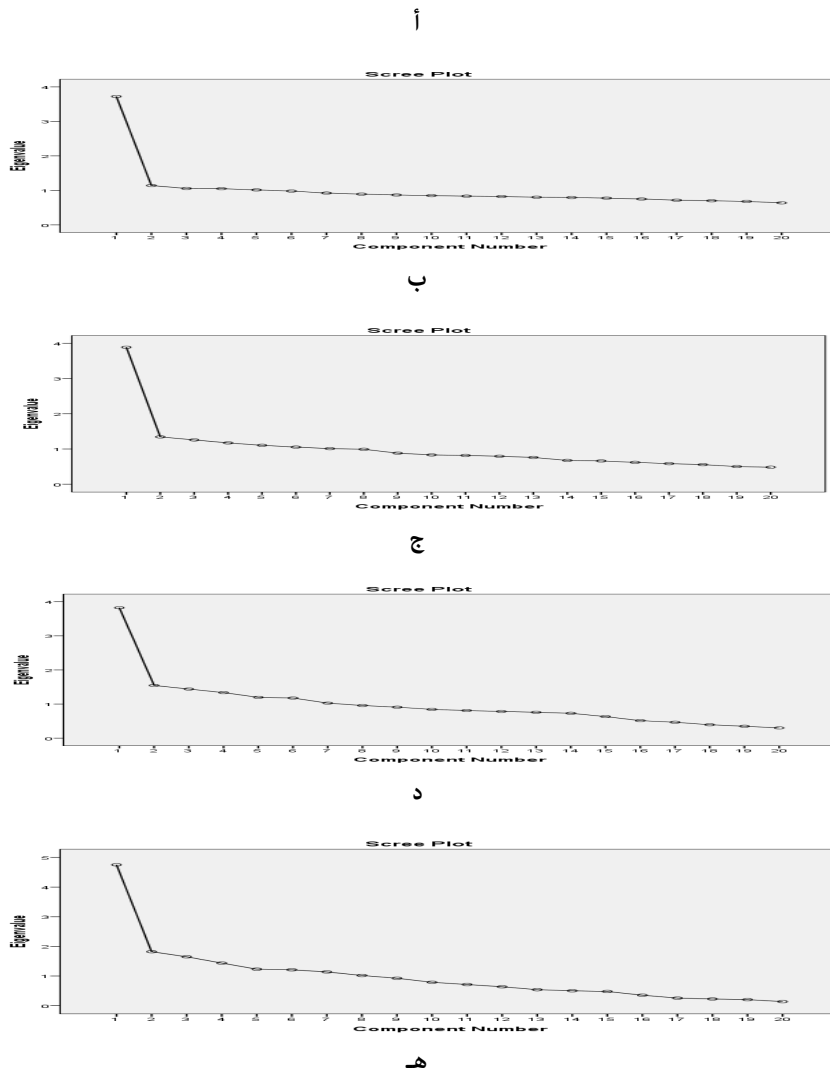
تم التحقق من افتراض أحادية البعد للاختبارات المستخدمة في هذه الدراسة بعدة مؤشرات اعتمدت على التحليل العاملي للمكونات الأساسية (Principal Components)، ومنها: نسبة التباين المفسر للعامل الأول، ونسبة الفرق بين الجذر الكامن الأول والجذر الكامن الثاني إلى الفرق بين الجذر الكامن الثاني والجذر الكامن الثالث، ونسبة الجذر الكامن الأول إلى الجذر الكامن الثاني. ودراسة العلاقة بين الجذر الكامن الثاني وبقية الجذور التي تليه بالإضافة إلى رسم بياني للجذور الكامنة للعوامل المكونة للاختبار (Scree Plot) باستخدام برمجية SPSS. حيث تبين من تحليل البيانات تحقق افتراض أحادية البعد لجميع الاختبارات. ويبين جدول ١ قيم الجذور الكامنة ونسبة التباين المفسر للاختبار المكون من ٢٠ فقرة والذي طبق على عينات مختلفة مكونة من (٥٠، ١٠٠، ٢٥٠، ٥٠٠، ١٠٠٠) مفحوص.

جدول ١ نتائج التحليل العاملي ل فقرات الاختبار ٢٠ ولعينات الدراسة الخمسة

حجم العينة	العامل	قيمة الجذر الكامن	التباين المفسر %	الجذر الكامن الأول - الجذر الكامن الثاني	الجذر الكامن الثاني - الجذر الكامن الثالث
١٠٠٠	الأول	٣,٨٥٦	٣٨,٣٦٢		
	الثاني	١,١٠٩	٥,٥٤٧		
	الثالث	١,٠٣٤	٥,١٦٩		
	الرابع	١,٠١٤	٥,٠٧١	٣٦,٦٢٧	٣,٤٧٧
٥٠٠	الأول	٣,٥٠٢	٣٦,٤٤٩		
	الثاني	١,٢٣٤	٦,١٢١		
	الثالث	١,١٢١	٥,٦٠٣		
	الرابع	١,٠٨٤	٥,٤١٩		
	الخامس	١,٠٦١	٥,٣٠٥		
	السادس	١,٠٠١	٥,٠٥٢		
	السابع	١,٠٠٥	٥,٠٢٦	٢٢,١١٧	٢,٨٦١
٢٥٠	الأول	٣,٨٨٩	٣٧,٠٧٧		
	الثاني	١,٣٤٤	٦,٧٢٠		
	الثالث	١,٢٥٨	٦,٢٩١		
	الرابع	١,١٧٣	٥,٨٦٥		
	الخامس	١,٠٥٠	٥,٥٢٦		
	السادس	١,٠٠٥	٥,٢٧٣		
	السابع	١,٠١١	٥,٠٠٥	٢٩,٥٩٣	٢,٨٩٤
١٠٠	الأول	٣,٨٢٤	٣١,٩١٤		
	الثاني	١,٥٤٦	٧,٧٢٨		
	الثالث	١,٤٤٠	٧,١٩٨		
	الرابع	١,٣٣٦	٦,٦٧٩		
	الخامس	١,١٩٧	٥,٩٨٦		
	السادس	١,١٧٧	٥,٨٨٤		
	السابع	١,٠٢٨	٥,١٣٩	٢١,٤٩١	٢,٤٧٤
٥٠	الأول	٤,٧٥٣	٢٩,٨٣٤		
	الثاني	١,٨٢١	٩,١٠٣		
	الثالث	١,٦٤٨	٨,٢٤٠		
	الرابع	١,٤٣٦	٧,١٨٢		
	الخامس	١,٢٣٠	٦,١٥٢		
	السادس	١,٢٠٩	٦,٠٤٧		
	السابع	١,١٣٧	٥,٦٨٥		
	الثامن	١,٠١٨	٥,٠٨٩	١٦,٩٤٨	٢,٦١٠

