

Power of IRT in GWAS: Successful QTL mapping of sum score phenotypes depends on interplay between risk allele frequency, variance explained by the risk allele, and test characteristics

Journal:	<i>Genetic Epidemiology</i>
Manuscript ID:	GenEpi-12-0129
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	13-Jul-2012
Complete List of Authors:	van den Berg, Stephanie; University of Twente, OMD Service, Susan; UCLA, Center for Neurobehavioral Genetics
Key Words:	IRT, measurement, statistical power

SCHOLARONE™
Manuscripts
Review

1
2
3
4
5
6 Power of IRT in GWAS: Successful QTL mapping of sum score phenotypes depends on interplay
7
8 between risk allele frequency, variance explained by the risk allele, and test characteristics
9
10

11
12
13
14
15 Stéphanie M. van den Berg¹ and Susan K. Service²
16

17
18
19 ¹ University of Twente, The Netherlands, Department of Research Methodology, Measurement
20
21 and Data Analysis
22

23
24 ² University of California, Los Angeles, Center for Neurobehavioral Genetics
25
26
27

28
29 Running title: QTL mapping of sum score phenotypes
30
31
32
33
34
35
36
37

38 Correspondence: Dr. Stéphanie M. van den Berg
39

40 University of Twente, Faculty of Behavioral Sciences
41

42 Department of Research Methodology, Measurement and Data Analysis (OMD)
43

44 P.O. Box 217, 7500 AE Enschede, The Netherlands
45

46 Telephone: +31 (0)53 489 2422 (office) 3616 (secr.)
47
48

49 E-mail: stephanie.vandenberg@utwente.nl
50
51
52
53
54
55
56
57
58
59
60

Abstract

As data from sequencing studies in humans accumulate, rare genetic variants influencing liability to disease and disorders are expected to be identified. Three simulation studies show that characteristics and properties of diagnostic instruments interact with risk allele frequency to affect the power to detect a QTL based on a test score derived from symptom counts or questionnaire items. Clinical tests, that is, tests that show a positively skewed phenotypic sum score distribution in the general population, are optimal to find rare risk alleles of large effect. Tests that show a negatively skewed sum score distribution are optimal to find rare protective alleles of large effect. For alleles of small effect, tests with normally distributed item parameters give best power for a wide range of allele frequencies. The item-response theory (IRT) framework can help understand why an existing measurement instrument has more power to detect risk alleles with either low or high frequency, or both kinds.

Key words: Item-Response Theory; IRT; measurement; statistical power; extreme samples design; case-control design; population sample design

Introduction

Many diagnostic instruments for a disorder consist of symptom counts. Often the disorder can be seen as the extreme tail of a continuous liability trait: the higher the liability, the more likely a subject shows certain symptoms. High liability persons will have many symptoms and on the basis of a diagnostic criterion are then labeled as affected with the disorder of interest. In some instances, GWAS studies are applied on symptom count data, rather than diagnosis, as this allows using more information [Van der Sluis et al., 2012].

Item-response theory (IRT) provides a formal statistical framework for modeling liability and diagnosis, and its application in the medical sciences is increasing [Reise & Waller, 2009]. IRT has also been successfully applied in genetics [Eaves et al., 2005 ; Van den Berg, Glas, & Boomsma 2007; Van den Berg et al., 2010]. It provides a useful framework for understanding the relationship between measurement problems and problems in detecting genetic variants. For example, using this IRT framework, Van der Sluis et al [2010] showed that ignored multidimensionality, measurement bias, and poor reliability can result in poor statistical power in QTL mapping studies.

Here we show how power is associated with test characteristics using different study designs, and how this association is moderated by allele frequency. We link the simulation results to the IRT concept of 'test information'. We start out with a brief introduction to IRT and so-called test information functions (TIFs). Next, we describe how this framework makes predictions about statistical power in QTL mapping. Three simulation studies demonstrate the intricate relationship between study design, allele frequency and the test information function.

Item-Response Theory

1
2
3 Item response theory (IRT) models item data as a function of both item characteristics as well as
4 person characteristics [Lord & Novick, 1968; Lord, 1980; Embretson & Reise, 2000]. An item can
5 be anything from a symptom that is scored in a diagnostic interview as being either present or
6 absent, or an item on a self-report questionnaire that can be answered with yes or no. Items do
7 not have to be dichotomous (i.e., yes/no, or 1/0), but for clarity of exposition we focus on
8 dichotomous items in our descriptions. In the Discussion we expand on alternative data types.
9
10
11
12
13
14
15

16
17 The one-parameter logistic IRT model for dichotomous items, or so-called Rasch model,
18 is
19
20
21

$$22 \quad P(X_{ij} = 1 | \theta_i, \beta_j) = \frac{1}{1 + \exp(\beta_j - \theta_i)} \quad \text{Equation 1}$$

23
24
25
26
27
28
29

30 where $P(X_{ij}=1)$ is the probability of a positive response for person i on item j (or the presence of
31 symptom j). Parameter θ_i is the person parameter for person i and can be thought of as person
32 i 's liability for a disorder. Parameter β_j is the item parameter for item j . In educational
33 measurement, this β parameter is usually called the difficulty parameter, where it refers to the
34 difficulty of a cognitive test item. A high β value indicates a difficult item, with an overall low
35 probability of making the item correct; only high ability individuals have a reasonable chance of
36 giving a correct answer. Analogously, for questionnaire and clinical data, a high value for β_j
37 indicates a high threshold for a positive response on a questionnaire, or a high threshold for a
38 symptom. A high β value therefore corresponds to low symptom prevalence, as the symptom is
39 not endorsed by many persons. As can be seen in Equation 1, when $\theta_i = \beta_j$, the probability of a
40 positive response is 50%. When $\theta_i > \beta_j$, the probability of a positive response is higher than the
41 probability of a negative response, and vice versa for $\theta_i < \beta_j$. In order to make the modeling
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 statistically identifiable, it is usually assumed that the population mean for parameter θ equals
4
5
6 0.

7
8 The logistic curve for $P(X_{ij}=1)$ in Eq. 1 is sigmoidal (see Figure 1) and has its maximum
9
10 slope at $\theta_i = \beta_j$; at this point the change in $P(X_{ij}=1)$ as a function of θ is at its maximum, rendering
11
12 maximum discrimination between those individuals with θ values below β (more likely to have a
13
14 negative response) and those individuals with θ values above β (more likely to have a positive
15
16 response). Thus, in Figure 1, the slope of the left curve is at its maximum at $\theta = -1.5$, whereas the
17
18 right curve has its maximum slope at $\theta = 1.5$. Discriminatory power of an item is at its lowest at
19
20 θ values very far removed from β ; for example, in the curve on the left in Figure 1 ($\beta = -1.5$),
21
22 individuals with high θ values will all have a probability of nearly one for a positive response; but
23
24 individuals with even higher θ values have about the same probability. All are very likely to have
25
26 the same positive response. An item j therefore yields very little information about individuals
27
28 and individual differences at θ levels far removed from β_j .

