

國立交通大學

統計學研究所

碩士論文

信賴區間與模擬研究
-對於穩定表現型的遺傳解釋比例

Confidence Interval and Simulation Studies for
the Proportion of Heritability Explained by
Endophenotypes

研究生：謝志強

指導教授：黃冠華 博士

中華民國九十五年六月

信賴區間與模擬研究
-對於穩定表現型的遺傳解釋比例
Confidence Interval and Simulation Studies for the Proportion of
Heritability Explained by Endophenotypes

研究生：謝志強
指導教授：黃冠華

Student : Chin-Chiang Hsieh
Advisor : Dr. Guan-Hua Huang

國立交通大學

統計學研究所

碩士論文

A Thesis

Submitted to Institute of Statistics

College of Science

Nation Chiao-Tung University

in partial Fulfillment of the Requirements

for the degree of Master

in

Statistics

June 2006

Hsinchu, Taiwan, Republic of China

中華民國九十五年六月

信賴區間與模擬研究 -對於穩定表現型的遺傳解釋比例

學生：謝志強

指導教授：黃冠華博士

國立交通大學統計學研究所

摘要

在生物學上，穩定表現型(endophenotype)和疾病有著相同的遺傳路徑，但穩定表現型卻比診斷上的表現型(phenotype)更為接近其相關的基因，這也顯示穩定表現型在複雜疾病上基因研究的重要性，在這篇報告裡，針對一個由穩定表現型所發展的指標，即穩定表現型的遺傳解釋比例，我們透過模擬提供其相關意義，同時，我們也提供此指標的信賴區間，藉此執行統計檢定和統計推論，除外，透過信賴區間和模擬的結果，建構一些準則，以幫助我們找尋有用的穩定表現型。

關鍵字：穩定表現型；遺傳率；基因分析

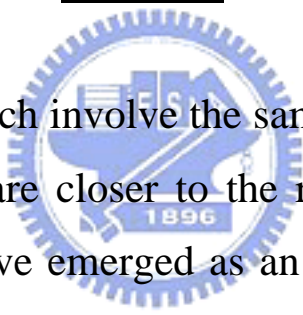
Confidence Interval and Simulation Studies for the Proportion of Heritability Explained by Endophenotypes

Student : Chin-Chiang Hsieh

Advisor : Dr.Guan-Hua Huang

Institute of Statistics
National Chiao Tung University

Abstract



Endophenotypes, which involve the same biological pathways as diseases but presumably are closer to the relevant gene action than diagnostic phenotypes, have emerged as an important concept in the genetic studies of complex diseases. In this report, we give some patterns about the developed index, the proportion of heritability explained (PHE) by the endophenotypes for validating endophenotypes. Besides, we provide a relevant confidence interval of PHE to perform a statistical test and to make some statistical inference. Using the relevant confidence interval of PHE, we construct some criteria to help us search a useful endophenotype.

KEY WORDS : endophenotype ; heritability ; genetic analysis

誌 謝

在交大統研所這2年的學習，實在讓我受益匪淺，非常感謝所上的所有老師，尤其是黃冠華老師，總是不厭其煩地指導我寫論文，讓我學會做研究時應該具備的精神與態度，在此，以一句「謝謝老師，您辛苦了」聊表心中無限的感激，同時，也要謝謝其他的口試委員(陳珍信老師、秋燕楓老師和洪志真老師)在口試時給我意見與建議，讓我的論文能更完整、更充實。

另外，也要謝謝一群很好的學長、同學與朋友，在我寫論文過程中，每遇到瓶頸和挫折，都能給我支持與鼓勵，特別是與我同個指導教授的秀慧，常幫我找程式指令和問題，讓我程式可以順利完成。

在此，僅以此篇論文獻給認識我的朋友們，謝謝你們。

謝志強 謹誌于
國立交通大學統計研究所
中華民國九十五年六月



CONTENTS

| | |
|---|-----|
| ABSTRACT(in Chinese) | i |
| ABSTRACT(in English) | ii |
| ACKNOWLEDGEMENTS(in Chinese) | iii |
| CONTENTS | iv |
| LIST OF TABLES | v |
| LIST OF FIGURES | vi |
| | |
| 1 Introduction..... | 1 |
| 2 Literature review..... | 3 |
| 2.1 Statistical validation of surrogate endpoints..... | 3 |
| 2.2 Statistical framework in genetic epidemiology..... | 4 |
| 3 Method..... | 9 |
| 3.1 Model..... | 9 |
| 3.2 Estimation..... | 13 |
| 3.3 Hypothesis Test..... | 19 |
| 4 Simulation studies..... | 19 |
| 4.1 Study design..... | 19 |
| 4.2 Result..... | 22 |
| 4.2.1 PHE..... | 22 |
| 4.2.2 The accuracy of the estimators of the standard error of PHE..... | 23 |
| 4.2.3 Test of PHE..... | 24 |
| 5 Discussion..... | 27 |
| Appendix..... | 30 |
| Reference..... | 34 |

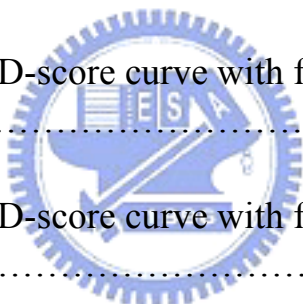
LIST OF TABLES

| | |
|---|----|
| Table 1. Genetic components of variance assuming mating..... | 8 |
| Table 2. The covariance components for relative pairs..... | 15 |
| Table 3. The derivative of covariance components for relative pairs.. | 17 |
| Table 4. Parameter setting under scenario II (1)..... | 21 |
| Table 5. Parameter setting under scenario II (2)..... | 21 |
| Table 6. Simulation results based on scenario I (1)..... | 37 |
| Table 7. Simulation results based on scenario I (2)..... | 38 |
| Table 8. Simulation results based on scenario I (3)..... | 39 |
| Table 9. Simulation results based on scenario I (4)..... | 40 |
| Table 10. Simulation results based on scenario II with $P > E$ (1)..... | 41 |
| Table 11. Simulation results based on scenario II with $P > E$ (2)..... | 42 |
| Table 12. Simulation results based on scenario II with $P > E$ (3)..... | 43 |
| Table 13. Simulation results based on scenario II with $P > E$ (4)..... | 44 |
| Table 14. Simulation results based on scenario II with $P < E$ (1)..... | 45 |
| Table 15. Simulation results based on scenario II with $P < E$ (2)..... | 46 |
| Table 16. Simulation results based on scenario II with $P < E$ (3)..... | 47 |
| Table 17. Simulation results based on scenario II with $P < E$ (4)..... | 58 |

LIST OF FIGURES

| | |
|---|----|
| Figure 1. A surrogate endpoint versus an endophenotype in the disease process..... | 49 |
| Figure 2. Two scenarios verified in the simulation studies..... | 50 |
| Figure 3. Scenario I histogram with family 200..... | 51 |
| Figure 4. Scenario I histogram with family 500..... | 52 |
| Figure 5. Scenario I histogram with family 200 & $\rho_\varepsilon=0.5$ | 53 |
| Figure 6. Scenario I histogram with family 500 & $\rho_\varepsilon=0.5$ | 54 |
| Figure 7. Scenario II histogram with family 200 & $P>E$ | 55 |
| Figure 8. Scenario II histogram with family 500 & $P>E$ | 56 |
| Figure 9. Scenario II histogram with family 200 & $P>E$ & $\rho_\varepsilon=0.5$... | 57 |
| Figure 10. Scenario II histogram with family 500 & $P>E$ & $\rho_\varepsilon=0.5$.. | 58 |
| Figure 11. Scenario II histogram with family 200 & $P<E$ | 59 |
| Figure 12. Scenario II histogram with family 500 & $P<E$ | 60 |
| Figure 13. Scenario II histogram with family 200 & $P<E$ & $\rho_\varepsilon=0.5$... | 61 |
| Figure 14. Scenario II histogram with family 500 & $P<E$ & $\rho_\varepsilon=0.5$... | 62 |
| Figure 15. Scenario I mean LOD-score curve with family 200..... | 63 |
| Figure 16. Scenario I mean LOD-score curve with family 500..... | 64 |
| Figure 17. Scenario I mean LOD-score curve with family 200 & $\rho_\varepsilon=0.5$ | 65 |

| | |
|---|----|
| Figure 18.Scenario I mean LOD-score curve with family 500 & $\rho_{\varepsilon}=0.5$ | 66 |
| Figure 19.Scenario I mean LOD-score curve with family 200 & $P>E$ | 67 |
| Figure 20.Scenario I mean LOD-score curve with family 500 & $P>E$ | 68 |
| Figure 21.Scenario I mean LOD-score curve with family 200 & $P>E$ & $\rho_{\varepsilon}=0.5$ | 69 |
| Figure 22.Scenario I mean LOD-score curve with family 500 & $P>E$ & $\rho_{\varepsilon}=0.5$ | 70 |
| Figure 23.Scenario I mean LOD-score curve with family 200 & $P<E$ | 71 |
| Figure 24.Scenario I mean LOD-score curve with family 500 & $P<E$ | 72 |
| Figure 25.Scenario I mean LOD-score curve with family 200 & $P<E$ & $\rho_{\varepsilon}=0.5$ | 73 |
| Figure 26.Scenario I mean LOD-score curve with family 500 & $P<E$ & $\rho_{\varepsilon}=0.5$ | 74 |



1 INTRODUCTION

In diseases with classic or Mendelian genetics as their distal causes, genotypes are usually indicative of phenotypes. However, this degree of genetic certainty does not exist for complex disease [Gottesman and Gould, 2003]. These “complex” diseases are influenced by multiple genes, environmental factors and their interactions on phenotypes. It leads the direct relationship between phenotype and genotype disrupted because that the same genotype may give rise to different phenotypes or the same phenotype may have arisen from different genotypes. To facilitate the identification of influential genetic markers of complex diseases, the endophenotype approach has been advocated. Other synonymies of endophenotype, such as intermediate phenotype, biological marker, sub-clinical trait, vulnerability marker, and phenotypic uncertainty, have been used interchangeably with slightly different implications. Gottesman and Gould [2003] provided a means of endophenotyps for identifying the “downstream” traits or facets of clinical phenotypes, as well as the “upstream” consequences of genes, and suggested the following five useful criteria for identification of endophenotypes:

1. The endophenotype is associated with illness in the population.
2. The endophenotype is heritable.
3. The endophenotype is primarily state-independent (manifests in an individual whether or not illness is active).
4. Within families, endophenotype and illness co-segregate.
5. The endophenotype found in affected family members is found in non-affected family members at a higher rate than in the general population.

Hence, the endophenotype is closer to the underlying gene than the phenotype in the course of disease’s natural history. Endophenotype-based genetic analysis is more likely to succeed than phenotype-based one in terms of search for the susceptibility genes; nevertheless, there are emerging needs of systematic statistical methods for endophenotype-based analysis.

On the other hand, surrogate endpoints have been frequently utilized in clinical research, especially in chronic diseases, when the primary endpoint is costly or time-consuming to obtain. A good deal of statistical research in the evaluation of surrogate endpoints have been undertaken for decades. Prentice [1989] presented a landmark definition of surrogate endpoints. Freedman et al. [1992] introduced “the proportion of treatment effect on the primary endpoint explained” (*PTE*) by the surrogate to supplement Prentice’s criteria.

Conceptually, an endophenotypes is a “downstream” biomarker for detection of heritable biological underpinning and a surrogate endpoint is an “upstream” biomarker for evaluation of treatment effect as illustrated in Figure 1. Noticeably, the causal pathway of intervention-surrogate endpoint-primary endpoint in surrogate analysis can be seen as an analogy of the pathway of genotype-endophenotype-phenotype in endophenotypes-based analysis. Both endophenotypes and surrogate endpoints lie in a biological pathway, but with two important differences: (i) the endophenotype is expected to be closer to the upstream genotype to increase the chance of identifying it, though the surrogate endpoint intends to substitute the downstream primary endpoint, and (ii) when the purpose of the study is to identify responsible genes for the phenotype, genotype information is usually unknown, whereas treatment in validating a surrogate is known.

Huang et al. [2005] defined an endophenotype to be “a trait for which a test of null hypothesis of no genetic heritability implies the corresponding null hypothesis based on the phenotype of interest” and developed a formal statistical methodology for accessing the utility of endophenotypes, motivated by the conditioning strategy used for surrogate endpoints commonly seen in clinical research. The methodology is especially useful for the situation where underlying genotype is unknown in that researchers use endophenotypes to increase opportunities of finding susceptible disease genes, not to verify whether a specific gene is the cause of disease. Similar to validating surrogate endpoints, various indices can be provided to use to validate endophenotypes. One of the indices is the proportion of heritability explained (*PHE*) by the endophenotype, similar to *PTE* introduced by Freedman et al. [1992].

Several authors had pointed out the major difficulty of using *PTE*: the confidence interval of *PTE* is generally too wide to convey any useful information. That is, the true *PTE* might

be anywhere from zero to well over 100% or be negative. To avoid the confidence interval of *PHE* far too wide to be of practical relevance, we provided a relevant confidence interval of *PHE*. Furthermore, for *PHE*, we perform a statistical test or to establish some criteria for determining whether there is an endophenotype. Also, extensive simulation studies were performed to verify the usefulness of *PHE*.

2 LITERATURE REVIEW

2.1 STATISTICAL VALIDATION OF SURROGATE ENDPOINTS

In most clinical researches, the primary endpoint is too difficult or costly or time-consuming to obtain, particularly in chronic diseases. It may force the investigators to use a substitute or “surrogate”, instead of true endpoint. Surrogate endpoints have been of clinical interest for decades, but it was not until Prentice published a seminal paper in 1989 that formal statistical investigation started. Prentice defined a surrogate endpoint to be “a response variable for which a test of null hypothesis of no relationship to the treatment groups under comparison is also a valid test of the corresponding null hypothesis based on the true (clinical) endpoint”. Prentice’s definition can be written as

$$f(S | X) = f(S) \iff f(T | X) = f(T) \quad (1)$$

where T denotes the status of a primary endpoint, S denotes the status of a surrogate endpoint, X is the treatment variable, $f(S)$ is the distribution of S, and $f(S|X)$ is the conditional distribution of S given X. Validation of Prentice’s definition involves the following two criteria:

$$f(T | S) \neq f(T) \quad \text{and} \quad f(T | S, X) = f(T | S) \quad (2)$$

[Prentice, 1989; Freedman et al., 1992; Buyse and Molenberghs, 1998]. The first criterion states that the surrogate endpoint must be correlated with the primary clinical endpoint, and the second criterion is that the surrogate endpoint should fully capture the treatment effect on the primary endpoint.

The surrogate endpoint described by Prentice mediates all of the effect of treatment on the primary endpoint, that is

$$X \rightarrow S \rightarrow T$$

A more complex, but more likely, situation arises when treatment has a direct effect on the primary endpoint that is not mediated through the surrogate [De Gruttola et al., 2001]:

$$X \begin{array}{c} \rightarrow S \rightarrow T \\ \searrow \nearrow \end{array}$$

Freedman et al.[1992] proposed to focus on the proportion of the treatment effect mediated through the surrogate. A good surrogate is one that explains a large proportion of that effect. The proposal can be made in the content of generalized linear models [McCullagh et al., 1989]. The net effect of X on T can be assessed through the regression coefficient β_T in the generalized linear model

$$g[E(T)] = \alpha_T + \beta_T X \quad (3)$$

where $g(\bullet)$ is the link function connecting the mean response and covariates, and the effect of X on T after inclusion of S is the regression coefficient β_{TS} in the following generalized linear model

$$g[E(T)] = \alpha_{TS} + \beta_{TS} X + \gamma_{TS} S \quad (4)$$

The proportion of the treatment effect (on the primary endpoint) explained (*PTE*) by the surrogate is given by

$$PTE = 1 - \frac{\beta_{TS}}{\beta_T} \quad (5)$$

The $100(1 - \alpha)\%$ confidence limits of *PTE* can be calculated using Fieller's theorem or the delta method [Buyse and Molenberghs, 1998]. Using Fieller's theorem is generally preferable the $100(1 - \alpha)\%$ confidence limits of PTE [Herson, 1975].

2.2 STATISTICAL FRAMEWORK IN GENETIC EPIDEMIOLOGY

First, consider a genetic locus defined by two alleles. If we assume two allelic variants, Q and q with frequencies of p_Q and $(1 - p_Q)$ at a given quantitative-trait locus(*QTL*), the

genotype-specific means are given by $\mu_{QQ} = \mu + a$, $\mu_{Qq} = \mu + d$, and $\mu_{qq} = \mu - a$. The genotypic mean values can be reparameterized in terms of $\mu'_{QQ} = a$, $\mu'_{Qq} = d$, and $\mu'_{qq} = -a$, so that the mean, μ' , is $p_Q^2 a + 2p_Q(1 - p_Q)d + (1 - p_Q)^2(-a)$ and the variance about the mean, σ_g^2 , is $p_Q^2 (\mu'_{QQ} - \mu')^2 + 2p_Q(1 - p_Q)(\mu'_{Qq} - \mu')^2 + (1 - p_Q)^2 (\mu'_{qq} - \mu')^2$. The variance about mean can be decomposed as

$$\sigma_g^2 = \sigma_a^2 + \sigma_d^2 \quad (6)$$

where

$$\begin{aligned} \sigma_a^2 &= 2p_Q(1 - p_Q) [p_Q \mu'_{QQ} + (1 - 2p_Q) \mu'_{Qq} - (1 - p_Q) \mu'_{qq}]^2 \\ &= 2p_Q(1 - p_Q) [a + d(1 - 2p_Q)]^2 \end{aligned} \quad (7)$$

is called the additive component of variance, and

$$\begin{aligned} \sigma_d^2 &= \{p_Q(1 - p_Q) [\mu'_{QQ} - 2\mu'_{Qq} + \mu'_{qq}]\}^2 \\ &= [2p_Q(1 - p_Q)d]^2 \end{aligned} \quad (8)$$

is called the dominance component of variance [Duggirala et al., 1997].

Let G_i and G_j represent the genotype of two individuals i and j . In general, under Hardy-Weinberg equilibrium and no inbreeding, the genetic covariance can be expressed as

$$cov(G_i, G_j) = 2\phi_{ij}\sigma_a^2 + \Delta_{ij}\sigma_d^2 \quad (9)$$

where, ϕ_{ij} , the coefficient of kinship, or coefficient of coancestry, is defined as the probability of randomly drawing a single allele in individual i that is identical by descent (*ibd*) to a single allele at the same locus randomly drawn from individual j , and Δ_{ij} , the fraternity coefficient, is defined as the probability that both alleles at a locus are shared *ibd* by individuals i and j [Duggirala et al., 1997].

After all, it is not very realistic. The involvement of several loci in the determination of the trait may be considered. Assume that there are m *QTLs* to influence the actual trait. If

the effects of single loci are independent, the covariance can be written as

$$\text{cov}(X_i, X_j) = 2\phi_{ij}\sigma_A^2 + \Delta_{ij}\sigma_D^2 \quad (10)$$

where X_i and X_j represent, respectively, the actual trait of individuals i and j , $\sigma_A^2 = \sum_{k=1}^m \sigma_{ak}^2$ is the total additive genetic variance, $\sigma_D^2 = \sum_{k=1}^m \sigma_{dk}^2$ is the total dominance genetic variance and σ_{ak}^2 and σ_{dk}^2 , are the additive and dominance genetic variance due to the k th locus, respectively [Iachine, 2004].

Besides, to describe the residual variation of the trait when the genotype is fixed, the so-called environmental effects may be introduced. Suppose the effects of genes and environment are additive. Under the additional assumption of independence between genotypic effects and environmental effects, the covariance can be written as

$$\text{cov}(X_i, X_j) = 2\phi_{ij}\sigma_A^2 + \Delta_{ij}\sigma_D^2 + \text{Var}(X_{E,ij}) \quad (11)$$

where $X_{E,ij}$ is environmental effect between individual i and individual j [Iachine, 2004].

In particular, this implies the following structure of the trait variance:

$$\text{Var}(X_i) = \sigma_A^2 + \sigma_D^2 + \text{Var}(X_{E,i}) \quad (12)$$

where $X_{E,i}$ is environmental effect of individual i . If we further assume, that we have the same the environmental variance $\text{Var}(X_E)$ and total variance σ^2 for all family members, the structure of the trait variance can be written as

$$\sigma^2 = \sigma_A^2 + \sigma_D^2 + \text{Var}(X_E) \quad (13)$$

A study point for many scientists investigating disease aetiology has often been to study the heritability of a particular trait. Formally, the heritability of a continuous trait is defined as the proportion of its total variance (σ^2) that is attributable to genetic factors in a particular population. Narrow-sense heritability is defined as σ_A^2/σ^2 and broad-sense heritability as

$(\sigma_A^2 + \sigma_D^2) / \sigma^2$. Usually, it is of interest to know the broad-sense heritability because its value can be used to predict the effect of searching for genes [Iachine, 2004; Burton and Tobin, 2003]. Let us decompose $Var(X_E)$ into σ_C^2 and σ_E^2 , where σ_C^2 is called the shared environmental component of variance and σ_E^2 is called the non-shared environmental component of variance, i.e.

$$\sigma^2 = \sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_E^2 \quad (14)$$

In some practical problems, it is often assumed that the dominance component of variance is negligible (i.e. $\sigma_D^2 = 0$), leading to the so-called ACE model.

As in mainstream epidemiology, many of the relevant models may helpfully be viewed as being generalized linear mixed models [Breslow and Clayton, 1993.]. Here, we will consider the structure of one such GLMM.

A general model with wide applicability may be written as

$$\begin{aligned} g(\mu_{ij}) &= \eta_{ij} = \alpha + \beta^T z_{ij} + \xi_{ij}, \\ Y_{ij} &\sim f(\mu_{ij}, \varpi) \\ \xi_{ij} &\sim N(0, [\sigma_A^2 + \sigma_D^2 + \sigma_C^2]) \\ cov(\xi_{ij}, \xi_{ik}) [j \neq k] &= 2\phi_{ij,ik}\sigma_A^2 + \Delta_{ij,ik}\sigma_D^2 + \lambda_{ij,ik}\sigma_C^2 \end{aligned} \quad (15)$$

where Y_{ij} is the observed phenotype in the j th member of the i th family, μ_{ij} is its expected value, and $f(\bullet)$ denotes an error distribution which may incorporate a nuisance parameter denoted ϖ [Burton and Tobin, 2003]. The expected value of the phenotype is predicted via a link function $g(\bullet)$ applied to a linear predictor (η_{ij}) comprising a baseline mean (α), a vector of observed covariates (z_{ij}), a corresponding vector of unknown regression parameters (β) and subject-specific random effects ξ_{ij} with an appropriate covariance structure. The components σ_A^2 , σ_D^2 and σ_C^2 represent, respectively, the variances arising from polygenic additive effects, polygenic dominance effects and shared environmental effects [Hopper, 2002]. The terms $\phi_{ij,ik}$ and $\Delta_{ij,ik}$ denote, respectively, the kinship coefficient and fraternity coefficient between individuals ij and ik . Table 1 details the $\phi_{ij,ik}$ and $\Delta_{ij,ik}$ values for selected relative pairs and

the total genetic variances that these imply [Burton and Tobin, 2003].

Table 1. Genetic components of variance assuming mating

| Relationship | ϕ | Δ | Genetic covariance |
|-------------------------|----------------|----------------|--|
| Same person | $\frac{1}{2}$ | 1 | $\sigma_A^2 + \sigma_D^2$ |
| Parent-child | $\frac{1}{4}$ | 0 | $\frac{1}{2}\sigma_A^2$ |
| Full sibling | $\frac{1}{4}$ | $\frac{1}{4}$ | $\frac{1}{2}\sigma_A^2 + \frac{1}{4}\sigma_D^2$ |
| Half sibling | $\frac{1}{8}$ | 0 | $\frac{1}{4}\sigma_A^2$ |
| Monozygous twins | $\frac{1}{2}$ | 1 | $\sigma_A^2 + \sigma_D^2$ |
| Grandparent-grandchild | $\frac{1}{8}$ | 0 | $\frac{1}{4}\sigma_A^2$ |
| Uncle/aunt-nephew/niece | $\frac{1}{8}$ | 0 | $\frac{1}{4}\sigma_A^2$ |
| First cousins | $\frac{1}{16}$ | 0 | $\frac{1}{8}\sigma_A^2$ |
| Double first cousins | $\frac{1}{8}$ | $\frac{1}{16}$ | $\frac{1}{4}\sigma_A^2 + \frac{1}{16}\sigma_D^2$ |
| Spoused | 0 | 0 | 0 |

In many situations, the elements, $\lambda_{ij,ik}$, is simply binary indicator denoting whether two individuals live together ($\lambda_{ij,ik} = 1$) or apart ($\lambda_{ij,ik} = 0$). However, the effect of shared environment may be modelled in a more sophisticated manner by adding some factors related to length of cohabitation and to time spent living apart [Hopper, 2002].

