

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلفة

إعداد/ دكتورة شيرين فاروق محمد طنطاوي

مدرس علم النفس الإحصائي بكلية الآداب – جامعة الفيوم

المستخلص:

هدفت الدراسة إلى البحث في تأثيرات طول الاختبار (١٠، ٢٠، ٣٠ مفردة) وطريقة التصحيح (نسبة التصحيح ثنائي الاستجابة ومتعدد الاستجابة) ونموذج تحليل المفردة (باستخدام مجموعتين من النماذج لتصحيح المفردات المختلفة: المجموعة الأولى (نموذج راش ثنائي الاستجابة Rasch ونموذج التقدير الجزئي متعدد الاستجابات PCM)، والمجموعة الثانية (النموذج ثنائي البارامتر ثنائي الاستجابات 2PL ونموذج الاستجابة المتدرجة متعدد الاستجابات GRM). على تقديرات القدرة ودالة معلومات الاختبار ومعايير التحسين للاختبارات ذات المفردات المختلفة. تكونت عينة الدراسة من ٣٤٠ طالباً وطالبة بالفرقة الثالثة من كلية الآداب جامعة الفيوم الذين يدرسون مقرر الاختبارات والمقاييس النفسية. ثم تم تحليل جميع الاستجابات للمفردات للتحقق من أول افتراضات نظرية الاستجابة للمفردة وهو افتراض أحادية البعد وذلك باستخدام طريقة التحليل العملي في برنامج SPSS وذلك بطريقة الارحجية القصوى. تم الوصول إلى عامل واحد لكل مجموعة من المفردات. وتحليل الاستجابات للاختبارات المختلفة تم استخدام مجموعتين من النماذج وهما المجموعة الأولى (نموذج راش ثنائي الاستجابة Rasch ونموذج التقدير الجزئي متعدد الاستجابات PCM)، والمجموعة الثانية (النموذج ثنائي البارامتر ثنائي الاستجابات 2PL ونموذج الاستجابة المتدرجة متعدد الاستجابات GRM). أشارت النتائج إلى أن كلاً من طول الاختبار ونسبة المفردات التي تم تصحيحها بشكل متعدد الاستجابة كان لها تأثير كبير على كمية معلومات الاختبار الناتجة عن الاختبارات ذات المفردات المختلفة. بشكل عام، أدت الاختبارات التي تم إجراؤها على ١٠٠٪ من المفردات التي تم تصحيحها بشكل متعدد الاستجابات إلى الحصول على أعلى معلومات شاملة. يبدو أن هذا ينطبق بشكل خاص على الممتحنين ذوي تقديرات القدرة المنخفضة. تمت المقارنة بين المجموعتين من النماذج من حيث قيم الخطأ

القياسية لتقديرات القدرة والثبات الهامشية ومؤشرات المطابقة (2LL-). وكانت الفروق بينهما واضحة في معدلات الخطأ المعيارية.

الكلمات المفتاحية: طول الاختبار، طريقة التصحيح، نموذج تحليل المفردة، تقديرات القدرة، دالة معلومات الاختبار، معايير التحسين، الاختبارات ذات المفردات المختلطة.

المقدمة

أخذ الاهتمام بالقياس شكل كبير طوال القرن العشرين وكذلك القرن الحالي، كما أصبح الاعتماد على المعلومات المشتقة من درجات الاختبار أهمية كبيرة في مجتمعنا المعاصر (Erguven,2014). حيث تعتمد جميع مستويات التعليم، من مرحلة رياض الأطفال إلى مرحلة الدراسات العليا، ومعظم إجراءات الترخيص المهنية والعديد من طرق التوظيف، اعتمادًا كبيرًا على أداء الاختبار لنشر الفرص وتعزيز المعايير المهنية وضمانها (Adegoke,2013). فجودة الاختبار لها أولوية كبيرة لمصممي الاختبارات، وكذلك الذين يعتمدون على درجات الاختبار في اتخاذ القرار. الآن أكثر من أي وقت مضى، حيث أن للاختبارات دور مهم في قياس القدرة وأن ما ينتج عنها من نتائج تعتبر مقاييس دقيقة لقدرة الممتحن. تركز المعايير المستخدمة لتحديد جودة الاختبار بشكل عام على مجالات تصميم الاختبار وتقنيات تحليل الاختبار وتفسير درجات الاختبار (Murphy,2021).

تتأثر جودة تصميم الاختبار بالعديد من العناصر بما في ذلك الشكل والطول وإجراءات التطبيق والبناء والصدق وطريقة وضع الدرجات. يتمثل العنصر الأكثر أهمية في المفردات المختارة لتكوين الاختبار (Oliveri and von Davier,2014). أكثر نوعي المفردات شيوعًا هما "الاستجابة المحددة selected response" و "الاستجابة المركبة constructed response". تتطلب مفردات الاستجابة المحددة من الممتحنين التعرف على إجابة واحدة واختيارها من مجموعة الاستجابات المقدمة (أي مفردات الاختيار من متعدد والمطابقة والصواب والخطأ). تحتوي هذه المفردات بشكل عام على مخطط تقدير الدرجات "ثنائي الاستجابة dichotomous" يسمح بنتيجتين محتملتين لكل استجابة للمفردة، "صحيحة" أو "غير صحيحة". في المقابل، تجمع مفردات الاستجابة المركبة معلومات حول المعرفة غير المكتملة التي قد يمتلكها الممتحن، من خلال مطالبة الممتحن بإنشاء استجابة لمفردة ما (على سبيل المثال، مفردات استجابة مفتوحة، وإجابة قصيرة، ومفردات مقالة). تمنح مفردات

الاستجابة المركبة "مصادقية جزئية" على مقياس هرمي للصواب باستخدام مقدر درجات المعروف باسم تقدير الدرجات "متعدد الاستجابات polytomous". (Von Davier et al.,2019)

بمجرد تطبيق الاختبار وجمع الاستجابات، توجد العديد من الأساليب التي يمكن استخدامها لتحليل استجابات المفردات. يشار إلى هذه الأساليب عمومًا باسم "النماذج الرياضية mathematical models" وتستخدم لتقدير قدرة الممتحن وخصائص المفردة. يشار إلى هذه العمليات الحسابية باسم "التقديرات estimates" بسبب خطأ القياس المتأصل المرتبط بجميع اختبارات القدرة المعرفية. ويؤدي تقليل هذا الخطأ إلى تحسين دقة القياس للاختبار وكذلك تقديرات خصائص المفردات (Cohen,2018).

توجد نماذج رياضية مختلفة لنظرية الاستجابة للمفردة (IRT) لتحليل استجابة المفردات بهدف تقدير قدرة الفرد وخصائص المفردة ومنها نموذج راش Rasch (Asriadi & Hadi,2021). تم تطبيق هذا النموذج على استجابات مسجلة بشكل ثنائي، ولكن توجد أيضًا إصدارات موسعة منه للتعامل مع البيانات متعددة الاستجابات وهو نموذج التقدير الجزئي PCM. ويطبق نموذج راش تقنية لوجستية، كما يحتفظ نموذج راش بمعلمات مفردة معينة ثابتة (Che Lah, Tasir & Jumaat,2021).

غالبًا ما تقارن نماذج نظرية الاستجابة للمفردة بشأن قدرتهم على تقدير كفاءات الممتحن ومعلمات المفردة عند تحليل الاستجابة للمفردات الثنائية أو المتعددة (Rahayu et al.,2021). ومع ذلك، لم تتم مقارنتها على نطاق واسع بشأن قدرتها على تحليل الاختبارات ذات المفردات المختلطة. تتوفر نماذج رياضية مختلفة في نظرية الاستجابة للمفردة لتحليل الاستجابة للمفردات التي تمتلك طرق تقدير مختلفة للدرجات، ولكن ليس ضمن الاختبارات ذات المفردات المختلطة.

مشكلة الدراسة

غالبًا ما يطبق المعلمون الاختبارات ذات المفردات المختلطة. يستفيد الطلاب عندما تتاح لهم الفرصة للاستجابة لمجموعة متنوعة من أنواع المفردات ومن المعروف أن الأشكال المختلفة للمفردات تختلف في نوع المعلومات التي يتم جمعها. في الأونة الأخيرة، تم اقتراح نظرية مفادها أن الاختبارات ذات المفردات المختلطة توفر

فرصة للاستفادة من المزايا الكامنة للأشكال المختلفة للمفردات في التحليل (Triwik,2021).

يتم تطبيق نماذج نظرية الاستجابة للمفردة في مواقف التقييم على نطاق واسع للمساعدة في تحليل الاختبار وإجراءات تعديل الاختبار (Stewart et al.,2018). ومع ذلك، يوجد تحد بين الأساليب المتاحة لتحليل أنماط الاستجابة من الاختبارات ذات المفردات المختلطة. معظم تقنيات القياس النفسي المتاحة مخصصة للاستخدام مع استجابات الممتحنين للمفردات ذات الدرجات المماثلة. حتى الآن، تم إجراء التحليل النفسي بشكل حصري تقريباً على مجموعات من مفردات الاختبار التي كانت من نفس نوع النتيجة (كلها ثنائية أو كلها متعددة). بالإضافة إلى ذلك، فإن معاملات ثبات وصدق نظرية القياس الكلاسيكية (Classical Test Theory ;CTT) تستند إلى أن جميع مفردات الاختبار هي المفردات ذات النسق الواحد، على سبيل المثال خيارات متعددة ولا تنطبق بسهولة على الاختبارات ذات المفردات المختلطة (Tseng& Wang,2021).

كما أن نماذج نظرية الاستجابة للمفردة لهما افتراضات يجب أن تليها البيانات قبل تطبيق النماذج واعتبار النتائج صحيحة.. يعتبر نموذج راش معلمة التمييز المشتركة عبر جميع المفردات صارماً ويصعب الوفاء به، والذي بدوره قد يؤدي إلى نتائج أقل صدقاً عند تحليل الاختبارات ذات المفردات المختلطة. إن من فوائد استخدام نموذج راش السماح بمقياس لوجستي مشترك لصعوبة المفردة وقدرة الشخص. غالباً ما تمت مقارنة نماذج نظرية الاستجابة للمفردة على تقديراتهم لمعاملات المفردة ومستويات كفاءة الممتحن، ولكن ليس على قدرتهما على تحليل الاختبارات ذات المفردات المختلطة (Wang, Drasgow and Liu,2016). يجب أيضاً أن يعتمد اختيار الطريقة المثلى لتحليل الاستجابات للاختبارات ذات المفردات المختلطة على تطبيق النموذج المناسب لنوع البيانات أو "مطابقة" النموذج وافتراضات النموذج للبيانات.

تفترض الدراسة الحالية أن نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة وطول الاختبار وطريقة تقدير الدرجات وطريقة التحليل تؤثر على معايير التحسين في الاختبارات ذات المفردات المختلطة. معايير التحسين، بدورها، لها تأثير ملموس عبر شكل المفردات المختلطة على ما يلي: ثبات تقديرات القدرة، ومعلومات المفردة، ومعلومات الاختبار، وتقديرات الكفاءة، ومعلومات المفردة.

تسعى الدراسة الحالية الى الإجابة على الأسئلة البحثية الآتية:

١. كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على تقديرات القدرة؟
٢. كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على دالة معلومات الاختبار (TIF)؟
٣. كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على التحسن الإجمالي للاختبار؟
٤. هل توجد فروق بين مجموعتي التصحيح: المجموعة الأولى (نموذج راش ثنائي الاستجابة Rasch ونموذج التقدير الجزئي متعدد الاستجابات PCM)، والمجموعة الثانية (النموذج ثنائي البارامتر ثنائي الاستجابات 2PL ونموذج الاستجابة المتدرجة متعدد الاستجابات GRM) في قدرتها على تحليل الاختبارات ذات المفردات المختلفة في قدرتها على تحليل الاختبارات ذات المفردات المختلفة؟

أهداف الدراسة

١. الكشف عن تأثير نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على تقديرات القدرة.
٢. معرفة تأثير نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على دالة معلومات الاختبار (TIF).
٣. معرفة تأثير نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على التحسن الإجمالي للاختبار.
٤. المقارنة بين مجموعتين باختلاف طريقة التصحيح: المجموعة الأولى (نموذج راش ثنائي الاستجابة Rasch ونموذج التقدير الجزئي متعدد الاستجابات PCM)، والمجموعة الثانية (النموذج ثنائي البارامتر ثنائي الاستجابات 2PL ونموذج الاستجابة المتدرجة متعدد الاستجابات GRM) في قدرتها على تحليل الاختبارات ذات المفردات المختلفة.

أهمية الدراسة:

تصميم الاختبار، وأساليب التحليل، وقضايا التحسين كلها عوامل حاسمة تؤثر على تقديرات القدرة، ومعلومات الاختبار، وفي النهاية ثبات تقديرات القدرة التي ينتجها أي اختبار. من المهم استكشاف ومقارنة:

(أ) أشكال الاختبار للاختبارات ذات المفردات المختلطة.

(ب) النماذج الرياضية المستخدمة لتحليل الاختبارات ذات المفردات المختلطة.

(ج) أساليب التحسين لتصميم وتحليل الاختبارات ذات المفردات المختلطة. لصالح أولئك الذين يخضعون للاختبارات، والذين يعدّون الاختبارات، وأولئك الذين يستخدمون درجات الاختبار لاتخاذ القرار، تحتاج هذه المجالات إلى مزيد من الاهتمام في مجال أساليب القياس النفسي.

أشكال الاختبار:

لقد أتاح التقييم ثنائي الاستجابة فرصة لتطبيق اختبارات أطول من أجل زيادة ثبات الدرجة مع الاستفادة من كفاءة خيارات وضع الدرجات بشكل آلي. في الأونة الأخيرة، تم التعرف على وضع الدرجات متعددة الاستجابة على أنها تقدم شكلاً أكثر دقة لقياس مستوى المهارة الذي يمتلكه الأفراد. مع فئات الاستجابة المتعددة وعتبات الفئات، يسمح وضع الدرجات متعددة الاستجابة بتقدير عدد أكبر من المعلمات. عادةً ما يحتوي الاختبار ذات المفردات المختلطة على نسبة كبيرة من المفردات ذات الدرجات ثنائية الاستجابة ونسبة صغيرة من المفردات ذات الدرجات متعددة الاستجابة. من المهم إجراء مزيد من البحوث حول مشكلات أشكال الاختبار التي تقدر قدرة الممتحن على أفضل وجه (Berger, 1998) مع تزايد شعبية الاختبارات ذات المفردات المختلطة، من المهم إجراء بحوث في مجال كيفية وضع الدرجات المجمعّة.

النماذج الرياضية: إن التحليل الإحصائي للاختبارات ذات المفردات المختلطة معقد. تم إجراء التحليل من خلال استخدام العديد من الإجراءات الحالية، ومع ذلك، فإن معظمها ينطوي على عملية فصل مجموعات أنواع المفردات ذات الدرجات المتشابهة، وتحليلها ثم إنشاء طريقة لدمج النتائج. تعتمد الاستنتاجات المتعلقة بالنماذج بشكل كبير على الأساليب المطبقة في التحليل والطريقة الذاتية في ترجيح النتائج والجمع بينها. لم يتم إجراء مقارنة واضحة بين هذه الأساليب في الاختبارات ذات المفردات المختلطة. لم يتم تأكيد طريقة لتحليل مجموعة كاملة من أنواع المفردات

المتنوعة على أنها متفوقة إلى الآن، وتعتبر معلمات المفردة من القضايا الحاسمة التي يجب تقييمها (Rezapour et al., 2021). من المهم ملاحظة أن كمية المعلومات المتاحة هي الأكبر بالنسبة للمفردات التي يتم وضع درجات لها على مقياس مستمر، على الرغم من أن البحث المتعلق بالعدد الأمثل لفئات الاستجابة للمفردات متعددة الاستجابة لم يكن حاسماً. تعد مساهمات معلومات الاختبار الشاملة للمفردات متعددة الاستجابة وثنائية الاستجابة، بالنسبة لبعضها البعض، قضية يجب البحث فيها ودراستها (Carlson, 1996).

أساليب التحسين:

قد يؤدي الجمع بين نوعين من المفردات في اختبار واحد إلى الاستفادة من المزايا التي يقدمها كل منهما، طالما أن أساليب التحليل المطبقة تستوعب الجمع بينهما. ركزت معظم الدراسات حول تحسين تصميم الاختبار وإجراءات التحليل التي أجريت في الماضي على الاختبارات ذات المفردات المماثلة. تعتبر افتراضات النموذج الإحصائي ومعايير التحسين لهذا النموذج هما المطلبان الرئيسيان لتصميم الاختبار الأمثل (Berger, 1998). إذا تم تطبيق طرق التحسين نحو وضع إرشادات صحيحة لاختيار المفردات، فيمكن إجراء اختبارات تقلل التباين في كل من معلمات المفردة وتقديرات القدرة (Berger, 1998). كما تعد نظريات التحسين وأشكال المفردات المختلطة كلاهما جديداً إلى حد ما وبالتالي لم يتم الجمع بينهما على نطاق واسع.

تعريف المصطلحات:

١. ثنائي الاستجابة Dichotomous:

مفردات الاختبار التي يتم تصحيحها باستخدام طريقة التصحيح الثنائية: إما صحيح (١) أو غير صحيح (٠).

٢. متعددة الاستجابة polychotomous:

مفردات الاختبار التي يتم تصحيحها على مقياس هرمي للصواب، مثل معيار تقييم الدرجات، حيث يمكن الحكم على كل إجابة على أنها أكثر صحة نسبياً أو أقل صحة على النحو الذي تحدده مجموعة من المعايير والمحتوى المحدد. يتم بعد ذلك تخصيص درجة مقياس، على سبيل المثال "٣" على مقياس يتراوح من ١ إلى ٤.

(Hambleton, Swaminathan and Rogers, 1991).