29
30 The information that an item j with known β_j gives about a trait θ is therefore a function
31
32 of θ . For the one-parameter model, the Fisher information from item j is

$$I_j(\theta) = P_j(\theta)(1 - P_j(\theta)) \quad \text{Equation 2}$$

33
34
35
36
37
38
39
40
41
42
43
44
45
46 where $P_j(\theta)$ is the probability of a positive response for item j for a person with liability $=\theta$. The
47
48 one-parameter model can be extended to a two-parameter model, where the slope of the item
49
50 response curve can vary over items, and by consequence the peaked-ness of the item
51
52 information functions, but for clarity of exposition we focus on the 1-parameter logistic IRT
53
54 model of Eq. 1. In the Discussion we expand to the 2-parameter model.
55
56
57
58
59
60

The information provided by all test items on trait θ combine additively to the overall shape of the test information function (TIF),

$$TIF(\theta) = \sum_{j=1}^k I_j(\theta) \quad \text{Equation 3}$$

and with $[TIF(\theta)]^{-1/2}$ we have a formula for the standard error of measurement (SEM; Lord 1980). Thus, at θ levels where test information is high, the SEM for the maximum likelihood estimate is small, so that we have high measurement precision. At those levels, it is possible to discriminate between individuals of different levels. With high information content, a small difference in θ level is associated with a difference in response pattern (more or fewer endorsed items, i.e. a different phenotypic sum score). At θ levels with little information, slight differences in θ do not result in different response patterns (sum scores), so there is no information to discriminate among individuals. Figure 2 shows a number of test information functions (top row), with varying distributions of the β values for the items. For instance, for a test with β values uniformly distributed between 1 and 3 (Test scenario 'Right': all items have low endorsement probabilities, therefore high β values), information reaches its maximum on the right-hand side of the distribution. Persons with high liability will therefore show variation in their sum scores, but not the persons with low liability.

Relationship between test information and distribution of sum scores

The test information function (TIF) is directly related to the shape of the distribution of sum scores. If the TIF is symmetric and centered around the average liability value of the population, the resulting distribution of sum scores will be symmetrical, too. If however the TIF has its

1
2
3 maximum at the right-hand side of the scale, that is, for above-average levels of θ , the expected
4 distribution of sum scores will be positively skewed. This is observed, for example, when clinical
5 tests are used in population samples: most of the items on clinical tests or diagnostic
6 instruments will have low response rates, as they relate to symptoms that less than half the
7 normal population are likely to exhibit. For such items, the probability of a positive response will
8 exceed 50% only for individuals with high values on the liability trait. The β values for such items
9 will therefore also be highly positively valued (Eq 1), with a resulting TIF that has its maximum at
10 a relatively high θ level. All individuals with average and below-average liability values will have
11 low to very low probability of symptoms, so many will have a sum score of 0 and there is no
12 further distinction possible among them. On the other hand, there will be quite some variation
13 in number of symptoms for above average trait values. As a result we see a skewed distribution
14 of symptom counts. So in short, for a given population with the average trait level defined as $\theta=$
15 0, we expect a positively skewed sum score distribution when most of the β values are positive,
16 a negatively skewed sum score distribution when most of the β values are negative, and more
17 symmetrical when the β values are scattered across the continuum. Figure 2 shows a number of
18 example TIFs and the typical sum score distributions that result from these (based on simulated
19 data for a 20-item test). The labels 'Right', 'Left', 'Uniform', and 'Normal' refer to the
20 distribution of the values of the β parameters in a test: 'Right' $\beta \sim U(1, 3)$, 'Left' $\beta \sim U(-3, -1)$,
21 'Uniform' $\beta \sim U(-3, 3)$ and 'Normal' $\beta \sim N(0, 1)$. Observe in Figure 2, that if the distribution of the
22 β parameters is uniform on the $[-3, 3]$ continuum ('Uniform'), the resulting sum score
23 distribution shows lower variance than if the β parameters are standard normally distributed
24 ('Normal').

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 Test information and power in different gene-mapping study designs
58
59
60

1
2
3 As with any mapping study, the power to find a QTL for a liability trait crucially depends on
4
5 having a data set where there is both genotypic variation and phenotypic variation. More
6
7 specifically, in the case of sum score phenotypes: there should be sum score variation among
8
9 those individuals that have different genotypes. If all heterozygotes Aa have the same sum score
10
11 or the same case-control status as homozygotes AA, there would be lower power to detect a
12
13 QTL than a situation where heterozygotes have a different average sum score than
14
15 homozygotes. If a risk allele has a strong phenotypic effect, one would expect that there would
16
17 be AA, Aa and aa individuals that have above average liability, assuming an additive model of
18
19 action of the A allele. The below average liability individuals will be mostly aa (see Figure 3). The
20
21 interesting genotypic variation therefore exists at the right-hand side of the liability scale. It is
22
23 then important to have good measurement resolution at that part of the scale: one would want
24
25 the three different genotypes to have different average sum scores. Thus, as a general rule, we
26
27 expect that the power to detect a QTL for a risk allele with frequency smaller than 0.5 and with a
28
29 large effect is best for tests that have a lot of test information at high trait levels. Conversely, we
30
31 expect that the power to detect a QTL for a risk allele with a high frequency (>.5; that is, rare
32
33 *protective* alleles) and of large effect is best for tests that have high levels of test information at
34
35 low trait levels. Following the same logic, studying the expected liability distributions of different
36
37 genotypes, we expect that alleles of small effect and/or alleles with a frequency of around 0.5
38
39 will generally require high levels of test information in the middle of the liability distribution.
40
41
42
43
44
45
46

47 This paper presents simulation studies that illustrate the tight relationship between test
48
49 information, the risk allele frequency of QTLs, the liability variance explained by a QTL, and the
50
51 statistical power to find such QTLs in a gene mapping study. Three different research designs are
52
53 studied: a population study where subjects are randomly sampled from the population (Study
54
55 1), an extreme samples design, where subjects are sampled from those with only very high or
56
57
58
59
60

1
2
3 very low sum scores (Study 2), and a case-control design (Study 3). We explore power in
4
5 situations where QTLs explain little liability variance (1%), and contrast these findings to the
6
7 situation where QTLs explain a large part of liability variance (25%). This contrast with 25%
8
9 variance explained serves to illustrate the main principles, rather than being a realistic setting
10
11 for gene mapping studies.
12
13

14 15 16 17 **Methods**

18
19
20
21
22 Number of test items and sample sizes were varied across simulation studies in order to
23
24 produce a wide range in power estimates to make patterns in the results as clear as possible.

25
26 Power was defined as the percentage of simulation replicates to be significant at the 0.01 alpha
27
28 level. The number of replicates for each simulation situation was 1,000 in all studies.
29
30

31 32 33 Simulation study 1. Population sample design

34
35 For each test information scenario ('Left', 'Right', 'Uniform' and 'Normal'), data sets were
36
37 simulated under the assumption of a SNP in perfect LD with a QTL that explains 1% of the
38
39 variance in a standard normally distributed quantitative liability trait θ under additive gene
40
41 action. That is, we simulated θ under a linear model where the expectation was equal to an
42
43 allelic effect of the risk allele times the number of risk alleles, and a normally distributed random
44
45 term with variance proportional to the variance due to the allelic effect. Theta values were then
46
47 rescaled to mean 0 and variance 1. The frequency of the risk allele (the allele that increases
48
49 liability) was varied: .001, .005, .01, .05, .10, .25, .50, .75, .90, .95, .99, .995 and .999. In each
50
51 replicated data set, genotypes for 1,000 individuals were simulated assuming Hardy-Weinberg
52
53 equilibrium. Theta (θ) values were standardized and used for simulating the data for a 20-item
54
55
56
57
58
59
60

1
2
3 test under a one-parameter logistic IRT model (Eq 1). For each individual, the item scores were
4
5 summed, and the sum score phenotypes were linearly regressed on the number of risk alleles to
6
7 study power.
8
9

10 Note that the proportion of variance in liability explained by the QTL was fixed to 1%,
11
12 but not the proportion of variance in the observed phenotypic sum scores. The proportion of
13
14 variance in the sum scores explained by the QTL is dependent on test information, which in turn
15
16 depends on the β parameter values for the items in the test. With a limited number of items,
17
18 the proportion of sum score variance explained is always lower than the proportion of liability
19
20 variance explained [Van den Berg, Glas, & Boomsma, 2007].
21
22
23