Furthermore, we are generally interested in examination of one or a few *QTLs* at a time. Having established the presence of genetic effects on the trait, we would like to investigate how much of this genetic variation can be attributed to genetic variation at a specific chromosome. That is, genetic effects are due to a specific locus and residual genetic effects. Assume that the quantitative trait X is influenced by the genetic loci L_1, L_2, \dots, L_m located on this chromosome. For example, if we are focusing on the analysis of the q th *QTL*, we can absorb the effects of all of the remaining *QTLs* in residual components of covariance. The covariance can be expressed as

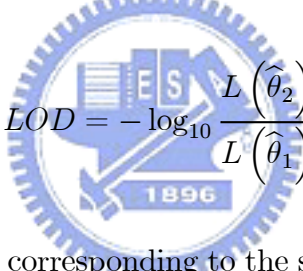
$$cov(X_i, X_j) = \pi_q \sigma_{aq}^2 + k_{2q} \sigma_{dq}^2 + 2\phi_{ij} \sigma_{A^*}^2 + \Delta_{ij} \sigma_{D^*}^2 + Var(X_{E,ij}) \quad (16)$$

where π_q is the probability of a random allele being *ibd* at the q th *QTL*, k_{2q} is the probability that both alleles at a locus are shared *ibd* at the q th *QTL*, $\sigma_{A^*}^2$ represents the residual additive genetic variance, $\sigma_{D^*}^2$ represents the residual dominance genetic variance, and $X_{E,ij}$ is environmental effect between individual i and individual j . For any given chromosome location, π and k_2 can be estimated from genetic marker data and information on the genetic map [Almasy and Blangero, 1998]. Similarly, it implies the following structure of the trait variance:

$$Var(X_i) = \sigma_{aq}^2 + \sigma_{dq}^2 + \sigma_{A^*}^2 + \sigma_{D^*}^2 + \sigma_C^2 + \sigma_E^2 \quad (17)$$

where σ_C^2 is environmental component of variance and σ_E^2 is non-shared environmental component of variance.

In linkage analysis, there is a tradition for using *LOD-score* from the null hypothesis H_0 of no linkage. This so-called *LOD-score* is defined as

$$LOD = -\log_{10} \frac{L(\hat{\theta}_2)}{L(\hat{\theta}_1)} \quad (18)$$


where $\hat{\theta}_2$ is the parameter estimate corresponding to the smaller model ($\sigma_A^2, \sigma_D^2, \sigma_C^2, \sigma_E^2$ be estimated in (14)) and $\hat{\theta}_1$ is the parameter estimate corresponding to the larger model ($\sigma_{aq}^2, \sigma_{dq}^2, \sigma_{A^*}^2, \sigma_{D^*}^2, \sigma_C^2, \sigma_E^2$ be estimated in (17)). Usually, the values of the *LOD-score* larger than 3 are interpreted as evidence of linkage.

3 Method

3.1 MODEL

Endophenotypes are useful for theorizing about clinical phenotypes and can mark the path between the genotype and the phenotype. Verification of existence of the pathway genotype-endophenotype-phenotype is the key of validating endophenotypes. Analogous to Prentice's definition [1989] that surrogate endpoint to be "a response variable for which a test of null hypothesis of no relationship to the treatment groups under comparison is also a valid test

of the corresponding null hypothesis based on the true (clinical) endpoint”, Huang et al. [2005] define an endophenotype to be “a trait for which a test of null hypothesis of no genetic heritability implies the corresponding null hypothesis based on the phenotype of interest”. More specifically, suppose P is the phenotype of interest, E is the selected endophenotype, and G represents an underlying genetic structure that fulfills the specified assumptions in calculating heritability, then the proposed definition is:

$$f(E | G) = f(E) \Rightarrow f(P | G) = f(P). \quad (19)$$

The definition has two important features [Huang et al. 2005]. First, “imply” replaces “if and only if” statement in Prentice’s definition of surrogate endpoints in avoidance of a problematic implication arisen in Begg and Leung [2000]. This change places endophenotype in higher upstream of the pathway from genotype to phenotype, instead of in the position that keeps the same distance with genotype as with phenotype. Second, genetic heritability is used as the measure of association with an underlying genetic structure. Heritability represents the proportion of variability attributable to genetic factors and can be obtained in a variance component approach [Hopper, 2002]. This is a perfect fit to our situation since it does not require knowledge of specific culprit genes and allows the likelihood of multiple gene influences.

The following is development of obtaining operational criteria of the proposal definition [Huang et al. 2005]. By definition, we have

$$f(P | G) = \int f(P, E | G) dE = \int f(P | E, G) f(E | G) dE \quad (20)$$

By (19), since $f(E | G) = f(E)$, we can obtain

$$f(P | G) = \int f(P | E, G) f(E) dE \quad (21)$$

If the condition

$$f(P | E, G) = f(P | E) \quad (22)$$

holds, then

$$f(P | G) = \int f(P | E) f(E) dE = f(P) \quad (23)$$

In pursuing a feasible approach, Huang et al. [2005] take (22) in a variance component model as the operational criterion for proposed endophenotype definition. It then requires heritability of phenotype becomes null, conditioning on candidate endophenotype, and implies genetic heritability of phenotype is captured by endophenotype.

Given a phenotype of continuous measurements, significance of (22) can be judged through the following variance component analysis for quantitative traits [Almasy and Blangero, 1998 and Huang et al. 2005]:

$$\begin{aligned} P_{ij} &= \alpha_H + \gamma_H E_{ij} + \tau_H Z_{ij} + G_{ij} + \epsilon_{ij}, \\ \epsilon_{ij} &\sim \text{Normal}(0, \sigma_R^2) \\ G_{ij} &\sim \text{Normal}(0, [\sigma_A^2 + \sigma_D^2 + \sigma_C^2]) \\ \text{cov}(G_{ij}, G_{ik}) &= 2\phi_{ij,ik}\sigma_A^2 + \Delta_{ij,ik}\sigma_D^2 + \lambda_{ij,ik}\sigma_C^2, \quad j \neq k \end{aligned} \quad (24)$$

where P_{ij} is the observed phenotype in the j th member of the i th family, E_{ij} is his/her corresponding specified endophenotype, Z_{ij} is his/her other covariates. ϵ_{ij} is the residual error term representing the effect of non-family factors. G_{ij} is the random effect for the underlying genetic structure. The components σ_A^2 , σ_D^2 and σ_C^2 represent the variance arising from polygenic additive effects, polygenic dominance effects and shared environmental effects, respectively. The (broad sense) heritability of P_{ij} , conditional on E_{ij} is

$$h = \frac{\sigma_A^2 + \sigma_D^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2} \quad (25)$$

The significance of rejecting the hypothesis $h = 0$ indicates the fulfillment of (22).

For a discrete phenotype of ordinal scale, the liability threshold model can be used in the preceding variance component setting^{[13][14]}. The model postulates the existence of an unobserved continuous trait (i.e., liability L_{ij}), and a set of thresholds t_1, t_2, \dots, t_{K-1} that

partition the liability distribution into intervals corresponding to distinct phenotypic states:

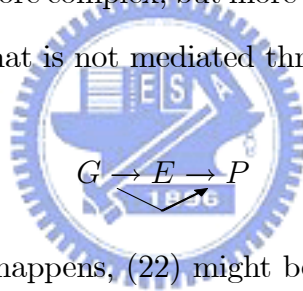
$$P_{ij} = \begin{cases} 1, & \text{if } L_{ij} < t_1 \\ 2, & \text{if } t_1 < L_{ij} < t_2 \\ \vdots & \quad \quad \quad \vdots \\ K, & \text{if } t_{K-1} < L_{ij} \end{cases}$$

The liability L_{ij} is then assumed to follow the same distribution as the P_{ij} in model (24) and heritability can be obtained based on the liability.

The endophenotype described above mediates all of the effect of genotype on phenotype, that is

$$G \rightarrow E \rightarrow P$$

This situation rarely happens. A more complex, but more likely, situation arises when genotype has a direct effect on phenotype that is not mediated through endophenotype:



If the more complex situation happens, (22) might be difficult to be satisfied in practice. This situation arises for most diseases. Huang et al. [2005] have provided some indices to evaluate the validation of endophenotypes. One of the important indices is the proportion of heritability explained (PHE) by the endophenotype defined as

$$PHE = 1 - \frac{h}{h_{NE}} \quad (26)$$

where h_{NE} is the heritability calculated from the variance component analysis (24) without including the endophenotype E_{ij} with any other covariates. A good endophenotype is one that explains a large proportion of heritability, thus, the greater the PHE value, the more likely E_{ij} an endophenotype.

3.2 ESTIMATION

Variance component analysis (24) can be performed using the SOLAR computer package [Almasy and Blangero, 1998]. As a result, PHE (26) can be estimated, that the estimators by of h and h_{NE} were obtained from the results of using the SOLAR computer package. Hence, we will focus on deriving the confidence limits of PHE or the estimator of the standard deviation of PHE . First, we we redefine (25) as

$$h \equiv h_A^{(1)} + h_D^{(1)}$$

where

$$h_A^{(1)} = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2}, \quad h_D^{(1)} = \frac{\sigma_D^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2}$$

Similarly, we redefine

$$h_{NE} \equiv h_A^{(2)} + h_D^{(2)}$$

PHE being the ratio of two parameter, its confidence limits can be calculated using Fieller's theorem or the delta method [Buyse and Molenberghs, 1998]:

Method1(*Fieller's theorem* [Buyse and Molenberghs,1998])

Using Fieller's Theorem, the 100 (1 - α) % confidence limits of $\left(1 - \frac{h}{h_{NE}}\right)$ are given by

$$1 - \frac{A \pm \sqrt{A^2 - BC}}{B}$$

where

$$\begin{aligned} A &= h \cdot h_{NE} - Z_\alpha^2 Cov(h, h_{NE}) \\ &= h \cdot h_{NE} - Z_\alpha^2 Cov\left(h_A^{(1)} + h_D^{(1)}, h_A^{(2)} + h_D^{(2)}\right) \\ &= h \cdot h_{NE} - Z_\alpha^2 \left\{ Cov\left(h_A^{(1)}, h_A^{(2)}\right) + Cov\left(h_A^{(1)}, h_D^{(2)}\right) \right. \\ &\quad \left. + Cov\left(h_D^{(1)}, h_A^{(2)}\right) + Cov\left(h_D^{(1)}, h_D^{(2)}\right) \right\} \end{aligned}$$

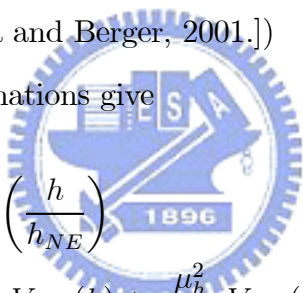
$$\begin{aligned}
B &= h_{NE}^2 - Z_\alpha^2 \text{Var}(h_{NE}) \\
&= h_{NE}^2 - Z_\alpha^2 \text{Var}(h_A^{(2)} + h_D^{(2)}) \\
&= h_{NE}^2 - Z_\alpha^2 \left\{ \text{Var}(h_A^{(2)}) + \text{Var}(h_D^{(2)}) + 2\text{Cov}(h_A^{(2)}, h_D^{(2)}) \right\}
\end{aligned}$$

$$\begin{aligned}
C &= h^2 - Z_\alpha^2 \text{Var}(h) \\
&= h^2 - Z_\alpha^2 \text{Var}(h_A^{(1)} + h_D^{(1)}) \\
&= h^2 - Z_\alpha^2 \left\{ \text{Var}(h_A^{(1)}) + \text{Var}(h_D^{(1)}) + 2\text{Cov}(h_A^{(1)}, h_D^{(1)}) \right\}
\end{aligned}$$

and Z_α is the $100 \times \left(1 - \frac{\alpha}{2}\right)$ percentile of the normal distribution (or, if sample numbers, n , were not large, of the student's t-distribution with $n-1$ degrees of freedom).

Method2(*delta method* [Casella and Berger, 2001.]

The first-order Taylor approximations give



$$\begin{aligned}
\text{Var}\left(1 - \frac{h}{h_{NE}}\right) &= \text{Var}\left(\frac{h}{h_{NE}}\right) \\
&\approx \frac{1}{\mu_{h_{NE}}^2} \text{Var}(h) + \frac{\mu_h^2}{\mu_{h_{NE}}^4} \text{Var}(h_{NE}) - 2\frac{\mu_h}{\mu_{h_{NE}}^3} \text{Cov}(h, h_{NE}) \\
&\approx \frac{1}{\mu_{h_{NE}}^2} \left\{ \text{Var}(h_A^{(1)}) + \text{Var}(h_D^{(1)}) + 2\text{Cov}(h_A^{(1)}, h_D^{(1)}) \right\} \\
&\quad + \frac{\mu_h^2}{\mu_{h_{NE}}^4} \left\{ \text{Var}(h_A^{(2)}) + \text{Var}(h_D^{(2)}) + 2\text{Cov}(h_A^{(2)}, h_D^{(2)}) \right\} \\
&\quad - 2\frac{\mu_h}{\mu_{h_{NE}}^3} \left\{ \text{Cov}(h_A^{(1)}, h_A^{(2)}) + \text{Cov}(h_A^{(1)}, h_D^{(2)}) \right. \\
&\quad \left. + \text{Cov}(h_D^{(1)}, h_A^{(2)}) + \text{Cov}(h_D^{(1)}, h_D^{(2)}) \right\}
\end{aligned}$$

We can use $\hat{h}_A^{(1)} + \hat{h}_D^{(1)}$ to estimate h in Method1 or μ_h in Method2 and use $\hat{h}_A^{(2)} + \hat{h}_D^{(2)}$ to estimate h_{NE} in Method1 or $\mu_{h_{NE}}$ in Method2. It is easy to estimate $\hat{h}_A^{(1)}$, $\hat{h}_D^{(1)}$, $\hat{h}_A^{(2)}$, and $\hat{h}_D^{(2)}$ by using the SOLAR computer package. But in both Method1 and Method2, we need $\text{Var}(\hat{h}_A^{(1)})$, $\text{Var}(\hat{h}_D^{(1)})$, $\text{Var}(\hat{h}_A^{(2)})$, $\text{Var}(\hat{h}_D^{(2)})$, $\text{Cov}(\hat{h}_A^{(1)}, \hat{h}_D^{(1)})$, $\text{Cov}(\hat{h}_A^{(2)}, \hat{h}_D^{(2)})$,

$Cov(\hat{h}_A^{(1)}, \hat{h}_A^{(2)})$, $Cov(\hat{h}_A^{(1)}, \hat{h}_D^{(2)})$, $Cov(\hat{h}_D^{(1)}, \hat{h}_A^{(2)})$, $Cov(\hat{h}_D^{(1)}, \hat{h}_D^{(2)})$ to estimate the remaining terms. Next, we will focus on deriving the estimator of the remaining terms.

Performing the above estimations involves $h_A^{(k)}$ and $h_D^{(k)}$, where $k = 1, 2$, that are related with σ_A^2 , σ_D^2 , σ_C^2 and σ_R^2 . To construct their relationship exactly, we let

$$h_1 = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2} (= h_A)$$

$$h_2 = \frac{\sigma_D^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2} (= h_D)$$

$$h_3 = \frac{\sigma_C^2}{\sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2}$$

$$h_4 = \sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_R^2$$

, i.e.

$$\sigma_A^2 = h_1 h_4, \sigma_D^2 = h_2 h_4, \sigma_C^2 = h_3 h_4, \sigma_R^2 = (1 - h_1 - h_2 - h_3) h_4$$

In other words, we make the 1-1 transformation between h_i s and σ_A^2 , σ_D^2 , σ_C^2 and σ_R^2 .

The following table shows the covariance components for relative pairs (Table 2):

Table 2. The covariance components for relative pairs

| Relationship | Covariance | V=Covariance after transformation |
|-------------------------|--|---|
| Same person | $\sigma_A^2 + \sigma_D^2 + \lambda\sigma_C^2 + \sigma_R^2$ | $h_1 h_4 + h_2 h_4 + \lambda h_3 h_4 + (1 - h_1 - h_2 - h_3) h_4$ |
| Parent-child | $\frac{1}{2}\sigma_A^2 + \lambda\sigma_C^2$ | $\frac{1}{2}h_1 h_4 + \lambda h_3 h_4$ |
| Full sibling | $\frac{1}{2}\sigma_A^2 + \frac{1}{4}\sigma_D^2 + \lambda\sigma_C^2$ | $\frac{1}{2}h_1 h_4 + \frac{1}{4}h_2 h_4 + \lambda h_3 h_4$ |
| Half sibling | $\frac{1}{4}\sigma_A^2 + \lambda\sigma_C^2$ | $\frac{1}{4}h_1 h_4 + \lambda h_3 h_4$ |
| Monozygous twins | $\sigma_A^2 + \sigma_D^2 + \lambda\sigma_C^2$ | $h_1 h_4 + h_2 h_4 + \lambda h_3 h_4$ |
| Grandparent-grandchild | $\frac{1}{4}\sigma_A^2 + \lambda\sigma_C^2$ | $\frac{1}{4}h_1 h_4 + \lambda h_3 h_4$ |
| Uncle/aunt-nephew/niece | $\frac{1}{4}\sigma_A^2 + \lambda\sigma_C^2$ | $\frac{1}{4}h_1 h_4 + \lambda h_3 h_4$ |
| First cousins | $\frac{1}{8}\sigma_A^2 + \lambda\sigma_C^2$ | $\frac{1}{8}h_1 h_4 + \lambda h_3 h_4$ |
| Double first cousins | $\frac{1}{4}\sigma_A^2 + \frac{1}{16}\sigma_D^2 + \lambda\sigma_C^2$ | $\frac{1}{4}h_1 h_4 + \frac{1}{16}h_2 h_4 + \lambda h_3 h_4$ |
| Spoused | $\lambda\sigma_C^2$ | $\lambda h_3 h_4$ |

Theorem 1 Suppose two models are $P_{ij} = x_{ij}^{(1)}\beta^{(1)} + G_{ij}^{(1)} + \varepsilon_{ij}^{(1)}$ and $P_{ij} = x_{ij}^{(2)}\beta^{(2)} + G_{ij}^{(2)} + \varepsilon_{ij}^{(2)}$, respectively, where $\varepsilon_{ij}^{(t)} \sim N\left(0, (\sigma_R^2)^{(t)}\right) \equiv N\left(0, \left(1 - h_1^{(t)} - h_2^{(t)} - h_3^{(t)}\right) h_4^{(t)}\right)$, $G_{ij}^{(t)} \sim N\left(0, (\sigma_A^2 + \sigma_D^2 + \sigma_C^2)^{(t)}\right) \equiv N\left(0, h_1^{(t)}h_4^{(t)} + h_2^{(t)}h_4^{(t)} + h_3^{(t)}h_4^{(t)}\right)$, and $\text{Cov}(G_{ij}, G_{ik}) [j \neq k] = (2\phi_{ij,ik}\sigma_A^2 + \Delta_{ij,ik}\sigma_D^2 + \lambda_{ij,ik}\sigma_C^2)^t \equiv 2\phi_{ij,ik}h_1^{(t)}h_4^{(t)} + \Delta_{ij,ik}h_2^{(t)}h_4^{(t)} + \lambda_{ij,ik}h_3^{(t)}h_4^{(t)}$. And P_{ij} is the observed value in the j th member of the i th family, x_{ij} is his/her corresponding covariate vector. Assumed R is the total number of family and there are n_i members in the i th family. Let $h^{(t)} = \left(h_1^{(t)}, h_2^{(t)}, h_3^{(t)}, h_4^{(t)}\right)$, then we have

$$\begin{aligned} & \text{Cov}\left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)}\right) \\ & \approx \left[\sum_{r=1}^R \left\{ \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \right. \right. \\ & \quad \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(k)} \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right\} \right] \\ & \times \left[\sum_{r=1}^R \left\{ \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right)' W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right\} \right] \\ & \times \left[\sum_{i=1}^R \left\{ \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right) \right. \right. \\ & \quad \left. \left. + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right\} \right] \end{aligned}$$

$$q = 1, 2, 3, 4 \quad q^* = 1, 2, 3, 4 \quad t = 1, 2 \quad t^* = 1, 2$$

where

$$S_r^{(t)} = \left(r_{r1}^{(t)} r_{r1}^{(t)}, r_{r1}^{(t)} r_{r2}^{(t)}, \dots, r_{r1}^{(t)} r_{rn_r}^{(t)}, \dots, r_{rn_r}^{(t)} r_{rn_r}^{(t)} \right)',$$

$$r_{rj}^{(t)} = P_{rj} - x_{rj}^{(t)}\beta^{(t)},$$

$$V_r^{(t)} = E\left(S_r^{(t)}; \beta^{(t)}, h^{(t)}\right) \text{ as given by Covariance after transformation in table I,}$$

$$W_{r \times r}^{(t)} = \begin{cases} 2\sigma_{ij}^{(t)2} & \text{for the } i, j \text{th and } l, m \text{th pairs} \\ \sigma_{il}^{(t)}\sigma_{im}^{(t)} + \sigma_{im}^{(t)}\sigma_{jl}^{(t)} & \text{for the } i, j \text{th and } l, m \text{th pairs} \end{cases},$$

$$\text{and } \frac{\partial W^{(t)}}{\partial h^{(t)}} = \begin{cases} 4\sigma_{ij} \frac{\partial \sigma_{ij}}{\partial h} & \text{for the } i, j\text{th and } l, m\text{th pairs} \\ \frac{\partial \sigma_{il}}{\partial h} \sigma_{jm} + \sigma_{il} \frac{\partial \sigma_{jm}}{\partial h} + \frac{\partial \sigma_{im}}{\partial h} \sigma_{jl} + \sigma_{im} \frac{\partial \sigma_{jl}}{\partial h} & \text{for the } i, j\text{th and } l, m\text{th pairs} \end{cases}$$

In the theorem, both $S_r^{(t)}$ and $V_r^{(t)}$ are vectors which their length are $\left[\left(\frac{n_r}{2}\right) + n_r\right]$, and $W_{r \times r}^{(t)}$ is a $\left[\left(\frac{n_r}{2}\right) + n_r\right] \times \left[\left(\frac{n_r}{2}\right) + n_r\right]$ matrix.

Proof. See Appendix. In the procedure, we have used Generalized Estimating Equations (GEE) method [Zeger and Liang, 1992; Amos, 1994], Taylor's expansion and some matrix operation [Harville, 1997]. ■

In our situation, $W^{(t)}$, $W^{(t^*)}$ and $\frac{\partial W^{(t)}}{\partial h_q^{(t)}}$ need to estimate. We estimate them with $\widehat{W}^{(t)}$, $\widehat{W}^{(t^*)}$ and $\frac{\partial \widehat{W}^{(t)}}{\partial h_q^{(t)}}$, where $\widehat{W}^{(t)}$, $\widehat{W}^{(t^*)}$ and $\frac{\partial \widehat{W}^{(t)}}{\partial h_q^{(t)}}$ are combination of \widehat{h}_1 , \widehat{h}_2 , \widehat{h}_3 and \widehat{h}_4 .

h_1 and h_2 are of our interest, so we only focus on the derivative of covariance components, related h_1 and h_2 ., for relative pairs. The following table shows the interested derivative of covariance components for relative pairs (Table 3):

Table 3. The derivative of covariance components for relative pairs

| Relationship | $\frac{\partial V}{\partial h_1}$ | $\frac{\partial V}{\partial h_2}$ | $\frac{\partial \tilde{V}}{\partial h_1}$ | $\frac{\partial \tilde{V}}{\partial h_2}$ |
|-------------------------|-----------------------------------|-----------------------------------|---|---|
| Same person | 0 | 0 | 0 | 0 |
| Parent-child | $\frac{1}{2}h_4$ | 0 | $\frac{1}{2}$ | 0 |
| Full sibling | $\frac{1}{2}h_4$ | $\frac{1}{4}h_4$ | $\frac{1}{2}$ | $\frac{1}{4}$ |
| Half sibling | $\frac{1}{4}h_4$ | 0 | $\frac{1}{4}$ | 0 |
| Monozygous twins | h_4 | h_4 | 1 | 1 |
| Grandparent-grandchild | $\frac{1}{4}h_4$ | 0 | $\frac{1}{4}$ | 0 |
| Uncle/aunt-nephew/niece | $\frac{1}{4}h_4$ | 0 | $\frac{1}{4}$ | 0 |
| First cousins | $\frac{1}{8}h_4$ | 0 | $\frac{1}{8}$ | 0 |
| Double first cousins | $\frac{1}{4}h_4$ | $\frac{1}{16}h_4$ | $\frac{1}{4}$ | $\frac{1}{16}$ |
| Spoused | 0 | 0 | 0 | 0 |

Corollary 2 Based on table 3, we can express the result of theorem 1 as follow:

$$\begin{aligned}
& Cov\left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)}\right) \\
& \approx \left[\sum_{r=1}^R \left\{ \widehat{h}_4^{(t)} \left(\frac{\partial \widetilde{V}_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \right. \right. \\
& \quad \left. \left. + \widehat{h}_4^{(t)} \left(\frac{\partial \widetilde{V}_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(k)} \left(\frac{\partial \widetilde{V}_r^{(t)}}{\partial h_q^{(t)}} \right) \widehat{h}_4^{(t)} \right\} \right] \\
& \times \left[\sum_{r=1}^R \left\{ \widehat{h}_4^{(t)} \left(\frac{\partial \widetilde{V}_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right)' W^{-1(t^*)} \left(\frac{\partial \widetilde{V}_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \widehat{h}_4^{(t^*)} \right\} \right] \\
& \times \left[\sum_{i=1}^R \left\{ \widehat{h}_4^{(t^*)} \left(\frac{\partial \widetilde{V}_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right) \right. \right. \\
& \quad \left. \left. + \widehat{h}_4^{(t^*)} \left(\frac{\partial \widetilde{V}_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(\frac{\partial \widetilde{V}_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \widehat{h}_4^{(t^*)} \right\} \right]
\end{aligned}$$

$$q = 1, 2, 3, 4 \quad q^* = 1, 2, 3, 4 \quad t = 1, 2 \quad t^* = 1, 2$$

In our case, $\left(\frac{\partial \widetilde{V}_r^{(t)}}{\partial h_q^{(t)}} \right) = \left(\frac{\partial \widetilde{V}_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)$ when $q = q^*$, where $t = 1, 2; t^* = 1, 2$.