٣. الاختبارات ذات المفردات المختلطة **Mixed-item format test**:

أداة أو اختبار للتقييم يحتوي على مزيج من نوعين من المفردات التي تتطلب طرق تصحيح مختلفة، على سبيل المثال، مفردات الاختيار من متعدد (ثنائي الاستجابة) والمفردات التي تتطلب طريقة تصحيح جزئية (متعددة الاستجابة) Swaminathan، (Rogers & Hambleton, 1991)

٤. التصميم الأمثل للاختبار **Optimal Test Design** :

توفر كل من افتراضات النموذج الإحصائي ومعايير التحسين لهذا النموذج هما المتطلبان الرئيسيان للتصميم الأمثل للاختبار (Berger, 1998)

٥. معايير التحسين للاختبارات ذات المفردات المختلطة **Optimality Criteria for the Mixed-Item Format Test**:

في الاختبار المركب، يتم تحقيق ذلك عندما:

(أ) يتوفر ثبات عالي للدرجات التي ينتجها كل اختبار فرعي (تتعلق عادةً بطول الاختبار).

(ب) الثبات الناتج عن الدرجات من الاختبار المركب يتجاوز الثبات الناتج عن أي من الاختبارين الفرعيين وحده.

(ج) إمكانية تطبيق الاختبار المركب (Downing, 2002).

الإطار النظري:

أشكال الاختبار وإجراءات التصميم:

القدرات المعرفية، التي يشار إليها عادة باسم القدرات "الكامنة" أو "السمات"، لا يمكن قياسها مباشرة، لذلك، تم تصميم الاختبارات وطرق التحليل للمساعدة في تقدير القدرات التي يمتلكها الممتحنون. تعتمد دقة هذه التقديرات على تصميم وبناء الاختبار والتحليل والتفسير. تتضمن ممارسات الاختبار الشائعة إدارة الاختبارات الكتابية المكونة من أنواع مختلفة من العناصر بما في ذلك الاختيار من متعدد والمقال والمطابقة والاستجابة المفتوحة والإجابة القصيرة. كما تزداد شعبية تقييمات الأداء التي تتطلب عادةً من الأفراد تنفيذ مهام محددة. بدأ البعض في النظر في هذه التقييمات

لتوفير مقاييس "أصلية" للقدرة نظرًا لشكلها. يعتبر هذا النهج في التقييم جديدًا إلى حد ما ويحتاج إلى عمالة كثيفة إلى حد ما لإدارته، وبالتالي فهو ليس شائعًا في الممارسة (Kline,2005).

يمكن تحديد جودة أداة التقييم بعدة طرق. يعد البناء والتصميم الدقيق لأداة الاختبار نفسها أمرًا بالغ الأهمية لتحديد درجة الثقة في النتيجة وصلاحيتها. تصميم أداة التقييم هو إجراء متعدد الخطوات يتكون من أربع خطوات رئيسية كما حددها رايت، ستون (Wright and Stone,1979). وهي كالتالي: تحديد المتغير المراد قياسه بالاختبار بوضوح؛ بناء المفردات التي تقيس هذا المتغير؛ إثبات أن المفردات هي مقياس دقيق للمتغير؛ وتقييم أنماط الاستجابة للاتساق المتوقع. هذه العملية عبارة عن عملية دائرية، فبمجرد تقييم أنماط الاستجابة، يمكن استخدام النتائج وذلك لمراجعة وتحسين المفردات وتصميم الاختبار للإدارات المستقبلية.

تتمثل الخطوة الأولى في تصميم أداة التقييم في التحديد الواضح للمتغير المراد قياسه بواسطة الاختبار. يعد تحديد المتغير ضروريًا لبناء مفردات ذات جودة تعكس السمة المرغوبة عند الإجابة بشكل صحيح. هذا يساعد على القضاء على الغموض في تصميم الاختبار. يعتمد ثبات وصدق الدرجات الناتجة على التحديد الواضح للمتغير المراد قياسه. من المهم أن يتم تصميم المفردات لقياس السمة المستهدفة بدقة. وهذا ما يسمى "أحادية البعد"، وهو مهم بشكل خاص عندما يتم دمج أشكال المفردات المختلفة في تقييم واحد. يجب معالجة فروق الأبعاد بين المفردات متعددة الاستجابة وثنائية الاستجابة قبل إجراء أي تحليلات (Carlson, 1996).

الخطوة الثانية، بناء المفردة، بهدف كتابة المفردة التي تقيس المتغير المحدد. عندما تكون مفردة الاختبار مبنية جيدًا ووجود مقاييس صادقة للسمة المستهدفة، فإن درجات الاختبار وأنماط استجابة الممتحن توفر دليلًا واضحًا فيما يتعلق بتقديرات مستوى السمات. تؤثر العديد من العوامل على جودة الاختبار ويحتاج الأشخاص الذين يصممون الاختبارات أو يتخذون قرارات تستند إلى درجات الاختبار إلى فهم واضح لهذه العوامل. بحيث تنتج الاختبارات درجات تعتبر تقديرات لسمات صادقة وأن يتم تفسير هذه الدرجات بعناية (Linn, 1990).

الخطوة الثالثة في تصميم الاختبار: إثبات أن المفردات التي تم إنشاؤها حديثًا هي مقياس دقيق للمتغير المستهدف. إذا تم تحديد المتغير المستهدف جيدًا وتم كتابة المفردات بعناية، فيجب إجراء اختبار ميداني مناسب للمفردات بحيث يمكن تحليل

الاستجابات. بمجرد الموافقة على مجموعة من المفردات، يتم دمجها بعد ذلك لإجراء اختبار. من أجل تأكيد صدق التقييم، يتم إجراء التحقيقات في المحتوى وبناء الصدق. تساعد التحليلات الارتباطية مع الدرجات من التقييمات الأخرى التي تقيس نفس المتغير المستهدف في دعم الصدق التقاربي لدرجات الاختبار (Marnat & Wright, 2016).

تتضمن الخطوة الرابعة في تصميم الاختبار إجراء تقييم سليم لأنماط الاستجابة لتحديد أي تناقضات وحلها ودعم ثبات درجات الاختبار، على سبيل المثال، يجب تصميم المفردات في الاختبار بحيث تكون أحادية البعد (أي أنها تقيس بناء واحد أو سمة واحدة). التحليل العاملي مفيد في تقييم مجموعة من مفردات الاختبار للأبعاد الأحادية. بالإضافة إلى ذلك، يمكن تقييم المفردات الفردية لمقدار المعلومات التي تجمعها. معلومات المفردات قابلة للقياس ويمكن تقييم كل مفردة لمساهمتها الفردية في إجمالي معلومات الاختبار. يعتبر الأعداد الجيدة المفردة ومعلومات الاختبار أمراً مهماً لإنشاء اختبار مصمم على النحو الأمثل (Kline, 2005).

نماذج التصحيح ثنائي الاستجابة:

يعد تطوير مفردات الاختيار من متعدد المخصصة للتقييم ثنائي الاستجابة مدخلا شائعاً في القياس. التصحيح ثنائي الاستجابة له مزايا وعيوب. وتتمثل إحدى المزايا في أن وقت الاستجابة المطلوب لتطبيق الاختبار يكون ضئيلاً نسبياً لكل مفردة بحيث يكون المطورون قادرين على التخطيط لتطبيق اختبارات أطول. يزيد طول الاختبار من ثبات ودقة تقديرات القدرة عن طريق تقليل مقدار الخطأ في القياس. في الثلاثينيات من القرن الماضي، أصبح من الممكن إكمال التصحيح الموضوعي لاستجابات الاختيار من متعدد في جزء بسيط من الوقت والتكلفة اللازمتين للتصحيح اليدوي. زاد هذا التقدم التكنولوجي من جاذبية طريقة الاختبار والتصحيح. انخفضت التكاليف المرتبطة بالحصول على هذه التكنولوجيا باستمرار وأصبح هذا النهج متاحاً الآن لمعظم معلمي الفصول الدراسية (صلاح الدين محمود علام، ٢٠١٨).

توجد مساوئ لهذا النهج التجريبي على الرغم من شعبيته المتزايدة. لا يعتبر استخدام مفردات الاختيار من متعدد طريقة تقييم "أصلية" من قبل البعض لأنه يؤكد فقط قدرة الممتحن على استدعاء المعلومات أو التعرف عليها بدلاً من إظهار المعرفة والقدرة. عيب آخر هو أن التخمين من قبل الممتحنين أمر لا مفر منه ويمكن أن يضيف متغير إزعاج كبير للتحليل وتقديرات القدرة (صلاح الدين محمود علام، ٢٠٠٥).

يعتمد تصحيح المفردات ثنائية الاستجابة على افتراض أن الشخص إما يعرف كل شيء أو لا يعرف شيئاً عن محتوى كل مفردة. لذلك، عندما يتم دمج المفردات المكتوبة بدرجات متفاوتة من الصعوبة لبناء اختبار، يتم استخدام عدد المفردات التي تمت الإجابة عليها بشكل صحيح لتحديد درجة الشخص. يكمن الخطر المرتبط بكل الدرجات ثنائية الاستجابة في أن الاستجابة الصحيحة لجميع المفردات في اختبار سهل للغاية يمكن أن تعني شيئاً مختلفاً تماماً عن الاستجابة بشكل صحيح لـ ٥٠٪ من المفردات في اختبار صعب للغاية. بشكل عام، تتجاوز نظرية الاستجابة للمفردة النسبة المئوية البسيطة للحساب الصحيح وتحلل الاستجابات والصعوبة النسبية لكل عنصر. ثم يتم استخدام هذه المعلومات لتحديد القدرة التقريبية لكل ممتحن (Kaplan, 2012).

نماذج التصحيح متعدد الاستجابة:

أوضحت تحليلات الاستجابات غير الصحيحة أن تحليل الاستجابات غير الصحيحة المختلفة قد أدى إلى معلومات مفيدة والتصحيح ثنائي الاستجابة يجعل هذه المعلومات غير قابلة للوصول من خلال تجميع جميع الاستجابات غير الصحيحة في فئة واحدة. يمكن تصحيح المفردات ذات فئات الاستجابة متعددة الاستجابة بشكل متعدد. يوفر تصحيح المفردات متعددة الاستجابة مقياساً لقياس المعرفة الجزئية لكل مفردة. يقترح وانج ولي (Wang and Lee, 1998) أن أفضل قياس لاكتساب المعرفة غير المكتملة هو بنية هرمية مفترضة توفرها استجابات للمفردات متعددة الاستجابة.

كانت البيانات المنشورة غير متسقة وغير حاسمة فيما يتعلق بتفوق أساليب نظرية الاستجابة للمفردة متعددة الاستجابة على أساليب نظرية الاستجابة للمفردة ثنائية الاستجابة. توجد علاقة مباشرة بين وعي الباحثين بعملية اكتساب المعرفة والحاجة إلى التقييم لتعكس القدرة بشكل أكثر دقة. للحصول على معلومات على مستوى المفردة حول قدرة الشخص على عكس النطاق المستمر لاكتساب المعرفة، يجب على كل مفردة قياس القدرة على مقياس مستمر. وتتراوح هذه المقاييس عموماً بين القيم المطلقة للخطأ والصحيح من خلال استخدام "مقياس التقييم" على سبيل المثال، وتتراوح في النقاط الممنوحة من "١" إلى "٤". تحدد نماذج وضع الدرجات مستويات الكفاءة أو القدرة بناءً على المعايير التي تم الوفاء بها على مستويات مختلفة من المعرفة أو المهارة المثبتة. في هذه الحالة، قد تمثل الدرجة "١" دليلاً محدوداً على القدرة على الكتابة بينما تمثل الدرجة "٤" قدرة كبيرة في الكتابة. يؤكد توثيق

المعرفة الجزئية أن المقاييس الدقيقة للمعرفة أو القدرة لم يعد من الممكن حصرها في استخدام المفردات ثنائية الاستجابة التقليدية (von Davier et al., 2019).

وصف كارلسون (Carlson, 1996) المفردات التي يتم تصحيحها بشكل متعدد الاستجابة بأنها فئات مرتبة متعددة تفترض وجود بنية هرمية في فئة الصعوبات. قد يظهر الممتحن معرفة جزئية بكل مفردة يتم تطبيقها، مما يسمح بجمع المزيد من المعلومات أكثر من المفردات ثنائية الاستجابة. يمكن تسليم مقياس القدرة، الذي يتطلب أربع مفردات ثنائية الاستجابة، من خلال مفردة واحدة متعددة الاستجابة بأربع فئات هرمية. تسمح هذه العملية بجمع كمية متساوية من البيانات من خلال استخدام عدد أقل من المفردات. لقد تبين أن الاختبارات المكونة بالكامل من المفردات التي يتم تصحيحها بشكل متعدد الاستجابة تجمع ما بين ضعفين وثلاثة أضعاف كمية المعلومات مثل الاختبارات التي تحتوي على مفردات ثنائية الاستجابة (Donoghue, 1994).

الجمع بين أنواع المفردات:

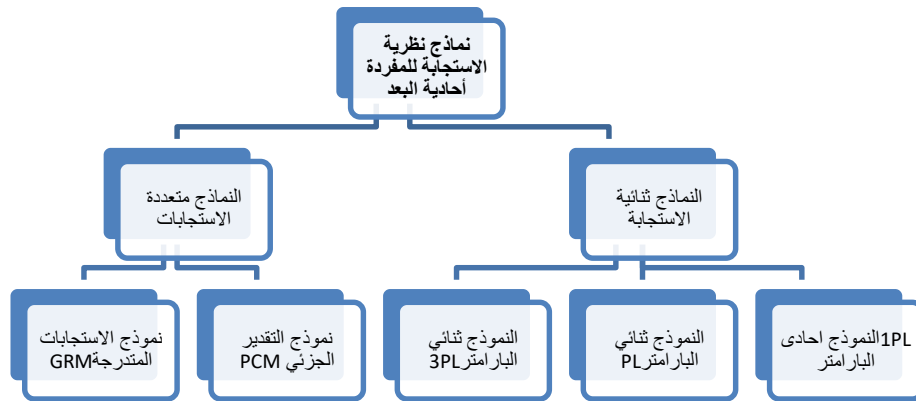
تم ذكر الاختبار ذات المفردات المختلطة كثيراً في الأبحاث الحديثة، وتحديداً أن المفردات ثنائية الاستجابة ومتعددة الاستجابة تكمل بعضها البعض عند دمجها في تقييم واحد مما يؤدي إلى تحسين درجة الثبات والصدق وتقليل التكلفة (Ericikan, 1998) تكمن الميزة الرئيسية لاستخدام نهج الاختبار ذات المفردات المختلطة لاختبار التصميم في فرصة جمع معلومات أكثر اتساقاً على جميع مستويات القدرة مع تحسين طول الاختبار. لقد وثقت العديد من الدراسات مزايا وعيوب كلا النوعين من المفردات، وقد يؤدي دمجها في اختبار واحد إلى التغلب على بعض العيوب الموثقة. درس وينر وتيسين (Wainer and Thissen, 1993) مساهمة المعلومات النسبية وطول الاختبار ووقت التطبيق للمفردات متعددة الاستجابة وثنائية الاستجابة، حيث اعتبروا المفردات متعددة الاستجابة على أنها "غير فعالة" استناداً إلى الوقت المستثمر في بناء المفردة، وتصحيح المفردات، والتطبيق والتكلفة مع وجود معلومات أقل في الدقيقة للتطبيق من المفردات ثنائية الاستجابة. ربما ساهمت نقطتان في بحثهم في هذا الاستنتاج: لقد قاما بتحليل بيانات استجابة الاختبار الميداني وانتهكوا افتراض الاستقلال المحلي مما سمح بمشاكل متعددة الأبعاد. بشكل عام، تعد مساهمات معلومات الاختبار الخاصة بالمفردات متعددة الاستجابة وثنائية الاستجابة بالنسبة لبعضها البعض مشكلة يجب

البحث فيها بشكل أكبر (Carlson, 1996). جمع البيانات وتصحيحها بشكل ثنائي الاستجابة هما طريقتان تقللان من كمية المعلومات المتاحة لتقدير قدرة الممتحن. توفر مفردات الاختبار التي تجمع المزيد من البيانات حول المعرفة والقدرة، والتي يتم تصحيحها على مقياس متعدد الاستجابة، معلومات أكثر دقة حول قدرة الممتحنين.

قد يؤدي تصحيح "الكل أو لا شيء" لمفردات الاختبار من متعدد إلى تحريف المعرفة الحقيقية أو قدرة الممتحنين. هذه الاستجابات التي تم تصحيحها في شكل ثنائي الاستجابة على أنها صحيحة بنسبة ١٠٠٪ أو ١٠٠٪ غير صحيحة تتجاهل نظرية اكتساب المعرفة الجزئية (Wang & Lee, 1998). تحد بيانات مفردات الاختبار التي لا يمكن تصحيحها وتحليلها إلا على المستوى ثنائي الاستجابة من المعلومات التي يمكن الحصول عليها من خلال التحليل. يتم اكتساب المعرفة أو القدرة الجزئية على طريق المعرفة الكاملة. لذلك، يجب أن يتم اكتساب المعرفة على نطاق مستمر وتراكمي كما تم قياسه بواسطة نماذج التصحيح متعددة الاستجابة. تمتلك المفردات متعددة الاستجابة فئات متعددة مرتبة تفترض بنية هرمية في فئة الصعوبات (Carlson, 1996).