يتبين من جدول ١ أن قيم الجذر الكامن للعامل الأول لجميع أحجام العينات في الاختبار المكون من ٢٠ فقرة يفسر أعلى من ٢٠% من التباين الكلي، وهذا يطابق ما اقترحه ريكا (Reckase, 1979) من أنه إذا فسّر العامل الأول على الأقل نسبة ٢٠% من التباين المفسر فذلك يدل على تحقق افتراض أحادية البعد. وبلغت نسبة الجذر الكامن الأول إلى الجذر الكامن الثاني أكبر من ٢ لجميع أحجام العينات في الاختبار المكون من ٢٠ فقرة، وهذا يطابق ما اقترحه هاتي (Hattie, 1985) من أن النسبة المرتفعة للجذر الكامن الأول إلى الجذر الكامن الثاني يعد دليلاً على افتراض أحادية البعد، وحدد ليناس (Linacre, 2008) تلك النسبة بحيث تكون أكبر أو تساوي القيمة ٢ للدلالة على أحادية البعد. كما بلغت قيمة الفرق بين الجذرين الكامنين الأول والثاني إلى الفرق بين الجذرين الكامنين الثاني والثالث أعلى من ٦ لجميع أحجام العينات للاختبار المكون من ٢٠ فقرة، وهذا يطابق ما اقترحه هاتي (Hattie, 1985) بأن نسبة طرح الجذر الكامن الثاني من الجذر الكامن الأول إلى طرح الجذر الكامن الثالث من الجذر الكامن الثاني مرتفعة، وجميعها أعلى من القيمة (٦) بما يفيد من تحقق افتراض أحادي البعد لجميع اختبارات الدراسة في العينات الخمسة. ويتعزز افتراض أحادية البعد من خلال اختبار فحص العوامل الذي يظهر في شكل ١. وتبين الأشكال (أ، ب، ج، د، هـ) في شكل ١ مخططات القيم الفارزة للجذور الكامنة (Scree Plot) التي أظهرت أحادية البعد للاختبار المكون من ٢٠ فقرة والذي طبق على أحجام عينات الدراسة.

الشكل (١) : مخطط قيم الجذور الكامنة في الاختبار ٢٠ فقرة والمطبق على عينات الدراسة الخمسة.



الاستقلال الموضوعي (Local Independence):

تم استخدام حزمة sirt في برمجية لغة R لفحص الاستقلال الموضوعي لفقرات الاختبارات والمفحوصين، من خلال مؤشر Q3 باستخدام حزمة المزاجية (Pairwise Package)، التي طورها العالم (Heine & Tarnai, 2015)، وبينت نتائج التحليل تحقق افتراض الاستقلال الموضوعي، وبين جدول ٢ نتائج اختبار الاستقلال الموضوعي لفقرات الاختبارات ولعينات الدراسة الخمسة.

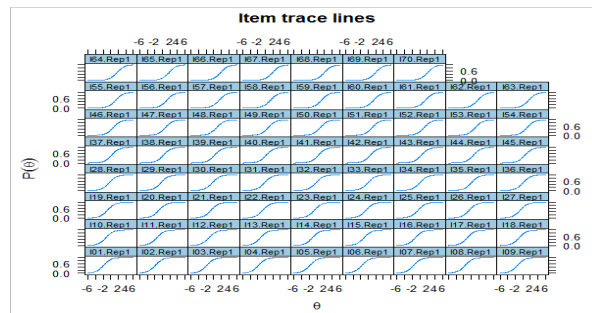
جدول ٢ نتائج اختبار الاستقلال الموضوعي لفقرات الاختبارات الأربعة ولعينات الدراسة الخمسة.

الوسط الحسابي لقيم Q3 للاختبار				
العينة	٢٠	٣٠	٤٠	٧٠
١٠٠٠	-0.029	-0.011	-0.005	0.007
٥٠٠	-0.019	-0.003	0.006	0.016
٢٥٠	-0.013	-0.007	0.016	0.029
١٠٠	0.024	0.025	0.045	0.054
٥٠	0.048	0.065	0.070	0.087

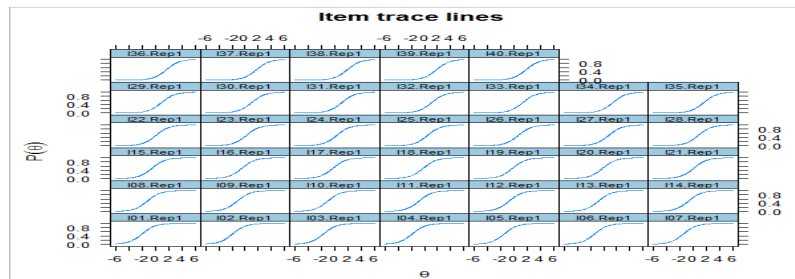
يلاحظ من جدول ٢ أن قيم الوسط الحسابي لـ Q3 لجميع الاختبارات الأربعة ولكل العينات الخمسة تراوحت ما بين (٠,٢٩، -٠,٠٨٧)، وبينت نتائج التحليل في جدول ٢ أن قيم معامل Q3 لم تتجاوز مستوى انتهاك افتراض الاستقلال الموضوعي الذي حدد ما بين (-٠,٢، ٠,٢) حسب ما اقترحه جين وثيسن (Chen & Thissen, 1997).

افتراض منحني خصائص الفقرة (ICC):
 للتحقق من افتراض اطراديه السمة الكامنة وفقاً للنموذج أحادي المعلمة في نظرية IRT، فقد تم رسم منحني خصائص الفقرة لكل فقرة من فقرات الاختبارات الأربعة التي طبقت على عينات الدراسة الخمسة باستخدام حزمة mirt في برمجة لغة R، حيث أظهرت أشكال منحنيات خصائص الفقرة أن زيادة قدرة المفحوص يقابلها زيادة احتمالية الاستجابة الصحيحة للمفحوص على الفقرة في جميع فقرات الاختبارات الأربعة المطبقة على جميع عينات الدراسة الخمسة، وبينت أشكال منحنيات خصائص الفقرة أنها لا تختلف باختلاف طول الاختبار وحجم العينة، وتبين الأشكال (أ، ب، ج، د) في شكل ٢ منحنيات خصائص الفقرة لاختبارات الدراسة الأربعة عند حجم العينة ١٠٠٠ مفحوص.

