



# THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Dependence of gene-by-environment interactions (GxE) on scaling

**Citation for published version:**

Murray, A, Molenaar, D, Johnson, W & Krueger, RF 2016, 'Dependence of gene-by-environment interactions (GxE) on scaling: Comparing the use of sum scores, transformed sum scores and IRT scores for the phenotype in tests of GxE interaction' Behavior Genetics, vol. 46, no. 4, pp. 552-572. DOI: 10.1007/s10519-016-9783-5

**Digital Object Identifier (DOI):**

[10.1007/s10519-016-9783-5](https://doi.org/10.1007/s10519-016-9783-5)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Behavior Genetics

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Behavior Genetics

## Dependence of gene-by-environment interactions (GxE) on scaling: Comparing the use of sum scores, transformed sum scores and IRT scores for the phenotype in tests of GxE

--Manuscript Draft--

<b>Manuscript Number:</b>	BEGE-D-15-00056R2
<b>Full Title:</b>	Dependence of gene-by-environment interactions (GxE) on scaling: Comparing the use of sum scores, transformed sum scores and IRT scores for the phenotype in tests of GxE
<b>Article Type:</b>	Original Research
<b>Keywords:</b>	gene-environment interaction; item response theory; transformation; scaling; skewness
<b>Corresponding Author:</b>	Aja Louise Murray Cambridge, UNITED KINGDOM
<b>Corresponding Author Secondary Information:</b>	
<b>Corresponding Author's Institution:</b>	
<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Aja Louise Murray
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Aja Louise Murray Dylan Molenaar Wendy Johnson Robert F Krueger
<b>Order of Authors Secondary Information:</b>	
<b>Funding Information:</b>	
<b>Abstract:</b>	<p>Estimates of gene-environment interactions (GxE) in behavior genetic models depend on how a phenotype is scaled. Inappropriately scaled phenotypes result in biased estimates of GxE and can sometimes even suggest GxE in the direction opposite to its true direction. Previously proposed solutions are mathematically complex, computationally demanding and may prove impractical for the substantive researcher. We, therefore, evaluated two simple-to-use alternatives: 1) straightforward non-linear transformation of sum scores and 2) factor scores from an appropriate item response theory (IRT) model. Within Purcell's (2002) GxM framework, both alternatives provided less biased parameter estimates, and improved false and true positive rates than using a raw sum score. These approaches are, therefore, recommended over using raw sum scores in tests of GxE. Circumstances under which IRT factor scores versus transformed sum scores should be preferred are discussed.</p>
<b>Response to Reviewers:</b>	<p>p.4,l.27-28: Non-normality could be due to many other reasons besides GxE or not capturing the full range of trait variation.</p> <p>Response: We have added a sentence noting this and provide an example of another reason for non-normality, namely, inadequate sampling from the population at the lowest or highest end of the relevant trait distribution (p.4).</p> <p>p.13: There is an alpha parameter in eq (3) and also in eq (2). Please correct.</p> <p>Response: We have replaced the alpha parameter in eq.3 and any references to it with 'a' to make a distinction between these two parameters.</p>

[Click here to view linked References](#)

Phenotype scaling in GxE

**Dependence of gene-by-environment interactions (GxE) on scaling: Comparing the use of sum scores, transformed sum scores and IRT scores for the phenotype in tests of GxE**

Aja Louise Murray<sup>1,2,3\*</sup>, Dylan Molenaar<sup>4</sup>, Wendy Johnson<sup>1,2</sup>, Robert F. Krueger<sup>5</sup>

<sup>1</sup> Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, UK

<sup>2</sup>Department of Psychology, University of Edinburgh, UK

<sup>3</sup>Violence Research Centre, University of Cambridge, UK

<sup>4</sup>Psychological Methods, University of Amsterdam, The Netherlands

<sup>5</sup>Department of Psychology, University of Minnesota- Twin Cities, US

\* Corresponding author: Institute of Criminology, Sidgwick Avenue, Cambridge, CB3 9DA.

Email: am2367@cam.ac.uk

**Abstract**

1  
2  
3 Estimates of gene-environment interactions (GxE) in behavior genetic models depend on how a  
4 phenotype is scaled. Inappropriately scaled phenotypes result in biased estimates of GxE and can  
5 sometimes even suggest GxE in the direction opposite to its true direction. Previously proposed  
6 solutions are mathematically complex, computationally demanding and may prove impractical for the  
7 substantive researcher. We, therefore, evaluated two simple-to-use alternatives: 1) straightforward  
8 non-linear transformation of sum scores and 2) factor scores from an appropriate item response theory  
9 (IRT) model. Within Purcell's (2002) GxM framework, both alternatives provided less biased  
10 parameter estimates, and improved false and true positive rates than using a raw sum score. These  
11 approaches are, therefore, recommended over using raw sum scores in tests of GxE. Circumstances  
12 under which IRT factor scores versus transformed sum scores should be preferred are discussed.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

26 **Keywords:** Gene-environment interaction; item response theory; transformation; scaling; skewness  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Dependence of gene-by-environment interactions (GxE) on scaling: Comparing the use of sum scores, transformed sum scores and IRT scores for the phenotype in tests of GxE**

Increasingly, theoretical perspectives on phenotypic development and expression are recognising that genes and environments transact in dynamic ways. Many posit some kind of gene-environment interaction (GxE) where GxE is defined as a differential response to environmental circumstances depending on genotype, or, a differential genetic expression depending on environment (Boomsa, & Martin, 2002; Eaves, Last, Marin & Jinks, 1977). GxE plays a central role in major theoretical models such as the diathesis-stress model, the differential susceptibility model, the vantage sensitivity model, and the bioecological model (Brofenbrenner & Ceci, 1994; Pluess & Belsky, 2013; Reiss, Leve & Neiderhiser, 2013; Rende & Plomin, 1992). The diathesis-stress model, for example, predicts that the genetic variance in a psychopathological trait is greater in more adverse environments whereas the bioecological model predicts that the genetic potential for a positive trait is realised to a greater extent in more stimulating, higher-quality environments (Asbury, Wachs & Plomin, 2005; Rende & Plomin, 1992). GxEs are also cited as mechanisms by which social factors regulate behavior, for example, in the idea that genetic influences on certain phenotypes are prevented from being expressed when there are stronger social norms or explicit prohibitions relating to those phenotypes (Shanahan & Hofer, 2005).

To keep pace with these theoretical developments, it has been necessary to develop statistical methodologies capable of modelling the more complex forms of interplay that they imply (e.g. Purcell, 2002). Despite the promise and widespread uptake of these methodologies, the ability to test theoretically implied GxE interactions is affected in practice by dependency of tests of interactions on the observed distributions or scales of the phenotypes (Eaves et al., 1977, 2002; Eaves, 2006; Mather & Jinks, 1971; Purcell, 2002; Schwabe & van den Berg, 2014).

The problem of dependency of GxE on phenotype scaling has been known since the time of R.A. Fisher, who noted that GxE interactions could be manipulated by re-scaling the variables involved. In fact, he went far as to advocate ‘transformations of scale’ to eliminate what he perceived

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

to be nuisance non-additivity (Tabery, 2008). This suggestion was controversial because he was recommending purging the same non-additivity that was and still is viewed by many substantive researchers as a meaningful clue as to the causal processes underlying phenotypic development. Since then, numerous methodological studies have further discussed and provided demonstrations of dependency of appearance of presence of GxE on scaling (Eaves et al., 1977; Martin, 2000; Molenaar, van der Sluis, Boomsma & Dolan, 2012; Purcell, 2002; Tucker-Drob, Harden & Turkheimer, 2009; van der Sluis, Dolan, Neale, Boomsma & Posthuma, 2006). In the section that follows we summarise and extend the key arguments of these authors.

The primary challenge in dependency of GxE on scaling concerns the multiplicity of possible causal structures that could underlie the same sample phenotypic distribution. Consider the case where the observed distribution of the phenotype is non-normal: a common occurrence in behavior genetic research, as well as psychological research in general (Beasley, Erickson & Allison, 2009; Miccerri, 1989). The problem is that when an observed phenotypic distribution is non-normal, this non-normality could reflect the presence of GxE, or it could simply be that the measurement instrument used has been unable to capture the full range of variation in the trait, leading to a skewed score distribution. A statistical test of GxE will not be able to distinguish among these possibilities easily.

The challenge of choosing between a ‘scaling’ and ‘GxE’ explanation for an apparent moderation effect is just one example of the broader challenge of selecting the correct model when a range of causal generating mechanisms could produce similar patterns in the observed data. **For example, non-normality could arise for a number of methodological reasons aside from improper scaling e.g., failing to adequately sample individual with the lowest or highest trait levels from the population.** In terms of theoretically important processes, GxE is also difficult to distinguish statistically from non-linear main effects of a moderator on a phenotype or from non-linear genetic or environmental influences on a phenotype (e.g. Rathouz et al., 2008; Zheng & Rathouz, 2015). However, there are good reasons to begin by attempting to rule out scaling as the alternative explanation for GxE effects. First, if improper scaling can account for apparent moderation effects,

1 there is no need to posit complex interactions between the etiological influences on a phenotype,  
2 whether this is GxE or some other form of interplay. At a scientific level, incorrectly accepting a  
3 ‘complex interplay’ explanation can lead to theories which lack parsimony and which when further  
4 pursued may lead to wasted research efforts. At a more practical level, falsely selecting a ‘GxE’  
5 explanation may foster the mis-impression that a candidate moderator is an important factor with  
6 respect to understanding variation in some phenotype, able to constrain or promote the expression of  
7 genetic liability, when in fact it is merely correlated with that phenotype.  
8  
9

10  
11  
12  
13  
14  
15  
16         Second, there is evidence that many phenotypic measures suffer from sub-optimal scaling.  
17  
18 Cases in point are measures of psychopathological constructs. These very commonly yield observed  
19 non-normal (positively skewed) distributions because majorities of participants score close to the low  
20 (non-pathological) ends of the measurement scales. It is often argued that these observed distributions  
21 are not necessarily appropriate representations of the population distributions of the phenotypes but  
22 arise as a result of the scales being developed with focus on the upper extremes of the traits (van den  
23 Oord, Pickles & Waldman, 2003; van den Oord, Simonoff, Eaves, Pickles, Silberg & Maes, 2000).  
24  
25 This argument is based on various pieces of evidence, including the apparent highly polygenic nature  
26 of many common psychopathological disorders (e.g. Wray et al. 2014); on the observed normal  
27 distributions obtained when special care is taken to measure ‘non-clinical’ levels of  
28 psychopathological traits (e.g. Baron-Cohen et al., 2001); and on statistical comparisons of models  
29 positing categorical versus dimensional models of psychopathological traits (e.g. Walton, Ormel &  
30 Krueger, 2011). None of these is definitive evidence that psychopathological traits are normally  
31 distributed in the population but together they suggest that this may be closer to the truth than the  
32 classical categorical models in which meaningful variation in psychopathological traits is restricted to  
33 a narrow, clinical range of trait values. Under this dimensional view, failure to observe a normal  
34 distribution for a trait may be a result of failing to measure that trait with items that have an  
35 appropriate range of difficulties to provide reliable coverage of the whole trait distribution.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56

57         Within an item response theory (IRT) framework, such a failure will be manifested as item  
58 difficulties that are tightly clustered at the high end of the range; a phenomenon observed in many  
59  
60  
61

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

psychometric studies of commonly used inventories of psychopathologies (Meijer & Egberink, 2012; Reise & Waller, 2009; Thomas, 2011). These scales have high discrimination in and around clinical cut-off points but poor discrimination in the healthy ranges. Thus, in a population-representative sample that would include predominantly subjects considered healthy, most participants completing such a test will endorse the lowest response options for most items, leading to a positively skewed score distribution and an apparent lack of individual differences at low levels of the phenotype due to the absence of items tapping this level.

If raw scores such as sums of items from scales affected in this way are used to represent phenotypes, they are likely to provide biased tests of GxE (Molenaar & Dolan, 2014; Schwabe & van den Berg, 2014). This is because GxE estimates depend on the degrees of individual differences in a phenotype at different levels of the moderator. Use of a scale that fails to these adequately at lower levels of the phenotype will tend to falsely indicate less variation at lower levels, when in fact this apparent observation is a function of weaker measurement at lower levels. The direction of the resulting bias in GxE depends on both skewness of the score and extent of correlation with the moderator. Positive skew combined with a positive moderator-phenotype correlation is liable to produce a positive interaction parameter, while negative skew combined with a positive moderator-phenotype correlation is liable to produce a negative interaction parameter. Thus moderation effects can arise even when there are no causal processes corresponding to our conceptual models of GxE influencing phenotypic development.