Corollary2 is almost same with Theorem1 in its form. But one of the advantage of Corollary2 is that the time of performing the program in the computer is less than Theorem1.

Now, let two models be

$$P_{ij} = \alpha_H + \gamma_H E_{ij} + \tau_H Z_{ij} + G_{ij} + \epsilon \equiv x_{ij}^{(1)} \beta^{(1)} + G_{ij}^{(1)} + \varepsilon_{ij}^{(1)}$$

and

$$P_{ij} = \alpha_H + \tau_H Z_{ij} + G_{ij} + \epsilon \equiv x_{ij}^{(2)} \beta^{(2)} + G_{ij}^{(2)} + \varepsilon_{ij}^{(2)}$$

Under the same assumptions, we can apply above Theorem1 or Corollary2 to compute some needful estimators for using Fieller's theorem or delta method. Moreover, we can obtain the confidence limits of PHE or the estimator of deviation of PHE to perform a statistical test or to establish some criteria for determining whether E is an endophenotype.

3.3 HYPOTHESIS TEST

For having more statistical meanings of PHE , we utilize the confidence interval to get more informations about PHE . We hope to find a value that it means there exist a useful endophenotype when PHE value is larger than the value. That is, do one-sided confidence interval, corresponding to such test,

$$\begin{cases} H_0 : PHE = a \\ H_1 : PHE > a \end{cases}$$

Under significance α , we reject H_0 if the lower bound of one-sided confidence interval of PHE , $\widehat{PHE} - Z_{1-\alpha} \times s.e.(\widehat{PHE})$, is larger than a . However, we get the $100(1-\alpha)\%$ confidence limits of PHE when Fieller's Theorem was used. So we take the lower bound of confidence limits of $100(1-2\alpha)\%$ as the lower bound of $100(1-\alpha)\%$ one-sided confidence interval. In our following simulation, we considered some different values of the cutpoint, a . Under $\alpha=0.05$, we calculated the proportion that $(\widehat{PHE} - 1.645 \times s.e.(\widehat{PHE})_{delta})$ is larger than 0, 0.25, 0.50 and 0.75 respectively, where $s.e.(\widehat{PHE})_{delta}$ is the estimator of $s.e.(\widehat{PHE})$ by using delta method, and the proportion that the lower bound of 95% one-sided confidence interval with using Fieller theorem is larger than 0, 0.25, 0.50 and 0.75 respectively. Based on different values of the cutpoint in our simulation, we hope to construct some criteria to help us validate useful endophenotypes.

4 SIMULATION STUDIES

4.1 STUDY DESIGN

The simulation studies evaluate the utility of the proposed index, PHE, under two different scenarios (Figure 2). In Scenario I, the disease gene has a direct effect on phenotype and endophenotype. Scenario II allows multiple disease genes. At the same time, we try to show the relationship between the PHE values and the LOD-score curve. The study design is as follows. There are five markers, each marker has five allele, each allele has population frequency

0.2, and they are on the same chromosome with each of the four intervals between adjacent markers being 10 cM. The disease gene is located at the midpoint of the second interval and has two alleles. The population frequency of most common allele was 0.9. With SIMULATE [Ott 2002], that is a computer program originally written by Joseph Terwilliger, the loci of the markers and the disease gene were simulated based on above description.

Our simulations assumed both endophenotype and phenotype to be continuous measurements. The quantitative trait y and genes that influence it were assumed to have a linear relation as described in Almasy and Blangero [1998]:

$$y = \mu + \sum_{i=1}^n \eta_i + \epsilon ,$$

where μ was the grand mean, η_i was the random effect of the i th disease gene, and ϵ represented a random non-family deviation. η_i and ϵ were assumed to be normally distributed and uncorrelated. For these simulations, dominance effects and shared environmental effects were not included, and therefore $var(\eta_i) = \sigma_{A_i}^2$. For scenario I, each of E (endophenotype) and P (phenotype) was generated to have the single-gene contribution from G (disease gene) simulated by SIMULATE. The non-family deviation of E (ϵ_E) and the non-family deviation of P (ϵ_P) were assumed to have a correlation ρ_ϵ . The multiple gene effect in scenario II included the action of gene $G1$ (disease gene) on E and P , the single-gene action of $G2$ on E and the single-gene action of $G3$ on P .

The simulated data contained either 200 or 500 unclear families, and two sibships were generated for each family. In scenario I, the heritability of P due to G was assumed to be 0.42, and the heritability of E due to G allowed being 0, 0.15, 0.42 or 0.74. The correlation between non-family deviations of E and P , ρ_ϵ , was 0, or 0.5. In scenario II, there are two situations under our consideration. One is that the total heritability of P is larger than the total heritability of E , the other is, on the contrary, the total heritability of P is smaller than

the total heritability of E . The parameter values were shown as the following tables:

| | situations | | | | | | |
|-------------------------------------|------------|------|------|------|------|------|------|
| the heritability of E due to $G1$ | 0 | 0.15 | 0.42 | 0.51 | 0.74 | 0.74 | 0.79 |
| the heritability of P due to $G1$ | 0.42 | 0.42 | 0.42 | 0.42 | 0.42 | 0.42 | 0.42 |
| $G2$ (other heritability of E) | 0.3 | 0.25 | 0.12 | 0.04 | 0.05 | 0.08 | 0.02 |
| $G3$ (other heritability of P) | 0.17 | 0.17 | 0.17 | 0.17 | 0.41 | 0.41 | 0.41 |
| the total heritability of E | 0.3 | 0.4 | 0.52 | 0.55 | 0.79 | 0.82 | 0.81 |
| the total heritability of P | 0.59 | 0.59 | 0.59 | 0.59 | 0.83 | 0.83 | 0.83 |

| | situations | | | | | |
|-------------------------------------|------------|------|------|------|------|------|
| the heritability of E due to $G1$ | 0 | 0.15 | 0.42 | 0.62 | 0.74 | 0.74 |
| the heritability of P due to $G1$ | 0.42 | 0.42 | 0.42 | 0.42 | 0.42 | 0.42 |
| $G2$ (other heritability of E) | 0.7 | 0.59 | 0.23 | 0.23 | 0.08 | 0.21 |
| $G3$ (other heritability of P) | 0.17 | 0.17 | 0.17 | 0.17 | 0.17 | 0.17 |
| the total heritability of E | 0.7 | 0.74 | 0.65 | 0.85 | 0.82 | 0.95 |
| the total heritability of P | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 |

The correlation between non-family deviations of E and P , ρ_e , was the same as scenario I. Two hundred replications were performed for each specified situation. For simplicity, we denote the coordinates, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_e)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G3$, and ρ_e means the correlation between non-family deviations of E and P . We can find scenario I is a special case of scenario II if $h(G1) = G$, $h(G2_E) = 0$, and $h(G3_P) = 0$ in the coordinates' expression, where G is a single-gene (disease gene) in scenario I.

The computer package SOLAR (Sequential Oligogenic Linkage Analysis Routines) [Blangero et al, 2004; Almasy and Blangero, 1998] was used. The SOLAR command “simqtl” was used to simulate the data following two scenarios. The variance component analysis (24) was performed using the SOLAR command “polymod”. Besides, we use the SOLAR command “multipoint” to create the LOD-score. Before using the SOLAR command “multipoint”, we must set chromosome information about our markers. We set 0cM, 10cM, 20cM, 30cM and 40cM as the positions of the markers in the chromosome respectively, that is, we hoped that there is a high LOD-score peak at 15cM to find the disease gene. Also, the estimates of the standard error of PHE was calculated by using R software. And we plot the mean LOD-score curve according as the results from 200 replications.

4.2 Result

Table 6-9 contain results under scenario I. Table 10-13 contain results under scenario II with the total heritability of $P >$ the total heritability of E and Table 14-17 contain results under scenario II with the total heritability of $P <$ the total heritability of E .

4.2.1 PHE

Table 6 and Table 7 contain results under the ideal causal relation (scenario I). The heritability of P due to G was fixed. The higher the heritability of E due to G , the lower the heritability of P conditional on E and the closer the PHE values to 1. No matter that we chose the correlation between non-family deviations of E and P is either 0 or 0.5, the trend is still kept.

Table 10 and Table 11 show the results when there exist multiple disease genes under scenario II with the total heritability of $P >$ the total heritability of E . When the heritability of P due to $G1$ were fixed as 0.42 and the heritability of P due to $G3$ were fixed as 0.17 or 0.41, the trend, that the higher the heritability of E due to $G1$, the higher the PHE values, is consistent with scenario I. Under scenario II with the total heritability of $P <$ the total heritability of E , Table 14 and Table 15 show a similar trend between the heritability of E due to $G1$ and PHE . However, we can find these values, the heritability of P due to $G3$ and

the heritability of E due to $G2$, also influence the PHE values. The higher the heritability of P due to $G3$ or the heritability of E due to $G2$, the lower the PHE values. Besides, the involvement of $\rho_\epsilon = 0.5$ leads the PHE values to be disrupted. That is, it reduces the efficiency to use the PHE values for searching a useful endophenotype.

4.2.2 THE ACCURACY OF THE ESTIMATORS OF THE STANDARD ERROR OF PHE

To check the accuracy of the estimators of the standard error of PHE calculated according to the delta method or the Fieller's theorem and our provided theorem or corollary, we compare the standard error of proportion of heritability explained by endophenotype(s.e) that was simulated with s.e(delta) and s.e(Fieller), where s.e(delta) is the mean of estimator of s.e by using delta method and s.e(Fieller) is the mean of the range of 95% confidence limits of PHE , used by Fieller method, divided by 2×1.96 . Table 6, Table 7, Table 10, Table 11, Table 14 and Table 15 contain these results under scenario I and scenario II. Let us regard the standard error of proportion of heritability explained by endophenotype(s.e), that was simulated, as the true standard deviation of proportion of heritability explained by endophenotype. We can find that, when the heritability of E due to the disease gene is lower than the heritability of P due to the shared gene, s.e(delta) and s.e(Fieller) tend to be overestimated. And s.e(delta) and s.e(Fieller) tend to be underestimated when the heritability of E due to the disease gene is higher than the heritability of P due to the shared gene. Also, we find that the relative error of the overestimators is larger than the relative error of the underestimators. But both the absolute error of the overestimators and the underestimators are small. That is, these estimators of the standard error of PHE are closer the true standard error of PHE . However, using these estimators calculated by either delta method or Fieller theorem don't have too wide confidence interval of PHE to make some statistical inferences. In other words, these estimators can be allowed.

4.2.3 TEST OF PHE

For using normal distribution to perform statistical tests or establish a confidence interval of PHE , we used Shapiro-Wilk statistic to test the normality of PHE . Table 6, Table 7, Table 10, Table 11, Table 14 and Table 15 also shows these p-values of using Shapiro-Wilk test under scenario I and scenario II. The histograms of PHE values under different situations are shown in Figure 3-14. Under scenario I, the normality of PHE doesn't hold in most situations. But the normality of PHE holds in most situations under scenario II. In other words, although using normal distribution is not good, it isn't too bad. Briefly, using normal distribution can be acceptable with a lower standard..

We first describe the information about the mean LOD-score curve under both scenario I and scenario II (Figure 15-26). The LOD-score in our simulation was found to be related to the number of families and the heritability of the trait due to the common disease gene, where the trait may be a phenotype or an endophenotype. When the heritability of the endophenotype due to the common disease gene is larger than the heritability of the phenotype due to the common disease gene, the endophenotype is useful to search the disease gene. We except that the heritability of the endophenotype due to the common disease gene isn't smaller than the heritability of the phenotype due to the disease gene. These results were consistent to the results from other papers [Almasy and Blangero, 1998; Williams et al., 1999].

Table 8 and Table 9 contain results under scenario I. At the same time, Figure 15, Figure 16, Figure 17 and Figure 18 show the mean LOD-score curve under scenario I. Based on these figures, when the heritability of P due G was assumed to be 0.42 and the heritability of E due to G allowed being 0 and 0.15, we find that using endophenotype to search for the disease gene is worse than using phenotype because the mean LOD-score of P was higher than the mean LOD-score of E . That is, we don't hope that these are endophenotypes. On the other hand, when the heritability of P due to G was assumed to be 0.42 and the heritability of E due to G was assumed to be 0.74, endophenotype-based genetic analysis is more likely to succeed than one in terms of search for the disease gene (i.e. the mean LOD-score of E is higher than the one of P). Besides, when the heritability of P due to G was assumed to be 0.42 and the heritability of E due to G was assumed to be 0.42, the phenotype-based effect

and the endophenotype-based effect are same. Altogether, when the heritability of P due to G was assumed to be 0.42 and the heritability of E due to G was assumed to be 0.42 or 0.74, the endophenotype-based effect isn't worse than the phenotype-based effect. As a result of above descriptions about mean LOD-score curve, based on Table A3 and Table A4, we view it endophenotype candidate if lower bound of 95% one-sided confidence interval is larger than 0.25 or 0.50. The criterion that lower bound of 95% one-sided confidence interval is larger than 0.50 can be seen as a stronger evidence and the criterion that lower bound of 95% one-sided confidence interval is larger than 0.25 is also a suitable frame of reference. With another viewpoint, using two cutpoints, 0.25 and 0.50, the power, that the probability of rejecting H_0 when H_1 holds, will exceed 0.7 or 0.8 except for the situation where the heritability of P due to G was assumed to be 0.42, the heritability of E due to G was assumed to be 0.42, ρ_ϵ was assumed to be 0, and cutpoint is set as 0.50. It implies that endophenotype-based effect isn't worse than the phenotype-based effect. If it is desired that there is a higher power such as 0.9, 0 may be an applicable cutpoint no matter ρ_ϵ was either 0 or 0.5. But it also leads the result that endophenotype-based effect is worse than the phenotype-based effect, happen, such as the situation where the heritability of P due G was assumed to be 0.42 and the heritability of E due to G allowed being 0.15.

In scenario II, on account of disrupted PHE values with the heritability of P due to $G3$ and the heritability of E due to $G2$, the criteria under scenario I may become improper. Based on Table 12, Table 13, Table 16 and Table 17, we downscale the standard of these criteria for searching the endophenotype successfully. The criterion that lower bound of 95% one-sided confidence interval is larger than 0.25 is still a suitable one. But many useful endophenotypes will be missed. So, we find that the criterion that lower bound of 95% one-sided confidence interval is larger than 0 should be seen as the criterion that search the potential candidate of endophenotype. Furthermore, if we want to let the higher power be kept for the goal that endophenotype-based effect isn't worse than the phenotype-based effect, considered cutpoint may be 0. However, if ρ_ϵ was assumed to be 0.5, the chosen cutpoint, 0, is not sufficient because of the lower power.

In summary, three criteria are provided as follows. The first criterion that lower bound of

95% one-sided confidence interval is larger than 0 is the potential evidence for searching the endophenotype. The second criterion that lower bound of 95% one-sided confidence interval is larger than 0.25 is the moderate evidence for searching the endophenotype. And the third criterion that lower bound of 95% one-sided confidence interval is larger than 0.50 is the stronger evidence for searching the endophenotype. However, you can choose some different criteria depended on the different goals of different cases or use lower bound of 95% one-sided confidence interval directly as the evidence for searching the endophenotype.

In another aspect, using the viewpoint of "power", we try to construct some steps to help us determine the desired endophenotype. The process of our construction is as follows. At the first step, check if ρ_ϵ is 0 because it brings different information about use of the PHE values. If it doesn't hold, we are careful with use of *PHE* values because there is a lower power of detecting the useful endophenotypes if ρ_ϵ is 0.5 even when the cutpoint is set as 0. That is, the involvement of $\rho_\epsilon \neq 0$ leads much uncertainty to use PHE values. Furthermore, if ρ_ϵ become larger, using the PHE values may loss much useful information of the endophenotypes. In other words, If the lower bound of 95% one-sided confidence interval isn't larger than 0 when ρ_ϵ is larger than 0, it doesn't imply that the endophenotype is helpless. If ρ_ϵ is 0, we will perform the second step.

At the second step, check if the lower bound of 95% one-sided confidence interval is larger than 0.25. If it holds, it implies two possibilities : (1) there is the single disease gene to lead a direct effect on phenotype and endophenotype such as Scenario I and endophenotype-based effect isn't worse than the phenotype-based effect; (2) it implies that both the influences of other genes on phenotype and endophenotype can be small, relative to the influences of the shared genes on phenotype and endophenotype such as Scenario II and endophenotype-based effect is better than the phenotype-based effect. If the lower bound of 95% one-sided confidence interval isn't larger than 0.25, we will proceed to perform the third step.

At the third step, check if the lower bound of 95% one-sided confidence interval is larger than 0. If it holds, there exists two possible situations : (1) there is the single disease gene to lead a direct effect on phenotype and endophenotype such as scenario I and endophenotype-based effect isn't better than the phenotype-based effect. It is out of our desire; (2) the

influence of other genes of either phenotype or endophenotype can be large relatively to the influence of the shared genes of either phenotype or endophenotype respectively such as scenario II and endophenotype-based effect isn't worse than the phenotype-based effect. If the lower bound of 95% one-sided confidence interval isn't larger than 0 when ρ_ϵ is 0, it means there is a high probability that it isn't a useful endophenotype. In sum, using three steps is helpful to search a useful endophenotype.

5 DISCUSSION

Based on definition of an endophenotype proposed by Huang et al. [2005], we have attempted to provide criteria that can be used to validate an endophenotype. Huang et al.,2005 had shown that the proposed index, *PHE*, is useful in validating endophenotypes. In our report, we use *PHE* proposed by Huang et al. [2005] as the index for evaluating endophenotypes to provide more clear informations, three criteria and three steps, through the one-sided confidence interval or the statistical test. However, we can find that the more the total numbers of family members, the more efficiency of detecting a useful endophenotype.

As discussed in corresponding index for validating surrogate endpoints such as *PTE*, confidence intervals of *PTE* can be calculated using Fieller's theorem [Buyse and Molenberghs, 1998], however, they are usually too wide to be useful. With our proposed theorem or corollary, we use Fieller's theorem or delta method to calculate confidence intervals of *PHE*. Our simulation results show that the estimators of standard error of *PHE* values' estimators are near "true" standard errors of these indices' estimators. That is, they are quite reasonable to avoid too wide confidence interval to be useful. However, although they may be overestimated or underestimated, they are helpful to detect the useful endophenotype easily. This is because that it tends to have a underestimator of standard error of *PHE* estimator for the good endophenotype and it leads the lower bound of 95% one-sided confidence interval to be easily larger than our set cutpoint. Otherwise, the lower bound of 95% one-sided confidence interval tend to be smaller than our set cutpoint for the useless endophenotype. In other words, it isn't too serious for using these overestimated or underestimated estimators of standard error of

PHE values' estimators to construct a reasonable one-sided confidence interval and to search a useful endophenotype.

Besides, our simulation results show that the multiple gene effect lowers *PHE* values to lead it confused for evaluating endophenotypes. We provide three criteria and three steps to help us understand the pattern of *PHE* values versus the relationship between endophenotype and phenotype. If you aren't interested in the relationship between *PHE* values and the heritabilities caused by different genes, the second step can be omitted. However, among three steps, we need to check that ρ_ϵ is 0. The SOLAR command "polygenic" can be used to calculate ρ_ϵ . If ρ_ϵ is near 0, we can view it 0 to use three criteria and three steps safely for searching a useful endophenotype. Furthermore, at the third step, we will face the situation that the influence of other genes of either phenotype or endophenotype can be large relatively to the influence of the shared genes of either phenotype or endophenotype respectively such as Scenario II. For the influence of other genes of phenotype or endophenotype, we can use linkage analysis to determine which heritability is relatively large. If the heritability of other genes of phenotype is relatively large to the heritability of the found disease gene of phenotype, it means that only using an endophenotype may be sufficient. We must to search more than one endophenotype to capture a complete feature of the specified phenotype. The following model can be tried to be considered.

$$P = \alpha_H + \gamma_{1H}E1 + \gamma_{2H}E2 + \tau_H Z + G + \epsilon,$$

where *E1* is assumed to being a found endophenotype and *E2* is assumed to being a new or interested endophenotype. And we calculate the *PHE* value, $1 - \frac{h_{E1E2}}{h_{NE}}$, directly and its lower bound of 95% one-sided confidence interval, where h_{E1E2} is the heritability calculated from the variance component analysis (24) including the endophenotypes, *E1* and *E2*, with any other covariates. To avoid to get same information or to find similar endophenotypes, we also calculate the partial proportion of heritability explained (*PPHE*) by the endophenotype defined as

$$PPHE = 1 - \frac{h_{E1E2}}{h_{E1}}$$

where h_{E1E2} is the heritability calculated from the variance component analysis (24) including the endophenotypes, $E1$ and $E2$, with any other covariates and h_{E1} is the heritability calculated from the variance component analysis (24) without including the endophenotype $E2$ with any other covariates. A good and new endophenotype is one that explains a large proportion of heritability given a found endophenotype $E1$, thus, the greater the $PPHE$ value, the more likely $E2$ an desired endophenotype.

In the future, to make it clear for using the PHE values, especially when $\rho_\epsilon \neq 0$, we should simulate with $\rho_\epsilon < 0$ and $\rho_\epsilon \gg 0$. The information of the PHE values involved with negative ρ_ϵ is a loss of our report. However, the much higher ρ_ϵ is considered to help us understand the efficiency of using the PHE values to detect a useful endophenotype clearly in a bad situation. If the power of using the PHE values to detect useful endophenotype candidates isn't too low when ρ_ϵ is a much larger value, PHE values will be very useful index to search a useful endophenotype to increase opportunities of finding susceptible disease genes.



Appendix:

Let model1: $P_{ij} = x'_{ij} \beta^{(1)} + G_{ij}^{(1)} + \varepsilon_{ij}^{(1)}$ and model2: $P_{ij} = x'_{ij} \beta^{(2)} + G_{ij}^{(2)} + \varepsilon_{ij}^{(2)}$, where

$$\varepsilon_{ij}^{(t)} \sim N\left(0, (\sigma_R^2)^{(t)}\right) \equiv N\left(0, \left(1 - h_1^{(t)} - h_2^{(t)} - h_3^{(t)}\right) h_4^{(t)}\right)$$

$$G_{ij}^{(t)} \sim N\left(0, (\sigma_A^2 + \sigma_D^2 + \sigma_C^2)^{(t)}\right) \equiv N\left(0, h_1^{(t)} h_4^{(t)} + h_2^{(t)} h_4^{(t)} + h_3^{(t)} h_4^{(t)}\right)$$

$$\text{Cov}(G_{ij}, G_{ik}) [j \neq k] = \left(2\phi_{ij,ik} \sigma_A^2 + \Delta_{ij,ik} \sigma_D^2 + \lambda_{ij,ik} \sigma_C^2\right)^t$$

$$\equiv 2\phi_{ij,ik} h_1^{(t)} h_4^{(t)} + \Delta_{ij,ik} h_2^{(t)} h_4^{(t)} + \lambda_{ij,ik} h_3^{(t)} h_4^{(t)}$$

$$t = 1, 2$$

By GEE [Zeger and Liang, 1992; Amos, 1994],

$$S_{\beta^{(t)}}\left(\beta^{(t)}, h^{(t)}\right) = \sum_{r=1}^R \left(\frac{\partial X_r^{(t)} \beta^{(t)}}{\partial \beta^{(t)}}\right)' \text{Cov}^{-1}(P_r) \left(P_r - X_r^{(t)} \beta^{(t)}\right) = 0$$

where $P_r = (P_{r1}, \dots, P_{rn_r})'$, and $X_r^{(t)} = (x_{r1}^{(t)}, \dots, x_{rn_r}^{(t)})'$.