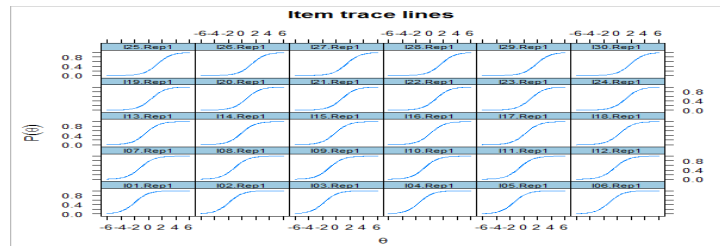
شكل (٢) منحنيات خصائص الفقرة لاختبارات الدراسة عند حجم العينة ١٠٠٠ مفحوص



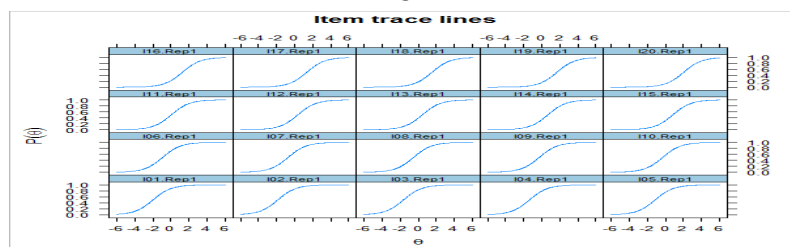
أ



ب



ج



د

تم التحقق من ثبات الاختبارات الأربعة التي طبقت على عينات الدراسة الخمسة من خلال حساب معامل أوميغا ومعامل الثبات التجريبي باستخدام برمجية لغة وبيبي جدول ٣ قيم معاملات الثبات لجميع الاختبارات في عينات الدراسة.

جدول ٣ قيم معامل ثبات أوميغا ومعامل الثبات التجريبي للاختبارات الأربعة في عينات الدراسة الخمسة

معامل ثبات أوميغا W و (معامل الثبات الامبريقي) للاختبار				
طول الاختبار				
حجم العينة	٢٠	٣٠	٤٠	٧٠
١٠٠٠	٠,٧٦ (٠,٨٠)*	٠,٨٢ (٠,٨٥)	٠,٨٨ (٠,٨٩)	٠,٩١ (٠,٩٢)
٥٠٠	٠,٨٠ (٠,٨١)	٠,٨٤ (٠,٨٦)	٠,٨٦ (٠,٨٧)	٠,٩٣ (٠,٩٣)
٢٥٠	٠,٧٣ (٠,٧٨)	٠,٨٥ (٠,٨٧)	٠,٨٦ (٠,٨٨)	٠,٩٣ (٠,٩٣)
١٠٠	٠,٨٠ (٠,٨٢)	٠,٧٦ (٠,٨٠)	٠,٩٠ (٠,٩١)	٠,٩١ (٠,٩٢)
٥٠	٠,٨٥ (٠,٨٥)	٠,٨٣ (٠,٨٣)	٠,٨٨ (٠,٨٩)	٠,٩٤ (٠,٩٤)

*: القيم داخل القوس تمثل معامل الثبات الامبريقي.

يتضح من جدول ٣ أن قيم معامل ثبات أوميغا قد تراوحت بين ٠,٧٣ و ٠,٩٤، ويتبين أيضاً أن قيم معاملات الثبات التجريبي قد تراوحت بين ٠,٨٠ و ٠,٩٤، وهذا يبين أن قيم معامل ثبات أوميغا والثبات التجريبي كانا مرتفعان لجميع اختبارات الدراسة والمطبقة على حجوم عينات مختلفة.

٢. تم تقدير معلمة صعوبة الفقرة بطريقة المزاجية باستخدام حزمة المزاجية (Pairwise Package)، وللحكم على دقة تقدير معلمة الصعوبة للفقرات تم استخدام مؤشر

الخطأ المعياري في التقدير باستخدام طريقة البوتستراب لفحص دقة تقدير معلمة صعوبة الفقرة باستخدام دالة pairSE في برمجية لغة R.

٣. تقدير معلمة قدرات الأفراد وفق النموذج أحادي المعلمة باستخدام طريقة الأرجحية الموزونة باستخدام حزمة (pairwise package).

نتائج الدراسة ومناقشتها

أولاً: النتائج الخاصة بسؤال الدراسة الأول الذي نص على: هل تختلف دقة تقدير معلمة صعوبة الفقرة باستخدام طريقة المزاجية باختلاف حجم عينة المفحوصين وطول الاختبار وفق النموذج أحادي المعلمة ؟

للإجابة عن السؤال الأول في الدراسة، تم تقدير أوساط قيم معلمة الصعوبة للفقرات باختلاف حجم عينة المفحوصين وطول الاختبار باستخدام طريقة المزاجية وفق النموذج أحادي المعلمة، وإيجاد قيم أوساط الجذر التربيعي لمعدل مربعات الأخطاء (RMSE) باستخدام طريقة البوتستراب في برمجية لغة R، وتم تحديد عدد عينة البوتستراب بحجم العينة الكلية (N). يبين جدول ٤ ملخص نتائج صعوبة الفقرات والخطأ المعياري في تقدير معلمة الصعوبة باختلاف حجم العينة وطول الاختبار.