In empirical studies a researcher is faced with the challenge of choosing the most appropriate scale for the measure used to capture the relevant phenotype. To the extent any phenotype actually has a latent dimensional distribution, it can be thought of as having some correspondingly dimensional scale of measurement, but for psychological constructs, we have little or knowledge of what these scales might be. Still, there are more or less appropriate choices given what is known about the underlying etiology of a trait, its distribution in the population, and the research question of interest (e.g. see Falconer & Mackay, 1996). The appropriate scale for a phenotypic measure cannot be selected based on its observed score distributions or other features of the data: it must be selected



## Phenotype scaling in GxE

1 based on conceptual knowledge and assumptions regarding the phenotype underlying the measures.  
2 Deviations of phenotypic distributions from expectations derived from these assumptions should be  
3  
4 cause for concern.  
5  
6

7           Compounding this challenge is the fact that most behavior genetic modelling approaches  
8  
9 require assumptions of multivariate normality<sup>1</sup> and that violations of those assumptions can lead to  
10  
11 incorrect inferences regarding the presence of GxE (van Hulle & Rathouz, 2015). With this in mind,  
12  
13 researchers have tended to deal with non-normal score distributions by employing straightforward  
14  
15 non-linear transformations intended to remove the non-normality. For positively skewed sum scores,  
16  
17 the log-transformation is popular (e.g Hicks, South, DiRago, Iacono, McGue, 2009; Johnson et al.  
18  
19 2010) but the square root transformation is also sometimes used (e.g. Distel, Middeldorp, Trull,  
20  
21 Derom, Willemsen & Boomsma, 2011). Given that the same approach is recommended to remove  
22  
23 GxE interactions that are artifacts of phenotypic scaling (e.g. see Falconer & MacKay, 1996 ch.17),  
24  
25 one might conclude that this also represents a solution to the problem of dependency of GxE on scale.  
26  
27 There are, however, at least two major reasons to question this. First, while there has been no  
28  
29 systematic simulation study evaluating their effectiveness in mitigating bias due to sub-optimal  
30  
31 scaling, Kang & Waller (2005) demonstrated that sum score transformations were only moderately  
32  
33 successful in reducing the tendency towards spurious phenotypic interactions in the context of  
34  
35 moderated multiple regression. Second, and more importantly: presence of GxE introduces non-  
36  
37 normality into the phenotypic distribution because it is by definition a relative expansion or  
38  
39 contraction of variance in the phenotype across levels of the moderator. This suggests that  
40  
41 transforming a non-normal score to normality could ‘transform away’ the very interaction effect of  
42  
43 potential interest.  
44  
45  
46  
47  
48

49           As another possible solution, some authors have suggested separating out scaling and GxE  
50  
51 sources of non-normality by modelling GxE using an explicit measurement model (the scaling part) in  
52  
53 combination with a biometric model (the GxE part). Essentially, the proposal is to model the scaling  
54  
55 properties of items to account for differences in informativeness of phenotypic estimates across levels  
56  
57 of the moderator. For example, if a scale has items that have difficulties that are clustered towards one  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

end of the scale, a psychometric model with potential to recognize this can be integrated into a broader biometric model so that these parameters can be freely estimated and reflected in the estimates of the biometric parameters. The particular choice of measurement model will vary from phenotype to phenotype and be dictated by expectations about the latent trait distribution and the item response format.

For continuous indicators, Molenaar et al. (2012) demonstrated the feasibility of this approach in a GxE model in which GxE was operationalised as heteroscedastic E or C variance across levels of A. They showed that when differences in item residual variances across phenotypic level were incorporated into a measurement model and combined with a test of GxE, biasing effects of poor scaling were substantially mitigated. Similarly, Tucker-Drob et al. (2009) suggested a procedure in which a factor model with quadratic factor loadings was estimated in one stage and then, in a second stage, the same measurement model (with parameters fixed to the values estimated from the first stage) was combined with Purcell's GxE model. Quadratic factor loadings allow for the relation between the items and latent phenotype to vary across levels of the phenotype: an effect that could otherwise be mis-attributed to GxE. However, truly continuous indicators are rare; therefore, Molenaar & Dolan (2014) and Schwabe & van den Berg (2014) proposed models for (ordered) categorical data that could be combined with a test of GxE. Again, using these models there was evidence of substantial reduction of bias in tests of GxE compared to using biometric models that did not explicitly model the scaling properties of the items used to measure the phenotype.

In spite of the potential utility of incorporating explicit measurement models for the phenotype into tests of GxE when an assumption about the underlying distribution of the genetic and environmental influences on the phenotype can be made, there have been very few studies taking this approach. One reason may be that the approach is mathematically complex and thus somewhat inaccessible for non-methodologists. There may also be a misconception that, because scores from these models will be highly correlated with sum scores, there would be essentially no benefit from using such models. It is not valid, however, to conclude that highly correlated measures will have the same properties in moderated models such as those that test for GxE. This is because correlations are

## Phenotype scaling in GxE

1 sensitive mainly to rank orders, which can be highly preserved even when distributional properties  
2 differ markedly. Distributional properties are particularly important in any situation involving any  
3 kind of nonlinearity such as that involved in interactions.  
4  
5

6  
7 Misconceptions aside, there are practical limitations to the various approaches discussed  
8 above, and it is not clear what the best approach might be. For example, the Schwabe & van den Berg  
9 (2014) approach requires assumption that IRT parameters are known, the Molenaar & Dolan (2014)  
10 approach is computationally intensive, and the approaches of Molenaar et al. (2012) and Tucker-Drob  
11 et al. (2009) require continuous indicators. Further, all were applied within the context of specific  
12 GxE models, potentially limiting their general applicability in practice.  
13  
14  
15  
16  
17  
18  
19  
20

21  
22 Given these potential practical limitations, another possibility is to use a two-step approach to  
23 estimating GxE. In this approach, an appropriate measurement model for the phenotype is estimated,  
24 factor scores are obtained from this model, and then in a separate stage, these factor scores are  
25 submitted to a biometric model to test GxE. The ‘two steps’ refer to the use of two separate models,  
26 and the approximation involved in using explicitly calculated factor scores to measure a variable  
27 conceptualized as latent. This is in contrast to the one-step approach described above in which the  
28 biometric and psychometric model are estimated together, in a single step.  
29  
30  
31  
32  
33  
34  
35  
36  
37

38 Although there has been no systematic study of this approach in GxE models, simulation  
39 studies have shown that a two-step approach works well in reducing bias due to scaling in phenotypic-  
40 level interactions in moderated multiple regression and factorial ANOVA (Embreston, 1996; Kang &  
41 Waller, 2005; Morse et al. 2012). For example, Kang & Waller (2005) showed that the tendency for  
42 spurious interactions to result from poor item scaling was substantially mitigated when IRT scores  
43 from a 2-parameter logistic model were utilised in place of sum scores. This strategy also proved  
44 more effective than a simple non-linear transformation of the score. Therefore, it is possible that a  
45 two-step approach could provide a compromise between the greater conceptual and computational  
46 simplicity of using a sum score and the effectiveness of IRT-based latent trait estimates in accounting  
47 for the scaling properties of items.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Based on the preceding argument, we compared a two-step approach to the currently most  
2 commonly used methods for handling observed non-normal phenotypes, that is, the raw sum scores  
3 and the transformed sum scores. We compared these three approaches using a statistical simulation  
4 study complemented by a real data example.  
5  
6  
7

### 8 9 **Modelling approach**

10 We based our analyses on the Purcellian GxM interaction (where the ‘M’ stands for measured  
11 environment) framework initially introduced by Purcell (2002) and subsequently extended and  
12 evaluated by others (Rathouz, van Hulle, Rodgers, Waldman & Lahey, 2008, van Hulle, Lahey &  
13 Rathouz, 2013; van Hulle & Rathouz, 2015; Zheng & Rathouz, 2015; Zheng, Van Hulle & Rathouz,  
14 2015). This framework is arguably the foremost in assessing theoretical hypotheses which predict  
15 moderation of genetic influences on a specific phenotype by a specific moderator because in addition  
16 to accommodating both gene-environment interaction and gene-environment correlation, it can also  
17 be used to evaluate a range of other forms of phenotype-moderator transactions (see Zheng &  
18 Rathouz, 2013). Uptake of the GxM modelling approach has been extensive; it has been employed to  
19 assess substantive hypotheses relating to a diversity of phenotypes including cognitive ability  
20 (Harden, Turkheimer & Loehlin, 2007), physical health (Johnson & Krueger, 2005), health behaviors  
21 (Timberlake et al., 2006), social relationships (South, Krueger, Johnson & Iacono, 2008), and  
22 psychopathological traits (South & Kruger, 2011). The popularity and influence of the approach is  
23 indicated by the fact that, at time of writing, the Purcell (2002) article has been cited almost 500  
24 times.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

47 We focussed on a form of the model that can be used to assess gene-by-measured  
48 environment interaction. The moderator (M) is modelled as:

$$49 M = a_M A_M + c_M C_M + e_M E_M$$

50  
51  
52  
53  
54  
55  
56 (1)

57 and the phenotype (P) as:

$$\begin{aligned}
P = & (a_C + \alpha_C M)A_M + (c_C + \gamma_C M)C_M + (e_C + \varepsilon_C M)E_M \\
& + (a_U + \alpha_U M)A_U + (c_U + \gamma_U M)C_U + (e_U + \varepsilon_U M)E_U,
\end{aligned}
\tag{2}$$

where  $A, C$  and  $E$  refer to mutually uncorrelated multivariate normally distributed (each with mean=0, variance=1) latent additive genetic, shared environmental and unshared environmental influences respectively,  $\alpha, \gamma$  and  $\varepsilon$  are moderation parameters that capture the moderation of  $A, C$  and  $E$  influences by  $M$ , with the subscripts  $c$  and  $u$  denoting ‘common’ (to P and M) and ‘unique’ (to P).

The parameter of interest is  $\alpha_U$  which captures moderation of the genetic influences on the phenotype that are not shared with the moderator. When this parameter is positive, genetic influences unique to the phenotype increase with the moderator and when it is negative, they decrease with the moderator.

### Simulation study

We evaluated the effect of poor scaling on estimates of  $\alpha_U$  using Eqs. 1 and 2 as our population biometric model, simulating poor scaling of the phenotype (explained below), and then estimating the model in Eqs. 1 and 2 using this poorly scaled phenotype. For our population biometric model, we used the following parameter magnitudes: For the moderator and phenotypic means we set  $\mu_M = \mu_P = 0$ ; for the latent genetic and environmental influences on the moderator and phenotype we set  $a_U = \sqrt{0.2}$ ,  $a_C = \sqrt{0.3}$ ,  $a_M = \sqrt{0.3}$ ;  $c_U = \sqrt{0.1}$ ,  $c_C = \sqrt{0.1}$ ,  $c_M = \sqrt{0.2}$ ;  $e_U = \sqrt{0.2}$ ,  $e_C = \sqrt{0.1}$ ,  $e_M = \sqrt{0.5}$ ; and for the moderation parameters we set  $\alpha_C = \gamma_C = \varepsilon_C = 0$  and varied the magnitude of  $\alpha_U, \gamma_U$  and  $\varepsilon_U$  across conditions. To explore how bias in  $\alpha_U$  was affected by direction of the skew of the observed score distribution and direction of the population interaction, we varied  $\alpha_U =$  to be -.15, 0, and .15 across conditions. In addition, as resolvability of the  $\alpha_U, \gamma_U$ , and  $\varepsilon_U$  parameters is often imperfect, we explored how the bias in  $\alpha_U$  is affected by whether  $\gamma_U$  and  $\varepsilon_U$  represented interactions in the same versus the opposite direction to that of  $\alpha_U$ . We did this by including a subset

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

of conditions in which  $\gamma_U$  and  $\varepsilon_U$  were specified to have the same sign as  $\alpha_U$  and a subset of conditions in which they were specified to have the opposite sign to  $\alpha_U$ . In both cases the absolute magnitudes of  $\gamma_U$  and  $\varepsilon_U$  were specified to be .20 and .08 respectively while  $\alpha_U$  was held constant at -.15. We chose these sets of conditions and corresponding parameter values with the goal of selecting realistic values based on our own experiences of working with empirical twin data and on other published studies. Because we could expect results to be broadly symmetrical for positive and negative skews and negative and positive interaction parameters, we did not implement a fully crossed simulation design, but focussed on models that were realistic and which covered key combinations of variables.

Together, this combination of population parameters resulted in a total of four population models, summarised in Tables 2 and 3. In each replication, we generated data for either 500 MZ and 500 DZ or 1000 MZ and 1000 DZ twins according to these models. To keep the model focussed on the question at hand, we did not consider sex differences.

### Observed item-level data generation

We generated item level data for twin 1 and twin 2 phenotypes using two different models that reflect common scaling practices. We did not manipulate the scaling of the moderator because - as in moderated multiple regression - the scaling of the predictor is far less critical with respect to the accuracy of estimates of interactions (e.g. van Hulle & Rathouz, 2015). First, we used a graded response model (GRM; Samejima, 1969) as the basis for linking the latent trait values for the phenotype (P) to observed item responses to give a set of conditions in which the scaling issues could be considered mild. These latent trait values were determined according to the GxM population models described in the previous section. We simulated these data using the catIrt package in R statistical software (Nydeck, 2014; R Core Team, 2014). Here, the items are considered in dichotomous steps, each characterised by a 2-parameter logistic model but with discriminations

constrained equal within items. Specifically, probability of a respondent  $i$  with level of the latent trait  $\theta_i$  having a response  $x_{ij}$  that falls at or above a given category ( $k = 1 \dots m_j$ ) is specified as:

$$P^*_{ijk} = P(x_{ij} \geq k | \theta_i, a_j, \beta_{jk}) = \frac{1}{1 + \exp[-a_j(\theta_i - \beta_{jk})]} \quad (3)$$

where  $a_j$  is the discrimination parameter of item  $j$  and  $\beta_{jk}$  is the category difficulty parameter of category  $k$  in item  $j$ . Note that  $\theta_i$  is identical to  $P$  in eq. 2.

We generated data for 20 items with  $a_j$  and  $\beta_{jk}$  parameters provided in Table 1. This gave items with five ordinal levels. The  $\beta_{jk}$  parameters were chosen to yield positively skewed item and sum score distributions that mimicked those commonly found in empirical research (e.g. Kang & Waller, 2005). To do this, we selected  $\beta_{jk}$  for successive response categories so that a disproportionate number of responses would fall into the first and second response categories. We also specified the  $\beta_{jk}$  parameters for a given category to show variability across the 20 items within our simulated test which is more realistic than setting them all equal. Discrimination parameters,  $a_j$ , were selected by randomly sampling from a uniform distribution with min=0.5 and max=2.5.