24 Under the test information scenario referred to as 'Left', the item parameters β were
25
26 randomly sampled from the uniform distribution $U(-3, -1)$, and under the scenario referred to as
27
28 'Right' using $U(1, 3)$. Under the scenario 'Uniform', the item parameters were sampled using $U(-$
29
30 $3, 3)$, and under the 'Normal' scenario, they were sampled from the standard normal
31
32 distribution, $N(0, 1)$. Figure 2 shows the corresponding test information functions.
33
34
35

36 Because low allele frequencies lead to very low counts for certain genotypes, the
37
38 assumptions underlying standard statistical testing might not be met. In practice, if only 1 or 2
39
40 individuals are homozygous for one allele, one might prefer an empirical p -value over the p -
41
42 value resulting from a standard t -test for the linear regression on genotype. Moreover, for 'Left'
43
44 and 'Right' test scenarios, the sum score distribution is very skewed, also violating the
45
46 assumptions of standard linear regression. Therefore, all p -values were empirically determined,
47
48 based on 400 permutations per data set.
49
50

51
52 In a second series of simulations for Study 1, the QTL explained 25% of the liability (θ)
53
54 variance. This value was chosen as an extreme contrast to the 1% simulations. Power was based
55
56
57
58
59
60

1
2
3 on simulating data for 50 individuals on a 10-item test. All other settings were the same as in the
4
5 first series.
6
7
8
9

10 Simulation study 2. Extreme samples design

11
12 The same test information scenarios were used as in Simulation Study 1. In a first series of
13
14 simulations, the variance of the liability explained by the QTL was fixed to 1%. In each test
15
16 information scenario and in each simulation, sum score data were simulated for 10,000 subjects
17
18 on 20 items. Two-hundred subjects were randomly sampled from those subjects with the lowest
19
20 observed sum score (with replacement), and 200 subjects were sampled from those with the
21
22 highest observed sum score. In each simulated data set, Fisher's exact test was performed on
23
24 the cross-tabulation of SNP genotype and 'extreme high'/'extreme low' status. One thousand
25
26 data sets were simulated for each condition of risk allele frequency: .005, .01, .05, 0.10, .25, .50,
27
28 .75, .90, .95, .99, and .995.
29
30
31
32

33 In a second series of simulations for Study 2, serving as a contrasting illustration, the
34
35 variance in liability explained by the QTL was fixed to 25%. All other settings were the same as in
36
37 the first series, except that now only 10 extremely low scoring and 10 extremely high scoring
38
39 subjects were sampled from the population of 10,000 subjects that had sum score data based
40
41 on 5 items.
42
43
44
45
46

47 Simulation study 3. Case-control design

48
49 The same test information scenarios were used as in Simulation Studies 1 and 2. In a first series
50
51 of simulations, the variance of the liability explained by the QTL was fixed to 1%. In each test
52
53 information scenario and in each simulation, sum score data were simulated for 10,000 subjects
54
55 on 20 items. Six-hundred controls were randomly sampled from those subjects with
56
57
58
59
60

1
2
3 subthreshold observed sum score (i.e., lower than the cut-off score for diagnosis; sampling with
4 replacement), and 600 cases were sampled from those with the cut-off sum score for diagnosis
5 or higher (with replacement). The cut-off score was defined as the 88th percentile based on the
6 observed sum scores in the simulated samples of 10,000 individuals (cf. Van der Sluis et al.,
7 2012). In each simulated data set, Fisher's exact test was performed on the cross-tabulation of
8 SNP genotype and case-control status. One thousand data sets were simulated for each
9 condition of risk allele frequency: .001, .005, .01, .05, .10, .25, .50, .75, .90, .95, .99, .995 and
10 .999.
11
12
13
14
15
16
17
18
19
20
21

22 In a second series of simulations for Study 3, serving as a contrasting illustration, the
23 variance in liability explained by the QTL was fixed to 25%. All other settings were the same as in
24 the first series, except that now 60 cases and 60 controls were sampled from the population of
25 10,000 subjects that had sum score data based on 20 items.
26
27
28
29
30
31
32

33 Results

34 Simulation study 1. Population sample design

35
36
37
38 Power to detect a QTL that accounts for a low percentage of trait variance is best detected
39 when the trait is measured using a test in which the β parameters are normally distributed,
40 except for extreme allele frequencies (Figure 4). Such a test with normally distributed β
41 parameters clearly outperforms a test where item parameters are uniformly distributed. Low
42 power is observed for detecting rare risk alleles for 'Left' tests, where item β parameter values
43 are all negative. The same low power is observed for rare protective alleles for 'Right' tests with
44 only positive items, that is, a clinical test that discriminates only among high-scoring individuals.
45
46
47
48 For very low frequencies, clinical tests seem optimal, outperforming the power of a test with
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 normally distributed item parameters. For very high frequencies (i.e., very rare protective
4 alleles), 'Left' tests outperform 'Normal' tests.
5
6

7
8 When the variance of liability explained is 25%, however, 'Left' and 'Right' type of tests
9
10 outperform 'Normal' tests, depending on allele frequency (Figure 4). For risk allele frequencies
11
12 lower than 0.25, 'Right' tests give best power; for risk allele frequencies higher than 0.75, 'Left'
13
14 tests give best power. For risk allele frequencies between 0.4 and 0.6, 'Normal' tests seem to
15
16 give best power.
17
18

19 20 21 Simulation study 2. Extreme samples design 22

23
24 When the percentage variance explained by the QTL is 1%, lowest power is observed for 'Left'
25
26 and 'Right' types of tests, irrespective of risk allele frequency (Figure 5). Overall best power is
27
28 shown by tests with normally distributed β parameters. In contrast, when percentage variance
29
30 explained by the QTL is 25%, the 'Left' and 'Right' types of tests outperform the other two: 'Left'
31
32 tests are optimal for risk allele frequencies above 0.5, and 'Right' tests are optimal for risk allele
33
34 frequencies below 0.5.
35
36

37 38 39 Simulation study 3. Case-control design 40

41
42 When the percentage variance explained by the QTL is 1%, lowest power is observed for 'Left'
43
44 types of tests, irrespective of risk allele frequency (Figure 6). Overall, 'Right' types of tests show
45
46 highest power, while tests with normally distributed β parameters are a close second. When
47
48 percentage variance explained by the QTL is 25%, power is again lowest for 'Left' types of tests.
49
50 In addition, there is a decrease in power with increasing frequency of the risk allele for allele
51
52 frequencies larger than 0.10.
53
54
55
56
57
58
59
60

Discussion

In educational measurement, the definition of the 'best test' is the test that minimizes measurement error over the target of measurement (Wright & Stone, 1979). For instance, when constructing a law school admission test, one would want to have maximum test information at the point on the scale where the threshold for sufficient aptitude is located; one would want to have confidence that the decision of pass/fail is reliable. There is less interest in reliably quantifying individual differences at the lower end of the scale. Translating this concept of 'optimal test design' to gene finding studies leads to constructing or searching for a test that discriminates best among those subjects with different genotypes for the type of QTL one hopes to find.