نماذج نظرية الاستجابة للمفردة أحادية البعد:

يوضح الشكل التالي نماذج نظرية الاستجابة للمفردة أحادية البعد كما يلي:



نماذج نظرية الاستجابة للمفردة أحادية البعد لصالح علام (٦٨: ٢٠٠٥)

نماذج نظرية الاستجابة للمفردة ثنائية الاستجابة لصيغة التصحيح:

هناك العديد من نماذج التحليل المختلفة المصنفة في نطاق نظرية الاستجابة للمفردة. تتضمن كل هذه النماذج مقياساً لقدرة الشخص (ثيناً) والاهتمام بمعلمات المفردات الأخرى المتعلقة بالقياس. تُستخدم النماذج اللوجستية ذات المعلمات الواحدة والثنائية والثلاثية (1PL, 2PL, 3PL) بشكل شائع لتحليل استجابات المفردات ثنائية الاستجابة في إطار نظرية الاستجابة للمفردة. تمت تسمية هذه النماذج على اسم عدد المعلمات التي لها حرية التغيير في التحليل. تتضمن هذه المعلمات مؤشر التمييز (أ) ومعامل صعوبة المفردة (ب) ومعامل التخمين (ج) (محمد طالب دبوس، ٢٠١٦).

يعد نموذج راش Rasch من ضمن نماذج نظرية الاستجابة للمفردة، تم تصميمه بواسطة جورج راش (Georg Rasch, 1960) لتحليل الاستجابة للمفردات ثنائية الاستجابة عن طريق تقدير قدرة الشخص وصعوبة المفردة بشكل منفصل. اشتق جورج راش هذا النهج لتحليل المفردات لغرض نمذجة سلوك الاختبار على وجه التحديد على مستوى المفردة ولتحليل الاستجابات المسجلة ثنائية الاستجابة. إن استخدام "الإحصائيات الكافية" في حساب تقديرات المفردة والأفراد يلغي الترابط بينهما. كما يطبق نموذج راش أسلوب لوجستي لتقدير معلمات المفردة وقدرات الشخص في قياسات لوغاريتمية نسبية. يتيح ذلك مقارنة قدرة الشخص وصعوبة المفردة على مقياس مشترك. كما يتضمن نموذج راش Rasch مقاييس لقدرة الشخص وصعوبة المفردة، مع الاحتفاظ بمعلمات المفردة الأخرى التمييز (١) على طول الاختبار والتخمين قيمة ثابتة وهي (٠) عبر جميع المفردات وهذا ما يجعله مختلفاً عن النموذج احادي البارامتر 1-PL حيث يعطى الأخير قيمة تقريبية لمعلم التمييز يضعها البرنامج المستخدم في التحليل لجميع المفردات وكذلك يثبت معلم التخمين لجميع المفردات. تعتبر الافتراضات التي يتطلبها راش Rasch صارمة ويصعب الوفاء بها، لذا يُعتقد أن مجموعات البيانات القليلة تلبى بشكل كاف افتراضات التحليل (Kaplan, 2012). يحتوي نموذج راش Rasch على ثلاث صفات تجعله أكثر جاذبية في الاستخدام من بعض النماذج الأخرى: سهولة الاستخدام بسبب عدد أقل من المعلمات، ومشكلات تقدير أقل بسبب عدد المعلمات الأقل، والموضوعية المحددة فيما يتعلق بتقدير المفردة ومعلمات القدرة، والذي كان سبب إنشائه (van der Linden, 2016).

هناك العديد من الافتراضات المتعلقة بنماذج نظرية الاستجابة للمفردة التي يجب التحقق منها في البيانات قبل تطبيق الإجراءات. يمكن أن يساعد المستوى الذي يتم تلبية هذه الافتراضات من خلال البيانات في تحديد النموذج الأنسب وإلى أي درجة

من الدقة يمكن تفسير النتائج. وهذه الافتراضات هي (أحادية البعد والاستقلال الموضوعي والسرعة ومنحنى خصائص المفردات (محمد محمد فتح الله، ٢٠١١).

الافتراض الأول أحادية البعد: Unidimensionality

يقصد بأحادية البعد أن تكون جميع المفردات الاختبارية متجانسة أي تقيس خاصية واحدة فقط وأن هناك قدرة كامنة وحيدة هي التي تفسر أداء الفرد في الاختبار، ويعد افتراض أحادية البعد أقوى افتراضات النظرية، فمن الصعب وجود قدرة واحدة تقيسها جميع مفردات اختبار ما بسبب وجود عوامل كثيرة مؤثرة في أداء الفرد على الاختبار كالتخمين الإيجابية الصحيحة وقلق الاختبار والقدرة على العمل بسرعة والدافع، حيث أن هذه العوامل وغيرها تؤثر في استجابة الفرد على الاختبار بالإضافة إلى القدرة الأساسية التي اعد الاختبار لقياسها ولكن بسبب توفر القدرة المسيطرة لدى جميع الأفراد والتي يعزى إليها الأداء بشكل عام وقلة تأثير العوامل الأخرى وعدم ثبات تأثيرها فإن القياس يرد إلى هذه القدرة الأساسية بينما تعد العوامل الأخرى أخطاء في تحليل نماذج الاستجابة للمفردة (Sinharay&Holland 2006,429) الثاني الافتراض الثاني الاستقلال الموضوعي Local

:independence

يعد افتراض الاستقلال الموضوعي من الافتراضات الأساسية في كلا من النظرية الكلاسيكية للاختبارات CTT ونظرية الاستجابة للمفردة IRT حيث تفترض النظرية الكلاسيكية أن أخطاء القياس مستقلة إحصائياً بين الفقرات المختلفة وبذلك تكون مستقلة في الاختبار ككل وكذلك في نظرية الاستجابة للمفردة يجب أن تكون الاستجابات لأي فقرتين مختلفتين مستقلين إحصائياً عن بعضهما للمفحوصين عند مستوى القدرة، أي أن استجابة المفحوص لمفردة ما لا تتأثر إيجاباً أو سلباً باستجاباته على مفردة أخرى (Lee, 2012)

فيعبر الاستقلال الموضوعي عن احتمال الإجابة الصحيحة للفرد على مفردة الاختبار يكون مستقلاً عن أجابته على أي مفردة أخرى في الاختبار، حيث أن نماذج منحنى خصائص المفردة Item characteristic curve اللوغاريتمية والمنحنى التراكمي الطبيعي يؤكد أن احتمال النجاح على المفردة يعتمد على بارامترات المفردة وقدرة المفحوص ولا شيء آخر (Mandal, 2002 , 36).

كما يقصد بهذا الافتراض أن استجابات الفرد للنبود في الاختبار مستقلة أي لا تؤثر استجابة الفرد لأحدى المفردات على استجاباته للمفردات الأخرى ويتضح هذا في:

تحرر القياس من توزيع العينة المستخدمة sample free وهذا يعنى ثبات تقدير كل من قدرة الفرد وصعوبة المفردة واستقرارهما بالرغم من اختلاف عينة الأفراد المستخدم في الاستجابة طالما أنها عينة ملائمة (Anastsi & warena , 2015,223).

الافتراض الثالث منحنى خاصية المفردة: (Item Characteristic Curve (ICC)

يعرف منحنى خاصية المفردة بأنه دالة رياضية تربط بين احتمال إجابة المفحوص على الفقرة إجابة صحيحة وبين السمة أو القدرة التي تقيسها المفردة أو يقيسها الاختبار الذي يحتوي على هذه الفقرة وتسمى هذه الدالة بمنحنى خاصية المفردة والتي تمثل بانحدار غير خطى لاحتمال الإجابة الصحيحة على المفردة بدلالة القدرة التي يتم قياسها بالاختبار ويعبر عنها بالصيغة التالية:

$$P_i(\theta) = F(\theta - b_i)$$

(θ) هي قدرة الفرد, (b) معلم صعوبة الفقرة (i). (أمينة كاظم, ١٩٨٨, ٤٤٠).

وتفترض نماذج الاستجابة للمفردة أحادية البعد وجود متصل Continuum للسمة المراد قياسها، وأنه يمكن تقدير احتمال إجابة فرد إجابة صحيحة على مفردة اختبار إذا علمنا موقعة على هذا المتصل، ويعبر عن ذلك رياضياً أن الاختبار يعد بمثابة راسم فوقى Mapping من عينة الأفراد المختبرين إلى متصل السمة الكامنة واحتمال الإجابة الصحيحة يكون دالة تزايدية لمواقع الأفراد على متصل السمة وبذلك يزداد احتمال توصل الفرد إلى الإجابة الصحيحة على مفردة الاختبار بازدياد مقدار السمة لديه.

ويتضح مفهوم المنحنى المميز للمفردة إذا قارنا بين النظريتين التقليدية CTT ونظرية الاستجابة للمفردة IRT فكلاهما يعتمد على فرضية مؤداها أن موقع الفرد على متصل سمة كامنة معينة يحكم احتمال أجابته إجابة صحيحة عن مفردة اختبار يقيس هذه السمة ولكنهما يختلفان في كيفية تحديد هذا الموقع. وارتباطه باحتمال الإجابة الصحيحة على المفردة أي أن الفرق بينهما يتعلق بشكل وخصائص الدالة التي تحدد هذا القيد فالنظرية الكلاسيكية تقترض أن منحنى الدالة المميزة للمفردة لا يختلف باختلاف مفردات الاختبار وتمثل هذه الدالة بالمنحنى الاعتدالي التراكمي

لمتغير عشوائي ويفترض في هذا المنحنى أن احتمال الإجابة الصحيحة عن مفردة اختبارية دالة تزايدية مطردة لموقع الفرد على متصل السمة الكامنة (صلاح علام, ٢٠١١, ٦٩٧).

الافتراض الرابع: التحرر من عامل السرعة في الإجابة Speedness:

أن نماذج الاختبارات التي تكون ملائمة لنماذج الاستجابة للمفردة لا يتم إعدادها تحت شرط السرعة أي أن إخفاق بعض الأفراد في الإجابة على بعض مفردات الاختبار يرجع إلى تجانس قدراتهم وليس إلى عامل السرعة في إجاباتهم. وهذا الافتراض لا يعلن عنه كأحد افتراضات نماذج الاستجابة للمفردة في معظم الأحوال, حيث أنه متضمن في الافتراض الخاص بأحادية البعد وذلك لأنه لو اعتبرت السرعة عاملاً مؤثراً في أداء الفرد على الاختبار ففي هذه الحالة يكون هناك سمتين تؤثران في الأداء وهما سرعة الأداء والسمة المقاسة من خلال محتوى الاختبار.

ويتم التحقق من هذا الافتراض عن طريق مقارنة تباين الفقرات المحذوفة مع تباين الفقرات التي تمت الإجابة عنها خطأ، أو من خلال مقارنة علامات الاختبار عند تحديد الوقت مع علامات نفس الاختبار عند ترك الوقت دون تحديد، أو من خلال التحقق النسبة المئوية للمفحوصين الذين أجابوا عن جميع الفقرات والتحقق من النسبة المئوية للمفحوصين الذين أكملوا ٧٥% من الاختبار، أو عدد الفقرات التي أجاب عنها ٨٠% من المفحوصين (أمينة كاظم وآخرون, ١٩٩٦; صلاح علام, ٢٠٠٥; مايا بركات, ٢٠١٠; ياسين الشواورة, ٢٠١٣, آمنه إبراهيم أبو حجر, ٢٠١١).

نماذج نظرية الاستجابة للمفردة متعددة الاستجابة لصيغة التصحيح:

(١) نموذج التقدير الجزئي Partial Credit Model (PCM):

طور ماسترز (Masters, 1982) ويعد هذا النموذج تعميم لنموذج راش للاستجابات الثنائية ليتناسب مع المفردات التي تتطلب عدد من الخطوات للاستجابة عليها، حيث يستخدم في تحليل مفردات الاختبارات، التي تشمل المفردات الاختبارية التي تتطلب مجموعة من الخطوات مثل المسائل الحسابية أو أسئلة المقال، وفيها لا تكون الاستجابة صحيحة أو خاطئة، إنما تقدر درجات للاستجابة الجزئية الصحيحة والذي يعبر عنه بالعلاقة الرياضية الآتية:

$$P_{ix}(\theta) = \frac{e^{\sum_{j=0}^x (\theta - \delta_{ij})}}{\sum_{i=0}^m (e^{\sum_{j=0}^x (\theta - \delta_{ij})})}$$

٢) نموذج الاستجابة المتدرجة (GRM) : Graded Response Model

يعتبر هذا النموذج امتداد للنموذج ثنائي المعلم والذي طوره ساميجاما (Samejima, 1969). ويعد هذا النموذج ملائماً لاستجابات المفردات مرتبة الفئات، ويستخدم في التنبؤ باستجابة الفرد المحتملة لفئة معينة والتي تتسق مع مستواه النمائي، حيث يهدف الى تحديد موضع هذه الفئة على متصل السمة الكامنة، ويتم تصنيف استجابات الافراد لكل مفردة وبشكل متدرج، ولحساب مقدار احتمال الاستجابة لكل فئة يتطلب خطوتين

منحنيات عددها يساوي عدد العتبات الفارقة باستخدام النموذج الثنائي كالتالي:

$$P_{ix}^*(\theta) = \frac{e^{a_i(\theta - \beta_{ij})}}{1 + e^{a_i(\theta - \beta_{ij})}} \quad \text{حيث } j = 1, 2, \dots, m$$

$P_{ix}^*(\theta)$ تسمى منحنيات مميزة حدية، boundary response curves، ويتطلب هذا النموذج تقدير منحنى مميز حدى لكل عتبة فارقة بين فئات الاستجابات، ويمكن تفسير قيم بارامترات (β_{ij}) على أنها تمثل مستوى السمة المطلوبة لكي تتخطى الاستجابة العتبة الفارقة (باحتمال قدره 50%). ، فالمفردة تعالج على أنها سلسلة من الفئات الثنائية ويستخدم النموذج ثنائي البارامتر في تقدير كل هذه الفئات الثنائية، تحت شرط أن ميل المنحنيات الحدية متساوي داخل كل مفردة، وبعد الحصول على هذه التقديرات يحسب الاحتمال الفعلي لفئات الاستجابات وفقاً للمعادلة التالية:

$$P_{ix}(\theta) = P_{ix}^* - P_{i(x+1)}^* \quad x = 1, 2, 3, 4$$

وتسمى بالمنحنيات الاستجابية للفئات Category response curve وهي تعبر عن احتمال استجابة الفرد في فئة معينة مشروطاً على مستوى السمة الخاصة به (صلاح، ٢٠٠٥، ٧٥).

تختلف النماذج متعددة الاستجابة عن النماذج ثنائية الاستجابة في أيهما أكثر ملائمة بناءً على استجابات الافراد. ويمكن تصحيح فئات الاستجابة متعددة الاستجابة بشكل

(E-mail of the correspondig: jcps@art.bsu.edu.eg)

هرمي وينتج عنها معلمات يشار إليها باسم "عتبات الفئة". لكل مفردة فئة $m + 1$ ، هناك معلمات حد m تفصل بين الفئات $a(k)$. كما هو الحال مع النماذج ثنائية الاستجابة، يمكن تحليل الاستجابة للمفردات لإنتاج تقديرات لقدرة الشخص ومعلمات المفردة. "مع الاهتمام المتزايد بـ "القياس النفسي"، يجب الاهتمام بنماذج نظرية الاستجابة للمفردة التي يمكنها التعامل مع التصحيح متعدد الاستجابة، نظرًا لأن القياس الحقيقي مرتبط باختبار الأداء وتصحيح أو وضع درجات لأداء الممتحن" (Hambleton, Swaminathan, & Rogers, 1991, p.153).

تحسين الاختبارات ذات المفردات المختلطة:

تصميم الاختبار:

بالنظر إلى التحول النموذجي في التحليل من تقدير مستويات السمات استنادًا إلى مجموع درجات الاختبار إلى الاستجابة للمفردات، فقد ذكر أنه يجب التخلي عن فكرة درجة الاختبار "الثبات" لصالح القياسات الأكثر ملاءمة لمعلومات المفردة والاختبار (Thissen, 1991). تم وصف معلومات المفردة والاختبار باستخدام معادلات فيشر Fisher's formulas بواسطة بيرنبوم Birnbaum الذي بدأ العمل في مجال تحليل المفردات في الخمسينيات من القرن الماضي. يعتمد تطبيق وظائف المعلومات في مراحل إنشاء الاختبار واختيار المفردة على مساهمة المفردة في معلومات الاختبار الإجمالية. عند فحص الكميات النسبية من المعلومات التي تم جمعها بواسطة أنواع المفردات المختلفة، اقترح بيرنبوم ترجيح مساهمات المفردات في مجموع درجات الاختبار بناءً على مساهماتها في إجمالي معلومات الاختبار. يمكن تنفيذ نهج الترجيح هذا من خلال استخدام مقياس درجات. طبق بيرنبوم أيضًا منهجه في الترجيح على النظريات المثلى لتسجيل درجات الاختبار (Nunnally & Bernstein, 1994).

تتضمن جهود تصميم الاختبار معالجة العوامل التالية: شكل أو صيغة المفردة واختيار المفردة وطول الاختبار وطرق تصحيح الدرجات وإجراءات تحسين الاختبار والجدول الزمنية والمخاوف المتعلقة بالميزانية. تشمل الجهود المبذولة لتحسين تصميم الاختبار كتابة المفردات التي تقيس بوضوح القدرة المستهدفة والتي تجمع الحد الأقصى من المعلومات بكفاءة عبر جميع مستويات القدرة. قد يكون شكل المفردات المختلطة هو التصميم الأمثل للاختبار لأنه قد يستفيد من مزايا وعيوب كلا النوعين من المفردات (Carlson, 1996). يعد تحقيق التوافق الجيد للنموذج المحدد

مع البيانات التي تم جمعها أمرًا بالغ الأهمية أيضًا إذا كانت درجات نمط الاستجابة ستحقق تقديرات القدرة ثيما المثلى (Thissen, Nelson and Swygert, 2001).