Second, we generated item-level data designed to be less favorable with respect to its scaling properties. Specifically, we used the same discrimination values but instead of using five ordinal levels, we used a 2PL model with only 2 ordinal levels (i.e., binary items), again selecting difficulty parameters such that disproportionate numbers of responses fell into the response category indicating a lower trait level. This gave us a set of conditions in which the scaling issues could be considered more serious than the polytomous case. Here, the model linking latent trait values to observed item level responses was:

$$P^*_{ijk}(x_{ij} = 1 | a_j, \beta_j) = \frac{1}{1 + \exp[-a_j(\theta_i - \beta_j)]}$$

1  
2  
3 The  $a_j$  and  $\beta_{jk}$  parameters used are provided in Table 1.  
4  
5

### 6 **True score**

7  
8  
9 As a control condition, we generated scores for the phenotype according to Eqs. 1 and 2 for  
10 without introducing any scaling issues. These scores can therefore be considered ‘true’ phenotypic  
11 scores. We considered these true phenotypic scores in order to provide a baseline against which we  
12 could compare the results. This is necessary because even in the absence of any scaling problems, it is  
13 likely that the GxM model will not perfectly recover all moderation parameters and because  
14 moderation parameters may be difficult to resolve from one another. For example, moderation of  
15 shared environmental influence may be to some extent mis-attributed to moderation of genetic  
16 influences.  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

### 27 **Sum score**

28  
29  
30 We created a sum score for the phenotype summing the scores from the 20 item responses  
31 generated as described above by Eq.s 1, 2, and 3 for the GRM and by summing the 20 item responses  
32 generated as described by Eq.s 1,2 and 4 for the 2PL . Examples of the resulting sum score  
33 distributions are shown in Figures 1 (polytomous) and 2 (binary). These sum score distributions  
34 exhibited positive skew, similar to that observed in many measures of psychopathological traits. In the  
35 binary case, this would be correspond to the kind of summed ‘presence versus absence’ symptom  
36 scores found in diagnostic data. Skew also depended on the direction of interaction in the population  
37 model, with positive interactions making score distributions more positively skewed and negative  
38 interactions making score distributions more negatively skewed. However, these effects were  
39 relatively minor in comparison to the effect of scaling on the phenotypic distribution.  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

### 54 **Transformed sum score**

55  
56  
57 We created transformed sum scores by  $\log_{10}$  transformations of the sum scores generated as  
58 described in the previous section. The  $\log_{10}$  transformation, the natural log transformation, and other  
59  
60  
61  
62  
63  
64  
65



1 similar kinds of transformation of the phenotype are commonly used in GxE models when the  
2 phenotype has a positively skewed distribution (e.g. Button et al. 2010; Hicks, Dirago, Iacono, &  
3 McGue, 2009; Hicks et al., 2009; Johnson, Kyvik, Mortensen, Skytthe, Batty & Deary, 2010;  
4 Silvetoinen et al. 2009; Tuvblad, Grann, & Lichtenstein, 2006). Transforming the sum scores gave  
5  
6 rise to approximately normal distributions (see Figures 3 and 4).  
7  
8  
9

### 10 **IRT scores**

11  
12 We obtained factor scores by fitting an IRT model to the item data and using the resulting  
13 item parameters to estimate IRT-based individual phenotype scores, usually referred to as ‘factor  
14 scores’ (Chalmers, 2012). To estimate item parameters for the polytomous items we fit graded  
15 response models and to estimate item parameters for the binary items we fit 2PL models. As we  
16 originally generated the data according to these models, we knew that these were the appropriate  
17 measurement models, however, in real applications this choice should be based on considerations of  
18 the response format of items and the likely form of relations between item responses and the latent  
19 phenotype. We then computed IRT-based estimates of the phenotypic level for each individual in the  
20 sample by combining information from their patterns of item scores with the estimated item  
21 parameters from fitting the graded response model. We used Expected a Posteriori (EAP) scores: a  
22 Bayesian approach based on finding the mean of a posterior distribution representing the likelihood of  
23 phenotypic scores given a response pattern (Embretson & Reise, 2000). The posterior distribution is  
24 computed by multiplying the prior distribution (likelihoods of phenotypic levels occurring in the  
25 population) by the likelihood of the observed response pattern given the phenotypic level (Embretson  
26 & Reise 2000). This method was selected among available factor score estimation approaches  
27 because it is easy to implement and available in most IRT software packages. In the context of the  
28 models used here in which the trait of interest was uni-dimensional and the sample size large, other  
29 commonly used scoring methods such as maximum a posteriori (MAP) scoring or maximum  
30 likelihood estimates (ML) should perform similarly to EAP. Unlike using sum scores as a proxy for  
31 the phenotype, this method takes into account the scaling properties of the items. For example, in an  
32 IRT model in which items differ in discrimination, each item’s contribution to the sum score will  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

depend on its discrimination. Estimating factor scores in this way gave phenotypic scores with an approximately normal distribution (see Figure 3).

### Summary of simulation conditions

The combination of GxE interaction parameters ( $\alpha_U = -.15$  vs  $0$  vs  $.15$ ), other interaction parameters ( $\gamma_U = .20$  and  $\epsilon_U = .08$  vs  $\gamma_U = -.20$  and  $\epsilon_U = -.08$ ), item response model (GRM vs 2PL) and score type (true, sum, transformed, IRT) resulted in 28 simulation conditions. These are outlined in Tables 2 and 3. We generated 100 datasets for each condition to give 100 replications per condition.

### Model fitting

To the 100 simulated datasets for each simulation condition (see Tables 2 and 3), we fit the GxM model described in Eqs. 1-2. We fit the models in Mx (Neale, Boker, Xie & Maes, 2006) using maximum likelihood estimation, making use of the script accompanying Purcell (2002) which the author made available on his website. All latent A,C and E variances and covariances were freely estimated,  $\alpha_C$ ,  $\gamma_C$ , and  $\epsilon_C$  were fixed to zero, and  $\alpha_U$ ,  $\gamma_U$  and  $\epsilon_U$  were freely estimated. In other words, the model we fit to each dataset was consistent with the true model. The main parameter of interest was  $\alpha_U$ , which captures the moderation of the additive genetic variance unique to the phenotype by M. Parameter bias was the difference between the population magnitude and the mean estimated value across the 100 replications within a condition. In addition, we conducted a likelihood ratio test (comparing a model in which  $\alpha_U$  was freely estimated to one in which it was constrained to zero) for each replication to evaluate the statistical significance (using  $\alpha = .05$ ) of the  $\alpha_U$  parameter. Based on these, we computed false positive and false negative rates across the 100 replications. False negative rate was defined as the proportion of replications in which  $\alpha_U$  was non-significant in the presence of a non-zero population parameter. False positive rate was defined as the proportion of replications in which  $\alpha_U$  was significant in the presence of a null population parameter or where  $\alpha_U$  was statistically significant but its value was in the opposite direction to its population value (e.g. negative sample value with a positive population value).

## Simulation Study Results

1 Simulation study results are provided in Tables 2 and 3. There was only one convergence  
2 failure across all the models fit; therefore, scaling of the phenotype did not seem to have a strong  
3 influence on model convergence. Both transforming to normality and using IRT scores provided  
4 overall improvement over using raw sum scores. Whether transformed or IRT scores performed better  
5 depended on the number of response options: IRT scores were superior for polytomous items but  
6 transformations to normality were superior for binary items. More specific results are discussed  
7 below.

### 16 Control conditions

18 Results for the control conditions are provided in the ‘true score’ rows of Table 2. In these  
19 conditions, the  $\alpha_U$  parameters were generally recovered well. There was a slight positive bias when  
20 the  $\alpha_U$  parameter was in the opposite direction to the other moderation parameters. This bias appeared  
21 to reflect the imperfect resolvability of  $\alpha_U$  from  $\gamma_U$  and  $\varepsilon_U$  because it was accompanied by a negative  
22 bias in these two parameters. Power to detect moderation of the genetic influences unique to the  
23 phenotype was also generally good, as indicated by the true positive rates of 75% and above. It was  
24 lowest in the condition in which  $\alpha_U$  was in the opposite direction to the other moderation parameters.  
25 The type 1 error rates fell short of nominal levels (i.e. 5%), staying at 0% across all population models  
26 at both sample sizes.

### 40 Sum scores conditions

42 Results using a poorly scaled sum score are provided in the ‘sum score’ rows of Tables 2 and  
43 3. In all of these conditions there was positive bias in the  $\alpha_U$  parameter. These biases are in the  
44 positive direction because the IRT parameters used to generate the data produced positively skewed  
45 sum scores when the true scores were approximately normally distributed. Had item parameters been  
46 selected to produce negatively skewed sum scores, negative biases would have occurred.

54 Positive  $\alpha_U$  bias was largest in conditions in which the true moderation parameter was in the  
55 opposite direction to the direction of skew (i.e. a negative or null population moderation parameter  
56 with a positively skewed score) and the other moderation parameters. Here the biasing effects of

1  
2 scaling and imperfect resolvability of the  $\alpha_U$  and  $\gamma_U$  parameters combined to give a larger overall  
3 positive bias. Bias was slightly worse when using binary rather than polytomous items.  
4

5 Both false and true positive rates varied considerably depending on the combination of skew  
6 and moderation direction. Power was lower when using binary items than when using ordered-  
7 categorical items and when analysing 1000 rather than 2000 twin pairs. Power was also, with the  
8 exception of the condition in which the scaling enhanced a positive moderation effect, quite poor.  
9  
10

11 False positive rates were also unacceptably high and far above nominal levels. For example,  
12 in the conditions in which there was no moderation effect; significant moderation was detected 54 and  
13 46% of the time using polytomous and binary items respectively. One notable result was that when  $\alpha_U$   
14 was negative and  $\gamma_U$  and  $\varepsilon_U$  were positive, detection of moderation using sum scores derived from  
15 summing binary items occurred *only* in the wrong direction. That is, while there were 13% false  
16 positives, there were no true positives at all. Collectively, these results suggest that moderation  
17 detected using sum scores suspected to depart from the distribution of the underlying phenotype  
18 should not be relied upon.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

### 32 **Transformed sum scores conditions**

33 Overall, the effect of transforming sum scores to normality was to reduce the bias in the GxE  
34 estimates. The effectiveness of the transformation varied considerably and for the most part some  
35 positive bias remained. The exception was that in the conditions in which a sum score was formed  
36 from binary items and in which  $\alpha_U$  was in the same direction as the other moderation parameters, the  
37 transformation over-corrected the scaling problems, leading to a negative bias in  $\alpha_U$ .  
38  
39  
40  
41  
42  
43  
44  
45  
46

47 In the conditions in which  $\alpha_U$  was negative, transforming sum scores improved but did not  
48 universally successfully recover all the statistical power lost by using inappropriately scaled sum  
49 scores. Again the conditions most affected were those in which  $\alpha_U$  was in the opposite direction to the  
50 other moderation effects. For example, the true positive rate dropped from 75% for the true scores to  
51 only 4% for the transformed sum scores when using either binary or polytomous items. However,  
52 transforming the sum scores to normality had the benefit of producing marked reductions in false  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62

1 positive rates. For example, when the population parameter was zero and N=2000 twin pairs, the  
2 false positive rate was only 23% when using a transformed sum score obtained from polytomous  
3 items, compared with 54% when using a raw sum score. The corresponding drop for the sum scores  
4 obtained from binary items was 46% to 0%.  
5  
6  
7

### 8 9 **IRT scores conditions**

10  
11  
12 Results using factor scores derived from the relevant IRT model are provided in the ‘IRT’  
13 rows of Tables 2 and 3. Like the transformed sum scores, these gave consistently less biased  $\alpha_U$   
14 parameter estimates than the raw sum scores. However, some positive bias remained in all cases,  
15 ranging from very mild (+.01) to substantial (+.16) and was again most pronounced when  $\alpha_U$  was in  
16 the opposite direction to the other moderation parameters. When considering smaller sample sizes, the  
17 IRT scores yielded less biased  $\alpha_U$  estimates than transformed sum scores for polytomous items;  
18 however, the opposite was true for binary items.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29  
30 Similar to transformed sum scores, IRT scores recovered some but not all of the statistical  
31 power lost by inappropriate scaling. Whether it yielded superior power to transforming sum scores  
32 depended on the directions of moderation parameters and whether binary or polytomous items were  
33 used. In general, IRT scores provided greater power when items were polytomous but transformed  
34 sum scores were superior in this respect with binary items. This suggests that IRT scores are  
35 advantageous primarily when trait-level indicators are rated at greater levels of detail.  
36  
37  
38  
39  
40  
41  
42

43  
44 IRT scores did not prevent scaling-related false positives and although they did bring the false  
45 positive rates down, these rates remained above nominal levels. Using polytomous items, IRT scores  
46 were more effective in reducing the false positive rates than transforming sum scores; however,  
47 transforming was more effective when using binary items.  
48  
49  
50  
51  
52

## 53 **Real Data Example**

### 54 55 56 **Participants**

1 We used data from the Minnesota Twin Registry (MTR), a comprehensive description of  
2 which can be found in Krueger & Johnson (2004). The full MTR includes data from twin pairs born in  
3 Minnesota in one of three year ranges. It includes 4307 twin pairs born between 1936 and 1955, 901  
4 twin pairs born between 1904 and 1943, and 391 male twin pairs born between 1961 and 1964.  
5  
6 Eligible participants were identified from birth records, located, and invited to participate via mail.  
7  
8 Additional incentives and invitations to participate were offered to those who did not initially respond.  
9  
10 Zygosity determination was by self-reported similarity in eye colour, hair colour, overall appearance,  
11 and the difficulties others had in distinguishing two members of a pair. Analysis of a sub-sample of 74  
12 twin pairs who underwent zygosity determination by serological analysis suggested that the self-  
13 report method had an estimated accuracy of 96%.  
14  
15  
16  
17  
18  
19  
20  
21  
22

23 Different subsets of the total MTR received different sets of measures. Data used in the  
24 current study were from 528 monozygotic twin pairs and 411 dizygotic twin pairs comprising 614  
25 males and 1264 females who had completed measures of both personality and leisure time interests.  
26  
27 The mean age of the sample was 37.11 (SD=7.8).  
28  
29  
30  
31

## 32 **Measures**

### 33 **Moderator**

34  
35  
36 As our moderator we used a composite of items from the Minnesota Leisure Time Interest  
37 Test (Lykken et al., 1990). The scale asks participants to rate the extent to which they would be  
38 interested in pursuing a given activity assuming no time, health, or financial constraints. Participants  
39 rated their interest on a 5-point scale from 1= 'No interest at all' to 5= 'I would certainly do this'. In  
40 total, 120 activities were rated, but we selected 6 items to form an 'Intellectual Interests' scale.  
41  
42 Selected items refer to the following activities: reading current non-fiction, taking a college course,  
43 reading literary classics, visiting galleries/museums/exhibitions, reading books/magazines or watching  
44 TV programs on science, and reading history/philosophy/biography. We checked that these items  
45 formed a reasonable uni-dimensional scale by fitting a single factor confirmatory factor model to the  
46 data from twin 1 of each twin pair. We used the Weighted Least Squares Means and Variances  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 (WLSMV) in estimator in *Mplus 7.0* (Muthén & Muthén, 2010) to account for the categorical item  
2 response format. The 6 items all showed standardised loadings of .50 or greater and yielded a good-  
3 fitting single factor model (RMSEA=.05, CFI=.99, TLI=.99, WRMR=0.56). We therefore used the  
4 unweighted sum score of these six items as our moderator variable. Cronbach's alpha of the scale was  
5  
6  
7  
8  
9 .63.