جدول ٤ أعلى قيمة وأدنى قيمة لمعلمة صعوبة الفقرة وقيم الوسط الحسابي للخطأ المعياري في تقدير صعوبة الفقرات باختلاف عينة المفحوصين وطول الاختبار

طول الاختبار												
٢٠			٣٠			٤٠			٧٠			حجم العينة
الوسط الحسابي للخطأ المعياري	أدنى قيمة	أعلى قيمة	الوسط الحسابي للخطأ المعياري	أدنى قيمة	أعلى قيمة	الوسط الحسابي للخطأ المعياري	أدنى قيمة	أعلى قيمة	الوسط الحسابي للخطأ المعياري	أدنى قيمة	أعلى قيمة	
٠,٠٠٢	-١,٩٧	٢,٠٧	٠,٠٠٢	-١,٩٨	١,٨٨	٠,٠٠٢	-٢,٠٨	١,٨٦	٠,٠٠٢	-٢,٠١	٢,١٠	١٠٠٠
٠,٠٠٥	-١,٩٤	٢,٠٧	٠,٠٠٥	-١,٩٨	١,٩٩	٠,٠٠٥	-١,٩٦	٢,٠٣	٠,٠٠٥	-٢,٠١	٢,٠١	٥٠٠
٠,٠١	-٢,٣٧	١,٩٥	٠,٠١	-٢,٠٣	١,٨٩	٠,٠١	-٢,٣٥	٢,٢٧	٠,٠١	-٢,١٧	٢,٢٣	٢٥٠
٠,٠٣	-٢	٣	٠,٠٣	-٢,٢١	٢,٤٤	٠,٠٣	-٢,٢١	٣,٢٢	٠,٠٣	-٢,٢٥	٢,٤٩	١٠٠
٠,٠٦	-٢,٥٩	٢,٣٤	٠,٠٥	-٢,١٩	١,٨٤	٠,٠٥	-٢,٤٤	٢,١٨	٠,٠٥	-٢,٤٤	٣,٦٢	٥٠

يلاحظ من جدول ٤ أن الوسط الحسابي للأخطاء المعيارية لتقديرات صعوبة الفقرات كانت صغيرة باختلاف حجم العينة وطول الاختبار، إذ تراوحت قيمة الوسط الحسابي بين ٠,٠٠٢ و ٠,٠٠٦، وهذا يدل على دقة تقدير صعوبة الفقرة باستخدام طريقة المزاجية وخاصة عند استخدام عينات صغيرة من طول الاختبار وحجم عينة المفحوصين، وأظهرت النتائج أن تقديرات معلمة الصعوبة باستخدام طريقة المزاجية كانت تعطي أخطاء معيارية قليلة في جميع اختبارات الدراسة والمطبقة على حجوم عينات مختلفة وقد تراوحت بين ٠,٠٠٢ و ٠,٠٥، ويفسر الباحثان هذه النتيجة بأن استخدام حزمة المزاجية الحديثة لتقدير معلمة الصعوبة والخطأ المعياري في تقديرها أدى إلى الحصول على نتائج دقيقة في تقدير معلمة صعوبة الفقرة وفق النموذج أحادي المعلمة وهذا يؤكد ما أشار إليه زانغ (Zhang, ٢٠١٤) إلى أن طريقة البوتستراب هي أفضل طريقة لتقدير الأخطاء المعيارية لتقديرات المعلمة وخاصة عند استخدام حجوم عينات صغيرة.

ويلاحظ من جدول ٤ ثبات قيمة الوسط الحسابي للخطأ المعياري في التقدير بتغير طول الاختبار وثبات حجم العينة، ويعزو الباحثان هذه النتيجة إلى أن تقدير معلمة صعوبة فقرة معينة لا تتأثر بأداء المفحوصين على فقرة أخرى وصعوبتها، وهذا يدل على جودة البيانات المولدة في هذه الدراسة ومطابقتها للنموذج أحادي المعلمة. وبين جدول ٥ نتائج تحليل التباين الثنائي لقيم الأخطاء المعيارية لتقديرات معلمة الصعوبة.

جدول ٥ نتائج تحليل التباين الثنائي لقيم الأخطاء المعيارية لتقديرات معلمة الصعوبة المقدرتها باختلاف حجم العينة وطول الاختبار

مصدر التباين	مجموع المربعات	درجات الحرية	متوسط مجموع المربعات	قيمة ف المحسوبة	الدلالة الإحصائية
طول الاختبار	٠,٠٠٠١٣	٣	٤,٣٣E-٥	٠,٨٢	٠,٤٩
حجم العينة	٠,٢٣٩٧٠	٤	٠,٠٦	١١٤٠,٦٨	٠,٠٠
طول الاختبار × حجم العينة	٠,٠٠٠٦٢	١٢	٥,١٧E-٥	٠,٩٨	٠,٤٦
الخطأ	٠,٠٤١٠٤	٧٨٠	٤,٣٦E-٥		
الكل	٠,٢٣٦٧٧٣	٧٩٩			

يلاحظ من جدول ٥ عدم وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات الصعوبة المقدرتها باستخدام طريقة المزاجية تعزى لمتغير طول الاختبار، والتفاعل بين حجم العينة وطول الاختبار؛ مما يعني عدم تأثر دقة تقدير صعوبة الفقرات بطول الاختبار والتفاعل بين حجم العينة وطول الاختبار، كما يبين جدول ٥ وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات الصعوبة تعزى لمتغير حجم العينة؛ أي أنه يوجد تأثيراً لمتغير حجم العينة في دقة تقدير معلمة الصعوبة. وبين جدول ٦ نتائج اختبار شافيه (Scheffe) للمقارنات البعدية لتحديد لصالح أي من حجومات العينات قد كانت الفروق بينها دالة إحصائياً.

جدول ٦ نتائج اختبار شافيه للمقارنات البعدية لمتغير حجم العينة

العينة أ	العينة ج	وسط الاختلاف بين العينتين	الدلالة الإحصائية
٥٠	١٠٠	٠,٢٩١٤٤	٠,٠٠٠
	٢٥٠	٠,٠٤٤٠٨٨	٠,٠٠٠
	٥٠٠	٠,٠٤٩٣١٣	٠,٠٠٠
	١٠٠٠	٠,٠٥١٩١٣	٠,٠٠٠
١٠٠	٢٥٠	٠,٠١٤٩٤٤	٠,٠٠٠
	٥٠٠	٠,٠٢٠١٦٩	٠,٠٠٠
	١٠٠٠	٠,٠٢٢٧٦٩	٠,٠٠٠
٢٥٠	٥٠٠	٠,٠٠٥٣٢٥	٠,٠٠٠
	١٠٠٠	٠,٠٠٧٨٢٥	٠,٠٠٠
٥٠٠	١٠٠٠	٠,٠٠٢٦	٠,٣٧

يبين جدول ٦ وجود فروق ذات دلالة إحصائية بين حجم العينة ٥٠ من جهة وحجوم العينات (١٠٠، ٢٥٠، ٥٠٠، ١٠٠٠) لمفحوص على التوالي، ولصالح حجومات العينات الأربعة، كما تبين وجود فروق ذات دلالة إحصائية بين حجم العينة ١٠٠ من جهة وحجوم العينات (٢٥٠، ٥٠٠، ١٠٠٠) لمفحوص على التوالي، ولصالح حجومات العينات الثلاثة، كما تبين وجود فروق ذات دلالة إحصائية بين حجم العينة ٢٥٠ من جهة وحجوم العينات (٥٠٠، ١٠٠٠) لمفحوص على التوالي، ولصالح حجم العينة ٥٠٠ و ١٠٠٠، كما تبين أيضاً وجود فروق ذات دلالة إحصائية بين حجم العينة ٥٠٠ من جهة وحجوم العينات (١٠٠٠) لمفحوص، ولصالح حجم العينة ١٠٠٠ لمفحوص.