معلومات المفردة والاختبار:

وجد كارلسون (Carlson, 1996) أن المفردات التي يتم تصحيحها بشكل متعدد الاستجابة تميل إلى توفير قدر أكبر نسبيًا من المعلومات مقارنة بالمفردات ثنائية الاستجابة وقدمت معلومات عبر نطاق أوسع من تقديرات الكفاءة مقارنة بالمفردات ثنائية الاستجابة. ووجد أيضًا أن المفردات ثنائية الاستجابة تبدو وكأنها تنتج منحني معلومات يتطابق بشكل كبير مع منحنى إتقان الممتحن أكثر من منحنى المفردات متعددة الاستجابة. بالنظر إلى الوقت المستغرق في كتابة المفردات متعددة الاستجابة والاستجابة لها وتصحيحها، يبدو أن معلوماتهم المتزايدة تكون مكلفة، لذلك عم كارلسون النتائج التي توصل إليها أن الاختبارات التي تجمع بين أنواع متعددة من المفردات ستستفيد من مزايا كل منها مع تقليل العيوب.

وجد دون وجو (Donoghue, 1994) أن المفردات متعددة الاستجابة قد تجمع ما بين ٢,١ إلى ٣,١ مرة من المعلومات مثل المفردات ثنائية الاستجابة. هذا هو الدعم للتوصية بترجيح المفردات متعددة الاستجابة الأثقل في التقييم. عندما تم تصحيح المفردات متعددة الاستجابة بشكل ثنائي، انخفض مقدار المعلومات التي تم جمعها بشكل كبير ولكنها كانت لا تزال أكبر من تلك التي تم إنتاجها باستخدام المفردات ثنائية الاستجابة بدقة. تستغرق المفردات متعددة الاستجابة مزيدًا من الوقت والمال لتطبيقها وتصحيحها.

معايير التقييم للتحسين:

كنتيجة للفوائد الموثقة المنسوبة إلى كل من المفردات ثنائية الاستجابة ومتعددة الاستجابة، من المهم إعطاء اعتبارات لتصميم الاختبارات التي تحتوي على مجموعة من هذه الأنواع من التصحيح أو وضع الدرجات. يقترح الاختبار ذو المفردات المختلطة من التحسين عند اختيار تصميم الاختبار وطرق التحليل في وقت واحد. يؤثر تصميم الاختبار وطرق التحليل أيضًا على الكمية المتاحة من معلومات المفردة والاختبار (Ercikan et al., 1998; Si, 2002). أشارت دراسة الأمثلة أو التحسين Optimality التي أجراها بليك (Plake, 1993) إلى أن أفضل النتائج التي يتم

الحصول عليها بواسطة مفردات الاختبار هي تلك "الأفضل" التي تتطابق مع قدرة الممتحن. من المتوقع أن تقدر المفردات التي تم تصحيحها بشكل متعدد القدرات بأعلى مستوى من الدقة.

لاحظ بيرجر Berger (١٩٩٨) أن بحث تصميم الاختبار الأمثل يعتمد على النموذج الإحصائي المفترض المطبق والمعيار الأمثل لهذا النموذج، مثل تلبية الافتراضات والبيانات المناسبة للنموذج. ستكون النتيجة المرجوة أكثر كفاءة في تقدير السمات والمعلومات (أي تباين أقل في التقديرات). هناك نوعان من التحديات التي لاحظها بيرجر في تحقيق التفاضل. يتعلق التحدي الأول بتصميم الاختبار الأمثل واعتماده على وظيفة معلومات فيشر. تعتمد هذه الوظيفة على قيم ثبات، والتي تكون غير معروفة في البداية في أي عملية تصميم اختبار. يمكن حل هذا الاعتماد باستخدام إجراء متسلسل محوسب ومعقد لتقدير المفردات التي يتم فيها تحديث تقديرات معلمة المفردة في كل خطوة. التحدي الثاني لتحقيق الأمثل هو تأثير عدد الفئات وعدد معلمات الفئة على إجراء تصميم الاختبار. لم يتم تحديد العدد الأمثل للفئات في الأدبيات. لاحظ بيرجر أيضاً أن هناك حاجة إلى مزيد من البحث فيما يتعلق باختيار المفردة الأمثل لتقدير القدرة الفعال للأفراد.

إجراءات الدراسة:

عينة الدراسة:

تكونت عينة الدراسة من ٣٤٠ طالبا وطالبة بالفرقة الثالثة من كلية الآداب جامعة الفيوم.

أداة الدراسة:

من أجل تحقيق هدف الدراسة، تم بناء اختبار، مكون من ٣٠ مفردة، يطلق عليه الاختبار ذي المفردات المختلطة، حيث يتكون من ٧٠٪ ثنائي الاستجابة و ٣٠٪ متعدد الاستجابة.

تحليل الاشكال المختلطة للمفردات:

تم تحليل جميع الاستجابات للمفردات من أجل أحادية البعد في برنامج SPSS باستخدام إجراء التحليل العاملي بطريقة الارحية القصوى. تم الوصول الى بناء عامل واحد لكل مجموعة من المفردات. لتحليل الاستجابات للمفردات ثنائية الاستجابة، تم استخدام نموذج راش Rasch من خلال تطبيق النموذج اللوجستي ذو المعلمة الواحدة. وتحليل الاستجابات للمفردات متعددة الاستجابة في المدخل القائم على نموذج راش Rasch، تم تطبيق نموذج التقدير الجزئي.

الثبات:

تم حساب الثبات باستخدام معامل ألفا كرونباخ للاتساق الداخلي Cronbach's α ، وجاءت قيمة معامل الثبات ($\alpha = 0.82$).

الصدق:

تم حساب صدق المحتوى للمقياس من قبل مجموعة من ١٠ خبراء، قيّموا مدى ملاءمة كل مفردة باستخدام مقياس Likert المكون من أربع نقاط (حيث يمثل ١ "غير ذي صلة" و ٤ يمثل "ذو صلة كبيرة")، وقدموا اقتراحاتهم وتعليقاتهم. تم الحكم على البنود الـ ٣٠ بأنها وثيقة الصلة نوعاً ما أو وثيقة الصلة بشكل كبير. تم حساب مؤشر صدق المحتوى على مستوى المفردة ($I-CVI = 0.90$).

التصميم:

تم استخدام تصميم تحليل التباين الأحادي ANOVA ذي التأثير الثابت $2 \times 3 \times 11$ مع أحد عشر شكلاً لتصحيح أو وضع درجات الاختبار (النسب المنظمة للمفردات ثنائية الاستجابة ومتعددة الاستجابة)، وثلاثة أطوال للاختبار (قصير = ١٠ مفردة، متوسط = ٢٠ مفردة وطويلاً = ٣٠ مفردة) ومجموعتين من النماذج الرياضية لنظرية الاستجابة للمفردة. يتكون أول أشكال التصحيح المنظم للاختبار الأحد عشر من ١٠٠٪ تصحيح متعدد الاستجابة للمفردة مع تصحيح ٠٪ ثنائي الاستجابة، متبوعاً بنسبة ٩٠٪ تصحيح متعدد الاستجابة مع ١٠٪ ثنائي الاستجابة، وهكذا حتى الوصول إلى التصحيح الحادي عشر الذي يتكون من ٠٪ تصحيح متعدد الاستجابة و ١٠٠٪ ثنائية الاستجابة. تم استخدام مجموعتين من النماذج لتصحيح المفردات المختلطة: المجموعة الأولى (نموذج راش ثنائي الاستجابة Rasch ونموذج التقدير

الجزئي متعدد الاستجابات (PCM)، والمجموعة الثانية (النموذج ثنائي البارامتر ثنائي الاستجابات 2PL ونموذج الاستجابة المتدرجة متعدد الاستجابات (GRM) نتج عن التصميم ثلاثي العوامل ٦٦ مجموعة بيانات كما هو موضح في الجدول ١.

الجدول (١) مجموعات العوامل لتصميم الدراسة

النسبة المئوية للمفردات ثنائية الاستجابة / متعددة الاستجابة												
0/100	10/90	20/80	30/70	40/60	50/50	60/40	70/30	80/20	90/10	100/0	نموذج التحليل	طول الاختبار
1	101	201	301	401	501	601	701	801	901	1001	المجموعة الأولى	10
1	101	201	301	401	501	601	701	801	901	1001	المجموعة الثانية	10
2	102	202	302	402	502	602	702	802	902	1002	المجموعة الأولى	20
2	102	202	302	402	502	602	702	802	902	1002	المجموعة الثانية	20
3	103	203	303	403	503	603	703	803	903	1003	المجموعة الأولى	30
3	103	203	303	403	503	603	703	803	903	1003	المجموعة الثانية	30

بناء مجموعات البيانات:

تم اختيار مجموعات الاستجابة المكونة من عشرة مفردات وعشرين مفردة بشكل عشوائي لبناء اختبارات متوسطة وقصيرة في الطول. استناداً إلى أطوال الاختبار التي تمت ملاحظتها خلال مراجعة الأدبيات والنظر في طول فترة تطبيق المفردات النموذجية التي تم تصحيحها باستخدام نموذج تقييم متعدد الاستجابة، تم اعتبار الاختبار المكون من ١٠ مفردات قصيراً، وكان الاختبار المكون من ٢٠ مفردة متوسط الطول وتم اعتبار الاختبار المكون من ٣٠ مفردة طويلاً. تم بعد ذلك تكرار مجموعات الاستجابة ١٠ و ٢٠ و ٣٠ مفردة ١١ مرة وتم اختيار نسبة محددة مسبقاً من المفردات بشكل عشوائي من كل اختبار وإعادة تصحيحها أو وضع درجات لها بشكل ثنائي الاستجابة. تم بعد ذلك تحديد أحد عشر شكلاً متناسباً عبر ثلاثة أطوال للاختبار من خلال النسبة المئوية للمفردات في الاختبار الذي تم تصحيحه بشكل متعدد الاستجابة وثنائي الاستجابة، على التوالي (١٠٠٪ / ٠٪، ٩٠٪ / ١٠٪، ٨٠٪ / ٢٠٪، ٧٠٪ / ٣٠٪، ٦٠٪ / ٤٠٪، ٥٠٪ / ٥٠٪، ٤٠٪ / ٦٠٪، ٣٠٪ / ٧٠٪، ٢٠٪ / ٨٠٪، ١٠٪ / ٩٠٪ و ٠٪ / ١٠٠٪). لإعادة تصحيح المفردات متعددة الاستجابة التي تراوحت من ١ إلى ٤، استند التقسيم ثنائي الاستجابة إلى معايير أن

الاستجابات في أعلى فئة استجابة متعددة الاستجابة (٤) فقط تم ترميزها بشكل ثنائي على أنها "صحيحة" بينما الاستجابات في جميع الفئات الدنيا (٢، ١، ٣) تم ترميزها على أنها "غير صحيحة". لذلك، تم الحصول على كل درجة ثنائية الاستجابة عن طريق إعادة ترميز ١ و ٢ و ٣ ك "٠" و ٤ ك "١" (٠، ١). نتج عن هذه المعالجات في البيانات أحد عشر اختباراً ذي مفردات مختلطة، يتكون كل منها من نسبة محددة مسبقاً من البيانات ثنائية الاستجابة ومتعددة الاستجابة. أسفرت هذه المجموعات من استجابات ذان الأشكال المختلطة عبر ثلاثة أطوال اختبار عن ٣٣ مجموعة بيانات من أشكال الاستجابة التي تم تحليلها باستخدام مجموعتين من النماذج لتصحيح المفردات المختلطة: المجموعة الأولى (نموذج راش للاستجابات الثنائية ونموذج التقدير الجزئي للاستجابات المتعددة)، والمجموعة الثانية (النموذج ثنائي البارامتر للاستجابات الثنائية ونموذج الاستجابة المتدرجة للاستجابات المتعددة).

إجراءات الدراسة:

قبل تطبيق الاختبار، تم إخطار المشاركين عن الهدف من الدراسة، وذكروا طوعية انهم يرغبون ويوافقون على الاشتراك في الدراسة. لضمان استجابة المشاركين للمفردات بأمانة، تم إخبارهم بعدم كتابة أسماءهم في ورقة الاختبار. كما تم إخبارهم أنه لا ينبغي أن يهتموا بأي شيء يتعلق بمشاركتهم في الدراسة وأن إجاباتهم تكون لأغراض بحثية فقط وستبقى سرية. تم إدخال جميع البيانات في ملف SPSS.

نتائج الدراسة:

افترضت الدراسة الحالية أن طريقة التصحيح (نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة) وطول الاختبار وطريقة التحليل تؤثر على معايير التحسين في الاختبارات ذات المفردات المختلطة. معايير التحسين، بدورها، لها تأثير ملموس عبر الاختبارات ذات المفردات المختلطة على ما يلي: الثبات، ومعلومات المفردة، ومعلومات الاختبار، وتقديرات القدرة، ومعلومات المفردة. وفيما يلي عرض للأسئلة البحثية على النحو الآتي:

السؤال الأول: كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على تقديرات الكفاءة؟

تم إجراء أحد عشر انوفا ANOVA ثلاثة X اثنين (واحد لكل مستوى نسبة) على ٦٦ مجموعة من تقديرات ثيتا Theta مع طول الاختبار والنموذج كمتغيرات مستقلة. تم

تقييم التأثيرات الرئيسية والتفاعلية للأهمية من خلال مقارنة قيم F المرصودة بقيمة F الحرجة. يتم تحديد قيمة F الحرجة من خلال درجات الحرية لكل تأثير ومستوى ألفا الذي تم تعيينه عند $\alpha = 0.05$. تمت مقارنة كل F تمت المرصودة بـ F الحرجة $(\infty, 1, 0.05) = 3.84$ لاختبار أهمية نموذج و F حرجة $(\infty, 2, 0.05) = 3.00$ لاختبار دلالة طول الاختبار والتفاعل بين النموذج وطول الاختبار. تم تقييم حجم التأثير أيضا.

حجم التأثير هو نسبة مجموع المربعات بين المجموعات والمجموع الإجمالي للمربعات، وفي هذه التحليلات تمثل نسبة التباين في تقديرات ثيتا التي يفسرها الفرق بين المجموعات. تم قياس تقديرات حجم التأثير باستخدام قيم التربيعية الجزئية وتم تقييمها بناءً على ما ذكره كوهين (1988) Cohen كمرجع عام لحجم التأثير لاختبارات F في أنوفا ANOVA (صغير = 0.1 ، متوسط = 0.25 وكبير = 0.4). تم حساب حسابات القدرة المرصودة أيضاً باستخدام SPSS وتم تضمينها في نتائج أنوفا ANOVA المجدولة.

كانت قيم F أنوفا المحسوبة لمستويات النسبة من 0% إلى 50% للدرجات متعددة الاستجابة غير دالة لجميع التأثيرات الرئيسية والتفاعلية. تراوحت قيم F هذه من 0.003 إلى 2.054 لجميع التأثيرات. كما كانت قيمة F لأنوفا ANOVA لمستوى نسبة واحدة (90% متعدد الاستجابة) دالة للتأثير الرئيسي للنموذج المرصود = 3.977 بحجم تأثير 0.001 وقوة تساوي 0.524 . أدت ANOVAs التي أجريت على مستويات نسبة 60% إلى 100% وضع درجات متعددة الاستجابة جميعها إلى قيم F دالة للتأثير الرئيسي لطول الاختبار. تراوحت قيم F هذه من F المرصودة = 6.813 إلى F المرصودة = 18.455 مع بقاء حجم التأثير مساوياً أو أقل 0.006 والقوة المتبقية تساوي أو تزيد عن 0.921 . لم يؤدي أي من أنوفا ANOVAs الاحد عشر الى قيمة F دالة لتأثير تفاعل طول النموذج حسب الاختبار. النتائج معروضة في الجداول 2 و 3 و 4 و 5.