## 10 11 **Phenotype**

12  
13  
14  
15 As our phenotypes we used personality scales from the 300-item Multidimensional  
16 Personality Questionnaire (MPQ; Tellegen & Waller, 2008). Participants were administered a version  
17 of the MPQ using a 2-point response scale. Items are phrased as statements to which participants  
18 answer 'True' or 'False' depending on whether they believe the statement describes their attitudes,  
19 opinions, interests or other characteristics.  
20  
21  
22  
23  
24  
25

26  
27 We selected two scales that yielded oppositely skewed scores. First, we used the negatively  
28 skewed 'Well-being' scale comprising 18 items. High scores on this scale are presumed to be  
29 indicative of a cheerful and happy disposition, feeling good about oneself, being optimistic, and  
30 enjoying an interesting and exciting life. Second, we used the positively skewed 'Aggression' scale  
31 comprising 18 items. High scores on this scale are presumed to be indicative of physical aggression,  
32 enjoyment of scenes of violence or upsetting or frightening others, victimisation of others for personal  
33 advantage, and vindictive and retaliatory tendencies.  
34  
35  
36  
37  
38  
39  
40  
41  
42

43 We varied how each phenotype was operationalised across conditions to mirror our  
44 simulation conditions. First, we used the raw sum score from each scale. Second, we used a  
45 transformation of the sum score that yielded an approximately normal distribution. Third, we used an  
46 IRT score for each scale. For this, we used a 2-parameter logistic model with a procedure otherwise  
47 identical to that described in the simulation study to estimate factor scores.  
48  
49  
50  
51  
52  
53

## 54 55 **Model fitting**

1 Model fitting broadly followed the procedure outlined in the simulation. However, because  
2 we were working with real data we did not know the true model and, therefore, relied on model fit  
3 comparisons to guide model selection. We first assessed whether it was possible to constrain  
4 moderation of the influences common to moderator and phenotype to zero without significant  
5 decrease in fit. We then attended to moderation of the influences unique to the phenotype. We present  
6 the parameter estimates from best-fitting model(s). In all cases, all latent A, C, and E variances and  
7 covariances were freely estimated.  
8  
9  
10  
11  
12  
13  
14  
15

## 16 **Real Data Example Results**

### 17 **Descriptive Statistics**

18  
19  
20  
21  
22 Descriptive statistics for the moderator and phenotypes are provided in Table 4. For the  
23 phenotypes, descriptive statistics are provided for sum scores, transformed sum scores and IRT  
24 scores. The Well-being sum score showed negative skew which was reduced considerably by a  
25 normalising transformation. The IRT factor scores for this phenotype showed a level of non-normality  
26 similar to the transformed sum score but slightly more negative. The correlation between Well-being  
27 and Intellectual interests was around  $r=-.18$  and practically unaffected by which phenotypic proxy was  
28 used. The correlations between the three kinds of scores derived from the Well-being items were all  
29  $>.97$ .  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

41 The Aggression sum score showed positive skewness. The transformation to normality  
42 produced scores with a near-normal distribution. The IRT factor scores for this phenotype also  
43 substantially reduced non-normality but these scores were more positively skewed than the  
44 transformed sum scores. The correlation between Aggression and Intellectual interests was around  
45  $r=-.12$  and practically identical across the three different kinds of phenotypic proxy. The correlations  
46 between the three kinds of scores derived from the Aggression items were also all  $>.97$ .  
47  
48  
49  
50  
51  
52  
53  
54

### 55 **GxM Model Fitting**



1 Fits for selected models for each phenotype and type of phenotypic score are provided in  
2 Tables 5 and 6. The parameter estimates from the best-fitting model for each phenotype across the  
3 three different phenotypic proxies (sum score, transformed sum score, and IRT score) are provided in  
4 Table 7.  
5  
6  
7

### 8 **Well-being**

9  
10 In the GxE models for Well-being, it was possible to constrain moderation of the common  
11 influences to zero without significant decrease in fit irrespective of whether a sum score, transformed  
12 sum score, or IRT score represented the phenotype. Therefore, this became the baseline model for all  
13 further model comparisons.  
14  
15  
16  
17  
18  
19  
20  
21

22 Using sum scores, model comparisons supported moderation of the genetic influences unique  
23 to the phenotype fairly unequivocally. Constraining this parameter to zero produced significant  
24 decreases in fit irrespective of whether moderation of the unique C and E influences on the phenotype  
25 were freely estimated or fixed to zero. Model fit comparisons suggested the latter model provided the  
26 best overall representation of the data: a conclusion on which there was agreement across all the  
27 information theoretical criteria examined. Thus, results suggested that the genetic influences unique to  
28 Well-being were smaller at higher levels of intellectual interests.  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

39 Using transformed sum scores, model fit comparisons suggested some moderation of unique  
40 influences for which moderation of the A influences unique to the phenotype best accounted.  
41 However, this result was not completely unequivocal: it was possible to constrain moderation of the A  
42 influences unique to the phenotype to zero without significant decrease in fit when moderation of the  
43 C and E influences were freely estimated but not when they were both fixed to zero. This further  
44 illustrates the lack of resolvability of  $\alpha_U$  and  $\gamma_U$  effects noted in the simulation study. The fact that GxE  
45 evidence was more marginal here was also reflected in the information theoretic fit criteria; for  
46 example, AIC was more negative for a model including  $\alpha_U$  while BIC was more negative for the  
47 nested model excluding this parameter. This is consistent with BIC having a larger parsimony penalty  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2 for these models. For these sets of comparisons, results suggested that the genetic variance unique to  
3 Well-being may be higher at higher levels of intellectual interests.  
4

5           When using IRT scores, results were highly similar to those for the transformed sum score in  
6 terms of fit differences and parameter magnitudes ( $\alpha_U$  was 0.04 when freely estimated but the other  
7 moderation parameters were fixed to zero). However, the difference in fit between the model in which  
8 moderation of all the unique A,C and E influences on the phenotype was fixed to zero and the model  
9 in which moderation of the unique A influences was freely estimated happened to fall just short of  
10 statistical significance. In addition, with the exception of saBIC, all information theoretic criteria were  
11 more positive for the model with  $\alpha_U$  than in the nested model excluding it. Therefore, there was  
12 technically no statistical evidence for GxE when using the IRT factor score, suggesting that the  
13 genetic influences unique to Well-being did not depend on level of intellectual Interests.  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

26           To summarise results from the Well-being scale, based on a naïve interpretation, all favoured  
27 different conclusions regarding the presence of GxE: GxE was in evidence using a sum score, was  
28 somewhat in evidence using a transformed sum score, and was not in evidence using an IRT score.  
29  
30 While the results in the latter two conditions were in actuality very similar, the fact that the statistical  
31 evidence lay on opposite sides of a statistical significance threshold and a naïve interpretation could  
32 lead to very different substantive conclusions in practice. Only the sum score condition appeared to  
33 show unambiguous support for GxE. This is consistent with the simulation conditions in which the  
34 presence of non-normality resulted in detection of GxE, irrespective of whether this non-normality  
35 was a result of moderation or poor scaling. The moderation observed using the sum score was in the  
36 direction expected for a negatively skewed sum score even when there was no true moderation. Thus,  
37 there would be reason to question the validity of the evidence for GxE observed in this real data  
38 example.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53

#### 54 **Aggression**

55  
56  
57           In all conditions, it was possible to constrain moderation of the influences common to  
58 moderator and phenotype to zero without significant drop in fit. From here, the best-fitting model  
59  
60  
61  
62  
63  
64  
65

## Phenotype scaling in GxE

1 using sum scores was one in which there was moderation of the unshared environmental influences on  
2 the phenotype captured by the  $\varepsilon_U$  parameter. Fixing  $\varepsilon_U$  to zero resulted in a significant deterioration in  
3 fit both when  $\alpha_U$  and  $\gamma_U$  were freely estimated and when they fixed to zero. Information theoretical  
4 criteria also unanimously supported the inclusion of  $\varepsilon_U$ . However, when this parameter was freely  
5 estimated, constraining moderation of neither shared environmental influences nor genetic influences  
6 on the phenotype resulted in statistically significant decrease in fit. Thus, using a sum score, there was  
7 evidence that only the unshared environmental influences unique to Aggression decreased with  
8 increasing Intellectual Interests. The direction of this moderation was in the opposite direction to the  
9 direction of the skew of the sum score. Given that the phenotype and moderator were negatively  
10 correlated, the moderation was in the direction consistent with the skew of the sum score.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

23 Using transformed sum scores, after constraining moderation of the influences common to  
24 moderator and phenotype to zero, the best-fitting model involved no moderation of the influences  
25 unique the phenotype. These could all be individually constrained to zero without significant decrease  
26 in fit, irrespective of whether moderation parameters for the other unique influences were also  
27 constrained or freely estimated. Based on information theoretic criteria, model fit was close between  
28 models including and excluding  $\varepsilon_U$ , but was - except according to AIC - better when it was excluded.  
29 Thus, on balance there was technically no evidence that the genetic or environmental influences on  
30 Aggression depended on level of Intellectual Interests.  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

42 Using IRT scores, after constraining moderation of the influences common to the moderator  
43 and phenotype to zero, there was some very weak support for moderation of the unshared  
44 environmental influences unique to the phenotype. Specifically, fixing moderation of unshared  
45 environmental influences unique to the phenotype to zero resulted in significant decrease in fit when  
46 all other moderation parameters were fixed to zero; however, the decrease in fit on constraining this  
47 parameter to zero was not statistically significant when moderation of the shared environmental and  
48 genetic influences unique to the phenotype was freely estimated. The best-fitting model according to  
49 BIC included no moderation, albeit by a small margin compared with one in which the moderation of  
50 the unshared environmental influences unique to the phenotype was freely estimated ( $\Delta\text{BIC}=0.99$ ).  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

However, AIC and saBIC favoured the model with moderation (DIC differed only in the 2<sup>nd</sup> decimal place between the two models). Considering these results together, the IRT factor score condition showed only very weak evidence for moderation intermediate between the results for the sum score (which showed evidence for moderation) and the transformed sum score (which showed no evidence for moderation). Again, the direction of moderation suggested smaller unshared environmental influences unique to Aggression at higher levels of Intellectual Interests.

## Discussion

It is well known that poorly scaled sum scores as phenotypic proxies in GxE tests can seriously bias tests of GxE. For example, using sets of items where the difficulty or location parameters are clustered near the high end of the phenotypic continuum can lead to positively skewed sum scores and, in turn, positively biased tests of GxE. In a simulation study, we assessed the extent to which this bias was mitigated by transforming non-normal sum scores to normality. We compared this to estimating phenotypic scores from an IRT model: a method that explicitly takes account of the scaling properties of items. Our results suggest that using IRT methods to provide formal models for the phenotype or appropriately transforming score distributions can provide much more accurate detection and quantification of GxE effects. Transformation may be preferred where there is insufficient information in the data (e.g. small sample size, small number of items, binary item response format) to provide good IRT latent trait estimates.