أشارت النتائج أنه يوجد فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات معلمة الصعوبة المقدرتها باستخدام طريقة المزاجية تعزى لمتغير حجم العينة لصالح حجومات العينات الكبيرة، مما يعني تأثر دقة تقدير معلمة صعوبة الفقرات بحجم عينة المفحوصين ولصالح العينات الكبيرة، كما أظهرت نتائج المقارنات البعدية باستخدام اختبار شافيه أنه كلما زاد حجم العينة قلت الأخطاء المعيارية في تقدير صعوبة الفقرة. ويمكن تفسير هذه النتائج بأن تقدير معلمة صعوبة الفقرة في نظرية IRT يحتاج إلى عدد كاف من المفحوصين فزيادة حجم العينة يزيد من دقة تقدير معلمة الصعوبة، حيث أن زيادة حجم عينة المفحوصين المطبق عليهم فقرات الاختبار يؤدي إلى تقديرات دقيقة لمعلمة صعوبة الفقرات وهذا ما أشار إليه هامبلتون (Hambleton, ١٩٨٩) بأن نظرية IRT تحتاج إلى حجومات عينات كبيرة للحصول على تقديرات دقيقة لمعلمة الفقرة وقدرة المفحوصين.

تتفق هذه النتائج مع نتائج دراسة ساهين وانيل (Sahin & Anil, 2017) التي أشارت إلى أن دقة التقدير لمعلمة صعوبة الفقرة تزداد بزيادة حجم عينة المفحوصين. كما واتفقت نتائج الدراسة مع نتائج دراسة عباينة (٢٠٠٤) حيث أشارت النتائج إلى أن دقة تقديرات معلمة الفقرة تزداد بزيادة حجم عينة المفحوصين، إذ أن متوسط الأخطاء المعيارية لتقديرات معلمة الفقرة تكون أقل ما يمكن عند استخدام حجم عينة ١٠٠٠ لمفحوص. كما وتتفق هذه النتائج مع نتائج دراسة هيونغ وآخرون (Huang et al, 2001) والتي أشارت إلى أن أخطاء التقدير لمعلمة الصعوبة ومعلمة التمييز تكون أكبر عندما يقل حجم عينة المفحوصين. وتتفق هذه النتائج مع نتيجة دراسة فينش وفرينش (Finch & French, 2019) التي أشارت أن بزيادة حجم عينة المفحوصين تقل الأخطاء المعيارية في تقدير معلمة الصعوبة، وإن زيادة طول الاختبار ليس بالضرورة أن يؤدي إلى دقة أفضل في تقدير معلمة الصعوبة.

كما أشارت النتائج إلى عدم وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات معلمة الصعوبة المقدر باستخدام طريقة المزاجية تعزى لمتغير طول الاختبار. ويعزو الباحثان هذه النتائج بأن تقدير معالم الفقرة في نظرية IRT لا يحتاج إلى عدد كبير من الفقرات، أي أن دقة تقدير معلمة الصعوبة لا تتأثر بطول الاختبار وهذا ما يؤكد ما رايت و ستون (Wright & Stone, 1979) إلى عدم وجود طول اختبار محدد يقدم أقل خطأ معياري في تقدير معالم الفقرة. واتفقت هذه النتيجة مع ما توصلت إليه دراسة فيتزباترك (Fitzpatrick, 2008) أن هناك اختلافات وعدم استقرار في معلمة الصعوبة عند استخدام اختبار طوله أقل من 15 فقرة، وأوصت الدراسة باستخدام اختبارات طولها أكبر من 15 فقرة لزيادة الاستقرار في تقدير معلمة صعوبة الفقرات.

كما وأشارت النتائج أيضاً إلى عدم وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات معلمة الصعوبة المقدر باستخدام طريقة المزاجية تعزى للتفاعل بين حجم العينة وطول الاختبار، مما يعني عدم تأثر دقة تقدير معلمة الصعوبة بالتفاعل بين حجم عينة المفحوصين وطول الاختبار. واختلفت هذه النتيجة مع نتيجة دراسة الدرايب (2001) بأن هناك فروقا إحصائية تعزى لتفاعل كل من حجم العينة وطول الاختبار على دقة تقدير معلمة صعوبة الفقرة. إن نتائج هذه الدراسة مشجعة لاستخدام طريقة المزاجية في تقدير معلمة الصعوبة مع حجوم عينات صغيرة؛ لأن الكثير من المجالات العملية وخاصة التربوية تتطلب استخدام حجوم عينات صغيرة، مما يجعل جمع عينات كبيرة بما يكفي لتقدير معالم نماذج نظرية IRT أمراً صعباً جداً. حيث أشارت نتائج هذه الدراسة إلى أن طريقة المزاجية تعطي دقة في تقدير معلمة الصعوبة عند استخدامها مع العينات الصغيرة والكبيرة على حد سواء.

ثانياً: النتائج الخاصة بسؤال الدراسة الثاني الذي نص على: هل تختلف دقة تقير معلمة القدرة باستخدام طريقة الأرجحية الموزونة باختلاف حجم عينة المفحوصين وطول الاختبار وفق النموذج أحادي المعلمة ؟

للإجابة عن السؤال الثاني في الدراسة، فقد قدرت قيم معلمة القدرة للمفحوصين باختلاف حجم عينة المفحوصين وطول الاختبار باستخدام طريقة تقدير الأرجحية الموزونة وفق النموذج أحادي المعلمة، وإيجاد قيم الأخطاء المعيارية في تقدير القدرة، وذلك باستخدام برمجية لغة R. ويبين جدول 7 ملخص نتائج قدرات المفحوصين والخطأ المعياري في تقدير القدرة باختلاف حجم العينة وطول الاختبار.