الجدول (2) تحليل التباين ثنائي الاتجاه ثابت التأثيرات على تقديرات ثيتا: النسبة 60% (النموذج، طول الاختبار)

المصدر	م.ع.	د.ج.	م.ع.	ف	قيمة P	مربع ايتا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
النموذج المصحح	5	33.662 (b)	6.732	7.869	0.001	0.007	39.347	1
التقاطع	1	43.157	43.157	50.446	0.001	0.008	50.446	1
النموذج	1	0.503	0.503	0.588	0.443	0.000	0.588	0.120
الطول	2	31.577	15.789	18.455	0.001	0.006	36.91	1
النموذج X الطول	2	1.582	0.791	0.925	0.397	0.000	1.849	0.211
الخطأ	5994	5127.947	0.856					
الكلية	6000	5204.766						
الكلية المصحح	5999	5161.609						

a محسوبة باستخدام الفا = 0,05 ، b التربيعية = 0,007 ، ر التربيعية المعدلة = 0,006

الجدول (3) تحليل التباين ثنائي الاتجاه ثابت التأثيرات على تقديرات ثيتا: النسبة ٧٠٪ (النموذج، طول الاختبار)

المصدر	م.ع.	د.ج.	م.ع.	ف	قيمة P	مربع ايتا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
النموذج المصحح	5	14.064 (b)	2.813	2.832	0.015	0.002	14.159	0.843
التقاطع	1	1.936	1.936	1.95	0.163	0.000	1.95	0.287
النموذج	1	0.018	0.018	0.018	0.893	0.000	0.018	0.052
الطول	2	13.535	6.767	6.813	0.001	0.002	13.626	0.921
النموذج X الطول	2	0.512	0.256	0.258	0.773	0.000	0.515	0.091
الخطأ	5994	5953.929	0.993					
الكلية	6000	5969.93						
الكلية المصحح	5999	5967.993						

a محسوبة باستخدام الفا = 0,05 ، b التربيعية = 0,002 ، ر التربيعية المعدلة = 0,002

الجدول (4) تحليل التباين ثنائي الاتجاه ثابت التأثيرات على تقديرات ثيتا: النسبة ٨٠٪ (النموذج، طول الاختبار)

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

المصدر	م.ع.	د.ج.	م.ع.	ف	قيمة P	مربع إبتا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
النموذج المصحح	17.553 (b)	5	3.511	3.553	0.003	0.003	17.763	0.923
التقاطع	2.6	1	2.6	2.631	0.105	0.000	2.631	0.368
النموذج	0.026	1	0.026	0.026	0.872	0.000	0.026	0.053
الطول	17.465	2	8.733	8.837	0.000	0.003	17.675	0.972
النموذج X الطول	0.062	2	0.031	0.031	0.969	0.000	0.063	0.055
الخطأ	5922.938	5994	0.988					
الكل	5943.091	6000						
الكل المصحح	5940.491	5999						

a محسوبة باستخدام الفا = 0.05 ، b ر التربيعية = 0.003 ، ر التربيعية المعدلة = 0.002

الجدول (5) تحليل التباين ثنائي الاتجاه ثابت التأثيرات على تقديرات ثيتا: النسبة ٩٠٪ (النموذج، طول الاختبار)

المصدر	م.ع.	د.ج.	م.ع.	ف	قيمة P	مربع إبتا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
النموذج المصحح	26.230 (b)	5	5.246	5.877	0	0.005	29.386	0.995
التقاطع	4.66	1	4.66	5.22	0.022	0.001	5.22	0.627
النموذج	3.55	1	3.55	3.977	0.046	0.001	3.977	0.514
الطول	21.055	2	10.527	11.794	0	0.004	23.588	0.995
النموذج X الطول	1.625	2	0.813	0.91	0.402	0	1.821	0.208
الخطأ	5350.28	5994	0.893					
الكل	5381.169	6000						
الكل المصحح	5376.51	5999						

a محسوبة باستخدام الفا = 0.05 ، b ر التربيعية = 0.005 ، ر التربيعية المعدلة = 0.004

الجدول (6) تحليل التباين ثنائي الاتجاه ثابت التأثيرات على تقديرات ثيتا: النسبة ١٠٠٪ (النموذج، طول الاختبار)

الفئة المصدرة (a)	المعلمة غير الرئيسية	مربع إينما الجزئي	P		ع.م		ع.م	المصدر
			قيمة	ف	د.ج	د.ج		
0.979	23.707	0.004	0	4.741	4.8	5	(b)24.000	
								النموذج المصحح
0.338	2.377	0	0.123	2.377	2.407	1	2.407	التقاطع
0.114	0.545	0	0.46	0.545	0.552	1	0.552	النموذج
0.993	22.491	0.004	0	11.245	11.385	2	22.769	الطول
0.104	0.671	0	0.715	0.335	0.34	2	0.679	النموذج X* الطول
					1.012	5994	6068.224	الخطأ
						6000	6094.631	الكلية
						5999	6092.225	الكلية المصحح

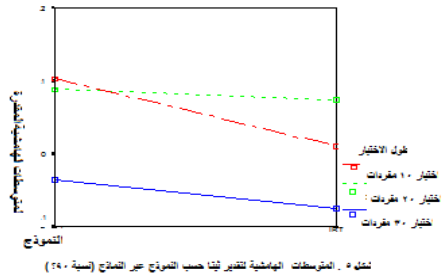
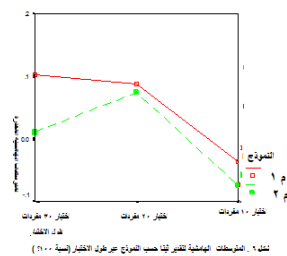
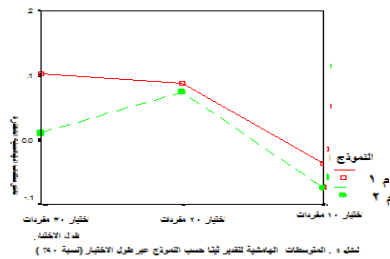
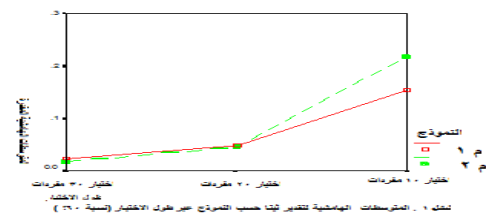
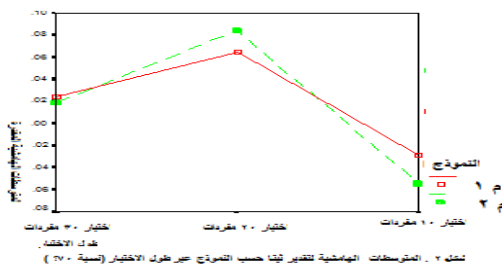
a محسوبة باستخدام الفا = 0,05 ، **b** التربيعية = 0,004 ، **r** التربيعية المعدلة = 0,003

أشارت هذه التحليلات إلى أنه مع زيادة نسبة المفردات متعددة الاستجابة إلى ما بعد ٥٠٪ من الاختبار، أثر طول الاختبار على تقديرات القدرة على الرغم من أن حجم التأثير لكل مقارنة كان ضئيلاً. تم إجراء مزيد من التحليلات للكشف عن تأثير المجموعات ضمن أطوال الاختبار ونماذج التحليل. الأشكال من ١-٦ تعرض بيانياً التأثيرات التي تركتها مجموعات العاملين two factors على تقديرات ثيتا theta estimates.

أجريت الاختبارات البعدية بعدما تبين دلالة "ف" F لتحديد اتجاه الفروق باستخدام اختبار شافيه Scheffe المتجانس لمجموعات فرعية عند مستوى ألفا يساوي 0,05. في نسبة ٦٠٪ من المفردات متعددة الاستجابة، كانت اختبارات الطول المكونة من ١٠ و ٢٠ مفردة متشابهة في المتوسطات الهامشية. ومع ذلك، كان أداء طول الاختبار المكون من ٣٠ مفردة مختلفاً بشكل دال عن الاختبارات الأقصر (الشكل ١). عند

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلفة

مستوى نسبة ٧٠٪ و ٨٠٪، كانت الاختبارات المكونة من ٢٠ مفردة و ٣٠ مفردة مختلفة بشكل دال (الشكل ٢ والشكل ٣). عند مستوى نسبة ٩٠٪، اختلف الاختبار المكون من ٣٠ مفردة اختلافاً دالاً في المتوسط الهامشي عن الاختبارين الأقصر (الشكل ٤) واختلف نمودجي المجموعتين (المجموعة الأولى والمجموعة الثانية) اختلافاً دالاً عبر أطوال الاختبار (الشكل ٥). عند مستوى النسبة ١٠٠٪، اختلفت النسب الهامشية لاختبار ٢٠ مفردة اختلافاً دالاً عن الاختبارين الآخرين (الشكل ٦).



السؤال الثاني: كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على وظيفة معلومات الاختبار (TIF)؟

تعرض الجداول ٦ و ٧ و ٨ معلومات الاختبار عبر مستويات ثيتا، ومجموع المعلومات وتحسين النسبة (بين نماذج التحليل) لكل مجموعة تصميم. توضح هذه الجداول أنه بالنسبة لجميع أطوال الاختبار، بشكل عام، تكون المعلومات أعلى في منتصف توزيع القدرة ثيتا والمعلومات عبر ثيتا وبشكل عام لها علاقة مباشرة بنسبة المفردات متعددة الاستجابة، ويحسن تحليل المجموعة الثانية (النموذج ثنائي البارامتر للاستجابات الثنائية ونموذج الاستجابة المتدرجة للاستجابات المتعددة). وكان إجمالي المعلومات أكثر للمجموعة الأولى (نموذج راش للاستجابات الثنائية ونموذج التقدير الجزئي للاستجابات المتعددة).

من المعلومات الواردة في الجداول 7 و 8 و 9، من الواضح أن الاختبارات التي تحتوي على نسبة أكبر من المفردات المصححة بشكل متعدد الاستجابة تجمع معلومات أكثر من تلك التي تم تصحيحها بشكل ثنائي الاستجابة. يُلاحظ هذا بشكل خاص في مستويات القدرة المنخفضة حيث كانت أعلى المعلومات دائمًا تقريبًا عند مستوى متعدد الاستجابة بنسبة ١٠٠٪. في مستويات القدرة الأعلى، تميل أعلى المعلومات إلى أن تكون أقل من ١٠٠٪. تعرض الأشكال ٧ و ٨ و ٩ مستويات المعلومات عبر توزيع ثيتا لاختبارات بطول ١٠ و ٢٠ و ٣٠ مفردة.

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

جدول (7) توزيع معلومات الاختبار وجمعها وتحسينها (اختبارات من ١٠ مفردات)

النموذج	التحسين %	معلومات على مدى مستويات فيتا										
		2.0	1.5	1.0	0.5	0.0	0.5-	1.0-	1.5-	2.0-		
المجموعة الثانية	0	5.6	6.5	6.5	6.7	5.6	4.5	3.5	3.5	1.8	11.5%	44.2
المجموعة الثانية	10	5.7	6.7	6.7	6.9	6.1	5.3	4.4	3.2	2.2	11.0%	47.2
المجموعة الثانية	20	5.8	6.8	6.9	7.4	7.3	6.9	5.7	3.7	2.3	11.9%	52.8
المجموعة الثانية	30	5.9	6.9	7.1	7.9	8.4	8.4	6.9	4.2	2.4	11.2%	58.1
المجموعة الثانية	40	5	6	7.7	9.1	9.3	8.9	7.1	4.3	2.5	4.0%	59.9
المجموعة الثانية	50	5	6	7.8	9.2	9.5	9.1	7.3	4.6	2.8	4.9%	61.3
المجموعة الثانية	60	5	6	7.8	9.2	9.5	9.2	7.7	5.2	3.5	3.2%	63.1
المجموعة الثانية	70	5	6	7.8	9.3	9.6	9.4	8	5.6	4.1	0.5%	64.8
المجموعة الثانية	80	4.7	6.1	8.2	9.8	10.1	9.8	8.2	5.7	4.1		66.7
المجموعة الثانية	90	4.4	5.8	8.1	9.9	10.4	10.1	8.5	6	4.3	1.8%	67.5
المجموعة الثانية	100	4.4	3.8	8	9.8	10.4	10.3	9	6.8	5.4		67.9
المجموعة الأولى	0	4.1	4.6	5.2	5.6	5.7	5.3	4.1	2.7	1.8		39.1
المجموعة الأولى	10	3.9	4.4	5	5.9	6.5	6.2	4.8	3.2	2.1		42
المجموعة الأولى	20	3.7	4.2	5	6.6	8.1	7.6	5.5	3.5	2.3		46.5
المجموعة الأولى	30	3.4	3.9	5.2	7.6	9.9	9	6.2	3.9	2.5		51.6
المجموعة الأولى	40	3.9	5.3	7	8.9	10.4	9.2	6.3	4	2.5		57.5
المجموعة الأولى	50	4	5.1	6.5	8.1	9.7	9.2	7.1	5.1	3.5		58.3
المجموعة الأولى	60	4	5.2	6.4	8	9.6	9.4	7.8	6.2	4.5		61.1
المجموعة الأولى	70	4.4	5.4	6.7	7.8	9.3	9.6	8.7	7.2	5.4		64.5
المجموعة الأولى	80	4.9	6.3	7.2	8.3	9.5	9.5	8.6	7.2	5.4	0.3%	66.9
المجموعة الأولى	90	4.5	6	7.4	8.6	9.6	9.5	8.4	7	5.3		66.3
المجموعة الأولى	100	4.5	6.1	7.5	8.6	9.5	9.5	9	8.1	6.3	1.7%	69.1

ملحوظة: القيم المكتوبة بخط عريض هي أعلى قيمة معلومات ضمن النموذج والنسبة والقدرة.

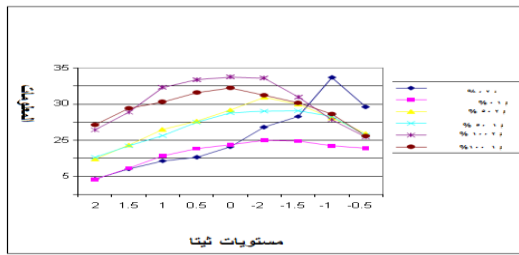
جدول (8) توزيع معلومات الاختبار وجمعها وتحسينها (اختبارات من ٢٠ مفردات)

التصنيف	مجموع العظومات	مدى مستويات بيتا										النسبة %	العلاج
		2.0	1.5	1.0	0.5	0.0	0.5-	1.0-	1.5-	2.0-			
%14.2	77.0	9.7	11.1	11.6	11.5	9.4	8	7	5.3	3.4	0	المجموعة الثانية	
%6.7	79.6	9.4	11.1	12.1	11.9	9.7	8.2	7.2	5.7	4.3	10	المجموعة الثانية	
%8.5	84.1	9.8	10.8	11.9	12.2	10.6	9.4	8.4	6.7	4.9	20	المجموعة الثانية	
%3.2	89.4	9.7	11.7	12.8	12.7	10.9	9.6	8.8	7.3	5.9	30	المجموعة الثانية	
%5.7	93.8	9.8	11.7	12.7	12.8	11.5	10.7	9.9	8.1	6.6	40	المجموعة الثانية	
%6.1	96.4	9.8	11.5	12.4	12.6	11.8	11.3	10.7	9	7.3	50	المجموعة الثانية	
%4.6	101.7	10.2	11.4	12.3	12.6	12.6	12.8	12.2	9.9	7.7	60	المجموعة الثانية	
%0.7	105.2	9.1	11.1	13	13.7	13.8	13.9	12.8	10.1	7.7	70	المجموعة الثانية	
%0.6	106.7	9.1	11.1	13	13.7	13.9	14	13.1	10.5	8.3	80	المجموعة الثانية	
%0.6	109.3	8.8	11.1	13.3	14.1	14.3	14.4	13.5	11	8.8	90	المجموعة الثانية	
%2.7	111.5	8.7	11.1	13.2	14.1	14.5	14.7	13.9	11.7	9.6	100	المجموعة الثانية	
	67.4	6.6	7.8	9.1	9.6	9.4	8.9	7.4	5.3	3.3	0	المجموعة الأولى	
	74.6	7.5	9.3	10.1	9.9	9.2	8.7	8	7	4.9	10	المجموعة الأولى	
	77.5	6.6	8.2	9.6	10.5	10.9	10.4	9	7.1	5	20	المجموعة الأولى	
	86.6	9.1	10.5	10.7	10.5	10.2	9.8	9.4	9.1	7.3	30	المجموعة الأولى	
	88.7	8	9.4	10.4	11.3	11.7	11	9.9	9.2	7.8	40	المجموعة الأولى	
	90.9	7.7	8.9	9.9	11.2	12.2	11.9	10.9	9.9	8.3	50	المجموعة الأولى	
	97.2	7.5	8.7	10.1	12.3	13.9	13.6	12	10.5	8.6	60	المجموعة الأولى	
	104.5	8	9.7	11.4	13.6	15.2	14.4	12.4	10.9	8.9	70	المجموعة الأولى	
	106.1	8	9.4	11	13.1	14.7	14.4	13.3	12.1	10.1	80	المجموعة الأولى	
	108.6	8.7	10.1	11.3	13	14.4	14.3	13.5	12.6	10.7	90	المجموعة الأولى	
	108.6	8.6	10.1	11.5	13.1	14.1	13.9	13.4	12.8	11.1	100	المجموعة الأولى	

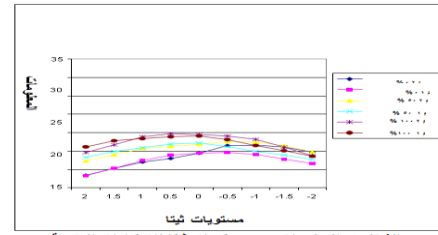
تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

جدول (9) توزيع معلومات الاختبار وجمعها وتحسينها (اختبارات من ٣٠ مفردات)

التصحيح الصبية	مجموع القطاعات	عدد مستويات بيتا										التصحيح %	المفردات
		2.0	1.5	1.0	0.5	0.0	0.5-	1.0-	1.5-	2.0-			
%34.7	141.2	24.3	32.4	21.6	18.7	13.2	10.3	9.3	7.1	4.3	0	المجموعة الثانية	
%25.0	142.8	22.7	27.7	20.3	18.8	15	12.9	11.8	8.6	5	10	المجموعة الثانية	
%13.0	147.3	20.6	21.2	20.8	20.8	17.6	15.7	14.1	10	6.5	20	المجموعة الثانية	
%7.1	151.4	18.9	20.4	21.6	22.1	18.9	16.6	14.7	10.7	7.5	30	المجموعة الثانية	
	163.9	16.7	20.2	23.6	25.4	21.8	18.6	16.3	12.2	9.1	40	المجموعة الثانية	
%4.9	175.2	16.9	21.2	25	26.9	23.4	20.3	18	13.7	9.8	50	المجموعة الثانية	
%11.6	199.9	16.5	23.1	35.1	28.2	29.7	21.4	19.3	14.9	11.7	60	المجموعة الثانية	
%4.1	194.7	16.7	20.9	24.9	27.7	26.2	24.5	22.7	17.6	13.5	70	المجموعة الثانية	
%6.8	224.1	16.4	20.9	26.7	31.6	31.7	30.8	28.4	21.3	16.3	80	المجموعة الثانية	
	209.8	15.6	22.1	25.5	28	30.7	28.4	23.2	20.2	16.1	90	المجموعة الثانية	
%5.7	230.4	15.9	20.7	27	32.2	32.5	31.8	29.6	22.8	17.9	100	المجموعة الثانية	
	104.8	12.8	13.5	14.8	15	13.8	12.7	10.7	7.3	4.2	0	المجموعة الأولى	
	114.2	12.7	13.8	15.4	16.2	16	14.9	12.2	8.3	4.7	10	المجموعة الأولى	
	130.3	12.8	15.2	17.3	18.9	18.9	17.1	13.4	9.8	6.9	20	المجموعة الأولى	
	141.4	14.1	17.8	19.9	20	19.4	17.1	13.6	10.9	8.6	30	المجموعة الأولى	
%0.5	164.7	15.6	21.2	23.1	22	20.5	17.8	15	13.2	16.3	40	المجموعة الأولى	
	167	16.4	21.6	23.2	23.1	22.6	20	16.3	13.5	10.3	50	المجموعة الأولى	
	179.2	16.4	21.8	23.6	24	23.9	21.2	17.9	16.6	13.8	60	المجموعة الأولى	
	187.1	15.8	20.6	22.7	24.3	25.9	24.1	20.5	18.3	14.9	70	المجموعة الأولى	
	209.8	15.6	22.1	25.5	28	30.7	28.4	23.2	20.2	16.1	80	المجموعة الأولى	
%2.4	214.8	16	21.6	24.3	26.5	29.5	28.7	25.7	23.6	18.9	90	المجموعة الأولى	
	217.9	16.2	22.3	25.4	27.5	29.5	28.2	25.6	23.9	19.3	100	المجموعة الأولى	