Based on our analyses, we can extend the arguments set out in the introduction in the following ways. First, we confirmed that biases in estimates of GxE can be introduced by phenotypic scaling that results in a sum score that fails to reflect the underlying distribution of the target phenotype. The nature of this bias is predictable: sum scores that are negatively skewed relative to their underlying phenotypic distribution will tend to produce negatively biased moderation parameters and sum scores that are positively skewed relative to their underlying phenotypic distribution will

1  
2 tend to produce positively biased moderation parameters. When there is no true moderation effect,  
3 this will often lead to unacceptably high false positive rates.  
4

5           These effects occur because non-normality due to poor scaling is not completely statistically  
6 distinguishable from non-normality due to presence of interaction. Where there is non-normality, the  
7 model is liable to attribute this to interaction; however, only when the observed phenotypic  
8 distribution reflects its population distribution will this estimate provide accurate quantification of  
9 GxE. Measuring the phenotype and capturing its population distribution as accurately as possible is,  
10 therefore, important in ensuring accurate assessment of GxE. When the raw score from an inventory  
11 fails to do this, there may be options for recovering this distribution via *post-hoc* manipulations of its  
12 measurement scale.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24           Our results showed, in particular, that transforming a score or using an IRT score in place of a  
25 non-normal sum score can be used to reduce in bias. We studied the case in which the *latent* genetic  
26 and environmental influences on the phenotype, absent the influence of the moderator could be  
27 assumed normally distributed in the population. This is a reasonable assumption in cases where there  
28 are a large number of small, independent effects on the phenotype. Here, a normal distribution of the  
29 joint effects of etiological contributors is predicted based on the central limit theorem. Under these  
30 conditions, using either a simple transformation or IRT scores reduced bias in GxE because they led  
31 to score distributions that better approximated the population distribution of the phenotype.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

43           In cases where there is no true moderation effect, using a phenotypic proxy that better reflects  
44 its population distribution than a sum score reduces false positive rates substantially. When the  
45 direction of the moderation is consistent with the direction of skew, either transforming to normality  
46 or using an IRT score will give close to unbiased parameter estimates and result in good power to  
47 detect the effect. However, in cases where moderation and skew are in opposite directions, these  
48 methods will under-estimate the effect and reduce power to detect GxE relative to situations in which  
49 the phenotype is not subject to scaling problems.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 We also provided a real data example from the Minnesota Twin Registry using two  
2 phenotypes with non-normal sum scores. Analysing the Well-being phenotype using (negatively  
3 skewed) sum scores yielded statistically and practically significant GxE whereas using IRT scores  
4 suggested no significant GxE. The transformed sum scores yielded evidence intermediate between  
5 these two outcomes. The direction of the GxE using sum scores was consistent with the direction of  
6 the skewness of the sum score. This suggests that the observed effect could be due to item scaling.  
7 Moreover, based on these results, researchers using sum scores rather than IRT scores could easily  
8 have been led to opposite substantive conclusions despite the high correlations between the raw and  
9 IRT scores.  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20

21 The Aggression phenotype did not yield evidence of GxE irrespective of whether (positively  
22 skewed) sum scores, transformed sum scores, or IRT scores were used. This shows that non-normal  
23 trait distributions will not automatically result in the appearance of GxE and that altering phenotypic  
24 distributions will not necessarily affect the GxE parameter. However, there was evidence for  
25 dependence of another moderation parameter on scaling: using sum scores and an IRT scores,  
26 negative moderation of the unshared environmental influences unique to the phenotype (captured by  
27 the  $\epsilon_U$  parameter) was detected. There was no such evidence using a transformed sum score. Taking  
28 into account the fact that the phenotype and moderator were negatively correlated, the  $\epsilon_U$  parameter  
29 was proportional to and in the direction consistent with the skew of the phenotypic proxy. That is, the  
30 parameter was most negative when the phenotypic proxy was strongly skewed (sum score), less  
31 negative when the phenotypic proxy was moderately positively skewed (IRT score) and effectively  
32 zero when the phenotypic proxy was only slightly positively skewed (transformed sum score). Thus,  
33 although we have focussed on the  $\alpha_U$  parameter because it is most often used to operationalise  
34 theoretical hypotheses, this example highlights the fact that the effects of scaling on GxE models are  
35 not confined to that one parameter.  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54

55 Our results reinforce the message that poorly scaled sum scores should be avoided in tests of  
56 GxE. Poorly scaled sum scores, in addition to producing high false positive rates, can yield results that  
57 suggest significant moderation in the opposite direction to the true moderation effect. Demonstrating  
58  
59  
60  
61  
62  
63  
64  
65

1 that sum scores are highly correlated with transformed sum scores or IRT scores for the same  
2 phenotype is thus not sufficient justification for using them in place of these better-performing  
3 methods. Because correlation coefficients are relatively unaffected by rank-preserving  
4 transformations, sum and functionally-transformed scores will show very high correlations, even  
5 when their distributions are markedly different. IRT scoring basically differentially weights the items  
6 or response options rather than weighting each one equivalently as does sum scoring, thus very  
7 closely preserving rank ordering. This was illustrated in our real data examples where, in spite of  
8 leading to diverging conclusions about the presence and strength of moderation effects, the three types  
9 of score were correlated with one another at  $>.97$ .

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21 The strategies of transforming sum scores to normality or using an IRT score did not suffer  
22 the limitations of poorly scaled sum scores to anywhere near the same extent; however, both resulted  
23 in tests that lacked statistical power when the moderation was in the opposite direction to skew and  
24 failed to control the type 1 error rate completely when GxE was not present. Overall, transforming  
25 non-normal sum scores to normality or using IRT scores will in many cases fail to address the biasing  
26 effects of poor scaling on GxE tests, especially when there is non-genetic moderation in the opposite  
27 direction to the genetic moderation. Therefore, evidence of GxE (or lack thereof) should be  
28 considered tentative even when obtained using transformed or IRT scores.

29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39 Although using IRT scores is more time consuming and technically demanding than  
40 transformations to normality, it may be worth the additional effort, especially when the raw scale  
41 items were rated using multiple response options. IRT scores can be estimated reasonably easily in a  
42 range of freely available software packages and have several practical and theoretical advantages over  
43 transformed sum scores. First, they are easily estimable in the presence of missing item data, or when  
44 respondents did not complete an identical set of items (Embretson & Reise, 2000). Second, the  
45 diversity of available IRT models means that many kinds of response formats, scale structures, or  
46 theories about how the latent trait relates to item responses can be accommodated. For example, a bi-  
47 factor model could be fit when it is desirable to partition general and specific trait variance captured  
48 by a set of items (Cai, Yang & Henson, 2011); if a scale has a categorical response format, a nominal

1 response model could be fit (Bock,1972); or if items follow an ideal point process an unfolding model  
2 can be fit (e.g. Chernyshenko, Stark, Drasgow & Roberts, 2007). All of these and other features can  
3 be easily dealt with in an IRT framework while posing significant problems or being simply  
4 impossible to take account of when using sum scores, both raw and transformed to normality.  
5  
6 Furthermore, while an IRT model can be chosen based on theoretical considerations, the choice of a  
7 transformation is somewhat arbitrary and usually driven by pragmatic considerations. The choice of  
8 an IRT model can be evaluated both overall and with respect to individual items using well-studied  
9 goodness-of-fit statistics and graphical checks. A beneficial side effect of this is that the process of  
10 fitting and evaluating IRT model(s) is likely to encourage explicit consideration of the assumptions  
11 that underpin the phenotypic proxy used. However, no analogous tests exist for transformations. More  
12 importantly, from a conceptual perspective, if the genetic and environmental influences on the  
13 phenotype in the absence of the influence of the moderator are normally distributed and there is true  
14 GxE in the population then the phenotype *should* show a non-normal distribution because GxE  
15 involves an expansion (or contraction) of the variance in a phenotype according to the levels of  
16 moderator. This expansion (or contraction) of variance shows up in the marginal distribution of the  
17 phenotype as non-normality that is commensurate with the GxE effect. Using a transformation to  
18 normality is, therefore, directly at odds with theoretical expectations when GxE is hypothesised. In  
19 IRT models, this is also a problem to some extent; however, the assumption of a normal latent  
20 distribution is not a necessity; where appropriate alternative prior distributions can be specified in a  
21 manner that is far more flexible than attempting to obtain that distribution through transformation of  
22 observed scores.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

47 The primary disadvantage of IRT scoring is practical: to be effective requires large sample  
48 sizes and ideally a large number of items with polytomous response formats. Where any of these is  
49 lacking, transformed sum scores may be more effective than IRT scores. This underlines the  
50 importance of assessing the empirical reliability of factor scores from IRT models, as one would for  
51 sum scores (see Culpepper, 2013). Unreliable IRT scores will not only be ineffective in addressing  
52 bias in GxE; they will also result in attenuated estimates of twin correlations and bias other model  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 parameters (van den Berg et al. 2007). Similarly, as the extent to which the accuracy of the scores as  
2 measures of the intended underlying dimension depends on the appropriateness of the IRT model, its  
3 specification should be carefully considered and its fit assessed empirically (see Embretson & Reise,  
4 2013).  
5  
6  
7

8  
9  
10 Where both approaches are limited is that the underlying liability distribution absent the  
11 influence of the moderator could be non-normal due to other moderators or the effects of rare but  
12 highly influential etiological factors that engender extreme effects. Analogous to the problem of  
13 distinguishing non-normality due to moderation versus poor scaling, it is not easy to disentangle non-  
14 normality due to the effect of a moderator of interest and non-normality due to other etiological  
15 factors without detailed a priori knowledge.  
16  
17  
18  
19  
20  
21  
22

23  
24 Further, the favourable performance of the IRT scores in the simulation study should be  
25 interpreted in light of the fact that they were estimated under idealised conditions. In practice their use  
26 is more complicated and may be less effective. For example, we fit graded response and 2-parameter  
27 logistic models to our polytomoyus and binary data respectively because we knew that these models  
28 had been used to generate the item responses. Thus, there was no risk of seriously mis-specifying the  
29 psychometric model. In reality, the appropriate model for the items will not be known in advance- it  
30 will have to be chosen on the basis of the item format and a hypothesis about how the latent trait is  
31 related to item responding and then tested for appropriateness. The lack of *a priori* knowledge about  
32 the appropriate IRT model for a given set of items increases the risk that the chosen model will be  
33 mis-specified in some important way. Further, parametric IRT models are also often poor fits to the  
34 very same kinds of data that prove problematic in GxE tests, such as those concerning  
35 psychopathological phenotypes. Less restrictive non-parametric IRT models are sometimes  
36 recommended as alternatives (Meijer & Baneke, 2004) but these methods do not allow estimation of  
37 factor scores for use in GxE tests. Finally, at a very pragmatic level, IRT models are only useful when  
38 item-level data are available, which is not always the case.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

In practice, it is worthwhile to compare results obtained using IRT scores with those obtained using raw and transformed sum scores. Comparison can highlight how sensitive results are to phenotypic scaling. Under some conditions, e.g. when the phenotype and moderator do not have strong association or the phenotypic distribution departs only slightly from its population distribution, scaling of the phenotype may make little difference to results. In addition, in rare cases where the phenotypic distribution is mis-specified in the IRT model used to estimate the scores but well approximated by the sum scores, the sum scores could, in principle, produce less biased results than the IRT scores. Even when the phenotypic distribution is correctly assumed to be normal, no non-linear transformation or IRT score estimation method guarantees a perfect reconstruction of the phenotypic distribution as it exists in the population. In fact, as argued above, the scores produced by a transformation to normality could be ‘too normal’ in the sense that in the presence of GxE non-normality of the phenotype would usually be expected. This is exactly what occurred in, for example, the condition of the simulation study in which all moderation parameters were positive in the population and in which a sum score from binary items was transformed to normality. Transforming to normality yielded a parameter estimate that was almost as negatively biased as the original estimate from using the sum score was positively biased. Moreover, the true positive rate dropped from 81% to 34% suggesting a significant drop in the power to detect GxE.

This result underscores the fact that near-normal observed score distributions should not always be the goal. Non-normal latent distributions would be expected when, for example, a phenotype is influenced by GxE processes (perhaps not related to the moderator of interest), when it is influenced by some genetic (or environmental) variants of disproportionately large effect, or when phenotypic expression is subject to a liability threshold. Without some knowledge of the etiology of the trait, the appropriate distribution to which to transform or to assume in an IRT model will not be obvious. For example, although empirical methods exist that attempt to determine a latent trait distribution and IRT parameters simultaneously (e.g. Woods, 2006), in practice the same patterns of item responses may be represented equally well by a range of combinations of distributions and IRT parameters (e.g. Pilkonis et al., 2011). There remains an important role of theoretical knowledge in

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

determining which of these combinations is the most biologically plausible. We believe that the continuing advances in characterising the etiologies of complex traits will increasingly serve to inform the reasonableness of distributional assumptions and measurement models for phenotypes in testing GxE. Although it was once necessary (at least in practical terms) to assume multivariate normality for parameter estimation, recent and continuing developments in statistical methodology mean that this is no longer the case. Rather, the primary limiting factor at present is the theoretical knowledge to guide the specification of an appropriate (implicit or explicit) measurement model, rather than the statistical models to operationalise it.

Finally, our results highlight some challenges with testing GxE even under optimal scaling conditions. In our control conditions, there was a slight negative bias in GxE estimates when this effect was in the opposite direction to moderation of shared and unshared environmental influences. In addition, although power to detect GxE was under optimal scaling, type 1 error rates were below nominal levels. This has also been observed in previous studies of the GxM model (van Hulle, Lahey & Rathouz, 2013) and suggests that nested model comparisons for the GxE provide conservative tests.

### **Limitations**

A limitation of the current study is that we did not directly compare the two-step IRT approach with a one-step approach presented here. A one-step approach has yet to be developed for testing of GxE within the Purcellian framework; however, it is possible to anticipate some of its disadvantages and advantages. First, the approach would share the limitation of the two-step approach that the true phenotypic distribution would not be known but assumed. Assuming a normal distribution for the phenotype when the true distribution is non-normal could, in principle, result in biased GxE tests in a similar way to using a poorly scaled sum score. It would also share the necessity to select an appropriate IRT model and freely estimate its parameters in a finite sample. A further disadvantage would be its statistical and computational complexity as compared to a two-step approach. However, an important advantage would be that the error-free latent trait could be decomposed directly and this is likely to result in less biased GxE tests. It would have the related

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

advantage that the IRT parameters would not have to be taken as given as they are in the second step of the two-step approach. Therefore, the imprecision in these parameter estimates could be appropriately taken account of. Further, and perhaps most importantly, a one-step approach is more appropriate from a conceptual perspective because it provides a much more direct operationalization of GxE hypotheses. In the two-step approach, a distribution for the phenotype is assumed in the first step; however, in tests of GxE it is important to distinguish between assumptions about the marginal distribution of the phenotype and the distribution of the underlying genetic and environmental influences absent the influence of the moderator. While the former would be expected to be non-normal because being subject to moderation skews the phenotypic distribution, the latter can usually be assumed normal. The two-step approach unfortunately conflates these distinct contributions because it specifies a distribution only for the latent phenotype. In addition, although we designed our simulation conditions to be as realistic as possible, we covered only a limited range of the possible conditions that could occur in the real world. Although the principles discussed are likely general, we conducted our analyses within specific GxE and IRT frameworks and used a limited range of parameter values. Similarly, while inclusion of a real data example is important to test conclusions from simulation studies in a more ecologically valid context, these too are limited by their specificity.