The correlation parameter h may be estimated by simultaneously solving

$$S_{\beta^{(t)}}\left(\beta^{(t)}, h^{(t)}\right) = 0$$

and

$$S_{h^{(t)}}\left(\beta^{(t)}, h^{(t)}\right) = \sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}}\right)' W^{-1(t)} \left(S_r^{(t)} - V_r^{(t)}\right) = 0$$

where

$$S_r^{(t)} = \left(r_{r1}^{(t)} r_{r1}^{(t)}, r_{r1}^{(t)} r_{r2}^{(t)}, \dots, r_{r1}^{(t)} r_{rn_r}^{(t)}, \dots, r_{rn_r}^{(t)} r_{rn_r}^{(t)}\right)'$$

$$r_{rj}^{(t)} = P_{rj} - x'_{rj} \beta^{(t)},$$

$$V_r^{(t)} = E\left(S_r^{(t)}; \beta^{(t)}, h^{(t)}\right) \text{ as given by Covariance after transformation in table I,}$$

$$\text{and } W_{r \times r}^{(t)} = \begin{cases} 2\sigma_{ij}^{(t)2} & \text{for the } i, j\text{th and } l, m\text{th pairs} \\ \sigma_{il}^{(t)}\sigma_{im}^{(t)} + \sigma_{im}^{(t)}\sigma_{jl}^{(t)} & \text{for the } i, j\text{th and } l, m\text{th pairs} \end{cases}$$

Since

$$S_{h^{(t)}}(\beta^{(t)}, h^{(t)}) = \sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' W^{-1(t)} (S_r^{(t)} - V_r^{(t)})$$

and

$$\frac{\partial^2 V_r^{(t)}}{\partial (h^{(t)})^2} = 0,$$

we have

$$\begin{aligned} & \frac{S_{h^{(t)}}(\beta^{(t)}, h^{(t)})}{\partial h^{(t)}} \\ &= \sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' \left(\frac{\partial W^{-1(t)}}{\partial h^{(t)}} \right) (S_r^{(t)} - V_r^{(t)}) + \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right) \right] \\ &= \sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' \left(-W^{-1(t)} \frac{\partial W^{(t)}}{\partial h^{(t)}} W^{-1(t)} \right) (S_r^{(t)} - V_r^{(t)}) + \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right) \right]. \end{aligned}$$

where

$$\frac{\partial W^{(t)}}{\partial h^{(t)}} = \begin{cases} 4\sigma_{ij} \frac{\partial \sigma_{ij}}{\partial h} & \text{for the } i, j\text{th and } l, m\text{th pairs} \\ \frac{\partial \sigma_{il}}{\partial h} \sigma_{jm} + \sigma_{il} \frac{\partial \sigma_{jm}}{\partial h} + \frac{\partial \sigma_{im}}{\partial h} \sigma_{jl} + \sigma_{im} \frac{\partial \sigma_{jl}}{\partial h} & \text{for the } i, j\text{th and } l, m\text{th pairs} \end{cases}$$

Using Taylor's expansion, we have

$$\begin{aligned} & \widehat{h}^{(k)} - h^{(k)} \\ &= \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' \left(-W^{-1(t)} \frac{\partial W^{(t)}}{\partial h^{(t)}} W^{-1(t)} \right) (S_r^{(t)} - V_r^{(t)}) \right. \right. \\ & \quad \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right) \right] \right)^{-1} \\ & \quad \times \left(\sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h^{(t)}} \right)' W^{-1(t)} (S_r^{(t)} - V_r^{(t)}) \right) \end{aligned}$$

According to above equation, we can obtain

$$\begin{aligned}
& Cov \left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)} \right) \\
&= Cov \left(\widehat{h}_q^{(t)} - h_q^{(t)}, \widehat{h}_{q^*}^{(t^*)} - h_{q^*}^{(t^*)} \right) \\
&= Cov \left\{ \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(-W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) (S_r^{(t)} - V_r^{(t)}) \right. \right. \right. \\
&\quad \left. \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)^{-1} \times \left(\sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} (S_r^{(t)} - V_r^{(t)}) \right), \right. \\
&\quad \left. \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(-W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) (S_r^{(t^*)} - V_r^{(t^*)}) \right. \right. \right. \\
&\quad \left. \left. \left. + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(-\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)^{-1} \times \left(\sum_{r=1}^R \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} (S_r^{(t^*)} - V_r^{(t^*)}) \right) \right\}
\end{aligned}$$

Note that

$$\left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(-W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) (S_r^{(t)} - V_r^{(t)}) + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)^{-1}$$

and

$$\left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(-W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) (S_r^{(t^*)} - V_r^{(t^*)}) + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(-\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)^{-1}$$

are 1×1 matrices.

Besides, for simplicity, we can replace them with

$$\left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(-W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) (\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)}) + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(-\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)$$

and

$$\left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(-W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) (\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)}) + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(-\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)$$

then

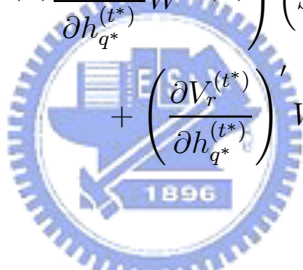
$$\begin{aligned}
& Cov\left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)}\right) \\
& \approx \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \right. \right. \\
& \quad \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)^{-1} \times \\
& \left[\sum_{r=1}^R Cov \left(\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} (S_r^{(t)} - V_r^{(t)}), \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} (S_r^{(t^*)} - V_r^{(t^*)}) \right) \right] \times \\
& \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right) \right. \right. \\
& \quad \left. \left. + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)^{-1}
\end{aligned}$$

Since $W^{(t)}$ and $W^{(t^*)}$ are symmetric matrices, $W^{-1(t)}$ and $W^{-1(t^*)}$ are also symmetric matrices. Above equation can be written as

$$\begin{aligned}
& Cov\left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)}\right) \\
& \approx \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \right. \right. \\
& \quad \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)^{-1} \times \\
& \left[\sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} Cov(S_r^{(t)} - V_r^{(t)}, S_r^{(t^*)} - V_r^{(t^*)}) W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \times \\
& \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right) \right. \right. \\
& \quad \left. \left. + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)^{-1}
\end{aligned}$$

We estimate $Cov(S_r^{(t)} - V_r^{(t)}, S_r^{(t^*)} - V_r^{(t^*)})$ with $(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)}) (\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)})$, then we

obtain

$$\begin{aligned}
& Cov\left(\widehat{h}_q^{(t)}, \widehat{h}_{q^*}^{(t^*)}\right) \\
& \approx \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' \left(W^{-1(t)} \frac{\partial W^{(t)}}{\partial h_q^{(t)}} W^{-1(t)} \right) \left(\widehat{S}_r^{(t)} - \widehat{V}_r^{(t)} \right) \right. \right. \\
& \qquad \qquad \qquad \left. \left. + \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right) \right] \right)^{-1} \times \\
& \left[\sum_{r=1}^R \left(\frac{\partial V_r^{(t)}}{\partial h_q^{(t)}} \right)' W^{-1(t)} (S_r^{(t)} - V_r^{(t)}) (S_r^{(t^*)} - V_r^{(t^*)}) W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \times \\
& \left(\sum_{r=1}^R \left[\left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' \left(W^{-1(t^*)} \frac{\partial W^{(t^*)}}{\partial h_{q^*}^{(t^*)}} W^{-1(t^*)} \right) \left(\widehat{S}_r^{(t^*)} - \widehat{V}_r^{(t^*)} \right) \right. \right. \\
& \qquad \qquad \qquad \left. \left. + \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right)' W^{-1(t^*)} \left(\frac{\partial V_r^{(t^*)}}{\partial h_{q^*}^{(t^*)}} \right) \right] \right)^{-1}
\end{aligned}$$


References

- [1] Almasy L, Blangero J. 1998. Multipoint quantitative-trait linkage analysis in general pedigrees. *Am J Hum Genet* 62:1198–1211.
- [2] Amos CI. 1994. Robust variance-components approach for assessing genetic linkage in pedigrees. *Am J Hum Genet* 54:535-543.
- [3] Beaty TH, Self SG, Liang KY, Connolly MA, Chase GA, Kwitrovich PO. 1985. Use of robust variance components of models to analyse triglyceride data in families. *Ann Hun Genet* 49:315-328.
- [4] Begg C, Leung DHY. 2000. On the use of surrogate end points in randomized trials (with comments). *JRSS A* 163:15–28.

- [5] Blangero J, Lange K, Dyer T, Almasy L, Göring H, Williams J, Charles Peterson C. 2004. SOLAR v.2.1.4. <http://www.sfbr.org/solar/index.html>. Southwest Foundation for Biomedical Research, San Antonio.
- [6] Breslow NE, Clayton DG. 1993. Approximate inference in generalized linear mixed models. *J Am Stat Assoc* 88:9–25.
- [7] Burton P.R., Tobin M.D., 2003, *Handbook of Statistical Genetics*, 2nd edition, Balding D.J., Bishop M. and Cannings C eds. John Wiley & Sons, Ltd, pp. 855-879.
- [8] Buyse M, Molenberghs G. 1998. Criteria for validation of surrogate endpoints in randomized experiments. *Biometrics* 54:1014–1029.
- [9] Buyse M, Molenberghs G, Buzykowski T, Renard D, Geys H. 2000. The validation of surrogate endpoints in meta-analyses of randomized experiments. *Biostatistics* 1:49–67.
- [10] Casella G, Berger RL. 2001. *Statistical Inference*, 2nd edition. pp. 240-245
- [11] Chen C, Wang H, Snapinn SM. 2003. Proportion of treatment effect (PTE) explained by a surrogate marker. *Stat Med* 22:3449-3459.
- [12] De Gruttola VG, Clax P, DeMets DL, Downing GJ, Ellenberg SS, Friedman L, Gail MH, Prentice R, Wittes J, Zeger SL. 2001. Considerations in the evaluation of surrogate endpoints in clinical trials: Summary of a National Institutes of Health workshop. *Control Clin Trials* 22:485–502.
- [13] Duggirala R, Williams JT, Williams-Blangero S, Blangero J. 1997. A variance component approach to dichotomous trait linkage analysis using a threshold model. *Genet Epidemiol* 14:987-992.
- [14] Falconer DS. 1989. *Introduction to Quantitative Genetics*, Third edn. John Wiley & Sons, New York.
- [15] Freedman LS, Graubard BI, Schatzkin A. 1992. Statistical validation of intermediate endpoints for chronic diseases. *Stat Med* 11:167–178.

- [16] Gottesman II, Gould TD. 2003. The endophenotype concept in psychiatry: etymology and strategic intentions. *Am J Psychiatry* 160:636–645.
- [17] Harville DA. 1997. Matrix algebra from a statistician's perspective. Springer, New York , pp285-331
- [18] Herson J. 1975. Fieller's theorem versus the delta method for significance intervals for ratios. *J Statist Comput Simul*:265-274
- [19] Hopper JL. 2002. In *Biostatistical Genetics and Genetic Epidemiology*, Elston,Olson and Palmer eds. Wiley, Chichester, pp. 371–372.
- [20] Hopper JL. 2002. In *Biostatistical Genetics and Genetic Epidemiology*, Elston,Olson and Palmer eds. Wiley, Chichester, pp. 778-788
- [21] Huang GH, Chen CH, Chen WJ. 2005 . Statistical Validation of Endophenotypes Using a Surrogate Endpoint Analytic Analogue with Application to Schizophrenia.
- [22] Iachine I. 2004. Statistical Methods in Genetic Epidemiology. <http://statmaster.sdu.dk/courses/st115>.
- [23] McCullagh P, Nelder JA. 1989. *Generalized Linear Models*, 2nd edition. Chapman and Hall,London.
- [24] Ott J. 2002. Documentation to the SIMULATE program. <http://linkage.rockefeller.edu/ott/simulate.htm>. Rockefeller University New York.
- [25] Prentice RL. 1989. Surrogate endpoints in clinical trials: definitions and operational criteria. *Stat Med* 8:431–440.
- [26] Willians JT, Eerdewegh PV, Almasy L, Blangero J. 1999. Joint Multipoint Analysis of Multivariate Qualitative and Quantitative Traits. I. Likelihood Formulation and Simulation Results.
- [27] Zeger SL, Liang KY. 1992. An overview of methods for the analysis of longitudinal data, *Stat Med* 11:1825-1839.

TABLE 6. Simulation results based on scenario I (1)

| <i>No. of families</i> | h_P^a | h_E^a | ρ_ϵ^b | h^c | PHE^d | $s.e^e$ | $s.e(\delta)^f$ | $s.e(Fieller)^g$ | $S.W - pvalue^h$ |
|------------------------|---------|---------|-------------------|-------|---------|---------|-----------------|------------------|------------------|
| 200 | 0.42 | 0 | 0 | 0.405 | -0.002 | 0.009 | 0.025 | 0.029 | < 0.001 |
| | | | 0.5 | 0.473 | -0.201 | 0.138 | 0.215 | 0.271 | < 0.001 |
| | 0.15 | 0 | 0 | 0.337 | 0.202 | 0.079 | 0.128 | 0.154 | < 0.001 |
| | | | 0.5 | 0.269 | 0.322 | 0.158 | 0.151 | 0.234 | 0.039 |
| | 0.42 | 0 | 0 | 0.183 | 0.562 | 0.138 | 0.107 | 0.204 | 0.698 |
| | | | 0.5 | 0.075 | 0.816 | 0.149 | 0.087 | 0.118 | < 0.001 |
| | 0.74 | 0 | 0 | 0.053 | 0.875 | 0.125 | 0.084 | 0.094 | < 0.001 |
| | | | 0.5 | 0.028 | 0.937 | 0.093 | 0.075 | 0.088 | < 0.001 |

^a h_P =heritability of P due to G; h_E = heritability of E due to G

^b ρ_ϵ =correlation between non-family deviations of E and P

^c h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(\delta)$ =mean of estimator of $s.e$ by delta method

^g $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times$ the range of confidence limits of PHE) by Fieller theorem

^h $S.W - pvalue$ =p value of using Shapiro-Wilk Test



TABLE 7. Simulation results based on scenario I (2)

| <i>No. of families</i> | h_P^a | h_E^a | ρ_ϵ^b | h^c | PHE^d | $s.e^e$ | $s.e(\delta)^f$ | $s.e(Fieller)^g$ | $S.W - pvalue^h$ |
|------------------------|---------|---------|-------------------|-------|---------|---------|-----------------|------------------|------------------|
| 500 | 0.42 | 0 | 0 | 0.422 | -0.0004 | 0.002 | 0.007 | 0.008 | < 0.001 |
| | | | 0.5 | 0.481 | -0.173 | 0.071 | 0.117 | 0.122 | < 0.001 |
| | 0.15 | 0 | 0 | 0.339 | 0.189 | 0.042 | 0.074 | 0.076 | 0.001 |
| | | | 0.5 | 0.282 | 0.331 | 0.081 | 0.084 | 0.088 | 0.282 |
| | 0.42 | 0 | 0 | 0.187 | 0.552 | 0.084 | 0.066 | 0.068 | 0.012 |
| | | | 0.5 | 0.076 | 0.817 | 0.092 | 0.050 | 0.052 | 0.003 |
| | 0.74 | 0 | 0 | 0.048 | 0.889 | 0.079 | 0.048 | 0.049 | < 0.001 |
| | | | 0.5 | 0.017 | 0.959 | 0.053 | 0.045 | 0.046 | < 0.001 |

^a h_P =heritability of P due to G; h_E = heritability of E due to G

^b ρ_ϵ =correlation between non-family deviations of E and P

^c h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(\delta)$ =mean of estimator of $s.e$ by delta method

^g $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times \text{the range of confidence limits of PHE})$ by Fieller theorem

^h $S.W - pvalue$ =p value of using Shapiro-Wilk Test



TABLE 8. Simulation results based on scenario I (3)

| <i>No. of families</i> | h_P^a | h_E^a | ρ_ϵ^b | delta method | | | | Fieller theorem | | | |
|------------------------|---------|---------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^c$ | $D0.25^c$ | $D0.50^c$ | $D0.75^c$ | $F0.00^d$ | $F0.25^d$ | $F0.50^d$ | $F0.75^d$ |
| 200 | 0.42 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0.01 | 0.005 | 0 | 0 |
| | 0.15 | 0 | 0.55 | 0 | 0 | 0 | 0.395 | 0.01 | 0.01 | 0.01 | |
| | | | 0.5 | 0.715 | 0.195 | 0.01 | 0 | 0.56 | 0.115 | 0.01 | 0 |
| | 0.42 | 0 | 0.99 | 0.815 | 0.255 | 0 | 0.95 | 0.71 | 0.19 | 0 | |
| | | | 0.5 | 0.995 | 0.98 | 0.825 | 0.365 | 0.99 | 0.945 | 0.8 | 0.34 |
| | 0.74 | 0 | 1 | 1 | 0.945 | 0.52 | 1 | 0.995 | 0.9 | 0.515 | |
| | | | 0.5 | 1 | 1 | 0.99 | 0.78 | 0.995 | 0.99 | 0.99 | 0.765 |

^a h_P =heritability of P due to G; h_E = heritability of E due to G

^b ρ_ϵ =correlation between non-family deviations of E and P

^c Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta})^e$ is larger than x ;

^d Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^e $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

TABLE 9. Simulation results based on scenario I (4)

| <i>No. of families</i> | h_P^a | h_E^a | ρ_ϵ^b | delta method | | | | Fieller theorem | | | |
|------------------------|---------|---------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^c$ | $D0.25^c$ | $D0.50^c$ | $D0.75^c$ | $F0.00^d$ | $F0.25^d$ | $F0.50^d$ | $F0.75^d$ |
| 500 | 0.42 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.15 | 0 | 0.935 | 0 | 0 | 0 | 0.89 | 0 | 0 | 0 | |
| | | | 0.5 | 0.975 | 0.28 | 0 | 0 | 0.945 | 0.24 | 0 | 0 |
| | 0.42 | 0 | 1 | 0.985 | 0.26 | 0.005 | 1 | 0.98 | 0.22 | 0.005 | |
| | | | 0.5 | 1 | 1 | 0.995 | 0.4 | 1 | 1 | 0.985 | 0.39 |
| | 0.74 | 0 | 1 | 1 | 1 | 0.74 | 1 | 1 | 1 | 0.725 | |
| | | | 0.5 | 1 | 1 | 1 | 0.975 | 1 | 1 | 1 | 0.965 |

^a h_P =heritability of P due to G; h_E = heritability of E due to G

^b ρ_ϵ =correlation between non-family deviations of E and P

^c Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta})^e$ is larger than x ;

^d Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^e $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

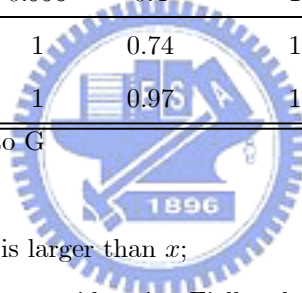


TABLE 10. Simulation results based on scenario II with P>E (1)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | h^c | PHE^d | $s.e^e$ | $s.e(delta)^f$ | $s.e(Fieller)^f$ | $S.W - pvalue^g$ |
|------------------------|---------------------|---------------------|-------------------|-------|---------|---------|----------------|------------------|------------------|
| 200 | 0/0.42 | 0.3/0.17 | 0 | 0.580 | -0.0009 | 0.0051 | 0.0116 | 0.0119 | < 0.001 |
| | | | 0.5 | 0.653 | -0.138 | 0.065 | 0.101 | 0.106 | < 0.001 |
| | 0.15/0.42 | 0.25/0.17 | 0 | 0.530 | 0.093 | 0.040 | 0.077 | 0.080 | 0.300 |
| | | | 0.5 | 0.581 | -0.004 | 0.095 | 0.113 | 0.118 | < 0.001 |
| | 0.42/0.42 | 0.12/0.17 | 0 | 0.424 | 0.273 | 0.087 | 0.089 | 0.093 | 0.004 |
| | | | 0.5 | 0.463 | 0.193 | 0.112 | 0.105 | 0.109 | 0.047 |
| | 0.51/0.42 | 0.04/0.17 | 0 | 0.380 | 0.344 | 0.101 | 0.087 | 0.090 | 0.124 |
| | | | 0.5 | 0.412 | 0.285 | 0.122 | 0.099 | 0.103 | 0.081 |
| | 0.74/0.42 | 0.05/0.41 | 0 | 0.674 | 0.181 | 0.057 | 0.053 | 0.054 | 0.033 |
| | | | 0.5 | 0.762 | 0.069 | 0.074 | 0.058 | 0.057 | 0.146 |
| | 0.74/0.42 | 0.08/0.41 | 0 | 0.682 | 0.174 | 0.069 | 0.053 | 0.053 | 0.777 |
| | | | 0.5 | 0.769 | 0.057 | 0.072 | 0.057 | 0.057 | 0.020 |
| | 0.79/0.42 | 0.08/0.41 | 0 | 0.660 | 0.191 | 0.063 | 0.055 | 0.056 | 0.537 |
| | | | 0.5 | 0.758 | 0.076 | 0.071 | 0.056 | 0.057 | 0.271 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P; h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(delta)$ =mean of estimator of $s.e$ by delta method; $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times$ the range of confidence limits of PHE) by Fieller theorem

^g $S.W - pvalue$ =p value of using Shapiro-Wilk Test

TABLE 11. Simulation results based on scenario II with P>E (2)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | h^c | PHE^d | $s.e^e$ | $s.e(delta)^f$ | $s.e(Fieller)^f$ | $S.W - pvalue^g$ |
|------------------------|---------------------|---------------------|-------------------|-------|---------|---------|----------------|------------------|------------------|
| 500 | 0/0.42 | 0.3/0.17 | 0 | 0.595 | -0.0003 | 0.0017 | 0.0039 | 0.0039 | < 0.001 |
| | | | 0.5 | 0.659 | -0.127 | 0.038 | 0.058 | 0.059 | < 0.001 |
| | 0.15/0.42 | 0.25/0.17 | 0 | 0.539 | 0.091 | 0.025 | 0.046 | 0.046 | < 0.001 |
| | | | 0.5 | 0.588 | -0.003 | 0.054 | 0.069 | 0.070 | 0.108 |
| | 0.42/0.42 | 0.12/0.17 | 0 | 0.432 | 0.267 | 0.051 | 0.055 | 0.056 | 0.367 |
| | | | 0.5 | 0.471 | 0.202 | 0.068 | 0.063 | 0.064 | 0.186 |
| | 0.51/0.42 | 0.04/0.17 | 0 | 0.388 | 0.344 | 0.053 | 0.053 | 0.054 | 0.084 |
| | | | 0.5 | 0.418 | 0.287 | 0.073 | 0.060 | 0.061 | 0.170 |
| | 0.74/0.42 | 0.05/0.41 | 0 | 0.672 | 0.185 | 0.038 | 0.034 | 0.034 | 0.805 |
| | | | 0.5 | 0.762 | 0.074 | 0.044 | 0.035 | 0.035 | 0.394 |
| | 0.74/0.42 | 0.08/0.41 | 0 | 0.681 | 0.175 | 0.038 | 0.033 | 0.034 | 0.495 |
| | | | 0.5 | 0.770 | 0.067 | 0.044 | 0.035 | 0.035 | 0.206 |
| | 0.79/0.42 | 0.08/0.41 | 0 | 0.664 | 0.192 | 0.041 | 0.034 | 0.034 | 0.681 |
| | | | 0.5 | 0.755 | 0.075 | 0.048 | 0.036 | 0.036 | 0.034 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P; h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(delta)$ =mean of estimator of $s.e$ by delta method; $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times \text{the range of confidence limits of PHE})$ by Fieller theorem

^g $S.W - pvalue$ =p value of using Shapiro-Wilk Test

TABLE 12. Simulation results based on scenario II with P>E (3)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | delta method | | | | Fieller theorem | | | |
|------------------------|---------------------|---------------------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^d$ | $D0.25^d$ | $D0.50^d$ | $D0.75^d$ | $F0.00^e$ | $F0.25^e$ | $F0.50^e$ | $F0.75^e$ |
| 200 | 0/0.42 | 0.3/0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.15/0.42 | 0.25/0.17 | 0 | 0.275 | 0 | 0 | 0 | 0.2 | 0 | 0 | 0 |
| | | | 0.5 | 0.04 | 0 | 0 | 0 | 0.02 | 0 | 0 | 0 |
| | 0.42/0.42 | 0.12/0.17 | 0 | 0.9 | 0.09 | 0 | 0 | 0.85 | 0.065 | 0 | 0 |
| | | | 0.5 | 0.585 | 0.025 | 0 | 0 | 0.515 | 0.025 | 0 | 0 |
| | 0.51/0.42 | 0.04/0.17 | 0 | 0.95 | 0.35 | 0.01 | 0 | 0.93 | 0.29 | 0.01 | 0 |
| | | | 0.5 | 0.805 | 0.2 | 0.01 | 0 | 0.755 | 0.16 | 0.01 | 0 |
| | 0.74/0.42 | 0.05/0.41 | 0 | 0.945 | 0.01 | 0 | 0 | 0.925 | 0.01 | 0 | 0 |
| | | | 0.5 | 0.41 | 0 | 0 | 0 | 0.38 | 0 | 0 | 0 |
| | 0.74/0.42 | 0.08/0.41 | 0 | 0.87 | 0.01 | 0 | 0 | 0.855 | 0.01 | 0 | 0 |
| | | | 0.5 | 0.35 | 0 | 0 | 0 | 0.335 | 0 | 0 | 0 |
| | 0.79/0.42 | 0.02/0.41 | 0 | 0.925 | 0.01 | 0 | 0 | 0.91 | 0.01 | 0 | 0 |
| | | | 0.5 | 0.415 | 0 | 0 | 0 | 0.39 | 0 | 0 | 0 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P;