Power to detect a QTL naturally depends both on variation in genotypes and variation in liability, but more crucially, it depends on whether the measurement tool that is used to assess liability discriminates among the genotypes. Firstly, power depends on how much genotypic variance there is and also where it is: as Figure 3 shows, allele frequency and variance explained determine in unison where and how much genetic variance there is: allele frequency determines the relative surface areas of the three histograms, whereas variance explained (actually, allelic effect) determines how far apart the means of the three histograms are. Together, allele frequency and variance explained determine whether most of the genotypic variation is among high liability individuals, low liability individuals or individuals of average liability. Secondly, the test information function (TIF) shows how well a test discriminates between individuals, and where on the scale on the scale it does so more clearly. Optimally therefore, the location of the maximum of the TIF on the liability scale should be exactly there where there is the most variation in genotypes. In the case of Figure 3, for example, a test that has a TIF that has its

1
2
3 maximum value on the far left-hand side of the distribution, with very low information content
4
5 on the right-hand side of the distribution, will result in extremely low power to find a QTL with
6
7 risk allele frequency 0.15, even when variance explained is 25%: sum scores would
8
9 phenotypically only discriminate among individuals that were homozygous for the protective
10
11 allele (i.e., they would vary in their sum score), whereas most heterozygotes and homozygotes
12
13 for the risk allele would have the same maximum sum score (cf. ceiling effect).
14
15

16
17 This is exactly what was seen in the simulations, both for a population design and an
18
19 extreme-samples design: even with highly influential risk alleles (25% variance explained), the
20
21 power to detect them can be dramatically low when the wrong phenotypic measurement tool is
22
23 used. In contrast, for case-control designs this effect was less dramatic.
24
25

26
27 In the case of QTLs that explain only little of the liability variance, we see all three
28
29 genotypes throughout the phenotypic range; all genotypes are scattered across the entire scale,
30
31 with all three genotype means close to 0. For maximum power, the information function should
32
33 therefore also be more spread out across the continuum, but also be high where most of the
34
35 individuals are. A test that is highly informative at scale locations where only few of the
36
37 individuals are located generates little statistical power. Our simulations clearly showed that the
38
39 optimal TIF for a wide range of allele frequencies and variance explained was where the TIF is
40
41 the result of item parameters that are standard normally distributed. As can be seen in Figure 2,
42
43 the corresponding sum score distribution then also shows the largest variance.
44
45

46
47 Even in a case-control design (Study 3), tests with normally distributed item parameters
48
49 performed very well, and approached the power of clinical tests. Nevertheless, a test that
50
51 maximizes the discrimination between cases and controls is the best tool for a case-control
52
53 study, independent of risk allele frequency or variance explained.
54
55
56
57
58
59
60

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

Thus, overall when looking for alleles of small effect, tests with normally distributed item parameters give best power for a wide range of allele frequencies, and for different study designs. Only in the case of a population sample design and very extreme allele frequencies (below 0.01, above 0.99), such tests are outperformed by 'Right' and 'Left' tests, respectively. In case-control designs, 'Right' types of tests are slightly more powerful than 'Normal' types of tests. In the test design phase, therefore, care should be devoted to choosing the optimal diagnostic test or self-report questionnaire. The test information function of the phenotyping instrument can be used as an aid for the identification of potential measurement problems. It can be plotted using standard IRT software packages such as Multilog [Thissen, Chen & Bock, 2003], or the ltm library in R [Rizopoulos, 2006] and compared to the test information functions in Figure 2. When interest is in finding QTLs that explain only a small amount of variance of a wide range of allele frequencies, the information function of the phenotyping instrument should ideally look like the 'Normal' information function in Figure 2.

As all results shown here are directly related to the shape of the Test Information Function (TIF), the results are easily extended to two-parameter IRT models (where items may have different factor loadings or discrimination parameters), and to IRT models for polytomous items, where instead of yes/no or absent/present, the data consist of scales with multiple categories like, for instance, 'never', 'sometimes', 'often', and 'always'. Such IRT models can be fitted using standard software after which the TIF can be plotted. These TIFs are directly comparable to those based on the 1-parameter IRT model for dichotomous data as used in the present simulations.

Usually in genome-wide association studies, SNPs with minor allele frequencies (MAF) lower than 1-2% are not included in analyses. In these simulations, we varied allele frequency in the population. Focus on low-frequency alleles is increasing with the increased interest in re-

1
2
3 sequencing studies. In case-control and extreme samples designs, sample frequencies of risk
4 alleles will be higher than their population values. Given the right design and the right sort of
5 phenotypic test, studies can be very powerful to detect allelic effects with low MAF, see for
6
7
8
9
10 example Figures 4 and 5, where high power was observed for MAFs below 1%.

11
12 The results reported here may also partly explain why relatively few QTLs of small effect
13 have been identified for disorders measured using clinical tests [e.g., Manolio et al., 2009]. Thus
14 far, the reasons mentioned in the literature for the lack of success include mainly limited sample
15 sizes [Sullivan, 2012], or reasons of genetic nature [e.g., Crow 2011]. Relatively little attention
16 has been paid to psychometric properties of phenotypic measures (but see Van der Sluis et al.,
17 2010, 2012). This study clearly shows that genetic studies could be much more successful if
18 researchers selected their phenotypic instruments to suit their study aim. Rather than using the
19 phenotypic measurement that is the most valid clinical tool for diagnosis, one should choose the
20 phenotypic measurement tool that gives highest power to find genetic variants that explain
21 variation in liability. Changing the phenotypic instrument might be much cheaper than ever-
22 increasing sample sizes with the same clinical test. Such a change might be as trivial as
23 rewording an item, such as for example changing “I like chocolate very much” into “I like
24 chocolate”. Making the item less extreme will increase prevalence of positive responses,
25 resulting in lower item β parameter values and therefore a maximum of the information
26 function lower on the liability scale, with a corresponding change in statistical power. Simulation
27 study 2 shows that power to find QTLs of small effect in an extreme samples design can be
28 dramatically increased when changing from a clinical test to a test that shows more variability in
29 low-liability individuals, with the same number of items.

30
31 In addition to carefully choosing a measurement instrument for the phenotyping, test
32 information, and therefore power, can be further increased by simply adding items to a scale
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

1
2
3 (see Eq 3). One approach is to combine data from different phenotypic measurement
4
5 instruments to optimize the test information function for a measure of interest. For example,
6
7 Van den Berg et al [submitted] optimized the test information function for a clinical measure for
8
9 schizotypy by adding information from a self-report instrument. This was done through the use
10
11 of a multidimensional IRT measurement model, which resulted in increased measurement
12
13 precision, especially at the lower end of the liability continuum. As shown by the present results,
14
15 this should increase overall power to find a QTL of small effect. In conclusion, Item-Response
16
17 Theory is a useful framework to identify potential strengths and limitations of an existing
18
19 measurement instrument based on symptom counts or questionnaire items, and to remedy
20
21 potential problems through more sophisticated modeling of multivariate phenotypic data sets.
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

Acknowledgments

This work was made possible in part due to a grant from the University of Twente Incentive Fund awarded to SMvdB. The authors thank Sophie van der Sluis and an anonymous reviewer for helpful comments.

References

- 1
2
3
4
5
6
7
8 Crow TJ. 2011. The missing genes: what happened to the heritability of psychiatric disorders?
9
10 Mol Psych 16:362–364.
11
12
13 Eaves L, Erkanli A, Silberg J, Angold A, Maes HH, Foley D. 2005. Application of Bayesian inference
14
15 using Gibbs sampling to item-response theory modelling of multi-symptom genetic data.
16
17 Behav Genet 35:765–780.
18
19
20 Embretson ES, Reise SP. 2000. Item Response Theory for Psychologists. Lawrence Erlbaum,
21
22 Mahwah, NJ.
23
24 Lord FM. 1980. Applications of item response theory to practical testing problems. Erlbaum,
25
26 Hillsdale, NJ.
27
28 Lord FM, Novick MR. 1968. Statistical theories of mental test scores. Addison-Wesley, Reading,
29
30 UK.
31
32
33 Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, McCarthy MI, Ramos EM,
34
35 Cardon LR, Chakravarti A, Cho JH, Guttmacher AE, Kong A, Kruglyak L, Mardis E, Rotimi
36
37 CN, Slatkin M, Valle D, Whittemore AS, Boehnke M, Clark AG, Eichler EE, Gibson G,
38
39 Haines JL, Mackay TFC, McCarroll SA, Visscher PM. 2009. Finding the missing heritability
40
41 of complex diseases. Nature 461:747-753.
42
43
44
45 Reise SP, Waller NG. (2009) Item Response Theory and clinical measurement. Ann Rev Clin
46
47 Psychol 5:27-48.
48
49 Rizopoulos D. 2006. Irm: An R package for latent variable modelling and Item Response Theory
50
51 analyses, J Stat Software 17:1-25.
52
53
54 Sullivan P. 2012. Don't give up on GWAS. Mol Psych 17:2–3.
55
56
57
58
59
60