جدول 7 أعلى قيمة وأدنى قيمة لمعلمة القدرة وقيم وسطها الحسابي وخطأ تقديرها باختلاف عينة المفحوصين وطول الاختبار

طول الاختبار																
20				30				40				70				حجم العينة
وسط الخطأ	وسط القدرة	أدنى قيمة	أعلى قيمة	وسط الخطأ	وسط القدرة	أدنى قيمة	أعلى قيمة	وسط الخطأ	وسط القدرة	أدنى قيمة	أعلى قيمة	وسط الخطأ	وسط القدرة	أدنى قيمة	أعلى قيمة	
0.56	0.01	-4.38	4.38	0.44	0	-3.58	3.56	0.39	-0.02	-3.88	3.87	0.29	-0.06	-3.26	3.26	1000
0.56	-0.05	-3.39	3.20	0.46	0.09	-4.79	4.78	0.39	0.06	-2.91	2.62	0.30	0.05	-3.28	3.83	500
0.55	0.08	-2.10	3.17	0.46	0.14	-3.59	4.74	0.39	-0.04	-2.62	2.93	0.30	0.06	-2.67	3.57	250
0.58	-0.03	-3.27	4.73	0.44	-0.06	-2.29	2.29	0.40	0.01	-2.66	3.45	0.30	-0.18	-2.38	1.90	100
0.62	-0.41	-4.59	2.24	0.45	0.17	-2.07	2.66	0.40	0.32	-1.99	2.44	0.30	0.10	-2.06	2.95	50

يلاحظ من جدول 7 أن قيمة الوسط الحسابي للأخطاء المعيارية تراوحت بين 0.29 و 0.62، باختلاف طول الاختبار وحجم عينة المفحوصين، كما أن أقل قيمة للخطأ المعياري عند حجم عينة 1000 مفحوص وطول اختبار 70 فقرة وأعلى قيمة عند حجم عينة 50 وطول اختبار 20، أي أن كلما زاد طول الاختبار قل الخطأ المعياري في تقدير معلمة القدرة. وقام الباحثان بإجراء تحليل تباين ثنائي Interaction ANOVA باستخدام برمجية (SPSS version 22) لقيم الأخطاء المعيارية لتقديرات معلمة قدرة المفحوصين المقدر باختلاف متغيري (حجم العينة، طول الاختبار)، ويبين جدول 8 نتائج تحليل التباين الثنائي.

جدول 8 نتائج تحليل التباين الثنائي لقيم الأخطاء المعيارية لتقديرات معلمة قدرة المفحوصين المقدر باختلاف حجم العينة وطول الاختبار

الدلالة الإحصائية	قيمة ف المحسوبة	متوسط مجموع المربعات	درجات الحرية	مجموع المربعات	مصدر التباين
0.00	2352.54	9.136	3	27.407	طول الاختبار
0.00	7.045	0.027	4	0.109	حجم العينة
0.00	5.56	0.022	12	0.259	طول الاختبار * حجم العينة
		0.004	7580	29.44	الخطأ
			7599	101.09	الكل

يلاحظ من جدول 8 وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات معلمة القدرة المقدر باستخدام طريقة الأرجحية الموزونة تعزى لمتغيري طول الاختبار وحجم العينة، والتفاعل بين حجم العينة وطول الاختبار. ولتحديد لصالح أي من حجوم العينات وطول الاختبار قد كانت الفروق بينها دالة إحصائياً، فقد قام الباحثان بإجراء مقارنات بعدية باستخدام اختبار شافيه (Scheffe). ويبين جدول 9 نتائج اختبار شافيه للمقارنات البعدية لمتغير حجم العينة.

جدول ٩ نتائج اختبار شافيه للمقارنات البعدية لمتغير حجم العينة

العينة أ	العينة ج	وسط الاختلاف بين العينتين	الدلالة الإحصائية
٥٠	١٠٠	٠,١٢٤	٠,٣٦١
	٢٥٠	٠,١٨٦	٠,٠٠٥
	٥٠٠	٠,١٦٧	٠,٠١١
	١٠٠٠	٠,٢٠٧	٠,٠٠٠
١٠٠	٢٥٠	٠,٠٦٢	٠,٥٨٥
	٥٠٠	٠,٠٤٣	٠,٨١٥
	١٠٠٠	٠,٠٨٣	٠,١٦٧
٢٥٠	٥٠٠	-٠,٠١٩	٠,٩٥٨
	١٠٠٠	٠,٠٢١	٠,٩٢٣
٥٠٠	١٠٠٠	٠,٠٠٤	٠,٢٣٠

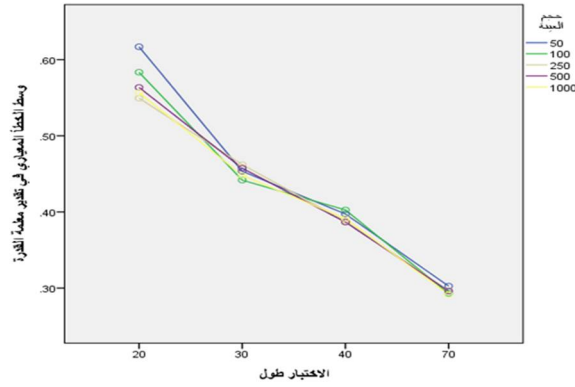
يبين جدول ٩ وجود فروق ذات دلالة إحصائية بين حجم العينة ٥٠ من جهة وحجم العينات (١٠٠، ٢٥٠، ٥٠٠، ١٠٠٠) مفحوص على التوالي، ولصالح حجوم العينات (٢٥٠، ٥٠٠، ١٠٠٠) مفحوص، كما تبين عدم وجود فروق ذات دلالة إحصائية بين حجم العينة ١٠٠ من جهة وحجم العينات (٥٠، ٢٥٠، ٥٠٠، ١٠٠٠) مفحوص على التوالي، كما تبين عدم وجود فروق ذات دلالة إحصائية بين حجم العينة ٢٥٠ من جهة وحجم العينات (٥٠٠، ١٠٠٠) مفحوص على التوالي، كما تبين أيضاً عدم وجود فروق ذات دلالة إحصائية بين حجم العينة ٥٠٠ من جهة وحجم العينة ١٠٠٠ مفحوص، مما يعني عدم تأثر دقة تقدير معلمة القدرة عند استخدام عينات ١٠٠ مفحوص فأكثر. ويمكن عزو هذه النتائج إلى أن دقة تقدير قدرات المفحوصين تتسم بالاستقرار عندما يتم استخدام عينات بأحجام كبيرة، كما تتأثر الدقة في تقدير قدرات المفحوصين بوضع مطابقة البيانات للنموذج اللوجستي أحادي المعلمة ويكون عندئذ من غير المؤكد الحصول على دقة أكبر عند زيادة حجوم عينات كبيرة. كما وتشير النتائج إلى أن استخدام عينات صغيرة ٥٠ مفحوص يؤثر سلباً على دقة تقدير قدراتهم وقد يستلزم زيادة دقة تقدير قدرات المفحوصين بشكل طفيف زيادة العينة المستخدمة في تقدير معلمة القدرة، وهذا ما أشار إليه هامبلتون وكوك (Hambleton & Cook, ١٩٨٣) إلى أن زيادة حجم العينة وطول الاختبار تزيد دقة التقدير لقدرة المفحوص. ويبين جدول ١٠ نتائج اختبار شافيه للمقارنات البعدية لمتغير طول الاختبار.