## Conclusions

Tests of GxE can be biased by inappropriate scaling of a phenotype, and reliance on raw scores that are suspected to mis-represent the underlying distribution of the target phenotype . Two potentially useful solutions are to transform sum scores to normality or to estimate IRT scores based on an appropriate model. Although these strategies will suffer low statistical power, they reduce the rate of spurious GxE detection and recover the correct direction of effects. Therefore, researchers can be more confident about the presence and direction of GxE when it is identified using one of these strategies than when using a raw sum score.

## Footnotes<sup>1</sup>

## Phenotype scaling in GxE

Purcell's GxM approach assumes a normal distribution for the phenotype conditional on the moderator; however, the presence of moderation will result in a skewed marginal distribution for the phenotype.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

References

- 1  
2  
3  
4  
5  
6 Asbury K, Wachs TD, Plomin R (2005) Environmental moderators of genetic influence on  
7  
8 verbal and nonverbal abilities in early childhood. *Intelligence* 33 (6): 643-661  
9  
10  
11 Baron-Cohen S, Wheelwright S, Skinner R, Martin J, Clubley E (2001) The autism-spectrum  
12  
13 quotient (AQ): Evidence from asperger syndrome/high-functioning autism, males and  
14  
15 females, scientists and mathematicians. *J Autism Dev Disord* 31(1): 5-17  
16  
17  
18  
19 Beasley TM, Erickson S, Allison, DB (2009) Rank-based inverse normal transformations are  
20  
21 increasingly used, but are they merited?. *Behav Genet* 39 (5) 580-595  
22  
23  
24 Boomsma DI, Martin NG (2002) Gene–environment interactions. In: D’haenen H, den Boer JA,  
25  
26 Willner P (eds) *Biological Psychiatry*. Wiley, New York, pp 181–187  
27  
28  
29  
30 Bock RD (1972) Estimating item parameters and latent ability when responses are scored in two or  
31  
32 more nominal categories. *Psychometrika* 37 (1): 29-51  
33  
34  
35 Bronfenbrenner U, Ceci SJ (1994) Nature-nuture reconceptualized in developmental perspective: A  
36  
37 bioecological model. *Psychol Rev* 101 (4): 568- 586  
38  
39  
40 Burt, SA, Klump KL (2009) The etiological moderation of aggressive and nonaggressive antisocial  
41  
42 behavior by age. *Twin Res* 12 (4): 343-350  
43  
44  
45 Button TM, Hewitt JK, Rhee SH, Corley RP, Stallings MC (2010) The moderating effect of  
46  
47 religiosity on the genetic variance of problem alcohol use. *Alcohol Clin Exp Res* 34 (9):1619-  
48  
49 1624  
50  
51  
52  
53 Cai L, Yang JS, Hansen M (2011) Generalized full-information item bifactor analysis. *Psychol*  
54  
55 *Methods* 16 (3): 221-248  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5 Chalmers RP (2012). mirt: A Multidimensional Item Response Theory Package for the R  
6 Environment. *J Stat Softw* 48 (60): 1-29  
7  
8  
9  
10 Chernyshenko OS, Stark S, Drasgow F, Roberts BW (2007). Constructing personality scales under the  
11 assumptions of an ideal point response process: toward increasing the flexibility of  
12 personality measures *Psych Assess* 19 (1): 88-106  
13  
14 Culpepper SA (2013) The reliability and precision of total scores and IRT estimates as a function of  
15 polytomous IRT parameters and latent trait distribution. *Appl Psych Meas* 37 (3): 201-225  
16  
17  
18 Distel MA, Middeldorp CM, Trull TJ, Derom CA, Willemsen G, Boomsma DI (2011) Life events  
19 and borderline personality features: the influence of gene–environment interaction and  
20 gene–environment correlation. *Psychol Med* 41 (4): 849-860  
21  
22  
23  
24  
25  
26 Eaves LJ, Last K, Martin NG, Jinks JL (1977) A progressive approach to non-additivity and  
27 genotype-environmental covariance in the analysis of human differences. *Br J Math Stat*  
28 *Psychol* 30 (1): 1-42  
29  
30  
31  
32  
33 Eaves LJ (2006) Genotype× environment interaction in psychopathology: fact or artifact?. *Twin Res*  
34 9 (1): 1-8  
35  
36  
37  
38  
39 Embretson SE (1996) Item response theory models and spurious interaction effects in factorial  
40 ANOVA designs. *Appl Psych Meas* 20 (3): 201-212.  
41  
42  
43  
44 Embretson SE, Reise SP (2000) Item response theory for psychologists. Psychology Press.  
45  
46  
47 Falconer DS, Mackay TF (1996) Introduction to quantitative genetics. Harlow. UK: Longman.  
48  
49  
50 Harden KP, Turkheimer E, Loehlin JC (2007) Genotype by environment interaction in adolescents'  
51 cognitive aptitude. *Behav Genet* 37 (2): 273-283  
52  
53  
54  
55  
56 Hicks BM, South SC, DiRago AC, Iacono WG, McGue M (2009) Environmental adversity and  
57 increasing genetic risk for externalizing disorders. *Arch Gen Psychiatry* 66 (6): 640-648  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- Hicks BM, DiRago AC, Iacono WG, McGue M (2009) Gene–environment interplay in internalizing disorders: consistent findings across six environmental risk factors. *J Child Psychol Psychiatry* 50 (10):1309-1317
- Johnson W, Krueger RF (2005) Genetic effects on physical health: lower at higher income levels. *Behav Genet* 35 (5): 579-590
- Johnson W, Kyvik KO, Mortensen EL, Skytthe A, Batty GD, Deary IJ (2011) Does education confer a culture of healthy behavior? Smoking and drinking patterns in Danish twins. *Am J Epidemiol* 173 (1): 55-63
- Kang SM, Waller NG (2005) Moderated multiple regression, spurious interaction effects, and IRT. *Appl Psychol Meas* 29(2): 87-105
- Krueger RF, Johnson W (2002) The Minnesota twin registry: current status and future directions. *Twin Res Hum Genet* 5 (5): 488-492
- Latvala A, Dick DM, Tuulio-Henriksson A, Suvisaari J, Viken RJ, Rose RJ, Kaprio J (2011) Genetic correlation and gene–environment interaction between alcohol problems and educational level in young adulthood. *J Stud Alcohol Drugs* 72 (2) 210-220
- Lykken DT, Bouchard TJ, McGue M, Tellegen A (1990) The Minnesota twin family registry: Some initial findings. *Acta Genet Med Gemellol* 39 (1) 35-70
- Martin N (2000) Gene-environment interaction and twin studies. In Spector, T., Sneider, H., MacGregor, A. (Eds). *Advances in twin and sib-pair analysis*. Greenwich Medical Media: London, UK.
- Mather K, Jinks JL (1971) *Biometrical genetics*. Biometrical Genetics (Ed. 2). Chapman and Hall: London, UK
- Meijer RR, Baneke JJ (2004) Analyzing psychopathology items: a case for nonparametric item response theory modeling. *Psychol Methods* 9 (3): 354-368



1  
2 Meijer RR, Egberink IJ (2012) Investigating invariant item ordering in personality and clinical scales  
3 some empirical findings and a discussion. *Educ Psychol Meas* 72 (4): 589-607  
4

5 Micceri T. (1989) The unicorn, the normal curve, and other improbable creatures. *Psychol Bull* 105  
6  
7 (1): 156-166  
8  
9

10 Molenaar D, Dolan CV (2014) Testing systematic genotype by environment interactions using  
11  
12 item level data. *Behav Genet* 44 (3): 212-231  
13  
14

15 Molenaar D, van der Sluis S, Boomsma DI, Dolan CV (2012) Detecting specific genotype by  
16  
17 environment interactions using marginal maximum likelihood estimation in the classical  
18  
19 twin design. *Behav Genet* 42 (3): 483-499  
20  
21  
22

23 Morse BJ, Johanson GA, Griffeth RW (2012) Using the graded response model to control  
24  
25 spurious interactions in moderated multiple regression. *Appl Psych Meas* 36 (2): 122-146  
26  
27

28 Muthén LK, Muthén BO (2010). *Mplus User's Guide: Statistical Analysis with Latent Variables:*  
29  
30 *User's Guide.* Muthén & Muthén  
31  
32

33 Neale MC, Boker SM, Xie G, Maes HH (2006) *Mx: statistical modeling*, 7th edn. VCU Department  
34  
35 of Psychiatry, Richmond  
36  
37  
38

39 Nydick SW (2014) *catIrt: An R Package for Simulating IRT-Based Computerized Adaptive Tests.*  
40  
41  
42 R package version 0.5-0. <http://CRAN.R-project.org/package=catIrt>  
43  
44

45 Pilkonis PA., Choi SW, Reise SP, Stover AM, Cella D (2011) Item banks for measuring emotional  
46  
47 distress from the Patient-Reported Outcomes Measurement Information System (PROMIS®):  
48  
49 depression, anxiety, and anger *Assess* 18 (3): 263-283  
50  
51

52 Pluess M, Belsky J (2013) Vantage sensitivity: Individual differences in response to positive  
53  
54 experiences. *Psychol Bull* 139 (4): 901-916  
55  
56

57 Purcell S (2002) Variance components models for gene–environment interaction in twin  
58  
59 analysis. *Twin Res* 5 (6): 554-571  
60  
61  
62

## Phenotype scaling in GxE

- 1  
2 R Core Team (2014) R: A language and environment for statistical computing. R Foundation for  
3 Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>  
4
- 5 Reiss D, Leve LD, Neiderhiser JM (2013) How genes and the social environment moderate  
6 each other. *Am J Public Health* 103 (S1): S111-S121  
7  
8
- 9  
10 Rathouz PJ, Van Hulle CA, Rodgers JL, Waldman ID, Lahey BB (2008) Specification, testing, and  
11 interpretation of gene-by-measured-environment interaction models in the presence of gene-  
12 environment correlation. *Behav Genet* 38 (3): 301-315  
13  
14  
15
- 16  
17 Rende R, Plomin R (1992) Diathesis-stress models of psychopathology: A quantitative genetic  
18 perspective. *Appl and Prev Psychol* 1 (4): 177-182  
19  
20  
21
- 22  
23 Samejima F (1969) Estimation of latent ability using a response pattern of graded  
24 scores. *Psychometrika Monograph Supplement* 34 (4): 100  
25  
26  
27
- 28  
29 Schwabe I, van den Berg SM (2014) Assessing genotype by environment interaction in case of  
30 heterogeneous measurement error. *Behav Genet* 44 (4): 394-406  
31  
32  
33
- 34  
35 Shanahan MJ, Hofer SM (2005) Social context in gene-environment interactions: Retrospect  
36 and prospect. *Journals Gerontol B Psychol Sci Soc Sci* 60 (Special Issue 1): 65-76  
37  
38
- 39  
40 Silventoinen K, Hasselbalch AL, Lallukka T, Bogl L, Pietiläinen KH, Heitmann BL, ... Kaprio J  
41 (2009) Modification effects of physical activity and protein intake on heritability of body size  
42 and composition. *Am J Clin Nutr* 90 (4): 1096-1103  
43  
44  
45
- 46  
47 South SC, Krueger RF, Johnson W, Iacono WG (2008). Adolescent personality moderates  
48 genetic and environmental influences on relationships with parents. *J Pers Soc Psychol* 94 (5):  
49 899-912  
50  
51  
52
- 53  
54 South SC, Krueger RF (2011). Genetic and environmental influences on internalizing  
55 psychopathology vary as a function of economic status. *Psychol Med* 41 (1): 107-117  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 South SC, Krueger RF (2014). Genetic strategies for probing conscientiousness and its relationship to  
2 aging Dev Psychol 50 (5): 1362- 1376  
3

4  
5 Tabery J (2008) RA Fisher, Lancelot Hogben, and the origin (s) of genotype–environment  
6  
7 interaction. J Hist Biol 41 (4): 717-761  
8  
9

10 Thomas ML (2011) The value of item response theory in clinical assessment: A review. Assessment  
11  
12 18 (3): 291-307  
13  
14

15  
16 Tellegen A, Waller NG (2008) Exploring personality through test construction: Development of the  
17  
18 Multidimensional Personality Questionnaire. In Boyle G, Matthews G, Saklofske DH. The  
19  
20 SAGE handbook of personality theory and assessment, volume 2: Personality  
21  
22 measurement and testing. Sage Publications: Thousand Oaks, CA, US.  
23  
24  
25

26  
27 Timberlake DS, Rhee SH, Haberstick BC, Hopfer C, Ehringer M, Lessem JM, ... Hewitt JK (2006)  
28  
29 The moderating effects of religiosity on the genetic and environmental determinants of  
30  
31 smoking initiation. Nicotine Tob Res 8 (1): 123-133  
32  
33

34 Tucker-Drob EM, Harden KP, Turkheimer E (2009) Combining nonlinear biometric and  
35  
36 psychometric models of cognitive abilities. Behav Genet 39 (5): 461-471  
37  
38

39 Tuvblad C, Grann M, Lichtenstein P (2006) Heritability for adolescent antisocial behavior  
40  
41 differs with socioeconomic status: gene–environment interaction. J Child Psychol  
42  
43 Psychiatry 47 (7): 734-743  
44  
45

46  
47 van den Berg SM, Glas CA, Boomsma DI (2007) Variance decomposition using an IRT  
48  
49 measurement model. Behav Genet 37(4): 604-616  
50  
51

52 van den Oord EJ, Pickles A, Waldman ID (2003) Normal variation and abnormality: an empirical  
53  
54 study of the liability distributions underlying depression and delinquency. J Child Psychol  
55  
56 Psychiatry 44 (2): 180-192  
57  
58  
59  
60  
61  
62