الشكل ٩. المعلومات عبر مستويات ثبات للاختبارات المكونة من ٣٠ مفردة (النسب ٢٠، ٢٥، و ٣٠) من



الشكل ٨. المعلومات عبر مستويات ثبات للاختبارات المكونة من ٢٠ مفردة (النسب ٢٠، ٢٥، و ٣٠) من

يعرض الجدول (٩) النسبة المئوية للزيادة في مستويات المعلومات عند زيادة طول الاختبار من ١٠ إلى ٢٠ مفردة لكلا النموذجين بثلاثة نسب كحد أقصى (٠٪، ٥٠٪، ١٠٠٪ تصحيح متعدد الاستجابة). يمكن ملاحظة أن زيادة طول الاختبار من ١٠ إلى ٢٠ مفردة، أو مضاعفة الطول، يزيد متوسط المعلومات عبر النماذج ونسبة المفردات متعددة الاستجابة بنسبة ٧١٪. يبدو أن هذا المعدل مستقر إلى حد ما عبر كلا النموذجين ونسب تصحيح مرتفعة. يعرض الجدول ١٠ النسب المئوية المحسوبة عند زيادة طول الاختبار من ٢٠ إلى ٣٠ مفردة. النسبة المئوية للزيادة في المتوسط هي ٨٢٪ مع أدنى قيمة لراش ٠٪ متعدد الاستجابة عند ٥٤٪ ونظرية الاستجابة للمفردة ١٠٠٪ متعدد الاستجابة الأعلى عند ١٠٤٪. يعرض الجدول ١١ النسب المئوية المحسوبة عند زيادة طول الاختبار من ١٠ إلى ٣٠ مفردة (ثلاثة أضعاف). تزيد هذه الزيادة في طول الاختبار من المعلومات الممكنة المتاحة بنسبة ٢٠٧٪ في المتوسط عبر جميع النماذج والنسب، وتتراوح من ١٦٧٪ لراش ٠٪ متعدد الاستجابة إلى ٢٥٤٪ لنظرية الاستجابة للمفردة ١٠٠٪ متعدد الاستجابة.

جدول (١٠). زيادة النسبة المئوية في المعلومات من اختبار ١٠ إلى ٢٠ مفردة

متوسط الزيادة										
	2	1.5	1	0.5	0	0.5-	1-	1.5-	2-	
%76	%73	%71	%78	%72	%68	%78	%100	%51	%89	IRT-0%
%74	%61	%70	%75	%71	%65	%68	%80	%96	%83	Rasch-0%
%71	%96	%92	%59	%37	%24	%24	%47	%96	%161	IRT-50%
%66	%93	%75	%52	%38	%26	%29	%54	%94	%137	Rasch-50%
%76	%98	%192	%65	%44	%39	%43	%54	%72	%78	IRT-100%
%60	%91	%66	%53	%52	%48	%46	%49	%58	%76	Rasch-100%
%71										الاجملي

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

جدول (١١). زيادة النسبة المئوية في المعلومات من اختبار ٢٠ إلى ٣٠ مفردة

متوسط الزيادة	2		1.5		1		0.5		0		0.5-		1-		1.5-		2-		
%73	%151	%192	%86	%63	%40	%29	%33	%34	%26	IRT-0%									
%54	%94	%73	%63	%56	%47	%43	%45	%38	%27	Rasch-0%									
%78	%72	%84	%102	%113	%98	%80	%68	%52	%34	IRT-50%									
%84	%113	%143	%134	%106	%85	%68	%50	%36	%24	Rasch-50%									
%104	%83	%86	%105	%128	%124	%116	%113	%95	%86	IRT-100%									
%100	%88	%121	%121	%110	%109	%103	%91	%87	%74	Rasch-100%									
%82										الاجملى									

جدول (١٢). زيادة النسبة المئوية في المعلومات من اختبار ١٠ إلى ٣٠ مفردة

متوسط الزيادة	2		1.5		1		0.5		0		0.5-		1-		1.5-		2-		
%202	%334	%398	%232	%179	%136	%129	%166	%103	%139	IRT-0%									
%167	%212	%193	%185	%168	%142	%140	%161	%170	%133	Rasch-0%									
%196	%238	%253	%221	%192	%146	%123	%147	%198	%250	IRT-50%									
%202	%310	%324	%257	%185	%133	%117	%130	%165	%194	Rasch-50%									
%254	%261	%445	%238	%229	%213	%209	%229	%235	%231	IRT-100%									
%220	%260	%266	%239	%220	%211	%197	%184	%195	%206	Rasch-100%									
%207										الاجملى									

للإجابة على السؤال البحثي الثاني، أشارت الحسابات إلى وجود علاقة مباشرة بين النسب الأعلى للمفردات متعددة الاستجابة ومقدار المعلومات المتاحة التي ينتجها الاختبار. أشارت جميع التحليلات إلى أن طول الاختبار يلعب دورًا رئيسيًا في تحديد مقدار المعلومات التي يتم جمعها من المفردات. يبدو هذا واضحًا لأن معلومات الاختبار الإجمالية هي مجموع كل قيم معلومات المفردة. مع وجود المزيد من المفردات، يتوفر المزيد من المعلومات.

السؤال الثالث: كيف تؤثر نسبة المفردات ثنائية الاستجابة ومتعددة الاستجابة عبر طول الاختبار ونماذج التحليل على التحسن الإجمالي للاختبار؟

تم استخدام ثلاثة معايير لتقييم أمثلية الاختبار عبر مجموعات التصميم في هذه الدراسة: الثبات الهامشي، تحليلات الأنوفا ANOVA لمعدلات الخطأ المعياري لكل مجموعة تصميم، ومستويات المعلومات. يعرض الجدول ١٣ قيم الثبات الهامشي عبر مجموعات التصميم. داخل كل اختبار طول يزداد الثبات عبر نسب المفردات متعددة الاستجابة. وتجدر الإشارة أيضًا إلى أنه عبر جميع أطوال الاختبار الثلاثة، زاد الثبات مع إطالة الاختبار. بالنظر إلى معاملات الثبات الهامشية عبر النماذج، فإن الاختلافات بين مجموعتي النماذج لا تكاد تذكر.

تتألف المعايير الثانية التي تم تحليلها لتحديد أمثلية الاختبار من ١١ انوفا 3×2 (واحد لكل مستوى نسبة) أجريت على ٦٦ مجموعة من معدلات الخطأ المعيارية بطول الاختبار والنموذج كمتغيرات مستقلة. تم تقييم التأثيرات الرئيسية والتفاعلية للأهمية من خلال مقارنة قيم F المرصودة بقيمة F الحرجة. يتم تحديد قيمة F الحرجة من خلال درجات الحرية لكل تأثير ومستوى ألفا الذي تم تعيينه عند $\alpha = 0.05$. تمت مقارنة كل F تمت مرصودة مع F حرجة $(\infty, 2, 0.05)$ $3.84 =$ لاختبار أهمية النموذج و F حرجة $(\infty, 2, 0.05)$ $3.00 =$ لاختبار أهمية طول الاختبار والتفاعل بين النموذج وطول الاختبار. تم تقييم حجم التأثير أيضاً. حجم التأثير هو نسبة مجموع المربعات بين المجموعات والمجموع الإجمالي للمربعات، وفي هذه التحليلات تمثل نسبة التباين في تقديرات الخطأ المعياري الموضح بالفرق بين المجموعات. وكما ورد سابقاً، تم قياس تقديرات حجم التأثير بقيم مربع إيتا الجزئية وتم تقييمها وفقاً لمعايير كوهين (صغير = ٠,١، متوسط = ٠,٢٥، وكبير = ٠,٤). تم حساب حسابات القدرة المرصودة أيضاً باستخدام SPSS وتم تضمينها في نتائج الانوفا ANOVA المجدولة

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

جدول (١٣). الثبات الهامشي عبر طول الاختبار ونسبة التصحيح والنموذج

المفردة	% الوزن	١م	2م	التقريب
	%0	0.78	0.78	0
	%10	0.81	0.81	0.01-
	%20	0.82	0.84	0.02-
	%30	0.84	0.85	0.03-
	%40	0.86	0.86	0
10	%50	0.86	0.87	0.01-
	%60	0.87	0.87	0
	%70	0.87	0.87	0
	%80	0.88	0.88	0
	%90	0.88	0.88	0
	%100	0.88	0.88	0
	%0	0.88	0.88	0
	%10	0.89	0.89	0
	%20	0.89	0.9	0.01-
	%30	0.9	0.9	0
	%40	0.9	0.91	0.01-
20	%50	0.91	0.91	0
	%60	0.91	0.92	0.01-
	%70	0.92	0.92	0
	%80	0.92	0.92	0
	%90	0.92	0.92	0
	%100	0.92	0.92	0

	%80	0.92	0.92	0
	%90	0.92	0.92	0
	%100	0.92	0.92	0
	%0	0.92	0.92	0
	%10	0.93	0.93	0
	%20	0.94	0.94	0
	%30	0.94	0.94	0
	%40	0.95	0.95	0
	%50	0.95	0.95	0
30	%60	0.95	0.96	0.01-
	%70	0.96	0.96	0
	%80	0.96	0.96	0
	%90	0.96	0.96	0
	%100	0.96	0.96	0

أدت الانوفا ANOVA المحسوبة لجميع مستويات النسبة الأحد عشر لثلاثة تأثيرات إلى ٣٠ قيمة F دالة للنموذج (١١)، وطول الاختبار (١١) وتفاعل طول النموذج على حدة (٨) مع الاستثناءات الثلاثة فقط: تأثيرات التفاعل لـ ٠ لم تكن نماذج النسبة ٠٪ و ٥٠٪ و ١٠٠٪ دالة. كانت نسب ٠٪ و ١٠٠٪ هي الاختبارات الوحيدة التي تم تحليلها والتي لم تكن اختبارات ذات مفردات مختلطة حيث تم تسجيل مستوى النسبة ٠٪ بشكل ثنائي الاستجابة بنسبة ١٠٠٪ وتم تصحيح نموذج النسبة ١٠٠٪ بشكل متعدد الاستجابة. تراوحت قيم F الدالة من تحليلات الأنوفا ANOVAs الأخرى من ١١,٠١٩ إلى ٩٥٠,٣١٩ للتأثير الرئيسي للنموذج، من ٢٢٦٥,٢١٣ إلى ٣٢٢١٦,٦٩٨ للتأثير الرئيسي لطول الاختبار، ومن ٤,١٣٩ إلى ٧٣٤٨,١١ لتأثير التفاعل بين النموذج وطول الاختبار. يتم عرض النتائج في الجداول من ١٤ إلى ٢٤.

تراوح حجم التأثير عبر جميع النسب للتأثير الرئيسي للنموذج من ٠,٠٠٢ إلى ٠,٠١٦ مع استثناءين: نسبة ٦٠٪ لها حجم تأثير ٠,١٣٧ ونسبة ٩٠٪ لها حجم تأثير ٠,٥٩١. تراوح حجم التأثير عبر جميع النسب للتأثير الرئيسي لطول الاختبار من ٠,٦٠٨ إلى ٠,٧٨١ مع نفس الاستثناءين: كان لنسبة ٦٠٪ حجم تأثير ٠,٤٣٠ ونسبة ٩٠٪ كان حجم تأثيرها ٠,٩١٥. تراوح حجم التأثير عبر جميع النسب لتأثير تفاعل النموذج حسب الطول من ٠,٠٠٠ إلى ٠,٠١٢ مع نفس الاستثناءين: نسبة ٦٠٪ كان لها حجم تأثير ٠,١٦٠ ونسبة ٩٠٪ لها حجم تأثير ٠,٧١٠.

ظلت القدرة عبر جميع النسب ثابتة بين ٠,٩١٣ و ١ للتأثير الرئيسي للنموذج وتساوي ١ لجميع التأثيرات الرئيسية بطول الاختبار. بالنسبة لتأثير التفاعل، تراوحت قوة تأثير التفاعل من ٠,٧٣٣ إلى ١ لثمانية من مستويات النسبة مع انخفاض مستويات النسبة ٠٪ و ٥٠٪ و ١٠٠٪ أقل بكثير من الباقي عند ٠,٣٧٥ و ٠,٣٤٣ و ٠,١٤٩ على التوالي.

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلفة

جدول (١٤). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٠.١٪ (النموذج، وطول الاختبار)

المصدر	مجموع المربعات	متوسط المربعات	P	مربع ابنا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
	درجات الحرية	ف	قيمة			
النموذج المصحح	(b)28.664	5.733	0	1990.358	0.624	9951.788
التقاطع	771.97	771.97	0	268017.996	0.978	268017.996
النموذج	0.059	0.059	0	20.515	0.003	20.515
الطول	28.595	14.297	0	4963.852	0.624	9927.703
النموذج * الطول	0.01	0.005	0.168	1.785	0.001	0.375
الخطأ	17.264	0.003				
الكلى	817.899	6000				
الكلى المصحح	45.929	5999				

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = ٠,٦٢٤ (ر التربيعية المعدلة = ٠,٦٢٤)

جدول (١٥). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٠.١٪ (النموذج، وطول الاختبار)

المصدر	مجموع المربعات	متوسط المربعات	P	مربع ابنا الجزئي	المعلمة غير الرئيسية	القوة المرصودة (a)
	درجات الحرية	ف	قيمة			
النموذج المصحح	(b)27.574	5.515	0	2426.215	0.669	12131.074
التقاطع	703.077	703.077	0	309315.367	0.981	309315.367
النموذج	0.031	0.031	0	13.484	0.002	13.484
الطول	27.525	13.762	0	6054.656	0.669	12109.312
النموذج * الطول	0.019	0.009	0.016	4.139	0.001	0.733
الخطأ	13.624	0.002				
الكلى	744.275	6000				
الكلى المصحح	41.198	5999				

(a) محسوبة باستخدام ألفا ٠,٠٥ ؛ b ر التربيعية = ٠,٦٦٩ (ر التربيعية المعدلة = ٠,٦٦٩)

جدول (١٦). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٢٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايٲا الجزئية	P قيمة	متوسط المربعات		درجات الحرية	مجموع المربعات	المصدر
				ف	ف			
1	11481.046	0.657	0	2296.209	4.703	5	(b)23.515	النموذج المصحح
								التقاطع
1	307302.91	0.981	0	307302.91	629.415	1	629.415	النموذج المصحح
1	41.979	0.007	0	41.979	0.086	1	0.086	الطول
1	11407.7	0.656	0	5703.85	11.683	2	23.365	النموذج * الطول
						0.032 2	0.064	
	0.999	31.368	0.005	0	15.684			الخطأ
						0.002 5994	12.277	
						6000	665.207	الكلية
						5999	35.792	الكلية المصحح

(a) محسوبة باستخدام ألفا ٠,٠٥؛ b ر التربيعية = ٠,٦٥٧ (ر التربيعية المعدلة = ٠,٦٥٧)

جدول (١٧). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٣٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايٲا الجزئية	P قيمة	متوسط المربعات		درجات الحرية	مجموع المربعات	المصدر
				ف	ف			
1	9382.93	0.61	0	1876.586	3.889	5	(b)19.445	النموذج المصحح
								التقاطع
1	278976.18	0.979	0	278976.18	578.143	1	578.143	النموذج المصحح
0.913	11.019	0.002	0.001	11.019	0.023	1	0.023	الطول
1	9299.424	0.608	0	4649.712	9.636	2	19.272	النموذج * الطول
						0.075 2	0.15	
	1	72.487	0.012	0	36.244			الخطأ
						0.002 5994	12.422	
						6000	610.01	الكلية
						5999	31.867	الكلية المصحح

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلفة

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.610 (ر التربيعية المعدلة 0.610)

جدول (١٨). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة 40% (النموذج، وطول الاختبار)

القوة المصدرة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	متوسط المربعات ف	درجات الحرية	مجموع المربعات	المصدر
1	10002.015	0.625	0	2000.403	3.791	5	النموذج المصحح (b) 18.955
							التقاطع
1	286552.14	0.98	0	286552.14	543.057	1	543.057
0.997	22.094	0.004	0	22.094	0.042	1	0.042
1	9968.824	0.625	0	4984.412	9.446	2	18.892
					0.011	2	0.021
	0.856	11.097	0.002	0.004	5.549		
					0.002	5994	11.359
						6000	573.372
						5999	30.315
							الخطأ الكلي المصحح

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.625 (ر التربيعية المعدلة 0.625)

جدول (١٩). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة 50% (النموذج، وطول الاختبار)

القوة المصدرة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	متوسط المربعات F	درجات الحرية	مجموع المربعات	المصدر
1	12409.319	0.674	0	2481.864	3.934	5	النموذج المصحح (b) 19.671
							التقاطع
1	333932.24	0.982	0	333932.24	529.345	1	529.345
1	63.74	0.011	0	63.74	0.101	1	0.101
1	12342.356	0.673	0	6171.178	9.782	2	19.565
					0.003	2	0.005
	0.343	3.223	0.001	0.2	1.611		
					0.002	5994	9.502
						6000	558.518
						5999	29.173
							الخطأ الكلي المصحح