جدول ١٠ نتائج اختبار شافيه للمقارنات البعدية لمتغير طول الاختبار

الاختبار أ	الاختبار ج	وسط الاختلاف بين العينتين	الدلالة الإحصائية
٢٠	٣٠	٠,١٠٨٣	٠,٠٠٠
	٤٠	٠,١٧٠١	٠,٠٠٠
	٧٠	٠,٢٦٦٧	٠,٠٠٠
٣٠	٤٠	٠,٠٦١٨	٠,٠٠٠
	٧٠	٠,١٥٨٤	٠,٠٠٠
٤٠	٧٠	٠,٠٩٦٦	٠,٠٠٠

يبين جدول ١٠ وجود فروق ذات دلالة إحصائية بين طول الاختبار ٢٠ فقرة من جهة وأطوال الاختبارات (٣٠، ٤٠، ٧٠) فقرة من جهة أخرى، ولصالح أطوال الاختبارات (٣٠، ٤٠، ٧٠) فقرة، كما تبين وجود فروق ذات دلالة إحصائية بين طول الاختبار ٣٠ فقرة من جهة وأطوال الاختبارات (٤٠، ٧٠) فقرة على التوالي، ولصالح أطوال الاختبارات (٤٠، ٧٠) فقرة، كما تبين وجود فروق ذات دلالة إحصائية بين طول الاختبار ٤٠ فقرة من جهة وطول الاختبار ٧٠ فقرة، ولصالح طول الاختبار ٧٠ فقرة، أي كلما زاد طول الاختبار تزداد دقة تقدير معلمة القدرة، مما يشير إلى زيادة طول الاختبار ساهم بشكل واضح في خفض الخطأ المعياري لتقديرات معلمة القدرة، ويعزو الباحثان هذه النتائج إلى أن زيادة طول الاختبار يؤدي إلى اجابة المفحوص على عدد كبير من الفقرات تؤدي إلى تقدير قدرته بشكل أفضل، وهذا ما يؤكد لورد (Lord, ١٩٨٠) إلى أن طول الاختبار يؤثر على جودة تقدير المعالم. تبين من نتائج التحليل على أن العامل الأهم في التأثير على دقة تقديرات معلمة القدرة وفق النموذج اللوجستي أحادي المعلمة هو طول الاختبار وليس حجم العينة عندما يكون حجم العينة متوسطاً (١٠٠، ٢٥٠) مفحوص أو كبيراً (٥٠٠، ١٠٠٠) مفحوص، حيث أشارت النتائج إلى أنه عند استخدام حجم عينة ١٠٠ مفحوص فأكثر فإن الوسط الحسابي لتقدير معلمة القدرة يبقى ثابتاً مع ثبات حجم العينة، وتتفق هذه النتيجة مع نتيجة دراسة عبابنة (٢٠٠٤) التي أشارت إلى أن دقة تقديرات معلمة القدرة تزداد بزيادة طول الاختبار. وأشارت النتائج أيضاً إلى وجود فروق ذات دلالة إحصائية عند مستوى الدلالة $\alpha=0.05$ بين الأخطاء المعيارية لتقديرات معلمة الصعوبة المقدره باستخدام طريقة الأرجحية الموزونة تعزى للتفاعل بين حجم العينة وطول الاختبار، إن نتائج هذا التفاعل يشير إلى أنه عند استخدام حجم عينة ١٠٠ مفحوص فأكثر فإن دقة تقدير معلمة القدرة تميل إلى الثبات، وهذا الثبات يزداد عندما يكون طول الاختبار أكبر من ٤٠ فقرة. وتتفق هذه النتيجة مع نتائج دراسة الدرايع (٢٠٠١) التي أشارت إلى وجود فروق ذات دلالة إحصائية لتفاعل كل من حجم عينة المفحوصين وطول الاختبار على دقة تقدير قدرة المفحوص. في حين لم تظهر النتائج في هذه الدراسة الحالية وجود فروق إحصائية في دقة تقدير معلمة القدرة تعزى لحجم عينة المفحوصين، ووجود فروق إحصائية في دقة تقدير معلمة القدرة تعزى لطول الاختبار لصالح طول الاختبار الكبير. واختلفت مع هذه الدراسات في أنها أظهرت بأن أفضل دقة لتقدير معلمة القدرة المفحوص عندما كان حجم العينة ١٠٠ مفحوص وطول الاختبار ٧٠ فقرة. ولتوضيح أثر التفاعل بين حجم العينة وطول الاختبار في دقة تقدير معلمة القدرة، فقد تم تمثيل هذا الأثر بيانياً كما في شكل ٣.

شكل (٣) الرسم البياني لأثر التفاعل بين حجم العينة وطول في دقة تقدير معلمة الصعوبة



التوصيات

في ضوء النتائج التي تم التوصل إليها في هذه الدراسة يمكن الخروج بالتوصيات:

١. نظراً لتمتع طريقة المزاجية بالدقة الكبيرة في تقدير معلمة الصعوبة وخاصة عند استخدامها مع العينات الصغيرة من المفحوصين وطول الاختبار، يوصي الباحثان بتقدير معلمة الصعوبة وفق النموذج ثنائي المعلمة والنموذج ثلاثي المعلمة.
٢. يوصي الباحثان بإجراء المزيد من الدراسات المستقبلية باستخدام بيانات حقيقية لتقدير معلمة صعوبة الفقرة.
٣. بسبب انخفاض الأخطاء المعيارية في تقدير معلمة الصعوبة عند استخدام طريقة المزاجية، ويوصي الباحثان بدراسة تهدف للمقارنة بين طريقة الأرجحية العظمى وطريقة بيز وطريقة المزاجية في دقة تقدير معالم الفقرة.
٤. يوصي الباحثان بدراسة تهدف للمقارنة بين البرمجيات الإحصائية (R, Bilog) لتقدير معلمة الصعوبة باستخدام طرق الأرجحية العظمى وطريقة بيز.