## Phenotype scaling in GxE

1 van den Oord EJ, Simonoff E, Eaves LJ, Pickles A, Silberg J, Maes H (2000) An evaluation of  
2 different approaches for behavior genetic analyses with psychiatric symptom scores. Behav  
3 Genet 30 (1): 1-18  
4  
5

6  
7 van der Sluis S, Dolan CV, Neale MC, Boomsma DI, Posthuma D (2006) Detecting genotype-  
8 environment interaction in monozygotic twin data: comparing the Jinks and Fulker test and a  
9 new test based on marginal maximum likelihood estimation. Twin Res Hum Genet 9 (3): 377-  
10 392  
11  
12  
13  
14  
15

16  
17 van der Sluis S, Posthuma D, Dolan CV (2012) A note on false positives and power in  $G \times E$   
18 modelling of twin data. Behav Genet 42 (1): 170-186  
19  
20  
21

22 van Hulle CA, Lahey BB, Rathouz PJ (2013) Operating characteristics of alternative statistical  
23 methods for detecting gene-by-measured environment interaction in the presence of gene-  
24 environment correlation in twin and sibling studies. Behav Genet 43 (1): 71-84  
25  
26  
27  
28  
29

30 Van Hulle CA, Rathouz PJ (2015) Operating characteristics of statistical methods for detecting gene-  
31 by- measured environment interaction in the presence of gene-environment correlation  
32 under violations of distributional assumptions. Twin Res Hum Genet 18(1):19-27  
33  
34  
35  
36  
37  
38  
39  
40  
41

42 Walton KE, Ormel J, Krueger RF (2011) The dimensional nature of externalizing behaviors in  
43 adolescence: Evidence from a direct comparison of categorical, dimensional, and hybrid  
44 models. J Abnorm Child Psychol, 39(4): 553-561  
45  
46  
47  
48  
49

50 Woods CM (2006). Ramsay-curve item response theory (RC-IRT) to detect and correct for nonnormal  
51 latent variables. Psychol Methods 11(3): 253-273  
52  
53

54 Wray NR, Lee SH, Mehta D, Vinkhuyzen AA, Dudbridge F, Middeldorp CM (2014) Research  
55 Review: Polygenic methods and their application to psychiatric traits.  
56 J Child Psychol Psychiatry 55(10): 1068-1087  
57  
58  
59  
60  
61

Phenotype scaling in GxE

1  
2 Zheng H, Rathouz P (2013) GxM: Maximum Likelihood Estimation for Gene-by-Measured  
3 Environment Interaction Models. R package version 1.0.  
4 <http://CRAN.Rproject.org/package=GxM>.  
5  
6

7 Zheng H, Rathouz PJ (2015) Fitting Procedures for Novel Gene-by-Measured Environment  
8 Interaction Models in Behavior Genetic Designs. Behav Genet 45(4): 467-479  
9

10  
11  
12  
13 Zheng H, Van Hulle CA, Rathouz PJ (2015) Comparing Alternative Biometric Models with and  
14 without Gene-by-Measured Environment Interaction in Behavior Genetic Designs:  
15 Statistical Operating Characteristics. Behav Genet, 45(4):480-491  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1: Parameter values for IRT models used to simulate item responses**

Item	Polytomous item parameters (GRM)					Binary item parameters (2PL)
	$a$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_4$	$\beta$
1	1.94	-0.27	0.84	2.23	2.74	1.29
2	1.93	-0.21	1.46	2.01	2.73	2.23
3	1.96	-0.11	1.50	2.38	2.82	0.67
4	2.13	-0.36	1.29	2.07	2.65	1.22
5	1.09	0.34	1.16	2.07	2.73	-0.03
6	1.13	-0.15	1.34	2.00	2.78	0.99
7	0.87	0.34	0.99	2.34	2.64	1.11
8	0.99	0.23	0.68	2.33	2.62	0.88
9	1.63	0.43	0.98	2.22	2.83	1.94
10	1.01	0.04	1.22	2.39	2.73	0.12
11	1.75	0.10	0.93	2.27	2.63	-0.33
12	0.80	0.01	0.67	2.20	2.75	0.89
13	0.67	0.37	1.49	2.42	2.67	0.45
14	1.91	0.13	0.89	2.29	2.92	1.01
15	1.06	0	1.29	2.09	2.96	2.20
16	0.55	0.50	0.76	2.32	2.81	2.03
17	1.88	-0.24	1.02	2.07	2.74	0.65
18	2.44	-0.40	0.80	2.09	2.86	1.00
19	0.90	-0.11	1.27	2.27	2.73	1.45
20	1.15	-0.24	0.65	2.17	2.73	1.20

*Note.*  $a$  is an item discrimination parameter,  $\beta_1$  -  $\beta_4$  and

$\beta$  are threshold parameters. The same  $a$  values were used in both the GRM- and 2PL-generated item responses. GRM=graded response model, 2PL= 2-parameter logistic model.

Phenotype scaling in GxE

**Table 2: Performance of sum score, transformed score and IRT score with polytomous items**

Score type	<i>N=1000 twin pairs</i>							<i>N=2000 twin pairs</i>						
	Population GxM values			Average $\alpha_U$	$\alpha_U$ Bias	$\alpha_U$ true positive rate	$\alpha_U$ false positive rate <sup>a</sup>	Average $\alpha_U$	$\alpha_U$ Bias	$\alpha_U$ true positive rate	$\alpha_U$ false positive rate <sup>a</sup>			
	$a_c$	$c_c$	$e_c$	$\alpha_U$	$\gamma_U$	$\epsilon_U$	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)		
True	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.15 (.04)	.00	98%	0%	.15 (.03)	+.00	100%	0%
True	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	-.12 (.05)	+.03	75%	0%	-.14 (.03)	+.01	97%	0%
True	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.00 (.03)	.00	N/A	0%	.00 (.02)	+.00	N/A	0%
True	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.15 (.05)	.00	96%	0%	-.15 (.04)	.00	96%	0%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.22 (.05)	+.07	94%	0%	.23 (.04)	+.08	98%	0%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	.03 (.08)	+.18	1%	8%	.02 (.05)	+.17	2%	8%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.14 (.07)	+.14	N/A	54%	.13(.05)	+.13	N/A	87%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.06 (.05)	+.09	15%	0%	-.05 (.03)	+.10	23%	0%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.16 (.03)	+.01	73%	0%	.16 (.03)	+.01	98%	0%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	-.02 (.05)	+.13	4%	0%	-.02 (.03)	+.13	8%	1%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.08 (.04)	+.08	N/A	23%	.08 (.03)	+.08	N/A	63%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.11 (.03)	+.04	68%	0%	-.11 (.02)	+.04	97%	0%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.16 (.04)	+.01	80%	0%	.16 (.03)	+.01	98%	0%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	-.06 (.05)	+.09	13%	0%	-.07 (.03)	+.08	50%	0%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.06 (.05)	+.06	N/A	16%	.05 (.03)	+.05	N/A	26%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.13 (.03)	+.02	79%	0%	-.12 (.02)	+.03	98%	0%

Phenotype scaling in GxE

<sup>a</sup>False positive defined as significant effect in opposite direction to population parameter or significant effect in any direction when population parameter is zero. True positive defined as significant effect in the correct direction.

**Table 3: Performance of sum score, transformed score and IRT score with binary items**

Score type	Population GxM values						<i>N=1000 twin pairs</i>				<i>N=2000 twin pairs</i>			
	$a_c$	$c_c$	$e_c$	$a_U$	$\gamma_U$	$\epsilon_U$	Average $a_U$ (SD)	$a_U$ Bias	$a_U$ true positive rate	$a_U$ false positive rate <sup>a</sup>	Average $a_U$ (SD)	$a_U$ Bias	$a_U$ true positive rate	$a_U$ false positive rate <sup>a</sup>
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.23 (.05)	+08	81%	0%	.22 (.04)	+07	97%	0%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	.05 (.09)	+20	0%	13%	.03 (.05)	+18	0%	11%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.14 (.07)	+14	N/A	46%	.14 (.05)	+14	N/A	79%
Sum	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.04 (.05)	+11	15%	0%	-.04 (.04)	+11	15%	0%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.09 (.04)	-.06	34%	0%	.10 (.03)	-.05	67%	0%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	-.04 (.05)	+11	4%	0%	-.05 (.03)	+10	13%	0%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.03 (.03)	+03	N/A	0%	.03 (.02)	+03	N/A	4%
Transformed	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.13 (.04)	+02	49%	0%	-.14 (.03)	+01	88%	0%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	.15	.20	.08	.17 (.03)	+02	79%	0%	.16 (.03)	+01	98%	0%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	.20	.08	.01 (.06)	+16	0%	2%	.00 (.04)	+15	1%	3%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	0	.20	.08	.09 (.04)	+09	N/A	32%	.09 (.03)	+09	N/A	59%
IRT	$\sqrt{.3}$	$\sqrt{.1}$	$\sqrt{.1}$	-.15	-.20	-.08	-.08 (.03)	+07	24%	0%	-.08 (.02)	+07	67%	0%

<sup>a</sup>False positive defined as significant effect in opposite direction to population parameter or significant effect in any direction when population parameter is zero. True positive defined as significant effect in the correct direction. Refer to Table 2 for results of control conditions.



**Table 4: Descriptive statistics for Well-being, Aggression and Intellectual Interests phenotypes**

Phenotypic proxy	N	N	Mean (SD)	Skew	Kurtosis	Correlation with moderator
	MZ pairs	DZ pairs				
<b>Intellectual Interests sum score</b>	528	411	13.32 (3.75)	0.13	-0.27	N/A
<b>Well-being sum score</b>	525	406 <sup>a</sup>	11.15 (2.21)	-1.06	0.71	.18
<b>Well-being sum score transformed</b>	525	406 <sup>a</sup>	0 (1)	-0.36	-0.90	.19
<b>Well-being IRT score</b>	528	411	0 (0.89)	-0.42	-0.32	.18
<b>Aggression sum score</b>	525	411	3.66 (3.21)	1.12	1.09	-.12
<b>Aggression sum score transformed</b>	525	411	0 (1)	0.23	-0.79	-.12
<b>Aggression IRT score</b>	528	411	-0.04 (0.86)	0.46	-0.40	-.13

<sup>a</sup>There were an additional 4 incomplete twin pairs for these measures which were included in the analysis.

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Phenotype scaling in GxE

Phenotype scaling in GxE

**Table 5: GxM model fits for Well-being phenotype**

Model (freely estimated parameters)	-2LL	df	BIC	AIC	saBIC	DIC
Sum score						
$a_C, c_C, e_C, \alpha_C, \gamma_C, \varepsilon_C, \alpha_U, \gamma_U, \varepsilon_U$	10204.50	3727	-7653.07	2750.50	-1734.73	-4228.18
$a_C, c_C, e_C, \alpha_U, \gamma_U, \varepsilon_U$	10204.78	3730	-7663.18	2744.80	-1740.08	-4235.54
$a_C, c_C, e_C, \gamma_U, \varepsilon_U$	10209.33	3731	-7664.34	2747.33	-1739.65	-4335.78
$a_C, c_C, e_C, \alpha_U$	10206.10	3732	-7669.38	2742.09	-1743.11	-4239.91
$a_C, c_C, e_C$	10222.75	3733	-7664.47	2756.75	-1736.61	-4234.08
Transformed sum score						
$a_C, c_C, e_C, \alpha_C, \gamma_C, \varepsilon_C, \alpha_U, \gamma_U, \varepsilon_U$	10214.25	3727	-7648.19	2760.25	-1729.85	-4223.30
$a_C, c_C, e_C, \alpha_U, \gamma_U, \varepsilon_U$	10214.92	3730	-7658.12	2754.92	-1735.02	-4230.48
$a_C, c_C, e_C, \gamma_U, \varepsilon_U$	10216.73	3731	-7660.64	2754.73	-1735.95	-4332.08
$a_C, c_C, e_C, \alpha_U$	10215.09	3732	-7664.88	2751.09	-1738.60	-4235.40
$a_C, c_C, e_C$	10219.96	3733	-7665.87	2753.96	-1738.00	-4235.47
IRT score						
$a_C, c_C, e_C, \alpha_C, \gamma_C, \varepsilon_C, \alpha_U, \gamma_U, \varepsilon_U$	9806.21	3739	-7893.28	2328.21	-1955.88	-4457.37
$a_C, c_C, e_C, \alpha_U, \gamma_U, \varepsilon_U$	9806.89	3742	-7903.21	2322.89	-1961.05	-4464.54
$a_C, c_C, e_C, \gamma_U, \varepsilon_U$	9808.22	3743	-7905.96	2322.22	-1962.22	-4466.38
$a_C, c_C, e_C, \alpha_U$	9807.08	3744	-7909.96	2319.09	-1964.62	-4469.45
$a_C, c_C, e_C$	9810.82	3745	-7911.51	2320.82	-1964.59	-4470.08