^d Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta}^f)$ is larger than x ;

^e Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^f $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

TABLE 13. Simulation results based on scenario II with P>E (4)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | delta method | | | | Fieller theorem | | | |
|------------------------|---------------------|---------------------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^d$ | $D0.25^d$ | $D0.50^d$ | $D0.75^d$ | $F0.00^e$ | $F0.25^e$ | $F0.50^e$ | $F0.75^e$ |
| 500 | 0/0.42 | 0.3/0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.15/0.42 | 0.25/0.17 | 0 | 0.71 | 0 | 0 | 0 | 0.665 | 0 | 0 | 0 |
| | | | 0.5 | 0.02 | 0 | 0 | 0 | 0.015 | 0 | 0 | 0 |
| | 0.42/0.42 | 0.12/0.17 | 0 | 1 | 0.09 | 0 | 0 | 1 | 0.085 | 0 | 0 |
| | | | 0.5 | 0.905 | 0.03 | 0 | 0 | 0.885 | 0.03 | 0 | 0 |
| | 0.51/0.42 | 0.04/0.17 | 0 | 1 | 0.53 | 0 | 0 | 1 | 0.5 | 0 | 0 |
| | | | 0.5 | 0.985 | 0.24 | 0 | 0 | 0.985 | 0.22 | 0 | 0 |
| | 0.74/0.42 | 0.05/0.41 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| | | | 0.5 | 0.6 | 0 | 0 | 0 | 0.585 | 0 | 0 | 0 |
| | 0.74/0.42 | 0.08/0.41 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| | | | 0.5 | 0.59 | 0 | 0 | 0 | 0.565 | 0 | 0 | 0 |
| | 0.79/0.42 | 0.02/0.41 | 0 | 0.995 | 0 | 0 | 0 | 0.995 | 0 | 0 | 0 |
| | | | 0.5 | 0.67 | 0 | 0 | 0 | 0.645 | 0 | 0 | 0 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P;

^d Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta}^f)$ is larger than x ;

^e Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^f $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

TABLE 14. Simulation results based on scenario II with P<E (1)

| <i>No. of families</i> | $h(G1_E)/h(G1_p)^a$ | $h(G2_E)/h(G3_p)^b$ | ρ_ϵ^c | h^c | PHE^d | $s.e^e$ | $s.e(delta)^f$ | $s.e(Fieller)^f$ | $S.W - pvalue^g$ |
|------------------------|---------------------|---------------------|-------------------|-------|----------|---------|----------------|------------------|------------------|
| 200 | 0/0.42 | 0.7/0.17 | 0 | 0.582 | -0.00009 | 0.0055 | 0.012 | 0.012 | < 0.001 |
| | | | 0.5 | 0.639 | -0.096 | 0.047 | 0.079 | 0.082 | < 0.001 |
| | 0.15/0.42 | 0.59/0.17 | 0 | 0.536 | 0.073 | 0.041 | 0.082 | 0.086 | < 0.001 |
| | | | 0.5 | 0.613 | -0.049 | 0.074 | 0.109 | 0.114 | 0.016 |
| | 0.42/0.42 | 0.23/0.17 | 0 | 0.434 | 0.243 | 0.074 | 0.093 | 0.097 | 0.555 |
| | | | 0.5 | 0.512 | 0.132 | 0.106 | 0.105 | 0.109 | < 0.001 |
| | 0.62/0.42 | 0.23/0.17 | 0 | 0.393 | 0.319 | 0.096 | 0.091 | 0.095 | 0.941 |
| | | | 0.5 | 0.477 | 0.182 | 0.129 | 0.103 | 0.108 | < 0.001 |
| | 0.74/0.42 | 0.08/0.17 | 0 | 0.329 | 0.426 | 0.109 | 0.086 | 0.089 | 0.002 |
| | | | 0.5 | 0.408 | 0.294 | 0.136 | 0.097 | 0.101 | 0.070 |
| | 0.74/0.42 | 0.21/0.17 | 0 | 0.381 | 0.347 | 0.108 | 0.089 | 0.092 | 0.043 |
| | | | 0.5 | 0.436 | 0.232 | 0.129 | 0.104 | 0.109 | 0.041 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_p)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_p)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P; h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(delta)$ =mean of estimator of $s.e$ by delta method; $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times \text{the range of confidence limits of PHE})$ by Fieller theorem

^g $S.W - pvalue$ =p value of using Shapiro-Wilk Test

TABLE 15. Simulation results based on scenario II with P<E (2)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | h^c | PHE^d | $s.e^e$ | $s.e(delta)^f$ | $s.e(Fieller)^f$ | $S.W - pvalue^g$ |
|------------------------|---------------------|---------------------|-------------------|-------|----------|---------|----------------|------------------|------------------|
| 500 | 0/0.42 | 0.7/0.17 | 0 | 0.589 | -0.00003 | 0.0019 | 0.0043 | 0.0043 | < 0.001 |
| | | | 0.5 | 0.647 | -0.091 | 0.028 | 0.046 | 0.046 | < 0.001 |
| | 0.15/0.42 | 0.59/0.17 | 0 | 0.553 | 0.069 | 0.025 | 0.047 | 0.048 | 0.170 |
| | | | 0.5 | 0.616 | -0.054 | 0.046 | 0.068 | 0.069 | 0.089 |
| | 0.42/0.42 | 0.23/0.17 | 0 | 0.446 | 0.243 | 0.049 | 0.056 | 0.057 | 0.990 |
| | | | 0.5 | 0.519 | 0.126 | 0.069 | 0.066 | 0.067 | 0.654 |
| | 0.62/0.42 | 0.23/0.17 | 0 | 0.405 | 0.313 | 0.058 | 0.056 | 0.057 | 0.249 |
| | | | 0.5 | 0.483 | 0.177 | 0.074 | 0.064 | 0.065 | 0.932 |
| | 0.74/0.42 | 0.08/0.17 | 0 | 0.337 | 0.431 | 0.069 | 0.051 | 0.052 | 0.730 |
| | | | 0.5 | 0.413 | 0.295 | 0.079 | 0.059 | 0.060 | 0.001 |
| | 0.74/0.42 | 0.21/0.17 | 0 | 0.388 | 0.340 | 0.065 | 0.056 | 0.056 | 0.980 |
| | | | 0.5 | 0.445 | 0.242 | 0.075 | 0.061 | 0.062 | 0.146 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P; h =mean of heritability of P, conditional on E

^d PHE =mean of proportion of heritability explained by endophenotype

^e $s.e$ =standard deviation of proportion of heritability explained by endophenotype

^f $s.e(delta)$ =mean of estimator of $s.e$ by delta method; $s.e(Fieller)$ =mean of $(\frac{1}{2 \times 1.96} \times \text{the range of confidence limits of PHE})$ by Fieller theorem

^g $S.W - pvalue$ =p value of using Shapiro-Wilk Test

TABLE 16. Simulation results based on scenario II with P<E (3)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | delta method | | | | Fieller theorem | | | |
|------------------------|---------------------|---------------------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^d$ | $D0.25^d$ | $D0.50^d$ | $D0.75^d$ | $F0.00^e$ | $F0.25^e$ | $F0.50^e$ | $F0.75^e$ |
| 200 | 0/0.42 | 0.7/0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.15/0.42 | 0.59/0.17 | 0 | 0.12 | 0 | 0 | 0 | 0.09 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.42/0.42 | 0.23/0.17 | 0 | 0.845 | 0.015 | 0 | 0 | 0.75 | 0.015 | 0 | 0 |
| | | | 0.5 | 0.415 | 0.02 | 0 | 0 | 0.315 | 0.015 | 0 | 0 |
| | 0.62/0.42 | 0.23/0.17 | 0 | 0.91 | 0.25 | 0 | 0 | 0.865 | 0.22 | 0 | 0 |
| | | | 0.5 | 0.585 | 0.03 | 0 | 0 | 0.545 | 0.03 | 0 | 0 |
| | 0.74/0.42 | 0.08/0.17 | 0 | 0.97 | 0.645 | 0.04 | 0 | 0.955 | 0.575 | 0.04 | 0 |
| | | | 0.5 | 0.805 | 0.225 | 0.015 | 0 | 0.775 | 0.23 | 0.01 | 0 |
| | 0.74/0.42 | 0.21/0.17 | 0 | 0.945 | 0.295 | 0.005 | 0 | 0.925 | 0.275 | 0.05 | 0 |
| | | | 0.5 | 0.67 | 0.12 | 0 | 0 | 0.61 | 0.085 | 0 | 0 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P;

^d Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta}^f)$ is larger than x ;

^e Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^f $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

TABLE 17. Simulation results based on scenario II with P<E (4)

| <i>No. of families</i> | $h(G1_E)/h(G1_P)^a$ | $h(G2_E)/h(G3_P)^b$ | ρ_ϵ^c | delta method | | | | Fieller theorem | | | |
|------------------------|---------------------|---------------------|-------------------|--------------|-----------|-----------|-----------|-----------------|-----------|-----------|-----------|
| | | | | $D0.00^d$ | $D0.25^d$ | $D0.50^d$ | $D0.75^d$ | $F0.00^e$ | $F0.25^e$ | $F0.50^e$ | $F0.75^e$ |
| 500 | 0/0.42 | 0.7/0.17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.15/0.42 | 0.59/0.17 | 0 | 0.4 | 0 | 0 | 0 | 0.32 | 0 | 0 | 0 |
| | | | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | 0.42/0.42 | 0.23/0.17 | 0 | 0.99 | 0.03 | 0 | 0 | 0.99 | 0.02 | 0 | 0 |
| | | | 0.5 | 0.575 | 0 | 0 | 0 | 0.54 | 0 | 0 | 0 |
| | 0.62/0.42 | 0.23/0.17 | 0 | 1 | 0.32 | 0 | 0 | 1 | 0.345 | 0 | 0 |
| | | | 0.5 | 0.805 | 0.015 | 0 | 0 | 0.76 | 0.01 | 0 | 0 |
| | 0.74/0.42 | 0.08/0.17 | 0 | 1 | 0.895 | 0.02 | 0 | 1 | 0.86 | 0.01 | 0 |
| | | | 0.5 | 0.971896 | 0.31 | 0 | 0 | 0.96 | 0.27 | 0 | 0 |
| | 0.74/0.42 | 0.21/0.17 | 0 | 1 | 0.515 | 0 | 0 | 1 | 0.47 | 0 | 0 |
| | | | 0.5 | 0.93 | 0.075 | 0 | 0 | 0.9 | 0.075 | 0 | 0 |

^a $h(G1_E)$ =heritability of E due to G1; $h(G1_P)$ = heritability of P due to G1;

^b $h(G2_E)$ =heritability of E due to G2; $h(G3_P)$ = heritability of P due to G3;

^c ρ_ϵ =correlation between non-family deviations of E and P;

^d Dx =the porportion that $(\widehat{PHE}-1.645 \times s.e(\widehat{PHE})_{delta}^f)$ is larger than x ;

^e Fx =the porportion that the lower 95% confidence limits at one side using Fieller theorem is larger than x ;

^f $s.e(\widehat{PHE})_{delta}$ =the estimator of $s.e$ by delta method

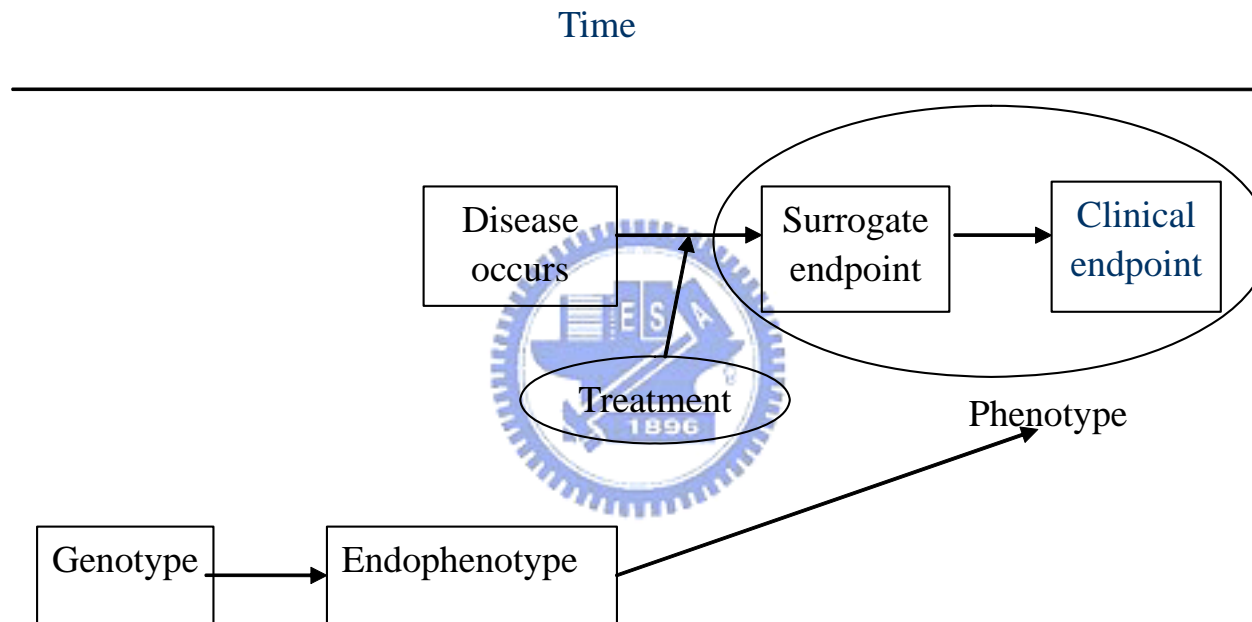


Figure 1: A surrogate endpoint versus an endophenotype in the disease process

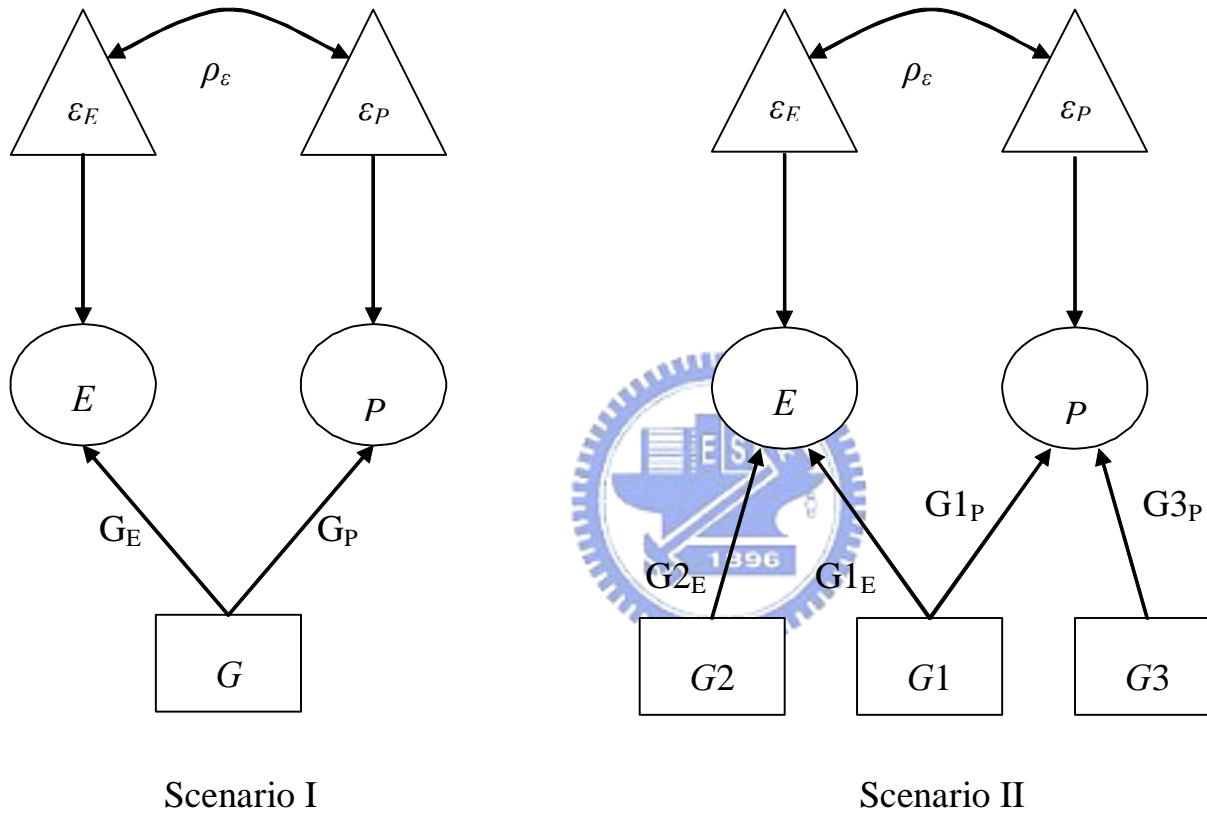


Figure 2: Two scenarios verified in the simulation studies: endophenotype (E), phenotype (P), underlying disease genes (G, G_1, G_2 and G_3), random non-family effects (ϵ_E and ϵ_P), $h(G'_E)$ means the heritability of E due to G' , $h(G'_P)$ means the heritability of P due to G' , and correlation between non-family effects (ρ_ϵ), where G' may be G, G_1, G_2 , or G_3 .

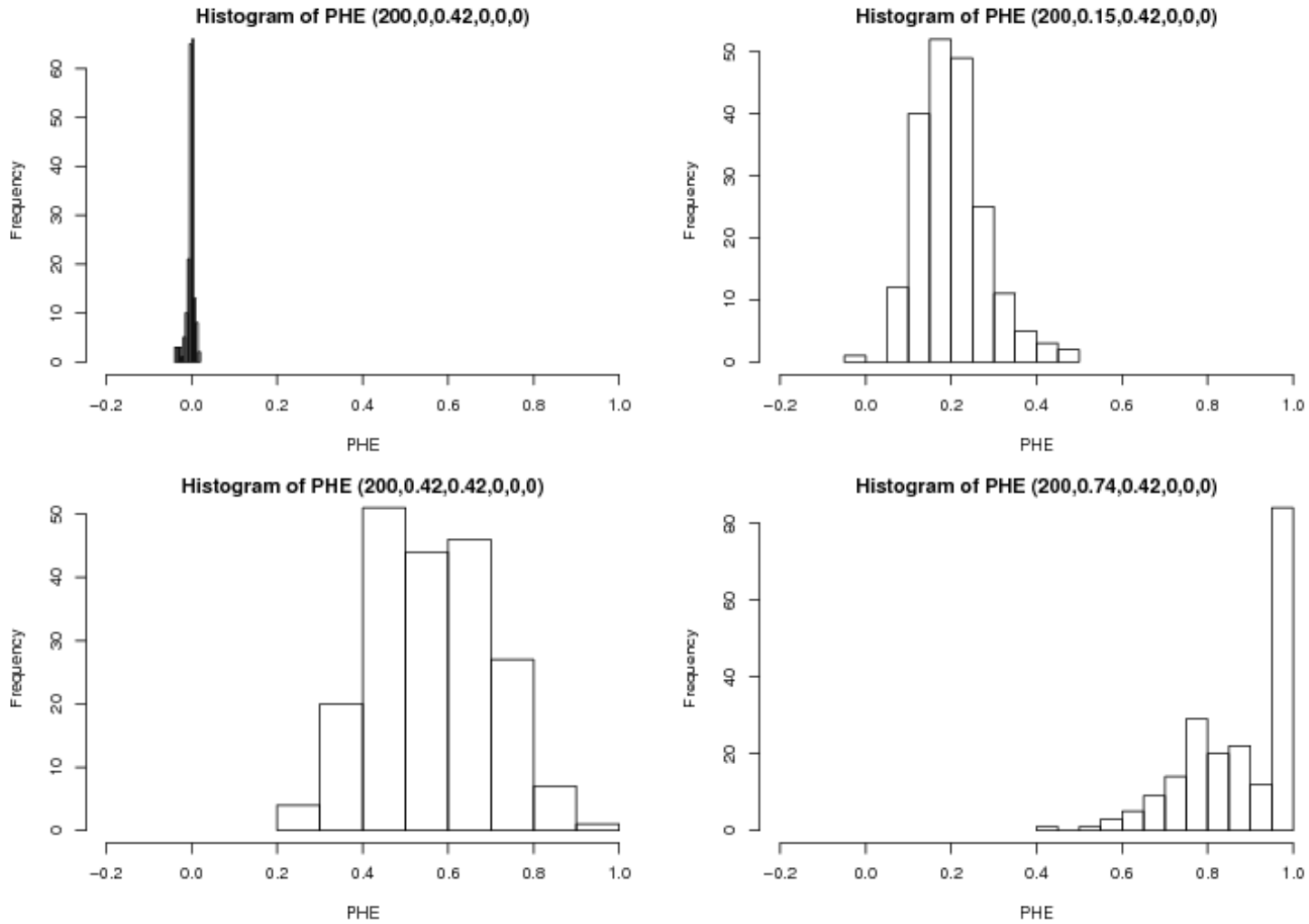


Figure 3: Scenario I histogram with family 200 The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

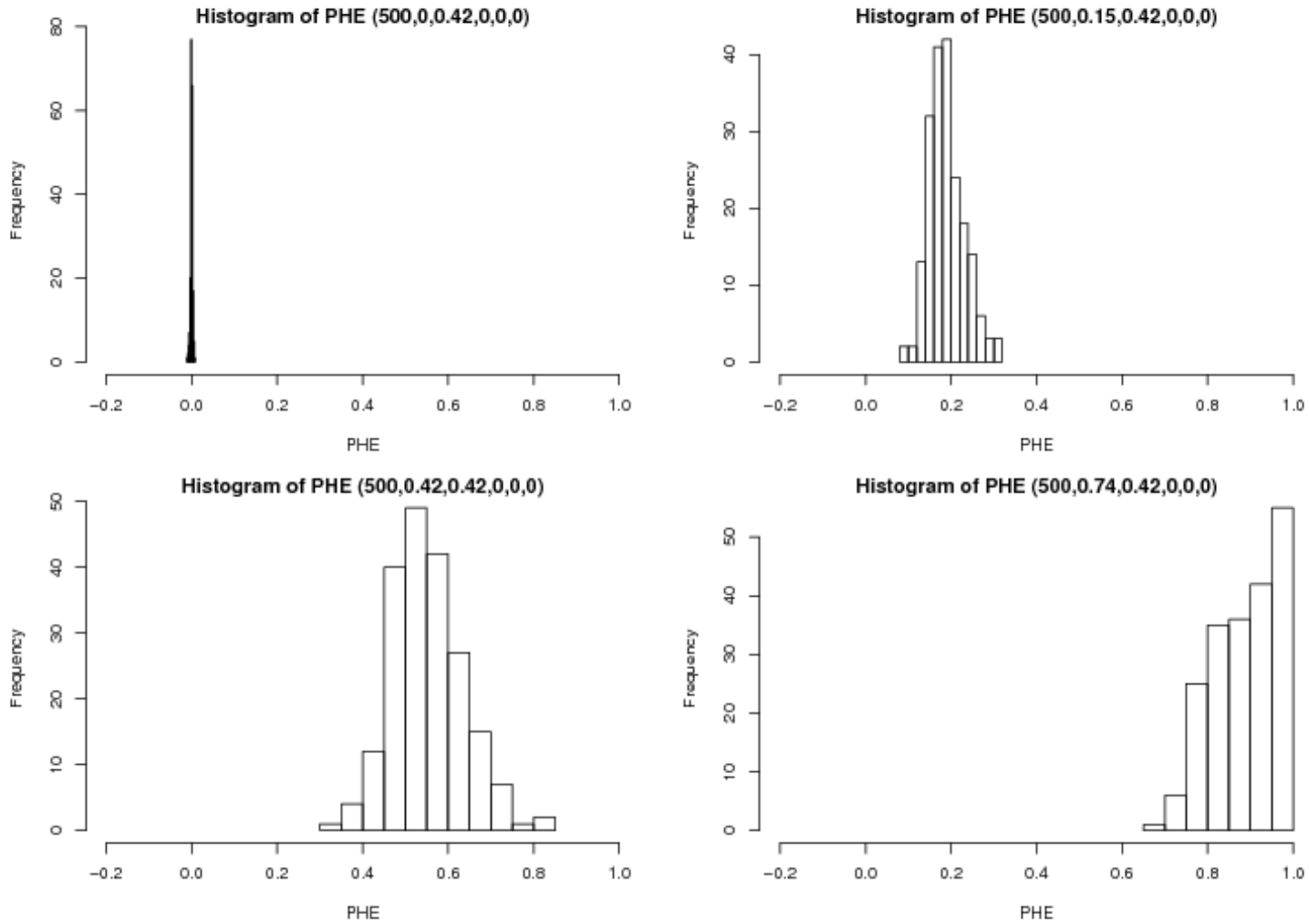


Figure 4: Scenario I histogram with family 500 The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