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.674 (ر التربيعية المعدلة = 0.674)

جدول (٢٠). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٦٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	متوسط المربعات	درجات الحرية	مجموع المربعات	المصدر
			ف				
1	6626.512	0.525	0	1325.302	2.823	5	التنوع المصحح (b)14.117
							التقاطع
1	263455.45	0.978	0	263455.45	561.275	1	561.275
1	950.319	0.137	0	950.319	2.025	1	2.025
1	4530.427	0.43	0	2265.213	4.826	2	9.652
						1.22 2	2.441
							النموذج * الطول
1	1145.766	0.16	0	572.883			
						0.002 5994	12.77
							الخطأ
						6000	588.162
						5999	26.887
							الكل المصحح

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.525 (ر التربيعية المعدلة = 0.525)

جدول (٢١). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٧٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	متوسط المربعات	درجات الحرية	مجموع المربعات	المصدر
			ف				
1	17012.977	0.739	0	3402.595	4.242	5	التنوع المصحح (b)21.211
							التقاطع
1	388599.48	0.985	0	388599.48	484.489	1	484.489
1	46.298	0.008	0	46.298	0.058	1	0.058
1	16954.014	0.739	0	8477.007	10.569	2	21.138
						0.008 2	0.016
							النموذج * الطول
0.9	12.666	0.002	0.002	6.333			
						0.001 5994	7.473
							الخطأ
						6000	513.174
						5999	28.684
							الكل المصحح

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

(a) محسوبة باستخدام ألفا $0,05$ ؛ b ر التربيعية = $0,739$ (ر التربيعية المعدلة $0,739 =$)

جدول (٢٢). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٨٠٪ (النموذج، وطول الاختبار)

المصدر	درجات الحرية	متوسط المربعات	P	مربع ايتا الجزئية	المعلمة غير الرئيسية	القوة المرصودة (a)
النموذج المصحح (b)	22.193	4.439	0	0.759	18898.362	1
التقاطع	1	462.381	0	0.985	393732.7	1
النموذج	1	0.019	0	0.003	16.27	0.981
الطول	2	11.082	0	0.759	18872.622	1
النموذج * الطول	2	0.006	0.009	0.002	9.471	0.793
الخطأ	5994	7.039				
الكلية	6000	491.614				
الكلية المصحح	5999	29.232				

(a) محسوبة باستخدام ألفا $0,05$ ؛ b ر التربيعية = $0,759$ (ر التربيعية المعدلة $0,759 =$)

جدول (٢٣). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ٩٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	ف	متوسط المربعات	درجات الحرية	مجموع المربعات	المصدر	
1	87779.444	0.936	0	17555.889	19.272	5	(b)96.361	النموذج المصحح	
								التقاطع	
1	529881.81	0.989	0	529881.81	581.684	1	581.684		
1	8649.715	0.591	0	8649.715	9.495	1	9.495	النموذج	
1	64433.396	0.915	0	32216.698	35.366	2	70.733	الطول	
						8.067	2	16.133	النموذج * الطول
1	14696.332	0.71	0	7348.166					
						0.001	5994	6.58	الخطأ
							6000	684.625	الكلى
							5999	102.941	الكلى المصحح

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.936 ، ر (التربيعية المعدلة = 0.936)

جدول (٢٤). تحليل التباين الثنائي للتأثيرات الثابتة ثنائية الاستجابة على الأخطاء المعيارية: النسبة ١٠٠٪ (النموذج، وطول الاختبار)

القوة المرصودة (a)	المعلمة غير الرئيسية	مربع ايتا الجزئية	P قيمة	ف	متوسط المربعات	درجات الحرية	مجموع المربعات	المصدر	
1	21439.27	0.782	0	4287.854	4.601	5	(b)23.003	النموذج المصحح	
								التقاطع	
1	416485.31	0.986	0	416485.31	446.861	1	446.861		
1	96.287	0.016	0	96.287	0.103	1	0.103	النموذج	
1	21341.806	0.781	0	10670.903	11.449	2	22.898	الطول	
						0.001	2	0.001	النموذج * الطول
	0.149	1.177	0	0.555	0.588				
						0.001	5994	6.431	الخطأ
							6000	476.295	الكلى
							5999	29.434	الكلى المصحح

(a) محسوبة باستخدام ألفا 0.05 ؛ b ر التربيعية = 0.782 ، ر (التربيعية المعدلة = 0.781)

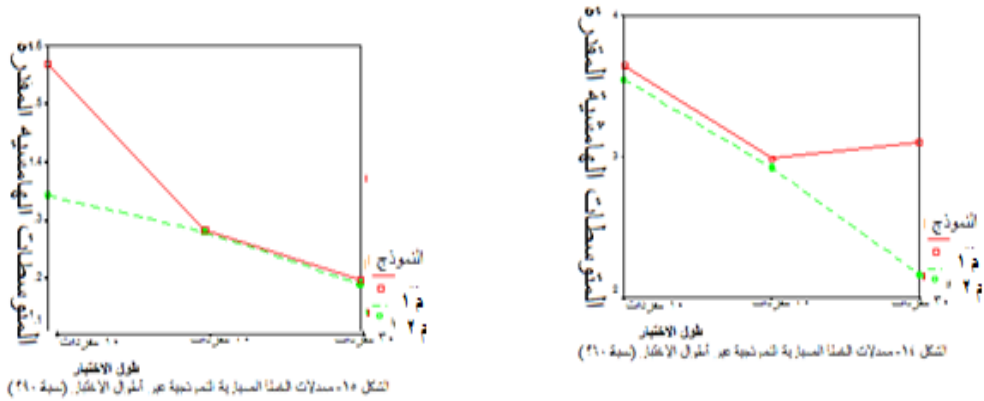
للإجابة على السؤال البحثي الثالث، أشار تحليل المعايير الأولى إلى أن الاختبارات المكونة من ٣٠ مفردة ذات النسب الأعلى من المفردات متعددة الاستجابة بغض النظر عن طول الاختبار أدت إلى معدلات ثبات هامشي أعلى. لا توجد فروق ذات دلالة إحصائية في حسابات الثبات المجموعتين من النماذج.

أشارت المعايير الثانية، تحليلات الأنوفا ANOVA لمعدلات الخطأ المعيارية، إلى أن التأثيرات الرئيسية لطول ونموذج الاختبار كانت دالة لجميع مجموعات التصميم و٨ من أصل ١١ اختبارًا للتفاعل. يوضح الشكلان ١٠ و ١١ العلاقة النموذجية بين العوامل لمعظم النماذج. أشارت تحليلات أنوفا ANOVAs إلى أن الاختبارات ذات المفردات المختلطة قد حصلت على معدلات خطأ معيارية لثبات والتي انخفضت مع زيادة طول الاختبار من ١٠ إلى ٢٠ ثم إلى ٣٠ مفردة. تظهر النتائج أيضًا أن معدلات الخطأ تميل إلى أن تكون مختلفة بشكل كبير بالنسبة للمجموعتين من نماذج التحليل في جميع أطوال الاختبار. ومع ذلك، يجب النظر إلى هذه النتيجة في ضوء حقيقة أن حجم التأثير كان ضئيلاً.

ومع ذلك، يجب النظر إلى هذه النتيجة في ضوء حقيقة أن حجم التأثير كان ضئيلاً. يعرض الشكلان ١٢ و ١٣ بياناً العلاقة بين ٦٠٪ و ٩٠٪ من المستويات النسبية التي اختلفت في الدلالة عن المستويات الأخرى عبر النماذج. يشير كلا الرقمين إلى أن معدل الخطأ لكل من المجموعتين من نماذج التحليل يختلف بشكل دال عن طول اختبار واحد. في كلتا الحالتين، كان أداء المجموعة الثانية أفضل من المجموعة الأولى لحزمة راش Rasch في تقليل معدلات الخطأ المعيارية لتقديرات ثبات. يوضح الجدول ٢٥ تقديرات حجم التأثير لجميع مستويات النسب. يُظهر النمط المعروف عبر تقديرات حجم التأثير زيادة في الأخطاء المعيارية. هذه الزيادة ناتجة عن طول الاختبار وزيادة في نسبة التصحيح متعدد الاستجابة.

كانت المعايير الثالثة التي تم تحليلها لتقييم الأمثلة هي مستوى المعلومات عبر أطوال الاختبار والنماذج ومستويات النسب. بالرجوع إلى الجداول ٦ و ٧ و ٨، لوحظ أنه بشكل عام، أدت نظرية الاستجابة للمفردة IRT إلى معلومات إجمالية أعلى عبر جميع مستويات النسب وحسنت كمية المعلومات التي تم جمعها مع زيادة طول الاختبار. أيضاً، تم إنتاج كميات أكبر من المعلومات مع زيادة نسبة المفردات متعددة الاستجابة.





جدول (٢٥). مقارنة حجم التأثير عبر نتائج الخطأ القياسية

تفاعل النموذج في الطول	الطول	النموذج	توليفة التصميم
0.001	0.624	0.003	%0
0.001	0.669	0.002	%10
0.005	0.656	0.007	%20
0.012	0.608	0.002	%30
0.002	0.625	0.004	%40
0.001	0.673	0.011	%50
*0.160	*0.430	*0.137	%60
0.002	0.739	0.008	%70
0.002	0.759	0.003	%80
*0.710	*0.915	*0.591	%90
0	0.781	0.016	%100

* تختلف أحجام التأثير للنموذج ٦٠٪ و ٩٠٪ عن نمط المخرجات الأخرى لأنها تعكس توزيعات بيانات الاختبار غير العادية

السؤال الرابع: هل توجد فروق بين مجموعتي التصحيح: أساليب المجموعة الأولى في التحليل (RASCH نموذج راش ونموذج التقدير الجزئي PCM) والمجموعة

(E-mail of the correspondig: jcps@art.bsu.edu.eg)

الثانية (2 PL) النموذج ثنائي البارامتر ونموذج الاستجابة المتدرجة (GRM) في قدرتها على تحليل الاختبارات ذات المفردات المختلطة؟

تم تقييم خمسة معايير لمقارنة أساليب المجموعة الأولى (نموذج راش للاستجابات الثنائية ونموذج التقدير الجزئي للاستجابات المتعددة)، والمجموعة الثانية (النموذج ثنائي البارامتر للاستجابات الثنائية ونموذج الاستجابة المتدرجة للاستجابات المتعددة) وهي افتراضات لكل نموذج (إحصائيات المطابقة)، وجودة تقدير الدرجة الحقيقية (تقديرات ثبات وحجم الأخطاء المعيارية)، والثبات والمعلومات.

البيانات أكثر ملاءمة للمجموعة الثانية (النموذج ثنائي البارامتر للاستجابات الثنائية ونموذج الاستجابة المتدرجة للاستجابات المتعددة). الذي سمح للتمييز بالتنوع. تم تقييم الإحصاء $2LL$ ، وهو تقدير لإحصاء مربع كاي، للنموذج الملائم للبيانات من خلال مقارنة الفرق χ^2 المرصودة مع χ^2 الحرجة. وكما هو واضح من الجدول ٢٦، توجد فروق دالة بين مطابقة البيانات لنماذج مجموعتين التحليل.

جدول (٢٦). سلبي ضعف نسبة احتمالية التسجيل (-LL٢)

المقررات	% البيولي	1م	L2	χ^2 الفرق	د.ح	p	مطابقه أفضل
		LL2-	LL2-				
10	%0	4943.5-	5017.5-	74	9	0.001	1م
10	%10	3573.1-	3650-	76.9	9	0.001	1م
10	%20	2355.4-	2423.6-	68.2	9	0.001	1م
10	%30	1203.4-	1247.5-	44.1	9	0.001	1م
10	%40	11.3	6.5-	17.8	9	0.05	2م
10	%50	714.5	677.2	37.3	9	0.001	2م
10	%60	1213.3	1194.7	18.6	9	0.05	2م
10	%70	1843.5	1848.1	4.6	9	0.9	مشابه
10	%80	3221.1	3230.9	9.8	9	0.5	مشابه
10	%90	4891.3	4875.6	15.7	9	0.1	مشابه
10	%100	5438.6	5419	19.6	9	0.05	2م
20	%0	2323.8	2156.1	167.7	19	0.001	1م
20	%10	3853	3713	140	19	0.001	2م

تأثيرات طول الاختبار وطريقة التصحيح ونموذج تحليل المفردة على تقديرات القدرة ومعايير التحسين للاختبارات ذات المفردات المختلطة

2م	0.001	19	179.7	6559.3	6739	%20	20
1م	0.001	19	81.1	7738	7819.1	%30	20
2م	0.001	19	163.3	9339.7	9503	%40	20
2م	0.001	19	143.9	11496.9	11640.8	%50	20
1م	0.001	19	112.8	13949.6	14062.4	%60	20
2م	0.001	19	44.9	16300	16344.9	%70	20
2م	0.001	19	50.5	17482.7	17533.2	%80	20
1م	0.02	19	34.7	19525.4	19560.1	%90	20
2م	0.001	19	89.5	21676.3	21765.8	%100	20
1م	0.001	29	318.5	6819.5	7138	%0	30
2م	0.001	29	297.5	10891.7	11189.2	%10	30
2م	0.001	29	254.3	13371	13625.3	%20	30
1م	0.001	29	216.3	16174.5	16390.8	%30	30
2م	0.001	29	138.1	18811.4	18949.5	%40	30
2م	0.001	29	197.6	22459.7	22657.3	%50	30
1م	0.001	29	156.4	24347.8	24504.2	%60	30
2م	0.001	29	116.6	27760.9	27877.5	%70	30
1م	0.001	29	1787.8	32970.7	31182.9	%80	30
2م	0.001	29	1803.3	31182.9	32986.2	%90	30
2م	0.001	29	96.5	36465	36561.5	%100	30

مناقشة النتائج:

تقديرات القدرة:

في هذه الدراسة، تمت مقارنة تقديرات ثيتا عبر نماذج تحليل نظرية الاستجابة للمفردة IRT عبر ٣٣ شكلاً مختلفاً من أشكال الاختبار (ثلاثة أطوال اختبار في أحد عشر نسباً ثابتة من المفردات متعددة الاستجابة وثنائية الاستجابة). تم تطبيق نماذج نظرية الاستجابة للمفردة IRT على كل مجموعة من مجموعات بيانات الاستجابة البالغ عددها ٣٣ لتوليد تقديرات القدرة. ثم تمت مقارنة تقديرات القدرة هذه باستخدام أسلوب تحليل التباين ANOVA العالمية ذات التأثيرات الثابتة. تمت مقارنة التأثيرات الرئيسية والتفاعلية لنموذج التحليل وطول الاختبار. لوحظت نتائج دالة للاختبارات

ذات النسب العالية من المفردات متعددة الاستجابة. تأثيرات النموذج والتفاعل، بشكل عام، لم تكن دالة.

توجد فروق في تقديرات القدرة عبر أطوال الاختبار عندما تم تصحيح ٦٠% إلى ١٠٠% من المفردات بشكل متعدد الاستجابة. يجب النظر إلى هذه الفروق في ضوء حقيقة أن حجم التأثير على تقديرات القدرة كان ضئيلاً، وعلى مدى المستويات الخمسة من النسبة، لم ينتج عن طول معين بشكل ثابت تقديرات مختلفة للقدرة. لذلك

بالنسبة للاختبارات ذات المفردات المختلطة، لا يمكن التوصل إلى استنتاجات قوية من هذه التحليلات بأن تقديرات القدرة تختلف بشكل دال بغض النظر عن طول الاختبار أو نموذج التحليل أو نسبة المفردات في الاختبار الذي تم تصحيحه بشكل متعدد الاستجابة. هذه النتيجة هامة عند أخذها في الاعتبار مع النتائج التي توصل إليها ساقى (Si,2002) بأن نماذج الاستجابة للمفردة تختلف بشكل دال في قدرتها على استعادة تقديرات القدرة الأصلية. إلا أن تحليلات ساقى Si لم تكن في الاختبارات ذات المفردات المختلطة ولكن على مجموعات استجابة ثنائية الاستجابة أو متعددة الاستجابة. كما أن نماذج نظرية الاستجابة للمفردة IRT التي تمت مقارنتها هنا أدت إلى تقديرات قدرة مماثلة لأن أشكال وضع الدرجات للاختبار كانت مختلطة.

معلومات الاختبار:

أظهرت هذه الدراسة أنه في الاختبارات ذات المفردات المختلطة، فإن طول الاختبار ونسبة المفردات التي تم تصحيحها بشكل متعدد الاستجابة لهما علاقات مباشرة مع كمية معلومات الاختبار المتاحة. بشكل عام، تتيح الاختبارات الأطول التي تحتوي على نسب أكبر من المفردات متعددة الاستجابة مزيداً من المعلومات عبر توزيع مستويات القدرة. هذا ملحوظ بشكل خاص في مستويات القدرة المنخفضة حيث كانت أعلى المعلومات دائماً تقريباً عند مستوى نسبة ١٠٠% للمفردات متعددة الاستجابة. بالنسبة للطلاب ذوي المستويات الأعلى من القدرة، تميل المعلومات الأعلى إلى أن تكون على مستويات متناسبة مع درجات متعددة الاستجابة أقل من ١٠٠%. نتجت أعلى معلومات اختبار إجمالية عن تصحيح متعدد الاستجابة بنسبة ١٠٠% عبر جميع أطوال الاختبار وعبر كلا نموذجي التحليل المطبقين في هذه الدراسة.