Phenotype scaling in GxE

**Table 6: GxM model fits for Aggression phenotype**

Model (freely estimated parameters)	-2LL	df	BIC	AIC	saBIC	DIC
Sum score						
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, a<sub>C</sub>, γ<sub>C</sub>, ε<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	10218.91	3732	-7662.97	2754.91	-1736.69	-4233.49
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	10222.38	3735	-7671.51	2752.38	-1740.46	-4239.27
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub></i>	10232.63	3736	-7669.80	2760.63	-1737.17	-4236.65
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, ε<sub>U</sub></i>	10224.28	3737	-7677.40	2750.28	-1743.18	-4243.33
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub></i>	10240.40	3738	-7672.76	2764.40	-1736.96	-4237.77
Transformed sum score						
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, a<sub>C</sub>, γ<sub>C</sub>, ε<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	10228.85	3732	-7658.00	2764.85	-1731.72	-4228.52
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	10232.34	3735	-7666.52	2762.34	-1735.48	-4234.29
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub></i>	10234.20	3736	-7669.01	2762.20	-1736.38	-4235.86
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, ε<sub>U</sub></i>	10234.73	3737	-7672.17	2760.73	-1737.96	-4238.10
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub></i>	10238.00	3738	-7673.96	2762.00	-1738.16	-4238.97
IRT score						
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, a<sub>C</sub>, γ<sub>C</sub>, ε<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	9676.16	3739	-7958.30	2198.16	-2020.91	-4522.39
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub>, ε<sub>U</sub></i>	9679.97	3742	-7966.67	2195.97	-2024.51	-4528.00
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, α<sub>U</sub>, γ<sub>U</sub></i>	9682.29	3743	-7968.93	2196.29	-2025.18	-4529.34
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub>, ε<sub>U</sub></i>	9682.21	3744	-7972.39	2194.21	-2027.06	-4531.88
<i>a<sub>C</sub>, c<sub>C</sub>, e<sub>C</sub></i>	9687.08	3745	-7973.38	2197.08	-2026.46	-4531.95

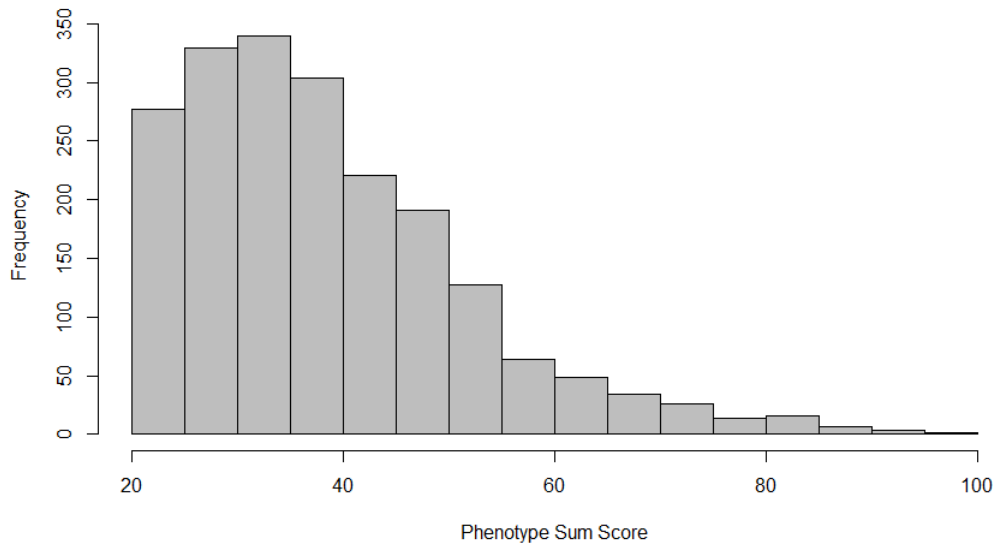
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## Phenotype scaling in GxE

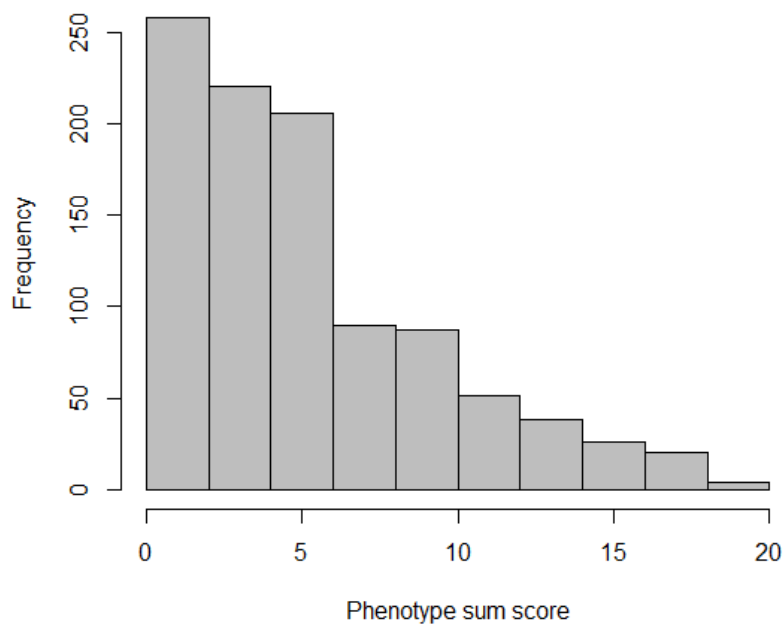
**Table 7: Parameter estimates from best-fitting models for Well-being and Aggression phenotypes**

Phenotype		GxM Parameter Estimates					
Phenotypic Proxy	Correlation with moderator	$\alpha_C$	$\alpha_U$	$\gamma_C$	$\gamma_U$	$\epsilon_C$	$\epsilon_U$
<b>Well-being</b>							
Sum score	.18	0 (fixed)	-.11	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)
Transformed sum score	.19	0 (fixed)	-.06	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)
IRT factor score	.18	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)
<b>Aggression</b>							
Sum score	-.12	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	-0.07
Transformed sum score	-.12	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)
IRT factor score	-.13	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	0 (fixed)	-0.03

**Figures**



**Figure 1**



**Figure 2**

Phenotype scaling in GxE

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

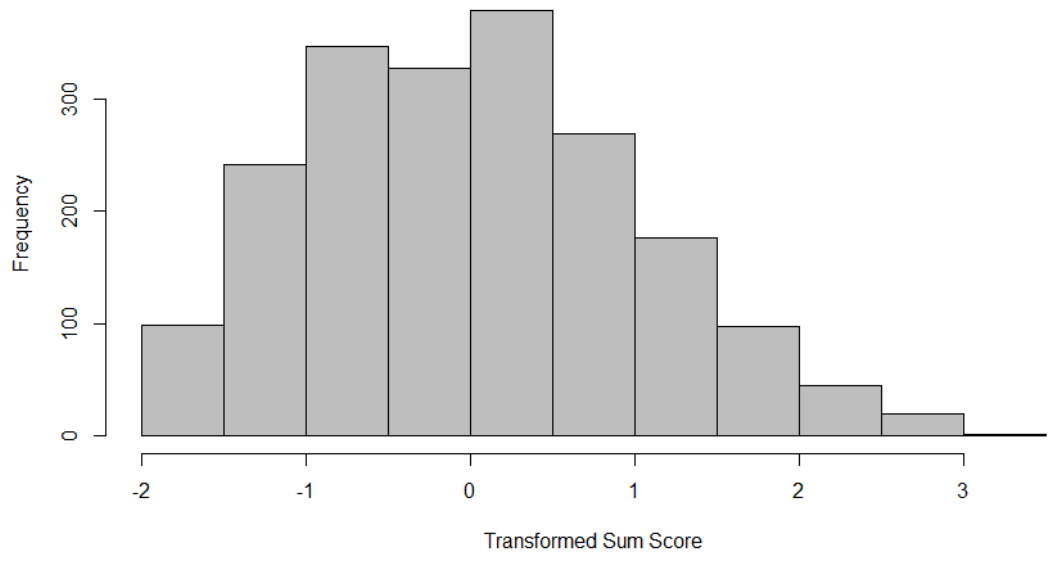


Figure 3

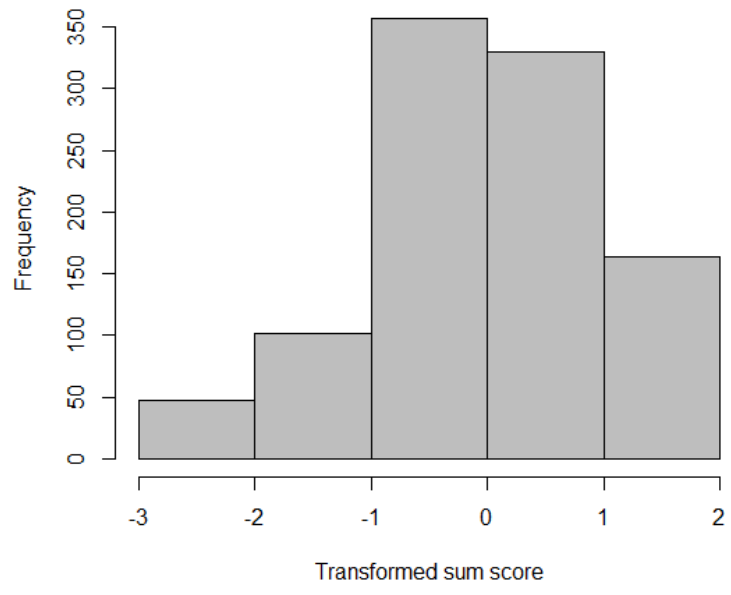
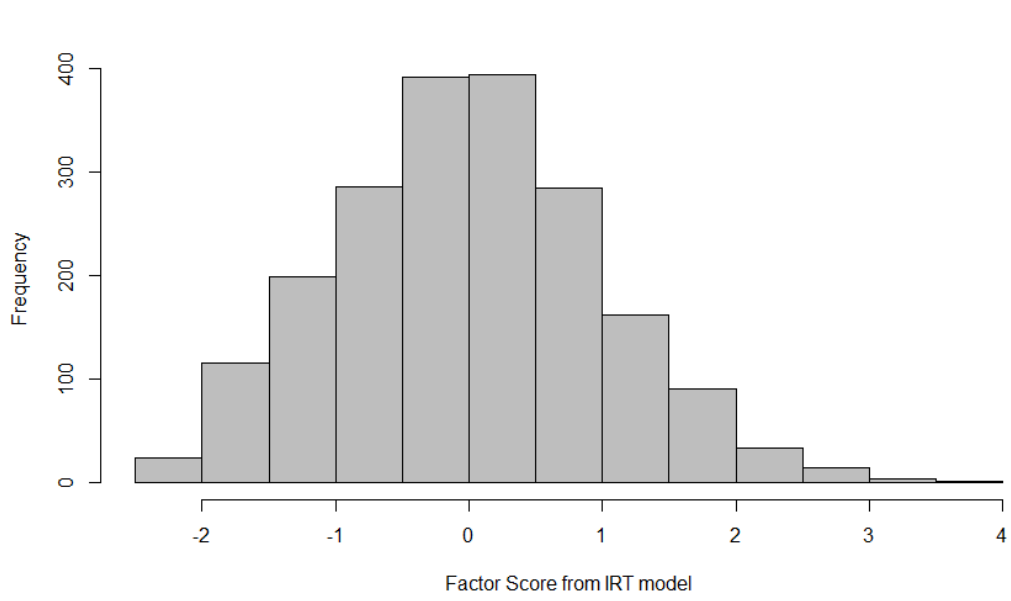


Figure 4

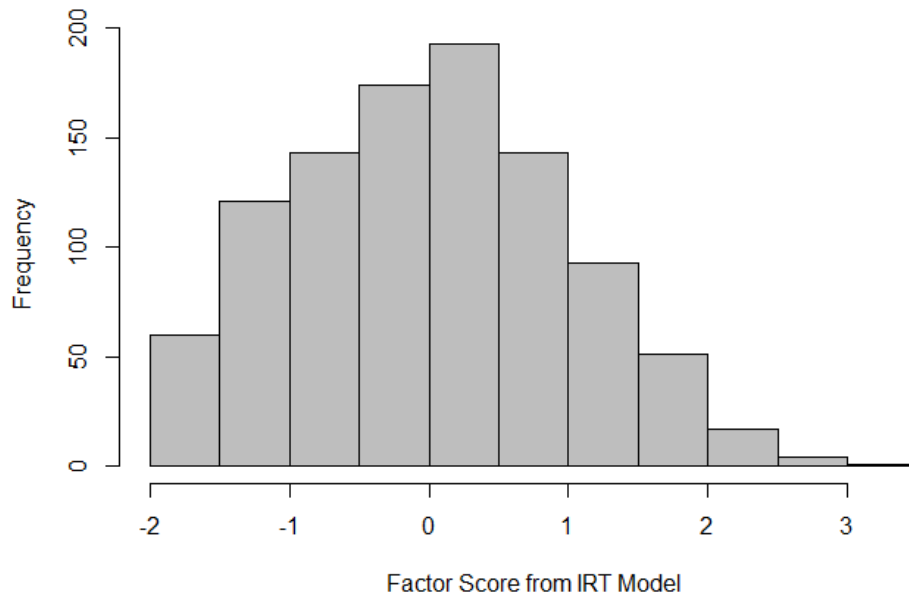


Phenotype scaling in GxE



23  
24  
25  
26  
27  
28  
29  
30  
31  
32

**Figure 5**



55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Figure 6**

**Figure Captions**

**Figure 1**

**Histogram showing the distribution of the sum score derived from generating item level data according to Eq. 3 with parameters in Table 1 (polytomous).**

**Figure 2**

**Histogram showing the distribution of the sum score derived from generating item level data according to Eq. 4 with parameters in Table 1 (binary).**

**Figure 3**

**Histogram showing the distribution of the transformed sum score derived from generating item level data according to Eq. 3 with parameters in Table 1 (polytomous) and then applying a  $\log_{10}$  transformation.**

**Figure 4**

**Histogram showing the distribution of the transformed sum score derived from generating item level data according to Eq. 4 with parameters in Table 1 (binary) and then applying a  $\log_{10}$  transformation.**

**Figure 5**

**Histogram showing the approximate distribution of factor scores derived from generating item level data according to Eq. 3 with parameters in Table 1 (polytomous), fitting a graded response model, and then obtaining factor scores based on this model.**

**Figure 6**

**Histogram showing the approximate distribution of factor scores derived from generating item level data according to Eq. 4 with parameters in Table 1 (binary), fitting a 2PL, and then obtaining factor scores based on this model.**