- 1
2
3 Thissen D, Chen WH, Bock RD. 2003. Multilog (version 7)[Computer software]. Scientific
4
5 Software International, Lincolnwood, IL.
6
7
8 Van den Berg SM, Fikse F, Arvelius P, Glas CAW, Strandberg E. 2010. Integrating phenotypic
9
10 measurement models with animal models. Proceedings of the 9th World Congress on
11
12 Genetics applied to Livestock Production, Leipzig, Germany, August 1–6 2010.
13
14
15 Van den Berg SM, Glas CAW, Boomsma DI. 2007. Variance decomposition using an IRT
16
17 measurement model. *Behav Genet* 37:604–616.
18
19
20 Van den Berg SM, Paap MCS, Derks EM, Genetic Risk and Outcome of Psychosis (GROUP)
21
22 investigators. Submitted. Using multidimensional modeling to combine self-report
23
24 symptoms with clinical judgment of schizotypy.
25
26
27 Van Leeuwen M, van den Berg SM, Boomsma DI. 2008. A twin-family study of general IQ.
28
29 *Learning and Individual Differences*, 18:76-88.
30
31
32 Van der Sluis S, Posthuma D, Nivard MG, Verhage M, Dolan CV. 2012, Epub ahead of print.
33
34 Power in GWAS: lifting the curse of the clinical cut-off. *Mol Psych* 22 May 2012, doi:
35
36 10.1038/mp.2012.65.
37
38
39 Van der Sluis S, Verhage M, Posthuma D, Dolan CV. 2010. Phenotypic complexity, measurement
40
41 bias, and poor phenotypic resolution contribute to the missing heritability problem in
42
43 genetic association studies. *Plos One* 5:e13929.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60 Wright BD, Stone MH. 1979. *Best Test Design-Rasch Measurement*. Chicago, IL: MESA Press.

Figure legends

Figure 1. Item response curves: The probability of a positive response plotted as a function of θ (see Equation 1), for two different values for β (left -1.5, right 1.5).

Figure 2. Test information functions for different types of tests, each consisting of 20 items, but different distributions of the items' β values (top row, see text for explanation of the labels.).

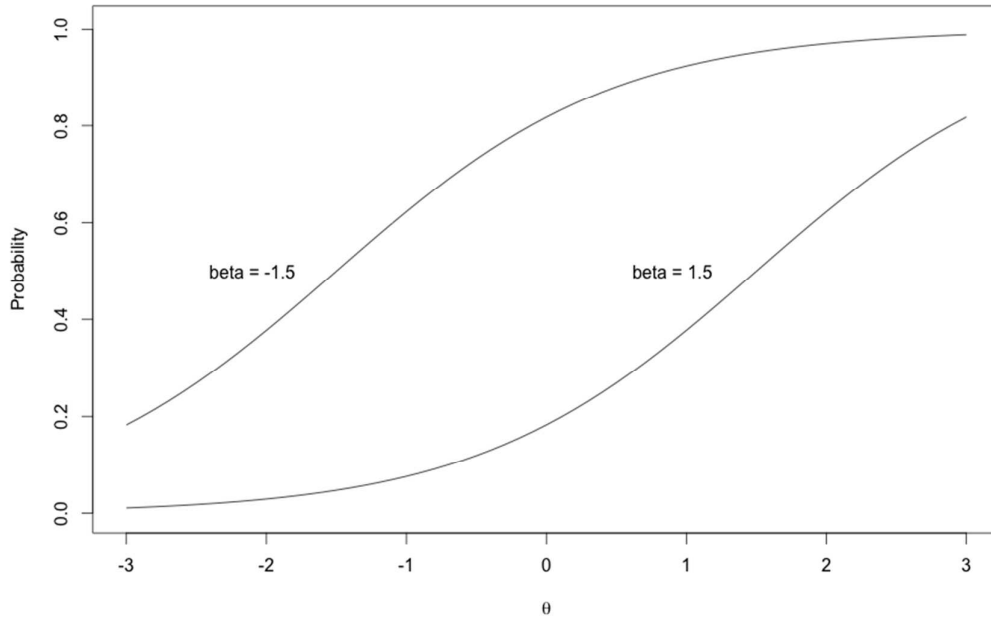
Bottom row shows corresponding distributions of sum scores when simulating data from 1000 individuals with ϑ values drawn from a standard normal distribution.

Figure 3. Superimposed histograms of simulated liability scores of aa genotypes (light gray), Aa genotypes, and AA genotypes (dark gray), with frequency of the A allele equal to 0.15 and variance explained equal to 25%, under Hardy-Weinberg and additive gene action.

Figure 4. Power to detect a QTL as a function of risk allele frequency and test characteristics, using a population sample design. Percentage liability variance explained by QTL is 1% (left panel), and 25% (right panel).

Figure 5. Power to detect a QTL as a function of risk allele frequency and test characteristics, using an extreme samples design. Percentage liability variance explained by QTL is 1% (left panel), and 25% (right panel).

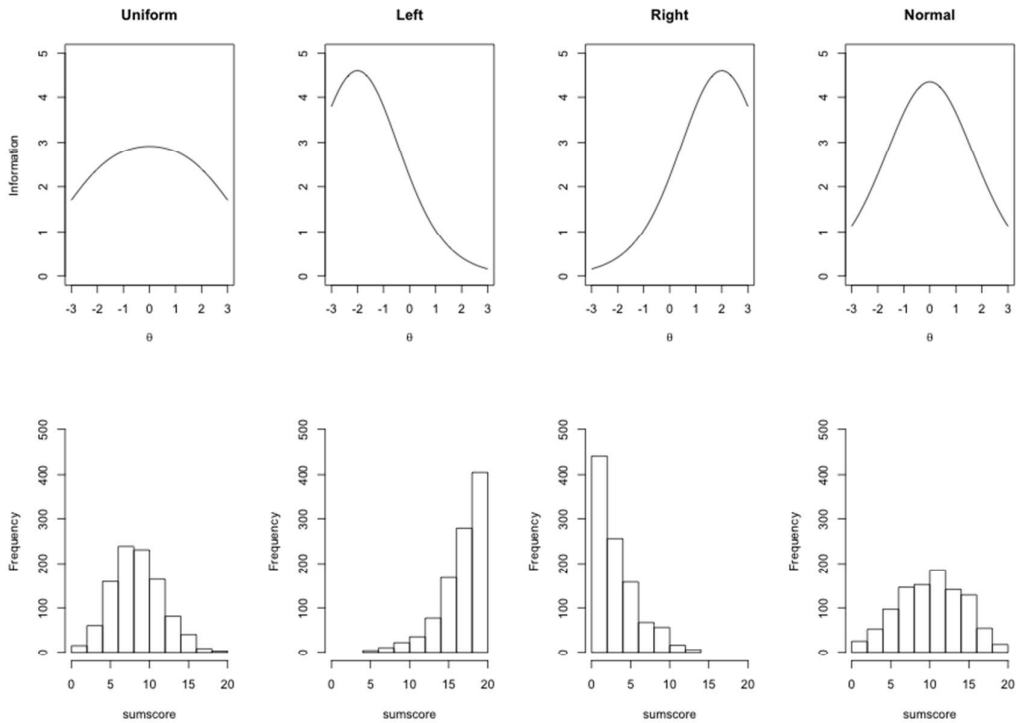
Figure 6. Power to detect a QTL as a function of risk allele frequency and test characteristics, using a case-control design. Percentage liability variance explained by QTL is 1% (left panel), and 25% (right panel).