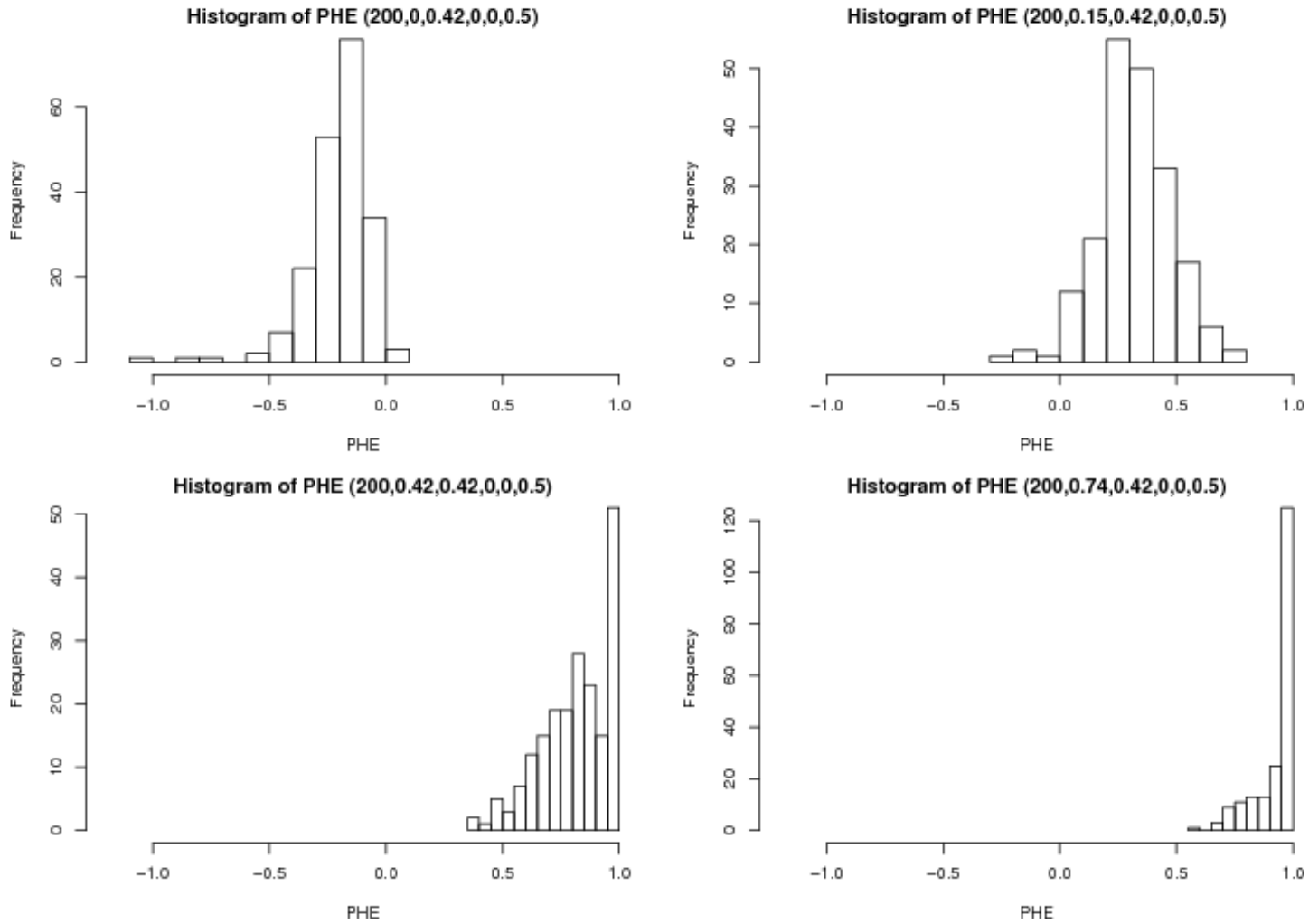


Figure 5: Scenario I histogram with family 200 & $\rho_\epsilon = 0.5$. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

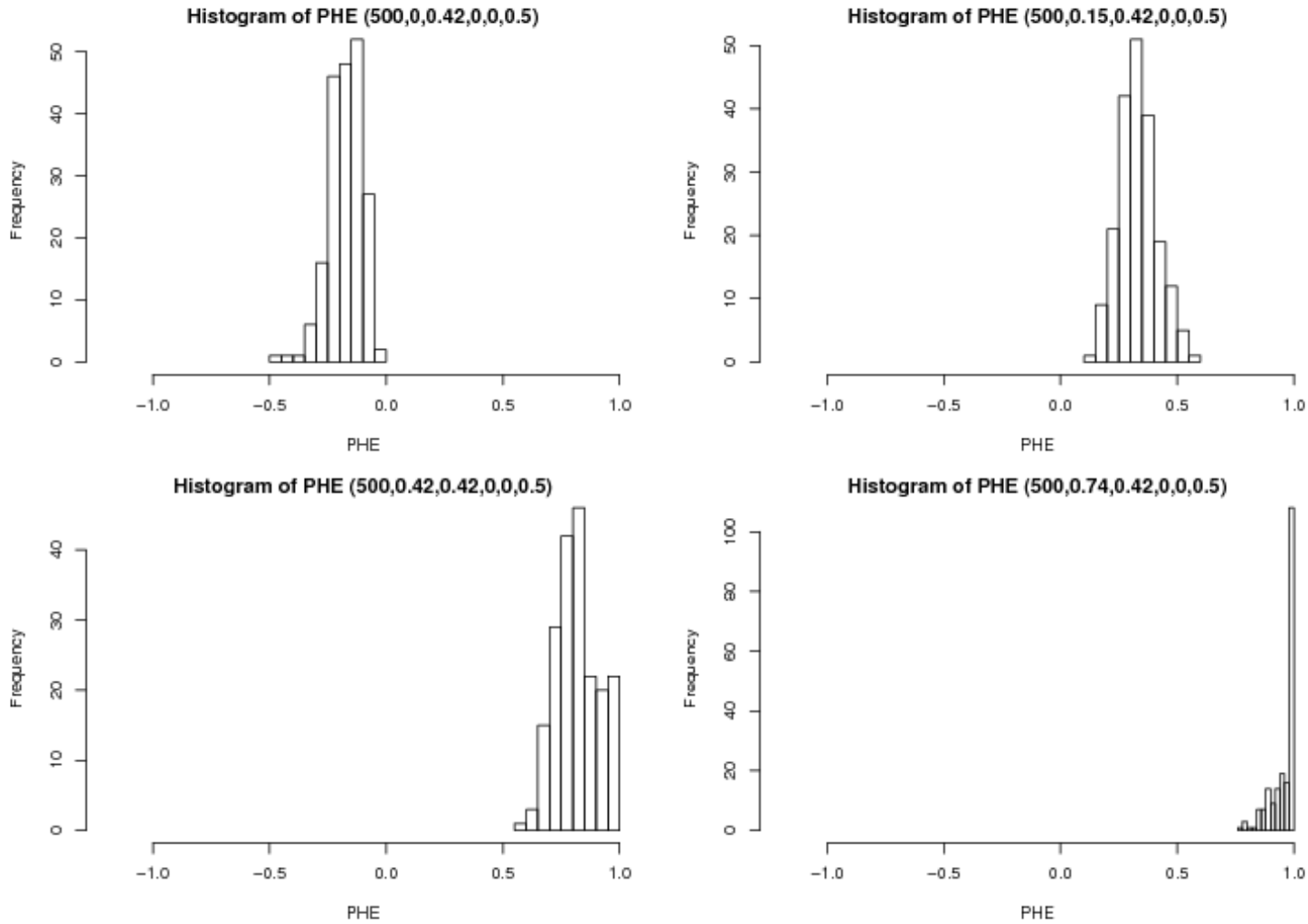


Figure 6: Scenario I histogram with family 500 & $\rho_\epsilon = 0.5$. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

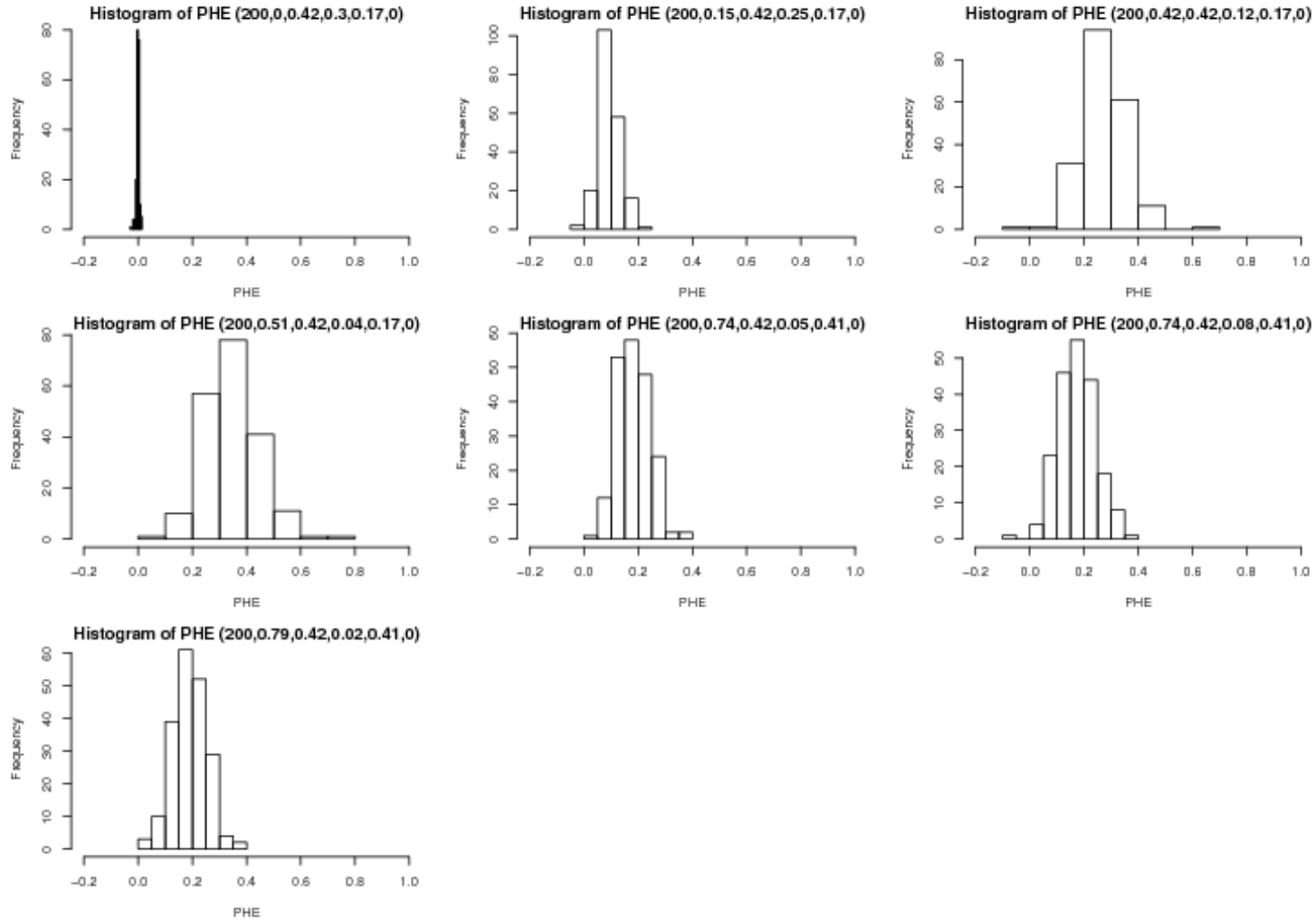


Figure 7: Scenario II histogram with family 200 & P>E The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

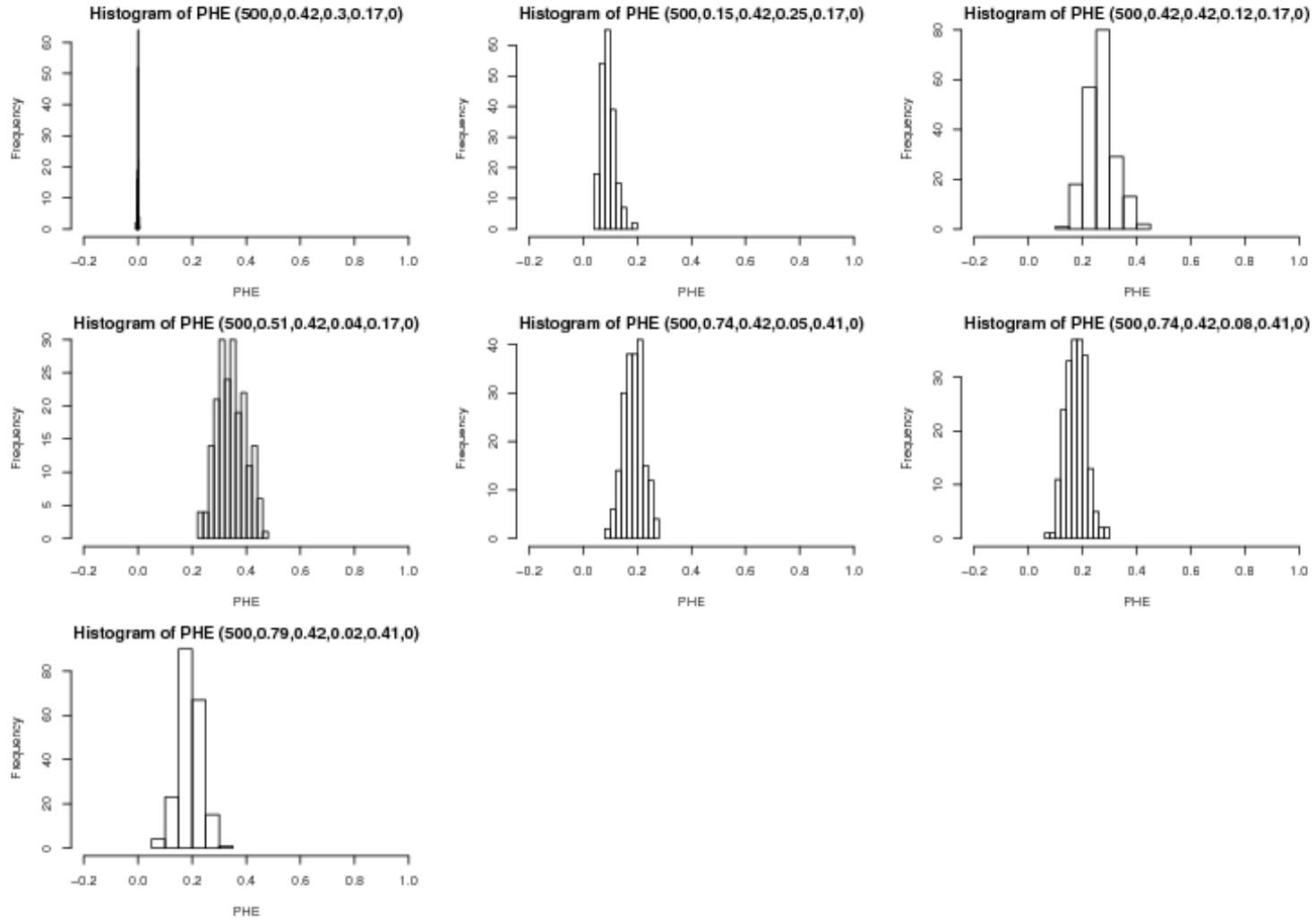


Figure 8: Scenario II histogram with family 500 & $P > E$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

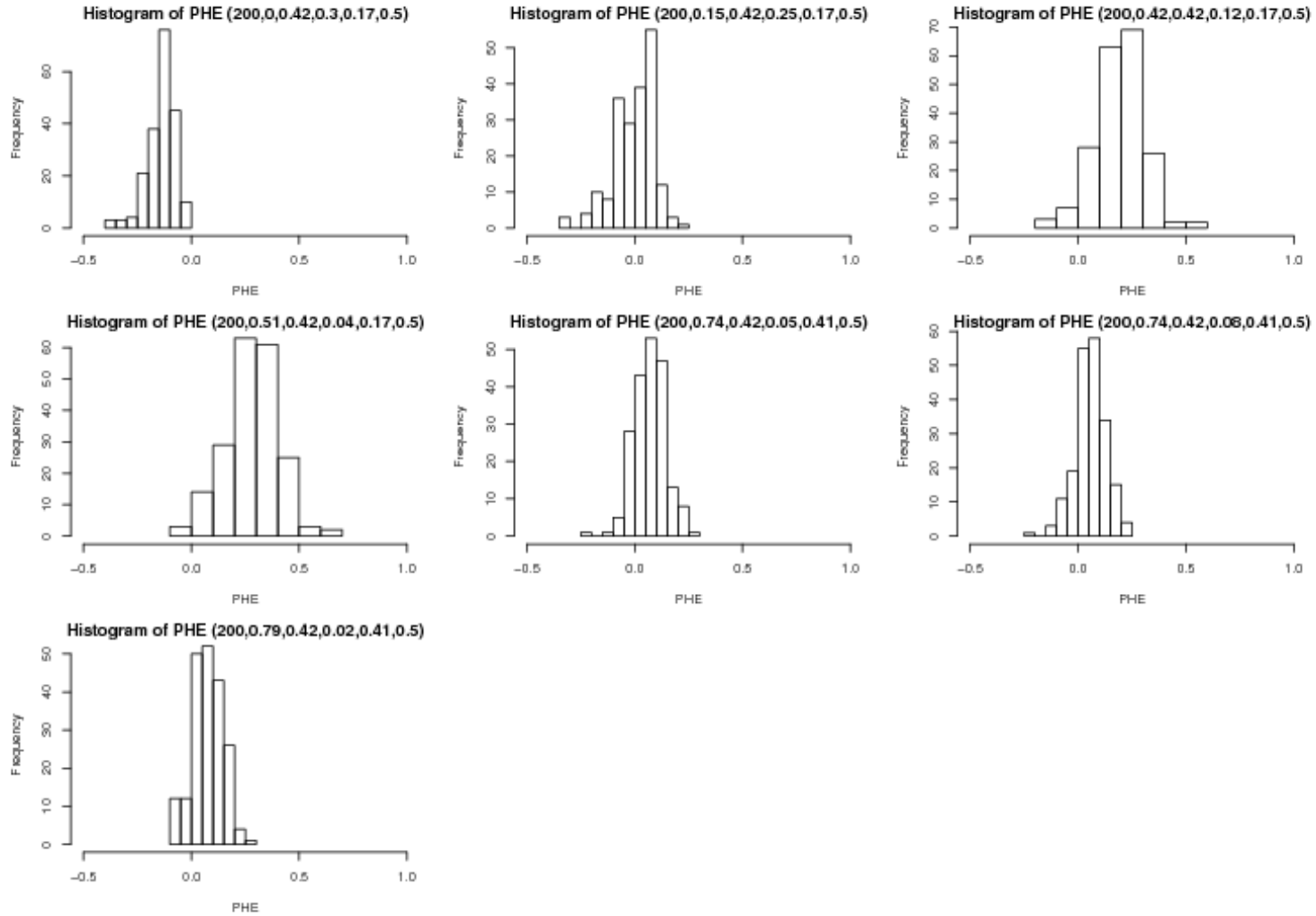


Figure 9: Scenario II histogram with family 200 & $P>E$ & $\rho_\epsilon = 0.5$ The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

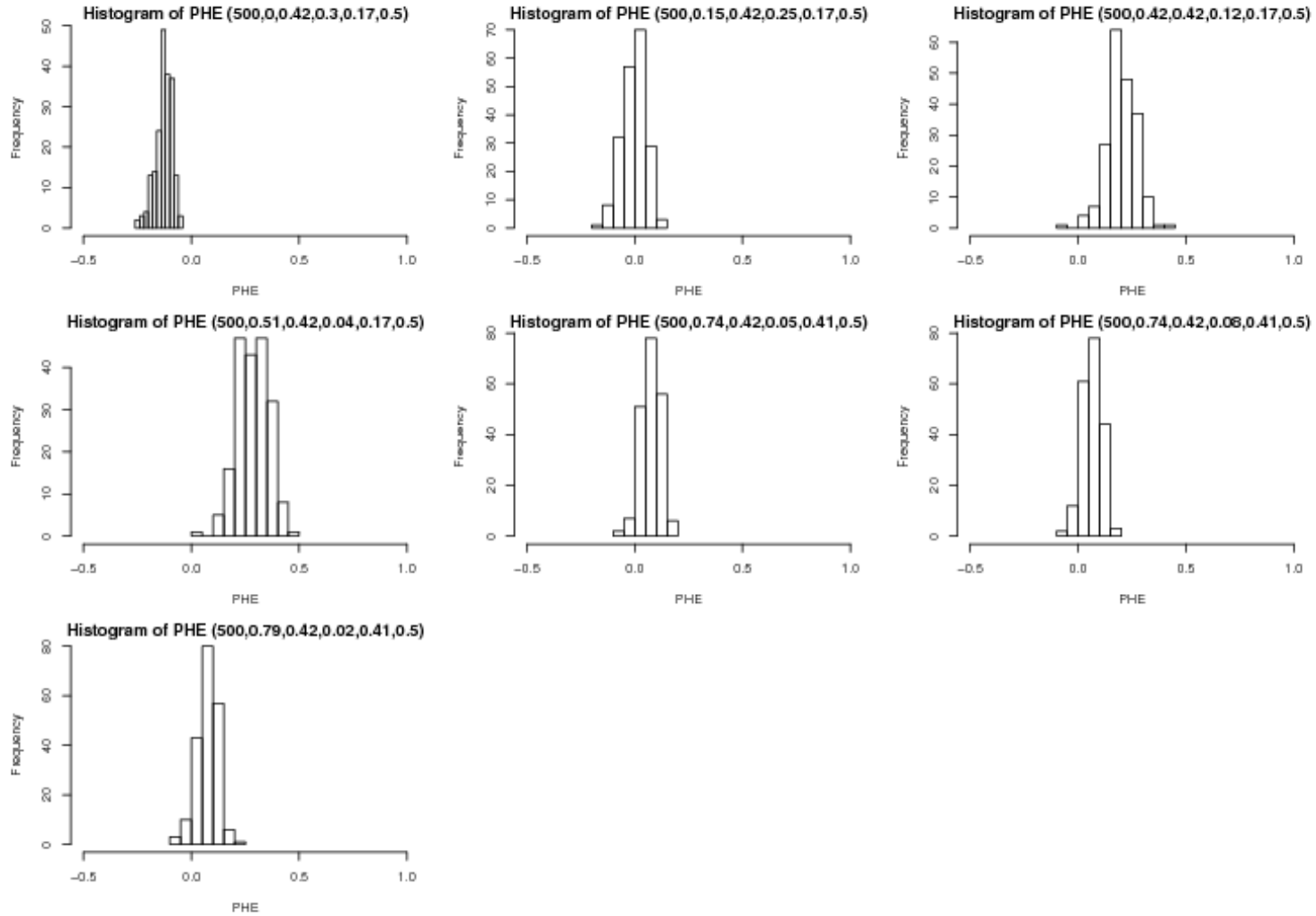


Figure 10: Scenario II histogram with family 500 & $P>E$ & $\rho_\epsilon = 0.5$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