- ◆ الدرابيع، ماهر. (٢٠٠١). فعالية النموذج اللوجستي ذي المعلمة الواحدة "نموذج راش" في دقة تقدير قدرة الفرد ومعامل صعوبة الفقرة باختلاف حجم العينة وطول الاختبار. مجلة دراسات-العلوم الإنسانية 28(١)، ١٩٧-٢٠٨.
- ◆ عبابنة، عماد. (٢٠٠٤). أثر حجم العينة وطريقة انتقائها وعدد الفقرات وطريقة انتقائها على دقة تقدير معالم الفقرة والقدرة لاختبار قدرة عقلية باستخدام نظرية استجابة الفقرة أطروحة دكتوراه غير منشورة جامعة عمان العربية للدراسات العليا، عمان.
- ◆ علام، صلاح الدين (٢٠٠٥). نماذج الاستجابة للمفردات الاختيارية أحادية البعد ومتعددة الأبعاد وتطبيقاتها في القياس النفسي والتربوي الطبعة الأولى دار الفكر العربي
- ◆ Barnes, L., & Wise, S. (1991). The utility of a modified one-parameter IRT model with small samples. *Applied Measurement in Education*, 4(2), 143-157.
- ◆ Birnbaum, A. (1969). Statistical theory for logistic mental test models with a prior distribution of ability. *Journal of Mathematical Psychology*, 6(2), 258-276.
- ◆ Chen, W., Lenderking, W., Jin, Y., Wyrwich, W., Gelhorn, H., & Revicki, A. (2014). Is Rasch model analysis applicable in small sample size pilot studies for assessing item characteristics? An example using PROMIS pain behaviour item bank data. *Quality of Life Research*, 23(2), 485-493.
- ◆ Chen, W. H., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, 22(3), 265-289.
- ◆ Choppin, B. (1985). A fully conditional estimation procedure for Rasch model parameters. *Evaluation in Education*, 9(1), 29-42.
- ◆ de Gruijter, D. & Van der Kamp, L. (2005). Statistical test theory for education and psychology. Graduate School of Education, University of Leiden.
- ◆ De la Torre, J., & Hong, Y. (2010). Parameter estimation with small sample size a higher-order IRT model approach. *Applied Psychological Measurement*, 34(4), 267-285.
- ◆ Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7, 1-26.
- ◆ Finch, H., & French, B. (2019). A comparison of estimation techniques for IRT models with small samples. *Applied Measurement in Education*, 32(2), 77-96.
- ◆ Fitzpatrick, A. (2008). The impact of anchor test configuration on student proficiency rates. *Educational Measurement: Issues and Practice*, 27(4), 34-40.
- ◆ Hambleton, R. (1989). Principles and Selected applications of item response theory, In Linn, Robert, L. (Ed.), *Educational Measurement* (3rd ed., PP. 147-201). American Council on Education, Macmillan Publishing Company.
- ◆ Hambleton, R., & Cook, L. (1983). Robustness of item response models and effects of test length and sample size on the precision of ability estimates. In D. J. Weiss & R. D. Bock (Eds.), *New horizons in testing: Latent trait test theory and computerized adaptive testing* (pp. 31-49). Academic Press.
- ◆ Hambleton, R. & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Kluwer, Nijhoff Publishing.
- ◆ Hambleton, R., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Sage Publications.
- ◆ Hattie, John. (1985). Methodology review: Assessing unidimensionality of tests and items. *Applied psychological measurement*, 9(2), 139-164.
- ◆ Heine, J., & Tarnai, C. (2015). Pairwise Rasch model item parameter recovery under sparse data conditions. *Psychological Test and Assessment Modeling*, 57(1), 3-36.

- ◆ Huang, C., Lohss, W., Lin, C., & Shin, D. (2001). Item calibrations for licensure tests with multiple specialty components. Submitted to division ID: *Educational Measurement, Psychometrics and Assessment*. Enabled Tiger a web – based Manuscripts Processing system. Michigan State University.
- ◆ Imon, A., & Masoom, A. (2005). Bootstrapping regression residuals. *Journal of the Korean Data and Information Science Society*, 16(3), 665-682.
- ◆ Jiang, S. Wang, C., & Weiss, D. (2016). Sample size requirements for estimation of item parameters in the multidimensional graded response model. *Frontiers in Psychology*, 7(109), 1-10.
- ◆ Lord, F. (1968). An analysis of the verbal scholastic aptitude test using Birnbaum's three-parameter logistic model. *Educational and Psychological Measurement*, 28(4), 989-1020.
- ◆ Lord, F. (1980). Applications of item response theory to practical testing problems. Lawrence Erlbaum Associates, publishers Hillsdale.
- ◆ Lord, F. (1986). Maximum likelihood and Bayesian parameter estimation in item response theory. *Journal of Educational Measurement*, 23(2), 157-162.
- ◆ Linacre, j. (2008). WINSTEPS Rasch measurement (3.63.2). Chicago, IL:MESA press.
- ◆ Rasch, G. (1960). Probabilistic models for some intelligence and attainment tests. Danish Institute of Educational Research.
- ◆ Reckase, M. D. (1979). Unifactor latent trait models applied to multifactor tests: Results and implications. *Journal of educational statistics*, 4(3), 207-230.
- ◆ Reeve, B., & Fayers, P. (2004). *Applications of item response theory (IRT) modelling for building and evaluating questionnaires measuring patient-reported outcome-s*. Bethesda MD: Advances in Health Outcomes Measurement. National Cancer Institute.
- ◆ Sahin, A., Anil, D., 2017. The effects of test length and sample size on item parameters in item response theory. *Educational Sciences: Theory and Practice*, 17 (1), 321-335.
- ◆ Warm, T. (1978). A Primer of item response theory. Coast Guard Institute.
- ◆ Warm, T. (1989). Weight likelihood estimation of ability in item response theory. *Psychometrika*, 54(3), 427-450.
- ◆ Wright, B., & Stone, M. (1979). Best test design. MESA Press.
- ◆ Zhang, G. (2014). Estimating standard errors in exploratory factor analysis. *Multivariate Behavioral Research*, 49(4), 339-353.