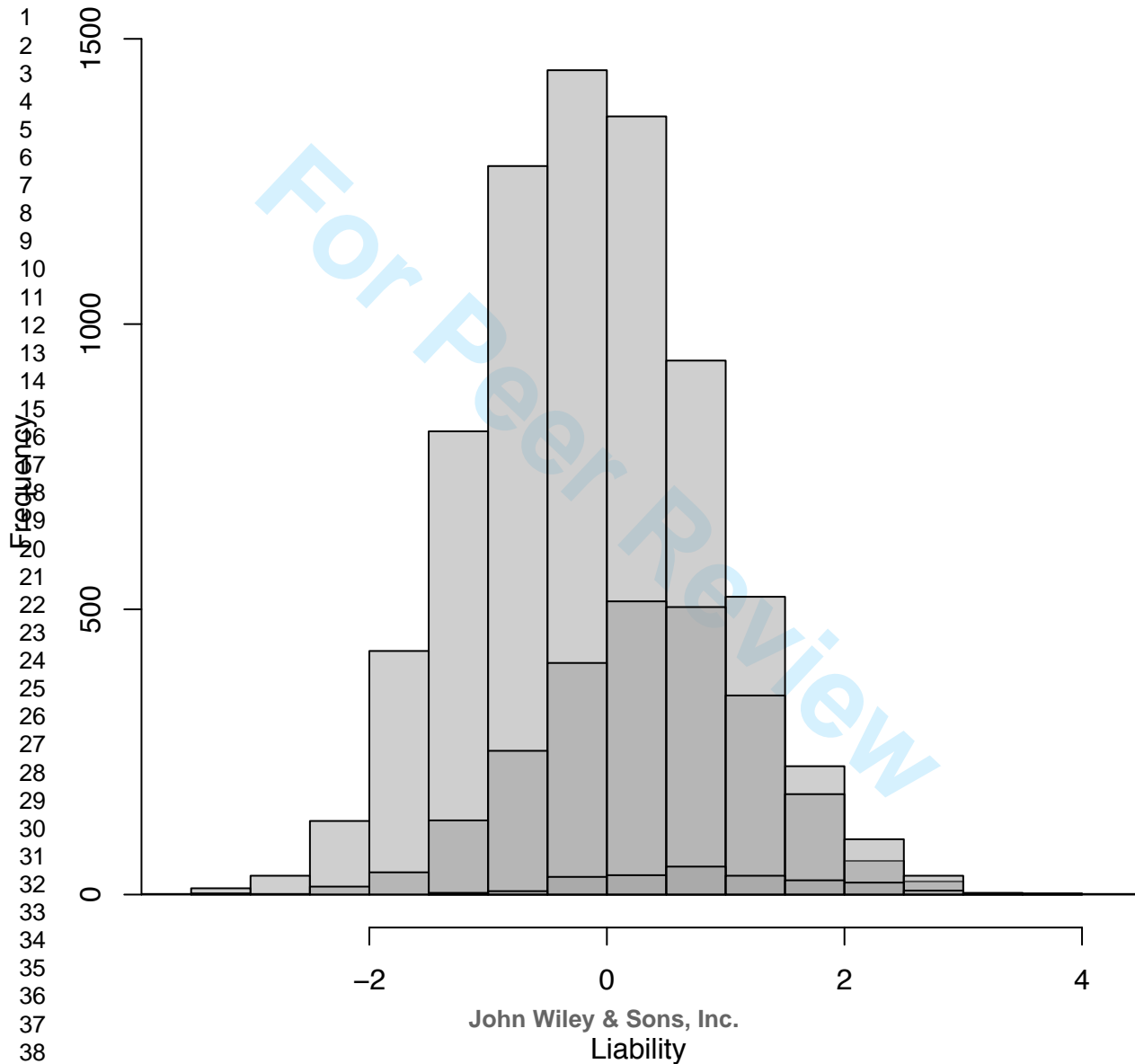
Peer Review

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