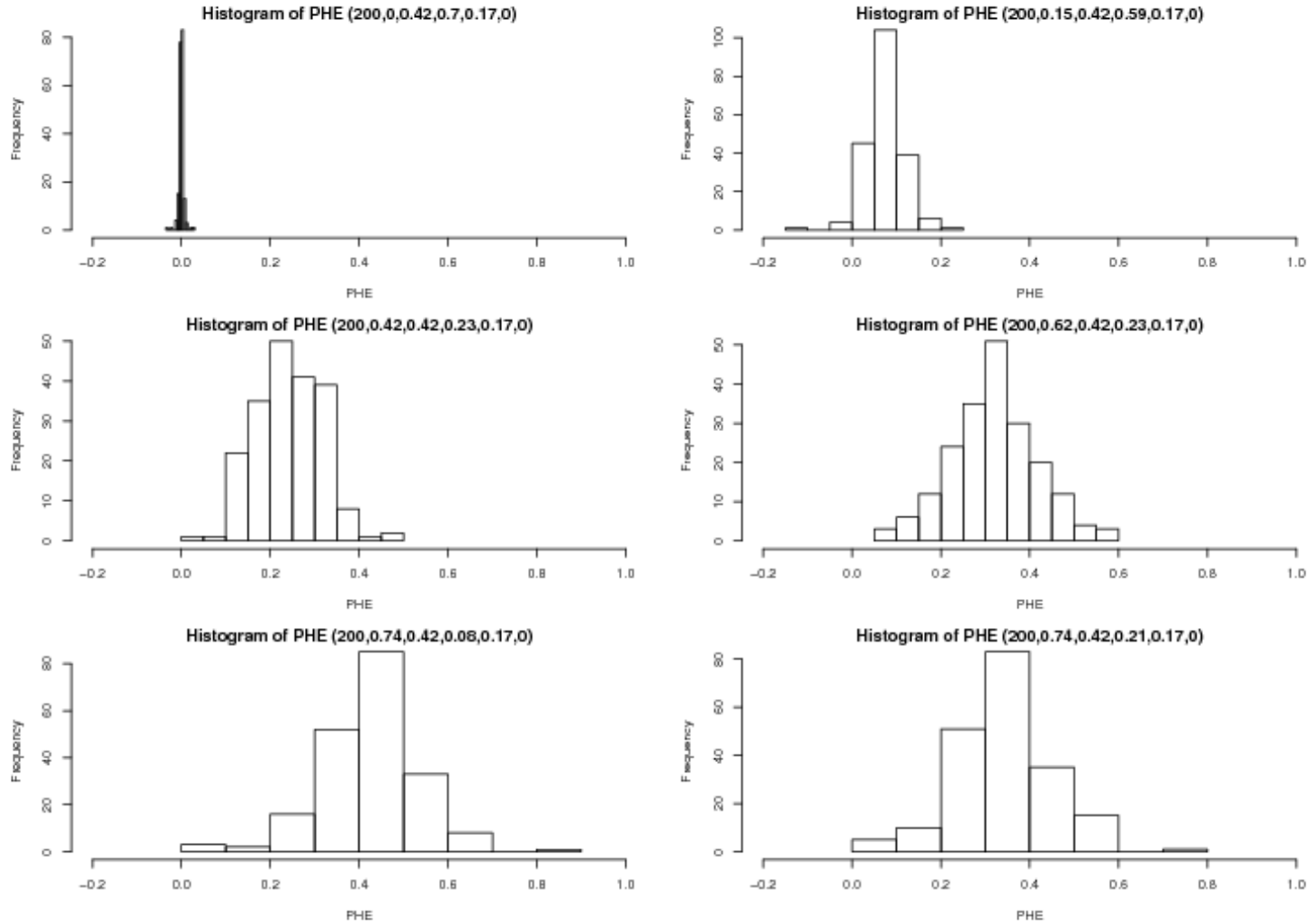


Figure 11: Scenario II histogram with family 200 & $P < E$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

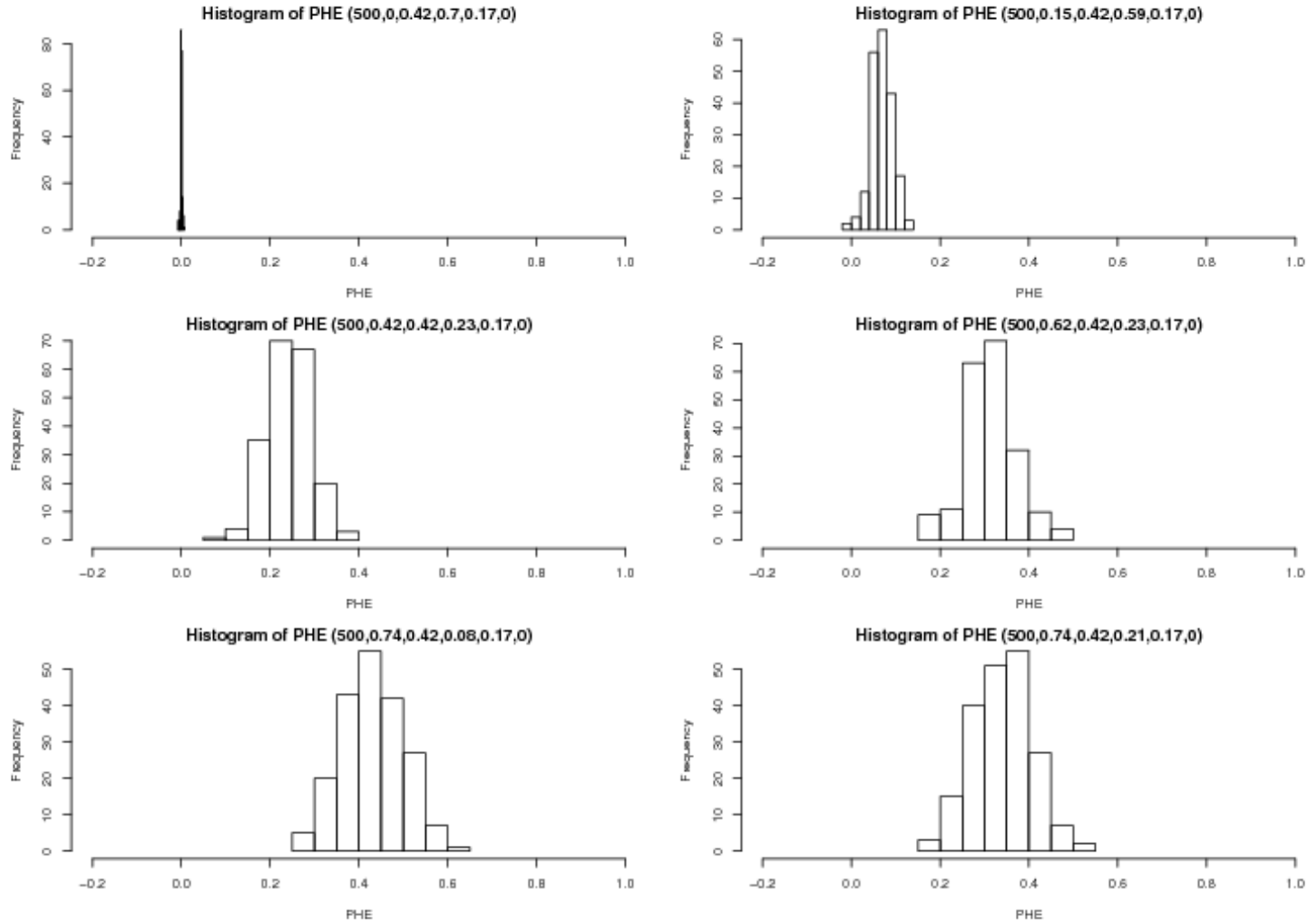


Figure 12: Scenario II histogram with family 500 & $P < E$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

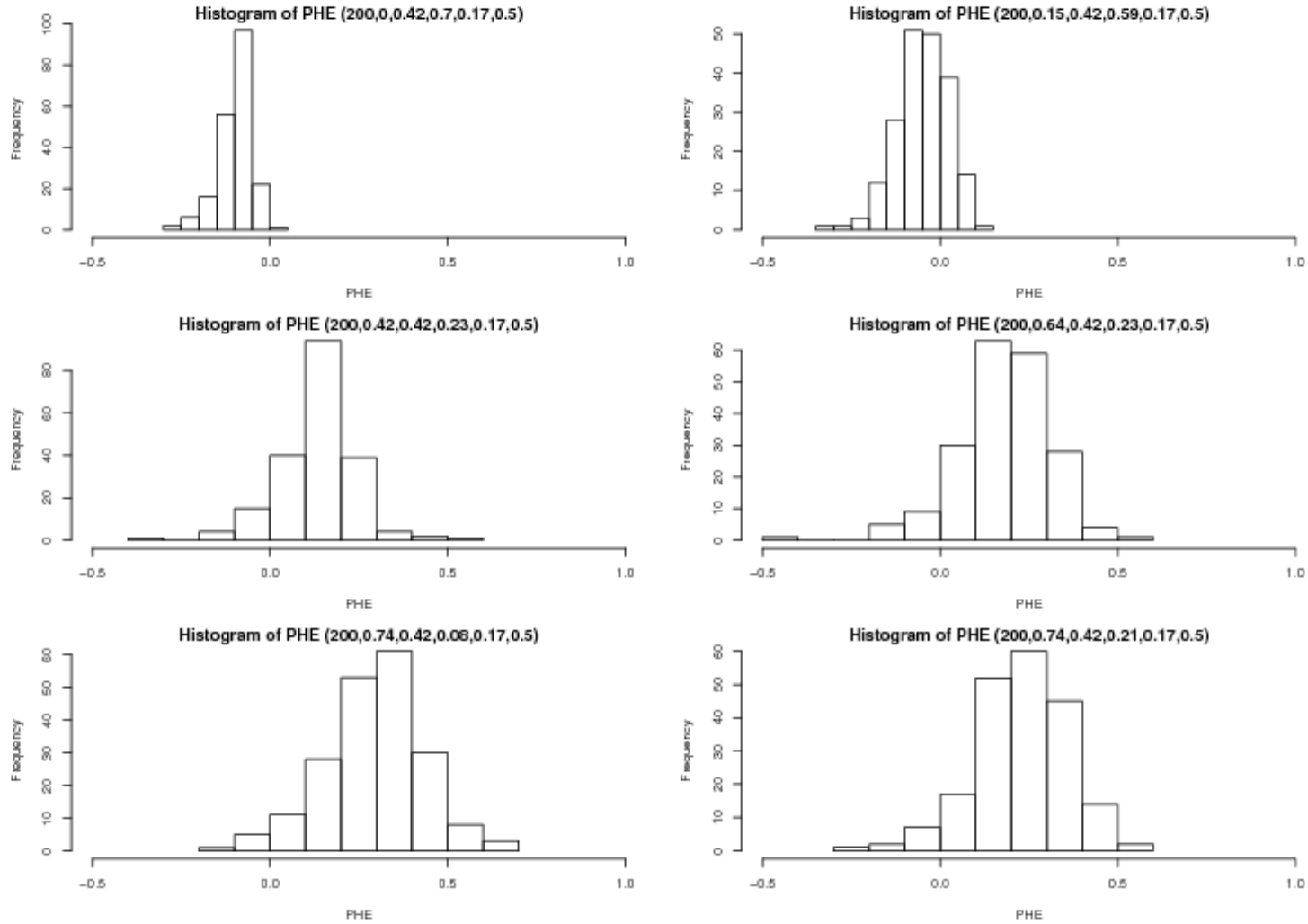


Figure 13: Scenario II histogram with family 200 & $P < E$ & $\rho_\epsilon = 0.5$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

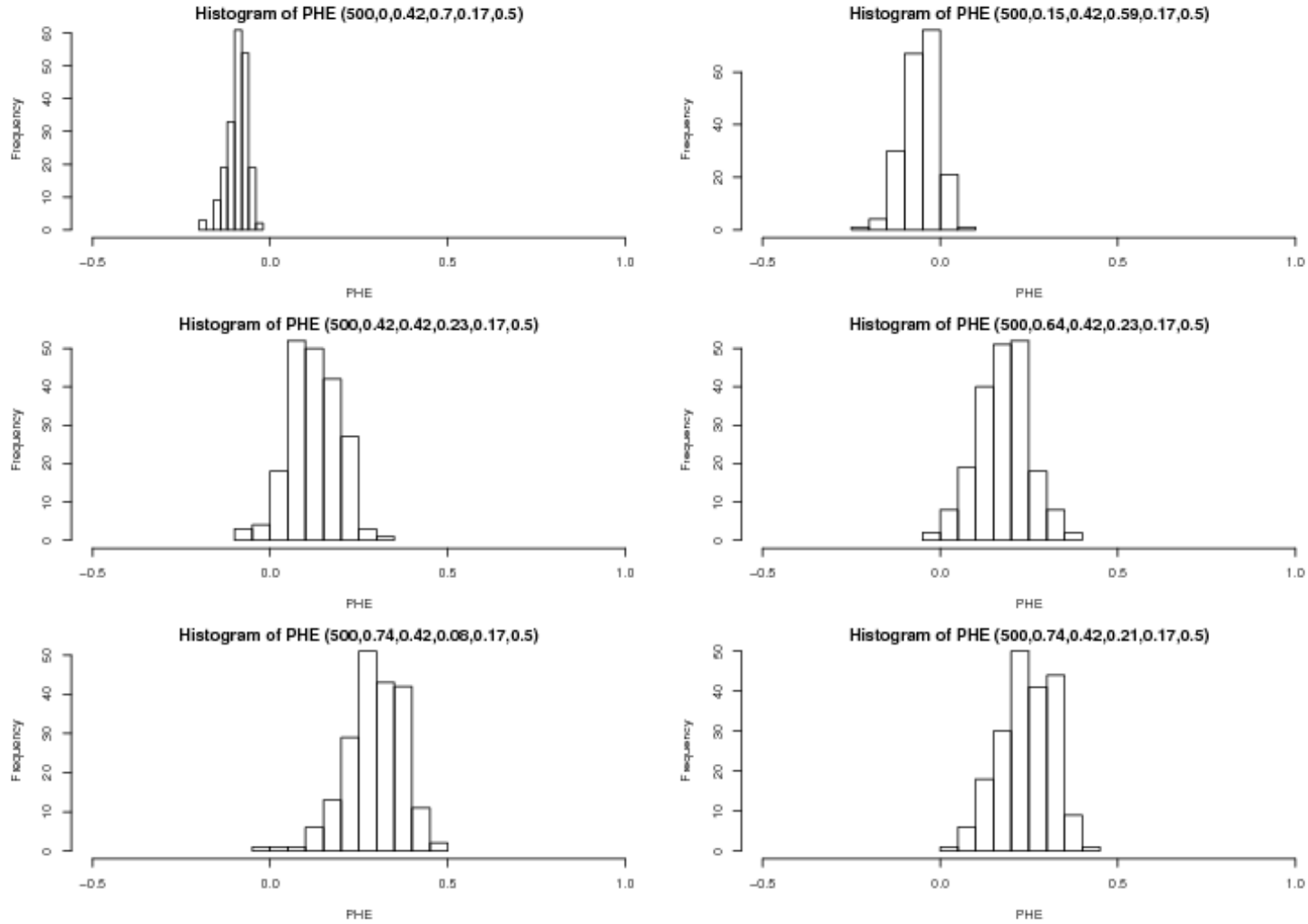


Figure 14: Scenario II histogram with family 500 & $P < E$ & $\rho_\epsilon = 0.5$. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

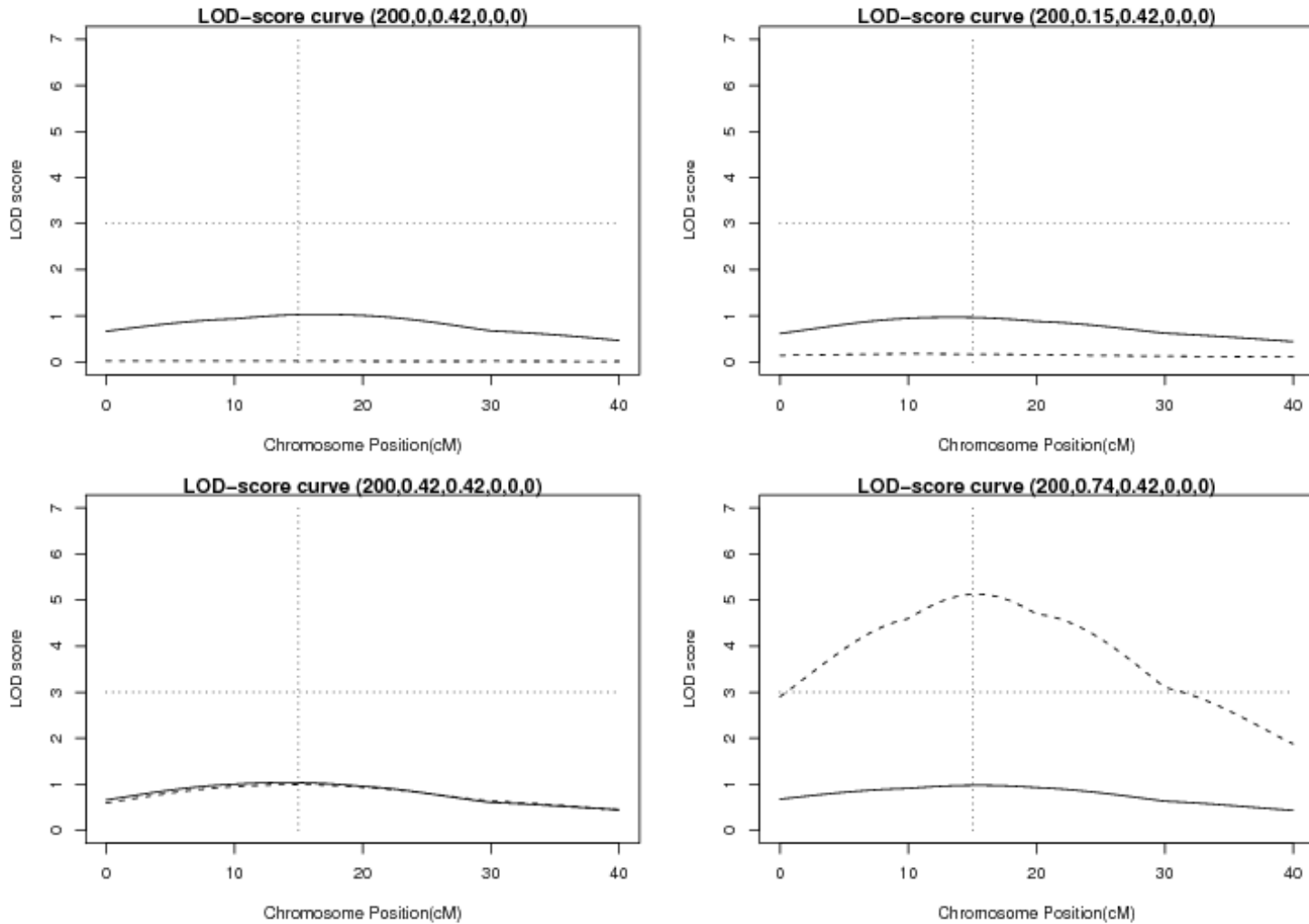


Figure 15: Scenario I mean LOD-score curve with family 200, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

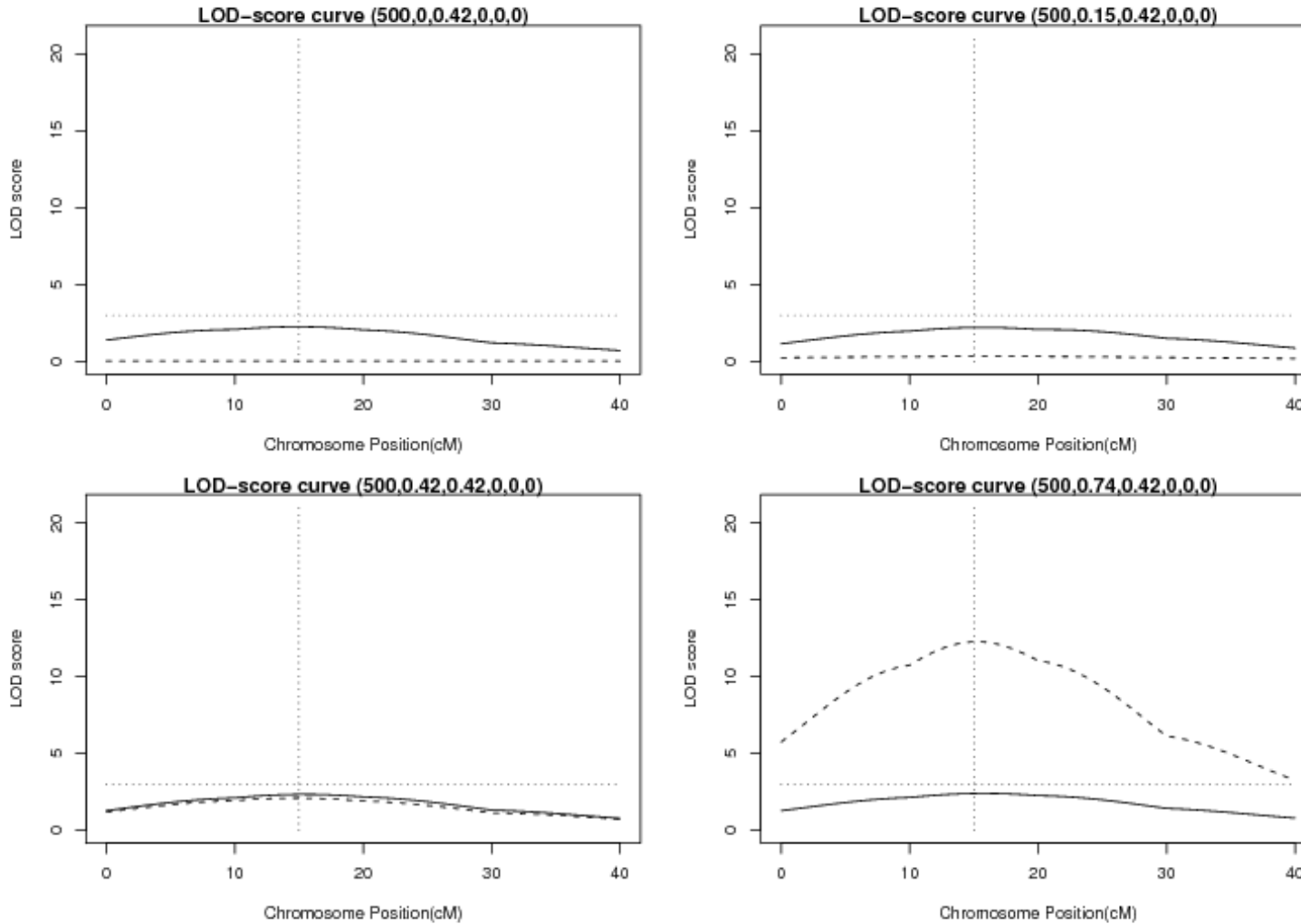


Figure 16: Scenario I mean LOD-score curve with family 500 , where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, $(fam, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

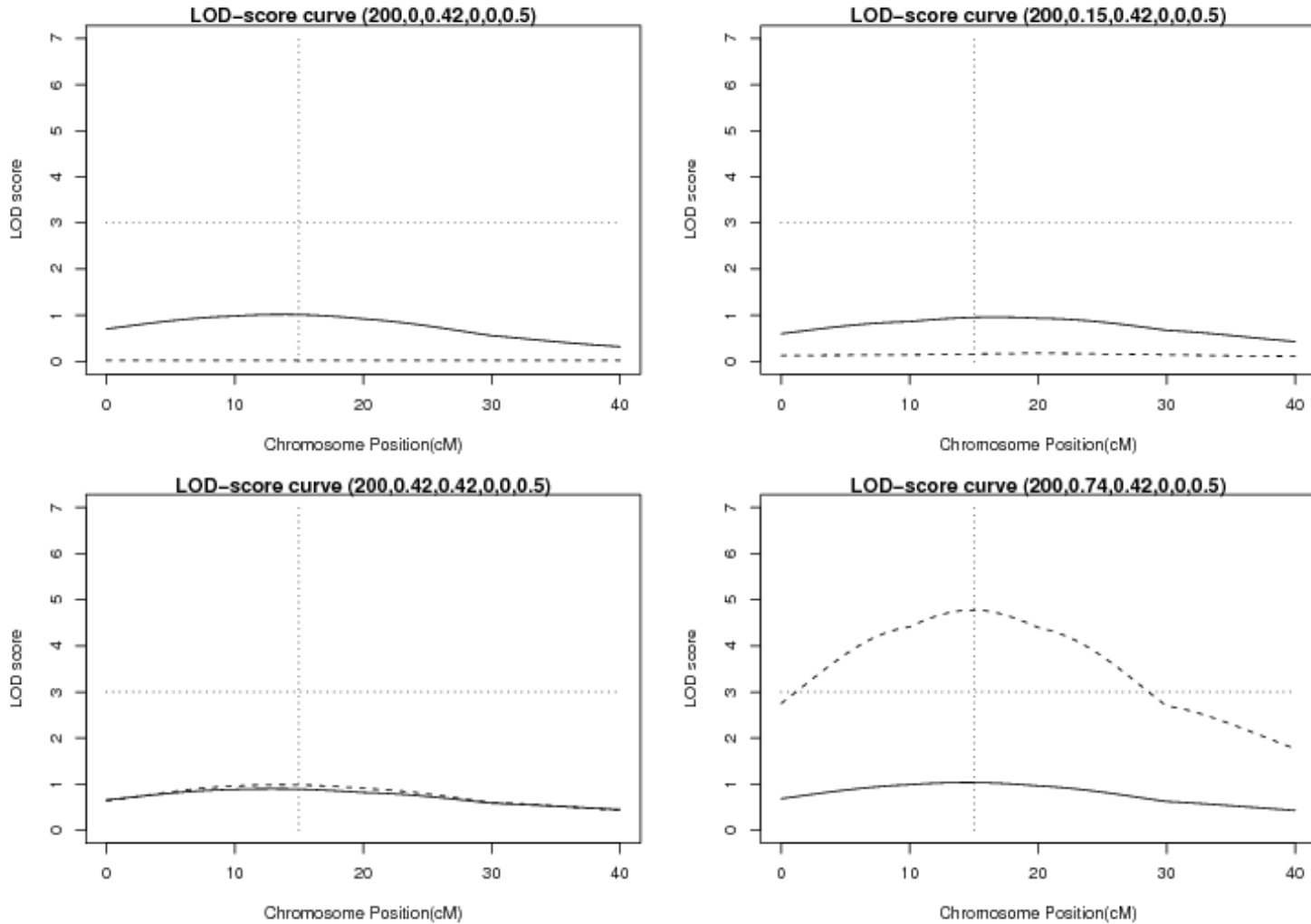


Figure 17: Scenario I mean LOD-score curve with family 200 & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