تشير الأدبيات (Ayanwale & Adeleke,2020) إلى أن المفردات متعددة الاستجابة تسمح للممتحنين بالحصول على تقدير جزئي لإثبات المعرفة الجزئية. تميل المفردات

ثنائية الاستجابة إلى تجاهل وجود معرفة جزئية بدلاً من تحديد جميع الاستجابات في سياق "الكل أو لا شيء" (Ayanwale, 2021). تسمح الاستجابات متعددة الاستجابة باستمرار القياس لمستويات القدرة، وبالتالي توفر فرصة لجمع المزيد من المعلومات من النطاق المتوسط للقدرة على كل مفردة. من الواضح أن هذا المدخل، بدوره، سيسمح بجمع المزيد من المعلومات. يجب أن يخدم الوعي بهذه النتيجة مصممي الاختبارات الذين يسعون جاهدين لقياس القدرة الحقيقية للممتحنين بشكل أكثر دقة من خلال تصميم مفردات اختبار أكثر فعالية.

بالنسبة للمفردات متعددة الاستجابة بشكل خاص، فإن التكاليف التي يتكبدها وقت تطبيق الاختبار، وكتابة المفردات، والاختبار الميداني للمفردات وتصحيح المفردات، جنباً إلى جنب مع التطبيق العملي لأطوال الاختبار، تعتبر مهمة في موازنة الرغبة في الحصول على مستويات معلومات عالية. ومع ذلك، فقد تم البحث في التطورات التكنولوجية الحديثة، مثل تصحيح المقالات الإلكترونية (Shermis & Berstein, 2002; Shermis, Koch, Page, Keith & Harrington, 2003) وتبين أنها فعالة بنفس القدر إن لم تكن أكثر فعالية من التصحيح اليدوي.

التحسين:

يعتمد تحسين الاختبار على تعظيم كل من افتراضات النموذج الإحصائي ومعايير التحسين لهذا النموذج، والتي تتضمن التوافق بين النموذج والبيانات. حدد الإحصاء الملائم الذي تم تحليله (-LL₂) أن البيانات تتلاءم مع النماذج المختارة من نظرية الاستجابة للمفردة IRT بشكل أفضل في اختبار 10 مفردات مع تصحيح متعدد الاستجابة بنسبة 40%. أظهرت نماذج المجموعة الثانية (النموذج ثنائي البارامتر ونموذج الاستجابة المتدرجة) مطابقة أفضل لمعظم مجموعات البيانات عند مقارنتها بنماذج المجموعة الأولى.

تم تحليل معدلات الخطأ المعيارية لتقديرات القدرة وفقاً لمعايير التحسين وتم الحصول على نتائج دالة عبر طول الاختبار بأحجام تأثير كبيرة. كما لوحظت التأثيرات على الأخطاء المعيارية من تأثيرات النموذج وتفاعل طول النموذج تلو الآخر لكن أحجام التأثير لكليهما كانت صغيرة. مع استثناءات قليلة، اختلفت معدلات الخطأ القياسية في جميع أطوال الاختبار وكلا النموذجين. انخفضت الأخطاء المعيارية مع زيادة طول الاختبار مما أظهر علاقة غير مباشرة أو عكسية. بالنسبة لهذه المقارنات، زاد حجم تأثير طول الاختبار على الخطأ المعياري مع زيادة نسبة

المفردات متعددة الاستجابة. كان هذا صحيحًا لجميع أطوال الاختبار الثلاثة وكانت قيم القدرة العالية موجودة أيضًا.

تبين زيادة نتائج الثبات عبر أطوال الاختبار ومستويات النسبة. حقق الاختبار المكون من ٣٠ مفردة عند التصحيح المتعدد الاستجابة بنسبة ١٠٠٪ أعلى معدلات الثبات عبر كلا النموذجين. تتوافق هذه النتيجة مع دراسة لو ووينج (Lau and Wang, 1998) اللذين وجدوا أن الجمع بين نوعين من المفردات في تقييم واحد يعزز ثبات التقييم. وتجدر الإشارة إلى أن مستويات المعلومات قدمت مقاييس أكثر دالة لتحسين الاختبار ذي المفردات المختلطة.

المقارنة بين المجموعتين من نماذج التحليل لنظرية الاستجابة للمفردة IRT:

بشكل عام، لم يتم ملاحظة فروق ذات دلالة إحصائية بين أداء المجموعتين على التحليلات في هذه الدراسة. بشكل عام عبر تحليلات المعلومات في هذه الدراسة، أدى استخدام تحليل المجموعة الثانية من النماذج في الاختبارات ذات المفردات المختلطة إلى تحسين طفيف في معلومات الاختبار المتاحة مقارنة بالمجموعة الأولى من النماذج عبر نطاق القدرات لمعظم أطوال الاختبار ومستويات النسبة للتصحيح متعدد الاستجابة. يبدو أيضًا أن المجموعة الثانية استفادت أكثر من الزيادة في طول الاختبار فيما يتعلق بالمعلومات التي تم جمعها. الاختلافات في الثبات وإحصائيات الملاءمة بين المجموعتين لا تذكر. وتتفق هذه النتيجة مع ما توصل إليها بيرجر (Berger (1998) تميل معدلات الخطأ المعيارية إلى أن تكون أعلى بالنسبة للمجموعة الأولى في جميع أطوال الاختبار التي تم تحليلها في هذه الدراسة.

بحوث مقترحة

لا توجد طريقة تحليل فائقة معروفة لتحليل الاختبار ذي المفردات المختلطة. بحثت هذه الدراسة في نماذج نظرية الاستجابة للمفردة IRT المستخدمة في الدراسة وكيفية مقارنتها عبر ثلاثة أطوال اختبار ومجموعة من أشكال المفردات المختلطة. كانت تحليلات هذه الدراسة كافية في إنتاج النتائج التي تناولت أسئلة البحث المطروحة هنا، ومع ذلك، هناك حاجة إلى مزيد من التحليل في المجالات التالية:

- النظر في صعوبة الاختبار عبر مجموعة متنوعة من خيارات التصحيح المختلطة،

- النظر في درجات اختبار الممتحن عبر مجموعة متنوعة من خيارات التصحيح المختلفة،
- تطوير طرق أكثر ملاءمة لتحليل ثبات ومطابقة النماذج الإحصائية للاختبارات ذات المفردات المختلفة،
- اختبار قوة النماذج عند انتهاك الافتراضات (السرعة والتخمين) في الاختبار ذات المفردات المختلفة،
- متابعة استخدام التصحيح الآلي للمفردات متعددة الاستجابة.
- البحث في المزيج المثالي بين المفردات ثنائية الاستجابة ومتعددة الاستجابة التي تعطي درجات اختبار تتسم بدرجة عالية من الثبات والصدق.
- بوجه عام، يجب إجراء مزيد من البحث في مجالات صعوبة الاختبار، ودرجات اختبار الممتحنين، والتصحيح الآلي للتقدير الجزئي جنبًا إلى جنب مع مقارنة مع المقاييس النفسية التقليدية الأخرى وكيفية معالجتها للتحديات المتعلقة بالاختبارات ذات المفردات المختلفة.

المراجع:

صلاح الدين محمود علام (٢٠٠٥). نماذج الاستجابة للمفردات الاختبارية احادية البعد ومتعددة الابعاد وتطبيقاتها في القياس النفسي والتربوي، الطبعة الاولى، القاهرة: دار الفكر العربي.

صلاح الدين محمود علام (٢٠١٨). الاختبارات والمقاييس التربوية والنفسية. الأردن، دار الفكر.

محمد طالب دبوس (٢٠١٦). استخدام نظرية الاستجابة للفقرة في بناء فقرات اختبار محكي المرجع في الرياضيات بفقرات ثنائية الاستجابة ومتعددة الاستجابة وفق النموذج اللوجستي ثنائي المعلم. مجلة جامعة النجاح للأبحاث - العلوم الإنسانية، مج ٣٠، ٧٤، ص ص ١٤٥٣ - ١٤٨٠

محمد محمد فتح الله (٢٠١١). تكامل مدخلي القياس محكي المرجع CRM ونظرية الاستجابة للمفردة IRT في تقييم فاعلية برنامج تدريبي لتنمية كفايات بناء

الاختبارات التحصيلية لدى معلمي التعليم قبل الجامعي. *مجلة التربية، جامعة الأزهر - كلية التربية، ع ١٤٦، ج ٦، ٤٥٩ - ٥٢٥*

Adegoke, B. A. (2013). Comparison of Item Statistics of Physics Achievement Test using Classical Test and Item Response Theory Frameworks. *Journal of Education and Practice*, 4(22), 87-96.

Asriadi, M., & Hadi, S. (2021). Implementation of item response theory at final exam test in physics learning: Rasch model study. *Proceedings of the 6th International Seminar on Science Education (ISSE 2020)*, 541(Isse 2020), 336–342. <https://doi.org/10.2991/assehr.k.210326.048>

Ayanwale, M.A. & Adeleke, J.O. (2020). Efficacy of Item Response Theory in the Validation and Score Ranking of Dichotomous Response Mathematics Achievement Test. *Bulg. J. Sci. Educ. Policy*, 14(2):260–285. Accessed: May 31, 2021. Available: <https://www.academia.edu/45182779/>

Ayanwale, M. A. (2021). Calibration of Polytomous Response Mathematics Achievement Test Using Generalized Partial Credit Model of Item Response Theory. *EDUCATUM Journal of Science, Mathematics and Technology*, 8(1), 57-69. <https://doi.org/10.37134/ejsmt.vol8.1.7.2021>

Berger, M. P. (1998). Optimal design of tests with dichotomous and polytomous items. *Applied Psychological Measurement*, 22(3) 248-258.

Bond, T. G., & Fox, C. M. (2015). *Applying the Rasch model: Fundamental measurement in the human sciences* (3rd). Routledge.

Bond, T.G., Yan, Z., & Heene, M. (2020). *Applying the Rasch Model: Fundamental Measurement in the Human Sciences* (4th Ed.). New York: Routledge.

Bozdağ, H. C., & Türkoğuz, S. (2021). A rasch model analysis of primary school students' conceptual understanding levels of the concept of light.

IOJPE: International Online Journal of Primary Education, 10(1), 160–179.

Calderón, C.; Beyle, C.; Véliz-García, O.; Bekios-Calfa, J. (2021). Psychometric properties of Addenbrooke's Cognitive Examination III (ACE-III): An item response theory approach. *PLoS ONE*, 16, e0251137.

Carlson, J. E. (1996, April). *Information provided by polytomous and dichotomous items on certain NAEP instruments*. Paper presented at the annual meeting of the American Educational Research Association, New York, NY.

Che Lah, N. H., Tasir, Z., & Jumaat, N. F. (2021). Applying alternative method to evaluate online problem-solving skill inventory (OPSI) using Rasch model analysis. *Educational Studies*, 00(00), 1–23. <https://doi.org/10.1080/03055698.2021.1874310>

Cohen, R. (2018) *Psychological Testing and Assessment: An Introduction to Tests and Measurement*. 9th edition, McGraw-Hill Education, 2 Penn Plaza, New York, NY 10121

Donoghue, J. R. (1994). An empirical examination of the IRT information of polytomously scored reading items under the generalized partial credit model. *Journal of Educational Measurement*, 31(4) 295-311.

Downing, S. M. (2002). Construct-irrelevant variance and flawed test questions: Do multiple-choice item-writing principles make any difference? *Academic Medicine*, 77(10), s103-104.

Ercikan, K., Schwarz, R. D., Julian, M. W., Burket, G. R., Weber, M. M. & Link, V. (1998). Calibration and scoring of tests with multiple-choice and constructed-response item types. *Journal of Educational Measurement*, 35(2), 137-154.

Erguven, M. (2014). Two Approaches to Psychometric Process: Classical Test Theory and Item Response Theory. *Journal of Education*, 2(2), 23-30.

- Hambleton, R. K. (1989). Principles and selected applications of item response theory. R.L. Linn (Ed.), *Educational Measurement: Third Edition*, (pp. 147-200). American Council on Education and Macmillan Publishing Company.
- Hambleton, R.K.; Swaminathan, H.; Rogers, H.J. (1991). **Fundamentals of Item Response Theory**; Sage: Thousand Oaks, CA, USA.
- Kaplan, R. (2012). **Psychological Testing: Principles, Applications, and Issues** 8th Edition. Wadsworth, Cengage Learning.
- Kline, T. J. B. (2005). *Psychological testing: A practical approach to design and evaluation*. Sage Publications, Inc.
- Linn, R.L. (1990). Has item response theory increased the validity of achievement test scores? *Applied Measurement in Education*, 3(2), 115-141.
- Marnat, G and Wright, M (2016). **Handbook of psychological assessment** Sixth Edition. John Wiley & Sons, Inc.
- Masers, G.N. (1982). A rasch model for partial credit scoring. *Psychometrika* 47(2), 149-174.
- Murphy, P. (2021). **Score Interpretation Guide for Educators**. State of New Jersey, Department of Education
- Nunnally, J., & Bernstein, I.H. (1994). **Psychometric Theory (3rd Ed.)**. New York: McGraw-Hill.
- Oliveri, M. E., and M. von Davier (2014). Toward increasing fairness in score scale calibrations employed in international large-scale assessments, *International Journal of Testing*, 14(1): 1-21.
- Plake, B. (1993). Applications of educational measurement: Is optimum optimal? *Educational Measurement: Issues and Practices*, 5-10.
- Rahayu, W., Putra, M. D. K., Rahmawati, Y., Hayat, B., & Koul, R. B. (2021). Validating an indonesian version of the what is happening in this class? (WIHIC) questionnaire using a multidimensional rasch

model. *International Journal of Instruction*, 14(2), 919–934.
<https://doi.org/10.29333/iji.2021.14252a>

Rezapour, M.; Cuccolo, K.; Veenstra, C.; Ferraro, F.R. (2021). An Item Response Theory to Analyze the Psychological Impacts of Rail-Transport Delay. *Sustainability*, 13, 6935. <https://doi.org/10.3390/su13126935>

Si, C. B. (2002). **Ability estimation under different item parameterization and scoring models**. Unpublished doctoral dissertation, University of North Texas.

Stemler, Steven E. and Naples, Adam (2021). Rasch Measurement v. Item Response Theory: Knowing When to Cross the Line," *Practical Assessment, Research, and Evaluation*: Vol. 26, Article 11. Available at: <https://scholarworks.umass.edu/pare/vol26/iss1/11>

Stewart, J., Zabriskie, C., Devore, S., & Stewart, G. (2018). Multidimensional item response theory and the force concept inventory. *Physical Review Physics Education Research*, 14(1), 10137. <https://doi.org/10.1103/PhysRevPhysEducRes.14.010137>

Thissen, D. (1991). **MULTILOG User's Guide: Multiple, Categorical item Analysis and Test Scoring using item Response Theory; version 6.0**. Scientific Software International Chicago: Illinois.

Thissen, D., Nelson, L. and Swygert, K. (2001). Item response theory applied to combinations of multiple-choice and constructed-response items- Approximation method for scaled scores. In Thissen, D. & Wainer, H. (Eds.), *Test Scoring*, (pp. 293-341). Mahwah, New Jersey: Lawrence Erlbaum Associates, Inc.

Triwik, J. (2021) Rasch analysis on item response theory: Review of model suitability. *AIP Conference Proceedings* **2326**, 020017 (2021); <https://doi.org/10.1063/5.0040305> Published Online: 08 February 2021

- Tseng, M. C., & Wang, W. C. (2021). The q-matrix anchored mixture rasch model. *Frontiers in Psychology*, 12(March), 1–9. <https://doi.org/10.3389/fpsyg.2021.564976>
- van der Linden, W.(2016). **Handbook of Item Response Theory** Volume 2: Statistical Tools. Chapman and Hall/CRC.
- Van der Linden, W. J. (Ed.). (2018). **Handbook of Item Response Theory, Three Volume Set**. CRC Press.
- von Davier, M., Yamamoto, K., Shin, H. J., Chen, H., Khorramdel, L., Weeks, J., Davis, S., Kong, N., & Kandathil, M., (2019). *Evaluating item response theory linking and model fit for data from PISA 2000–2012. Assessment in Education: Principles, Policy, and Practice*, 26(4), 466-488. doi:10.1080/0969594X.2019.1586642
- Wainer, H. & Thissen, D. (1993). Combining multiple-choice and constructed- response test scores: toward a marxist theory of test construction. *Applied Measurement in Education*, 6(2), 103-118.
- Wang W, Drasgow F and Liu L (2016) Classification Accuracy of Mixed Format Tests: A Bi-Factor Item Response Theory Approach. *Front. Psychol.* 7:270. doi: 10.3389/fpsyg.2016.00270
- Wright, B. D., & Stone, M. H. (1979). **Best test design: Rasch measurement**. Chicago, IL: Mesa Press.

Abstract

This study investigated the effects of test length, scoring schema and item analysis model on the ability estimates, test information levels and optimization criteria of mixed item format tests. SPSS software across ability estimates and standard errors of ability estimates using a 3 x 11 x 2 fixed factorial ANOVA was used to make the comparison between two sets of models (first group consists of Rasch model and Partial Credit Model, second group consists of Two Parameter Logistic Model and Graded Response Model of item analysis procedures. Effect sizes and power were reported for each procedure. Scheffe post hoc procedures were conducted on significant factors. Test information was analyzed and compared across the range of ability levels for all 66-design combinations. The results indicated that both test length and the proportion of items scored polytomously had a significant impact on the amount of test information produced by mixed item format tests. Generally, tests with 100% of the items scored polytomously produced the highest overall information. This seemed to be especially true for examinees with lower ability estimates. Optimality comparisons were made between IRT models procedures based on standard error rates for the ability estimates, marginal reliabilities and fit indices (-2LL). The only significant differences reported involved the standard error rates for both two sets procedures. This result must be viewed in light of the fact that the effect size reported was negligible.

Keywords: test length, scoring schema, item analysis model, the ability estimates, test information function, optimization criteria, mixed item format tests.