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



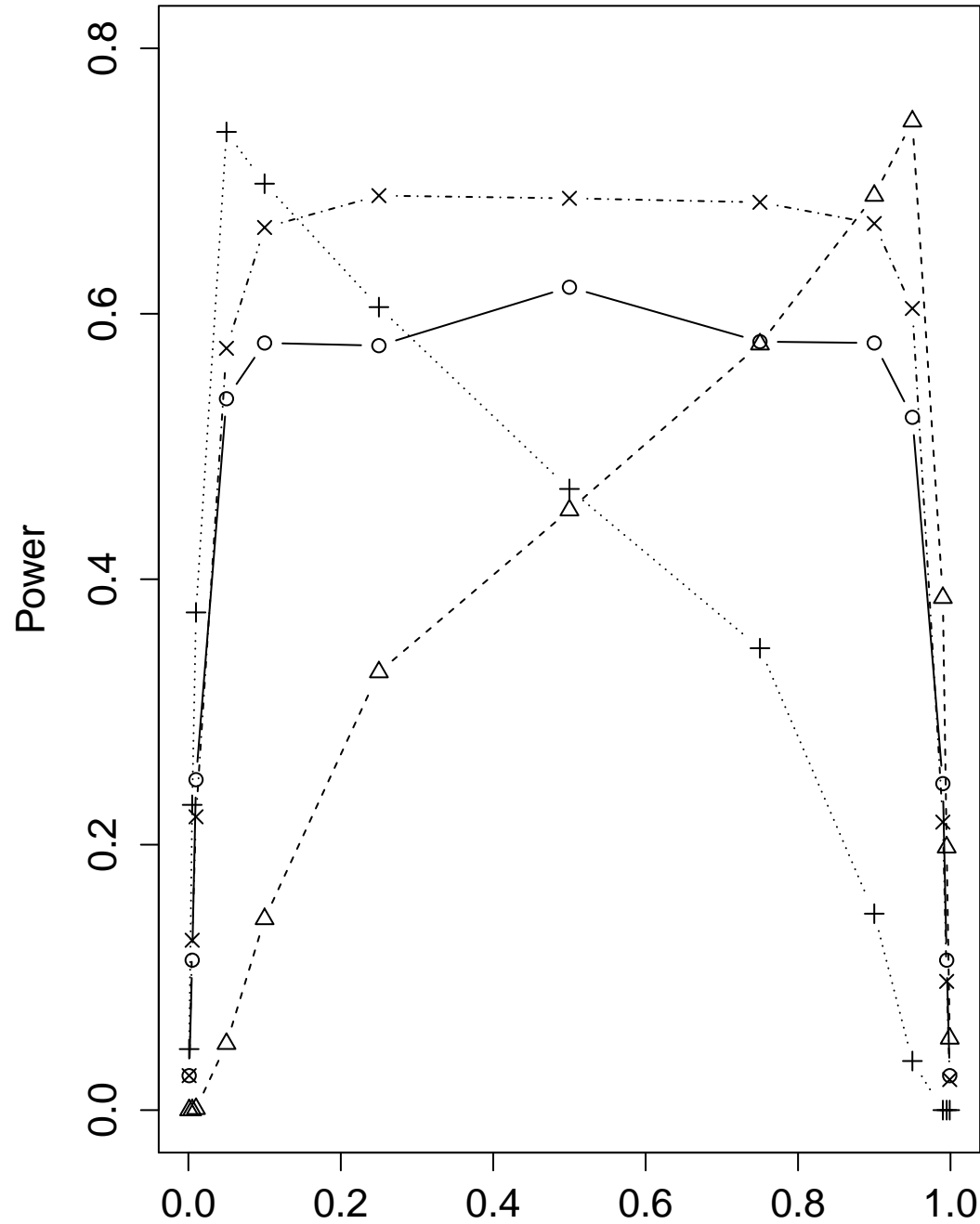
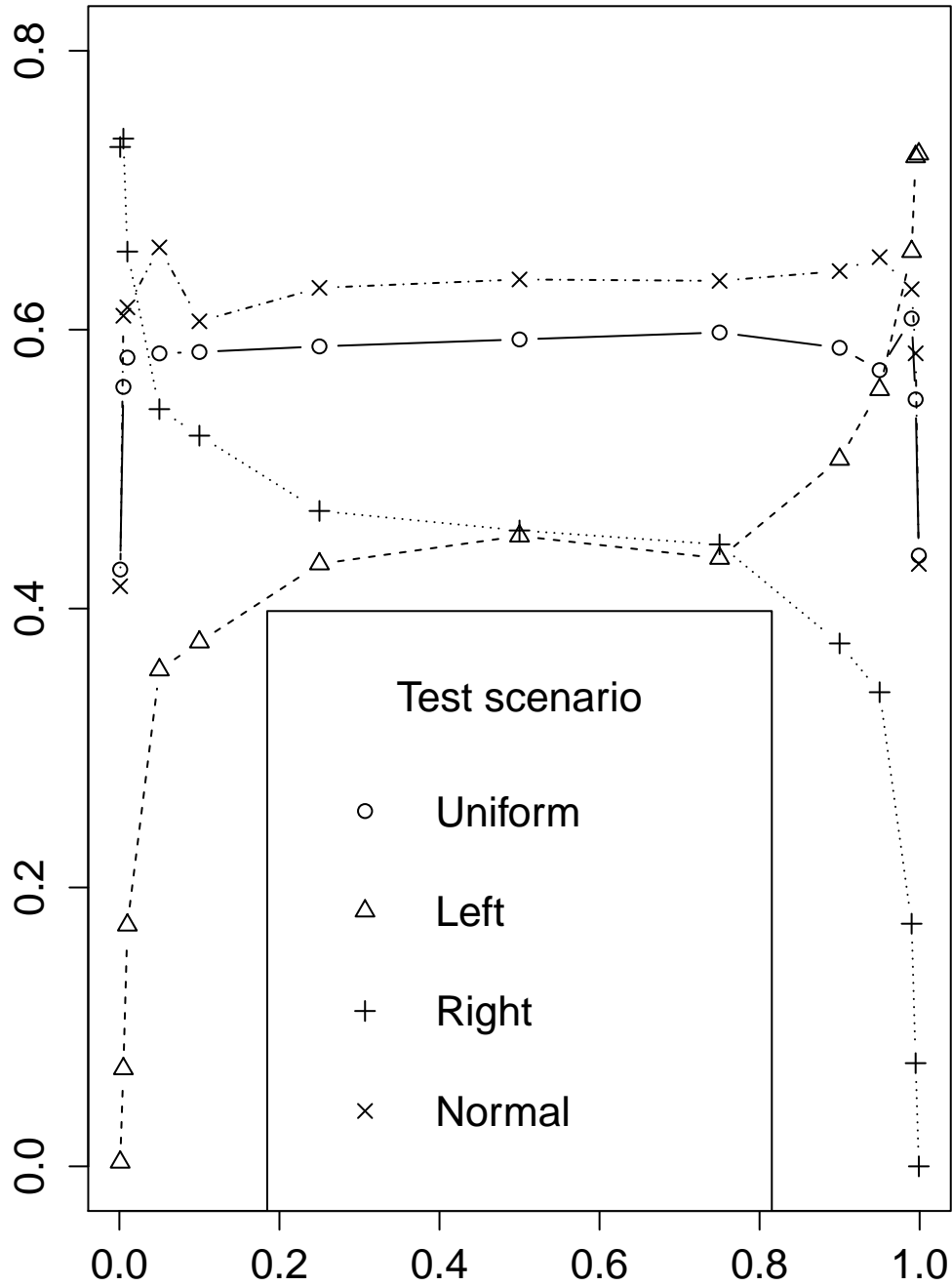
Peer Review

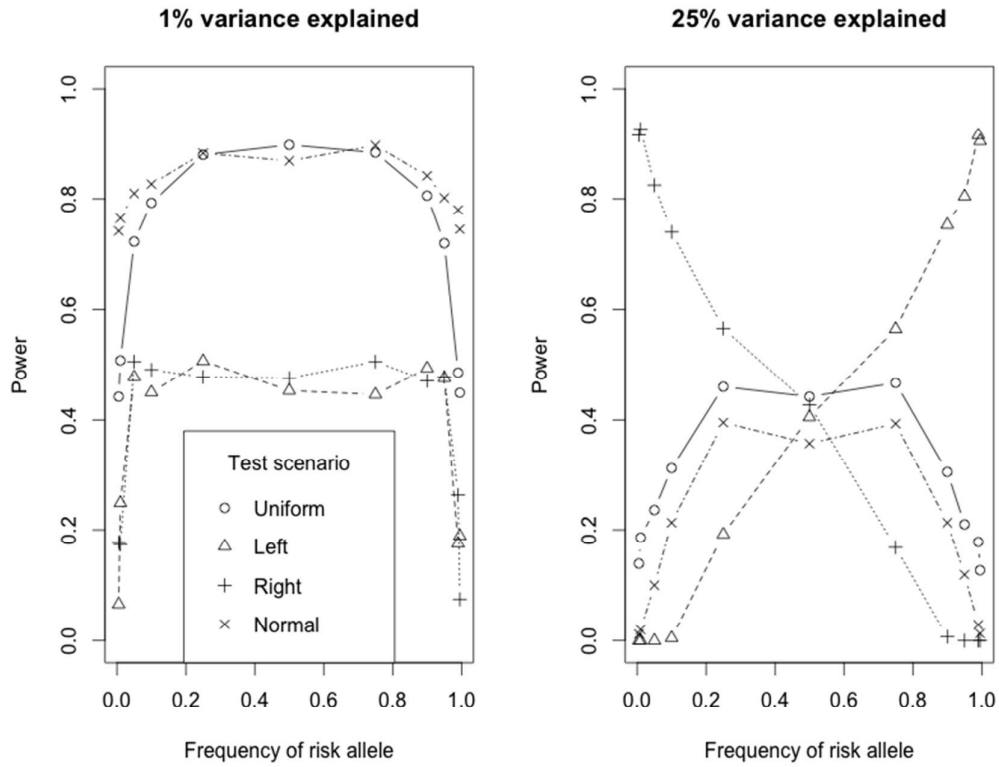


1% variance explained

25% variance explained

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

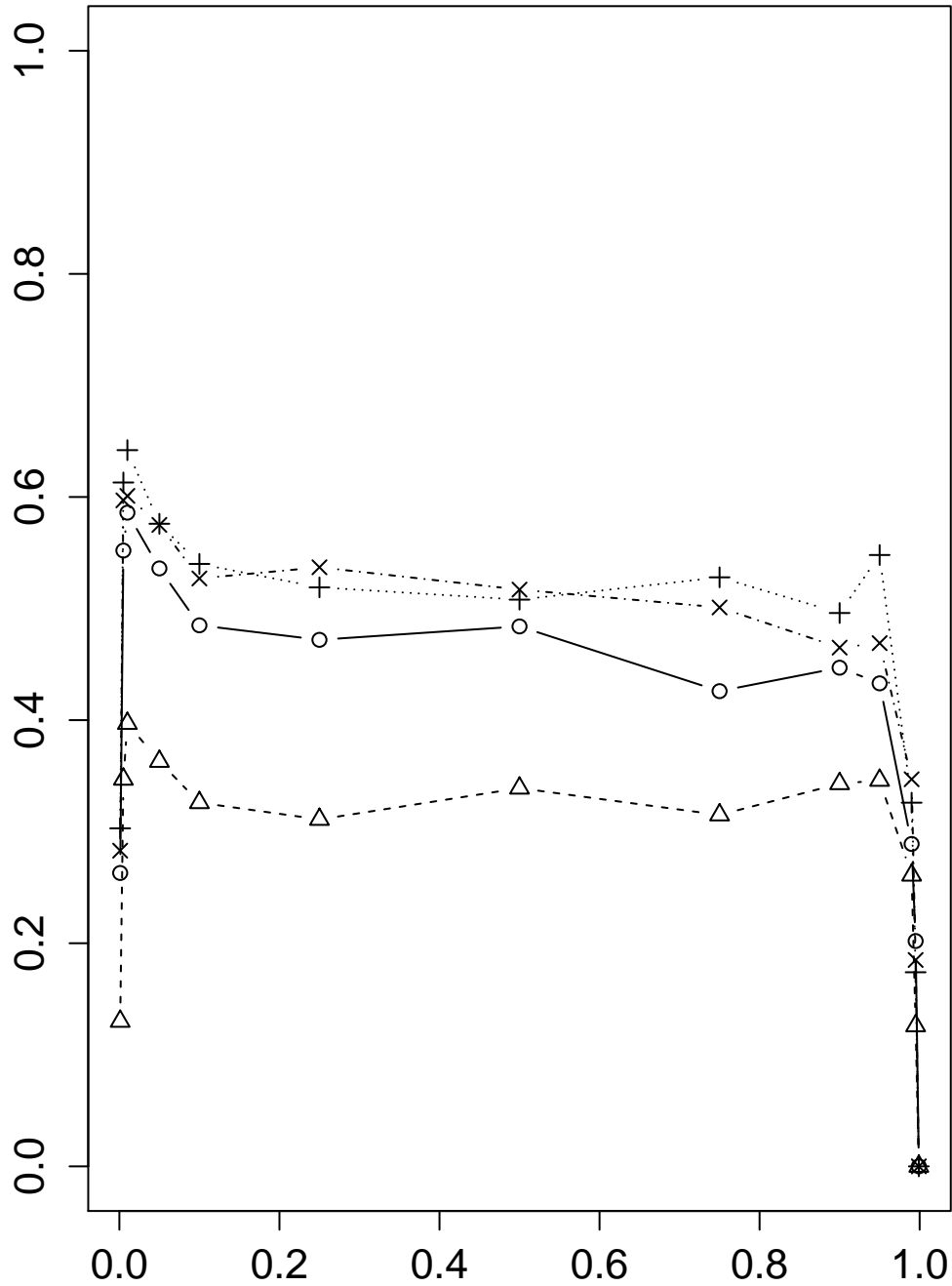




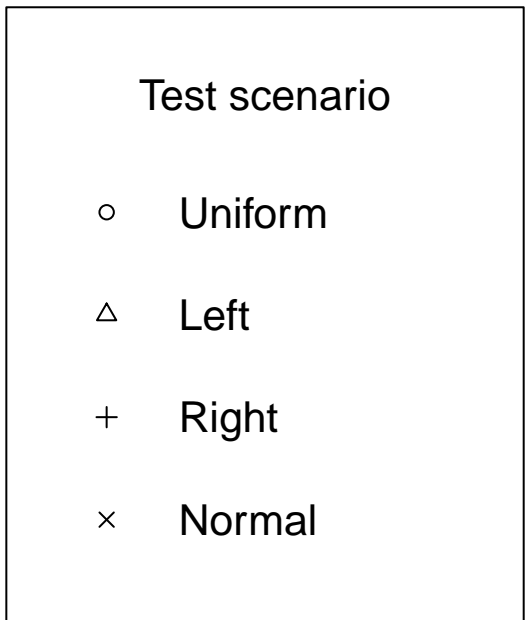
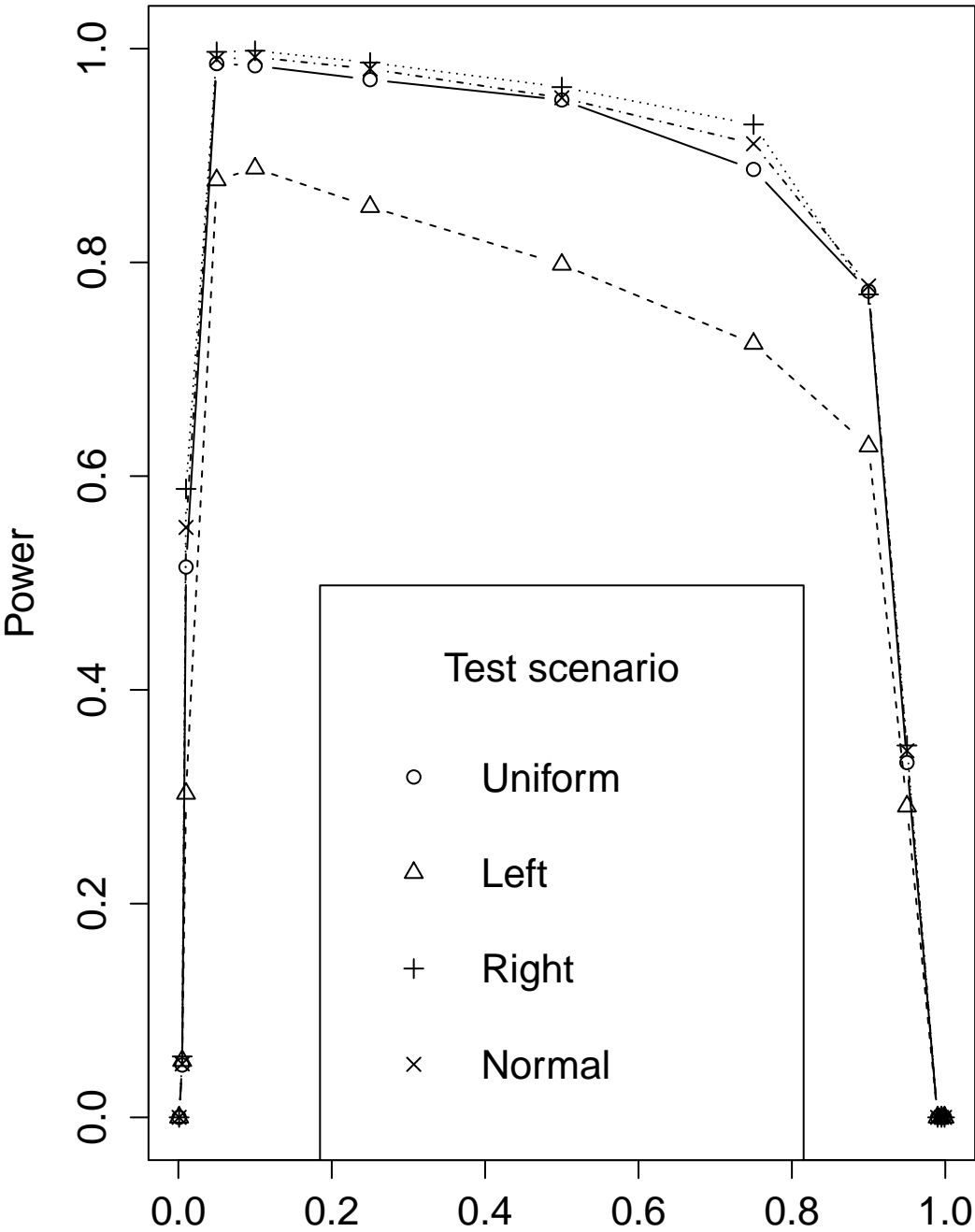
1% variance explained

25% variance explained

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



John Wiley & Sons, Inc.



Frequency of risk allele