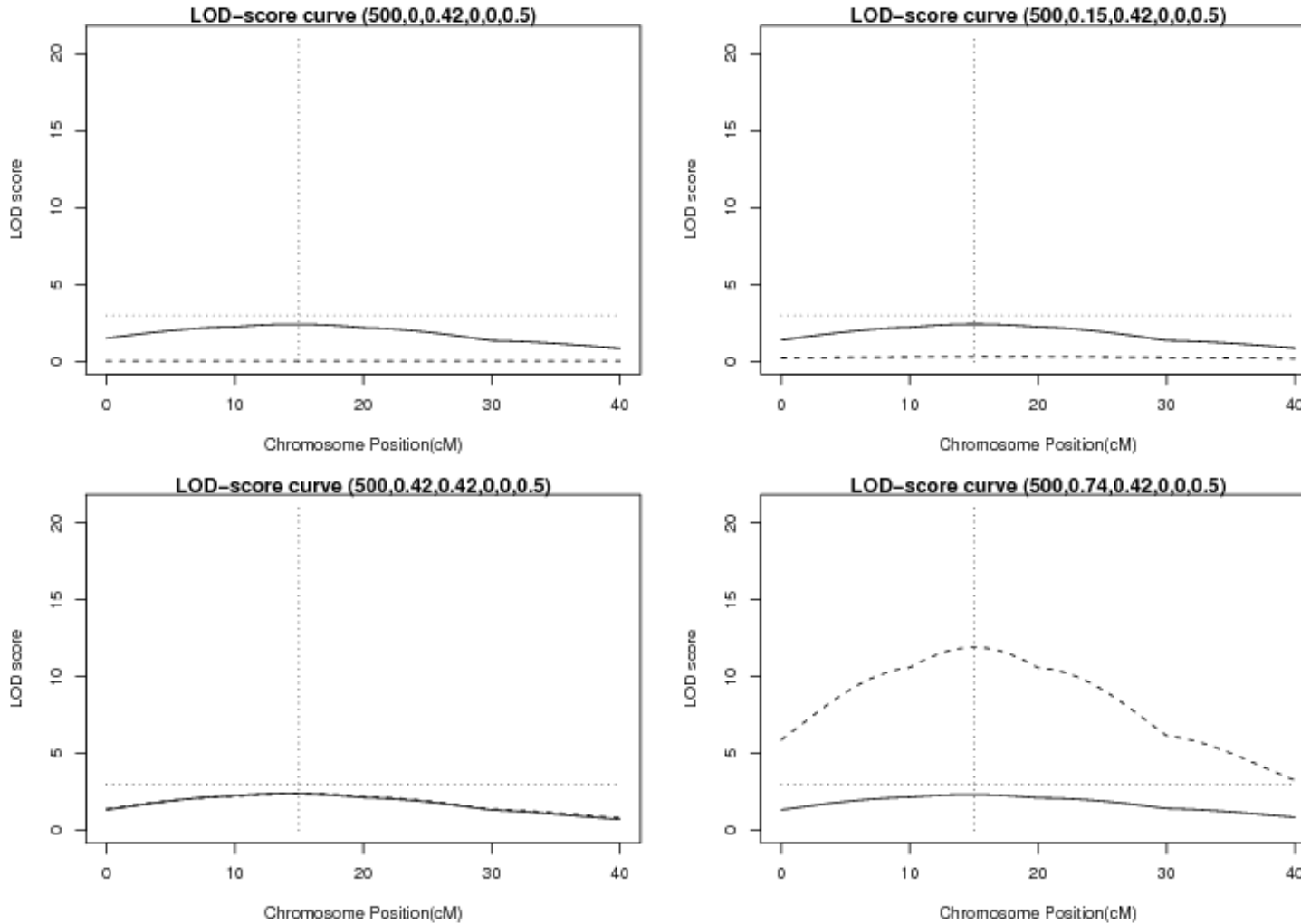


Figure 18: Scenario I mean LOD-score curve with family 500 & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

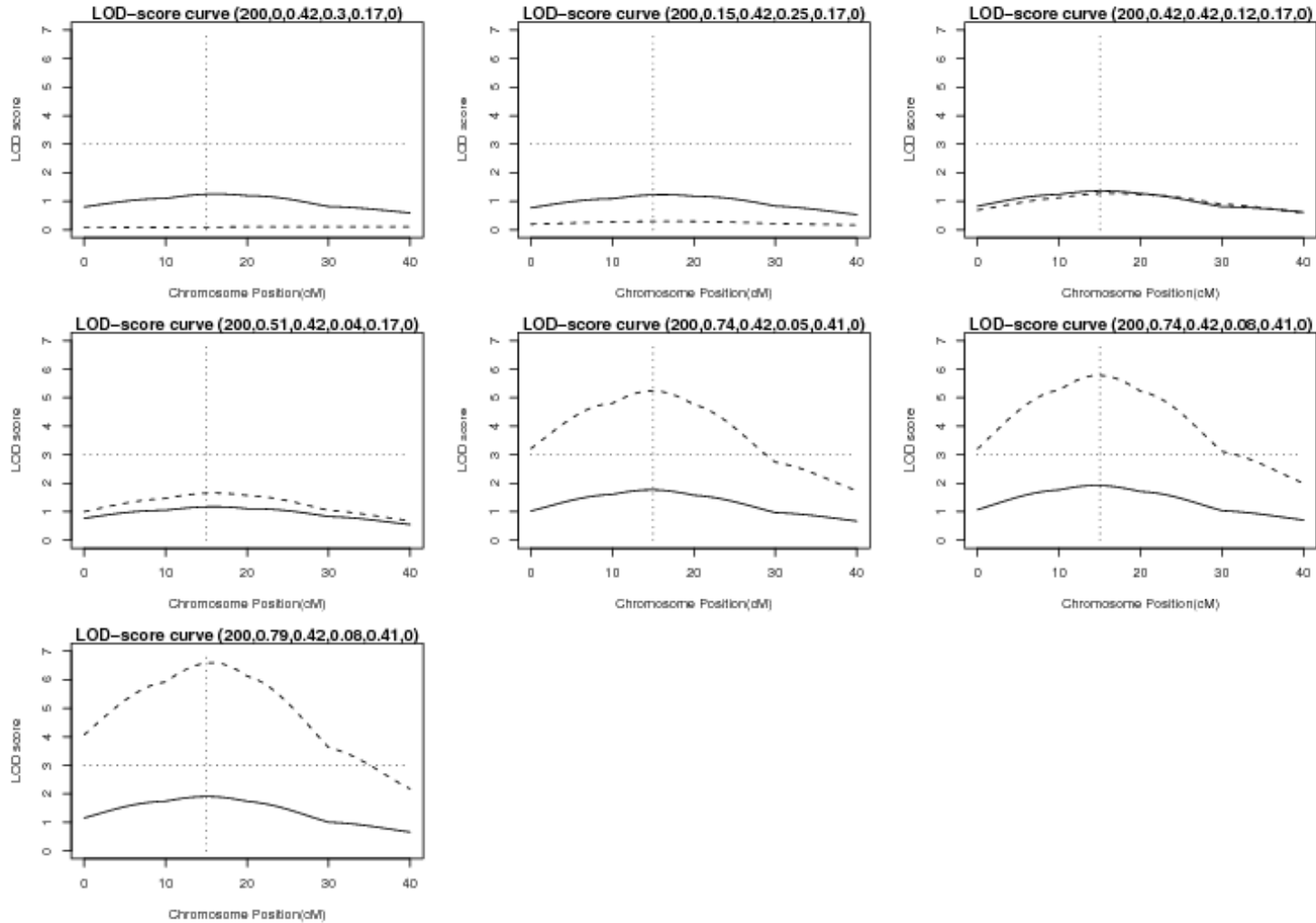


Figure 19: Scenario II mean LOD-score curve with family 200 & $P > E$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is $\text{LOD-score}=3$, and vertical dotted line is the position of disease gene. The title in each figure, $(\text{fam}, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

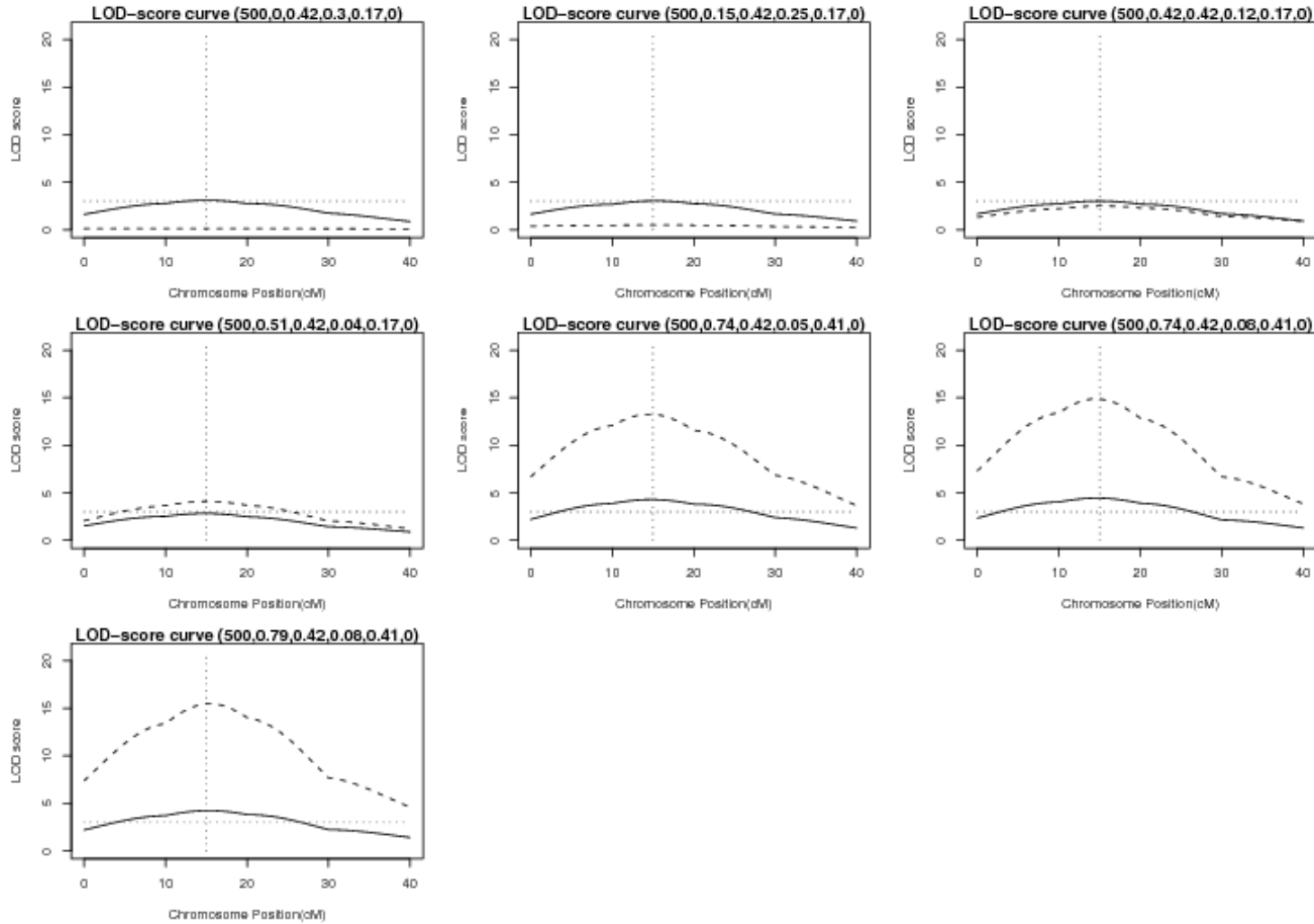


Figure 20: Scenario II mean LOD-score curve with family 500 & $P > E$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is $\text{LOD-score}=3$, and vertical dotted line is the position of disease gene. The title in each figure, $(\text{fam}, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

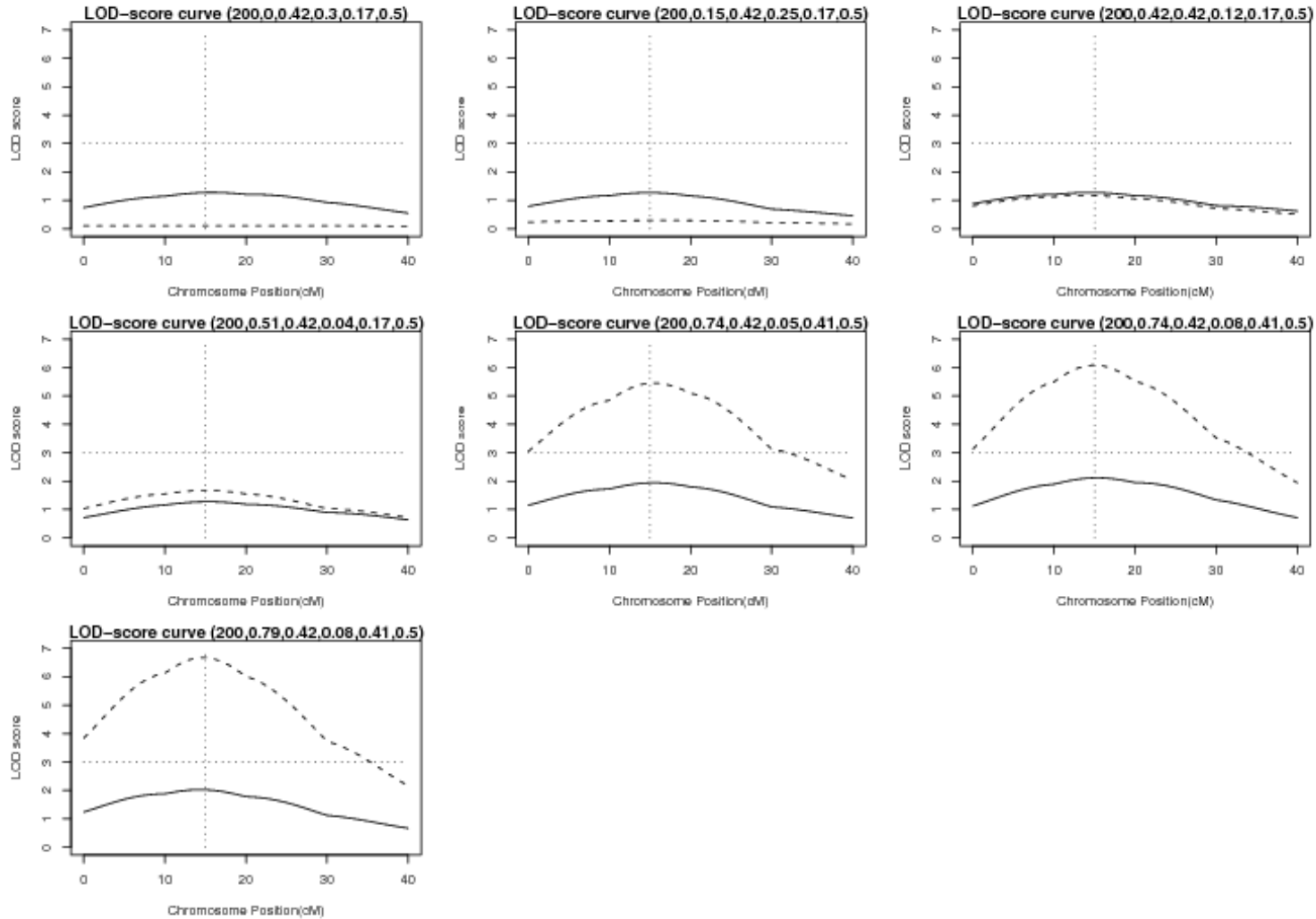


Figure 21: Scenario II mean LOD-score curve with family 200 & $P > E$ & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is $\text{LOD-score}=3$, and vertical dotted line is the position of disease gene. The title in each figure, $(\text{fam}, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

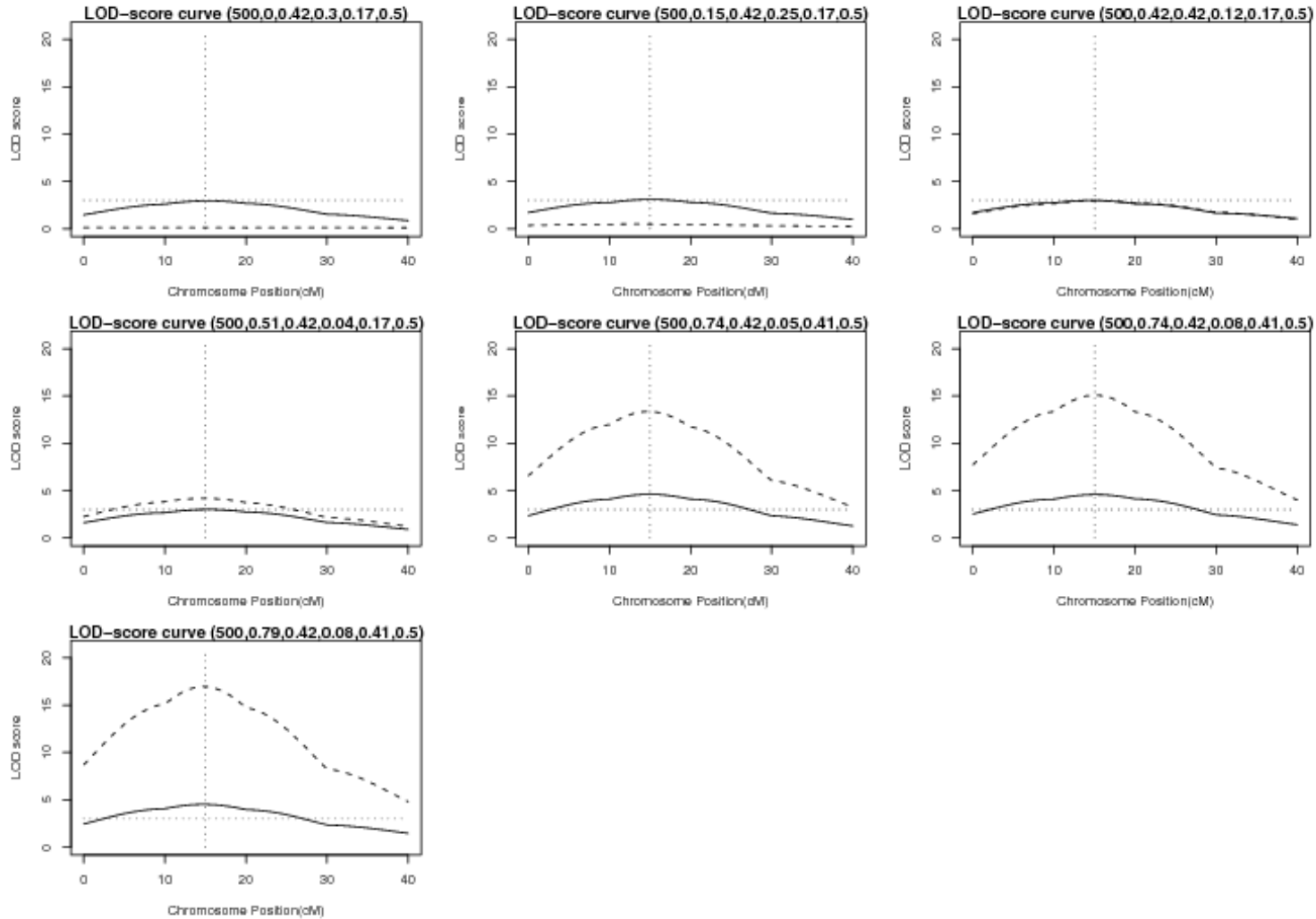


Figure 22: Scenario II mean LOD-score curve with family 500 & $P > E$ & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is $\text{LOD-score}=3$, and vertical dotted line is the position of disease gene. The title in each figure, $(\text{fam}, h(G1_E), h(G1_P), h(G2_E), h(G3_P), \rho_\epsilon)$, to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

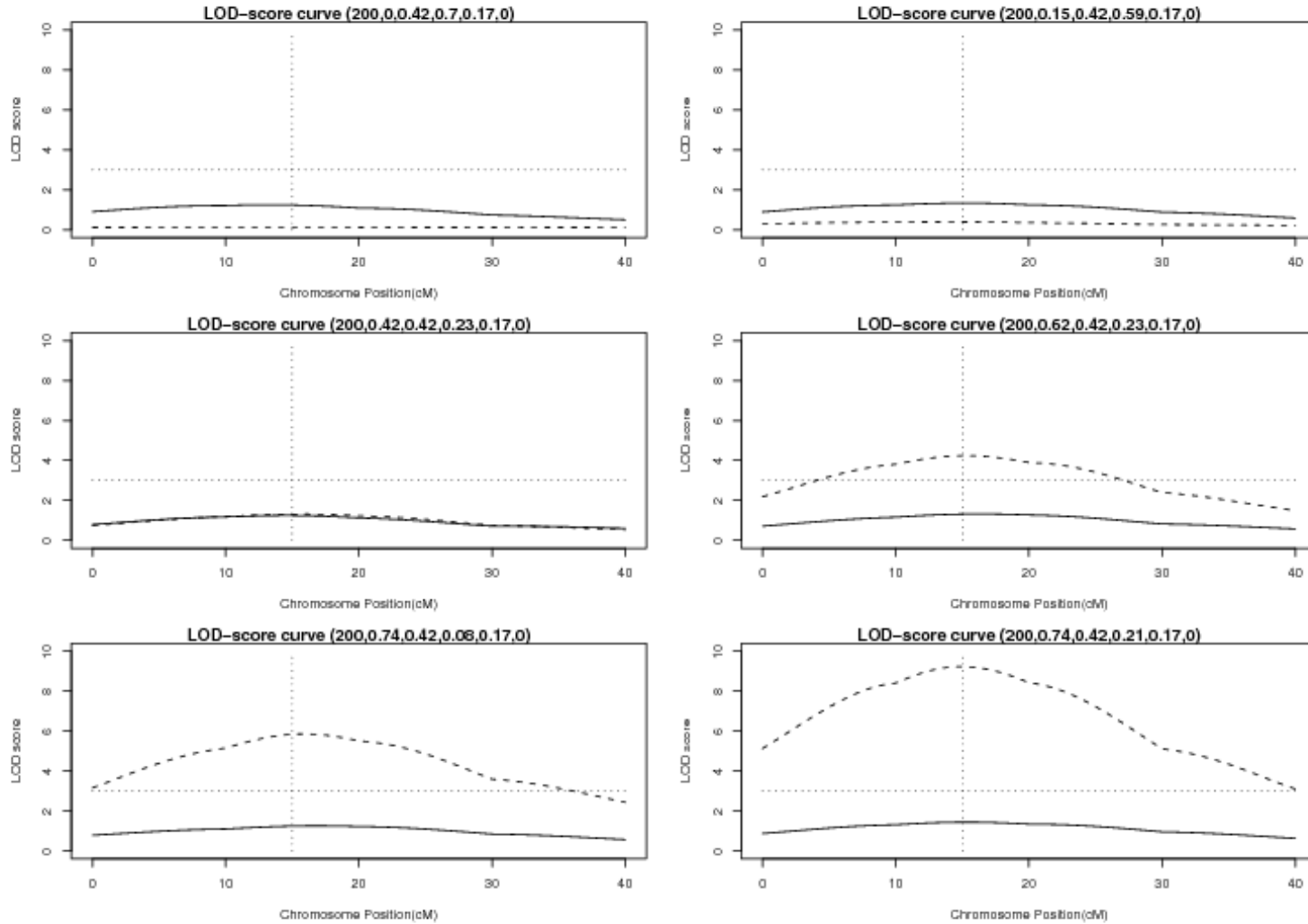


Figure 23: Scenario II mean LOD-score curve with family 200 & $P < E$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

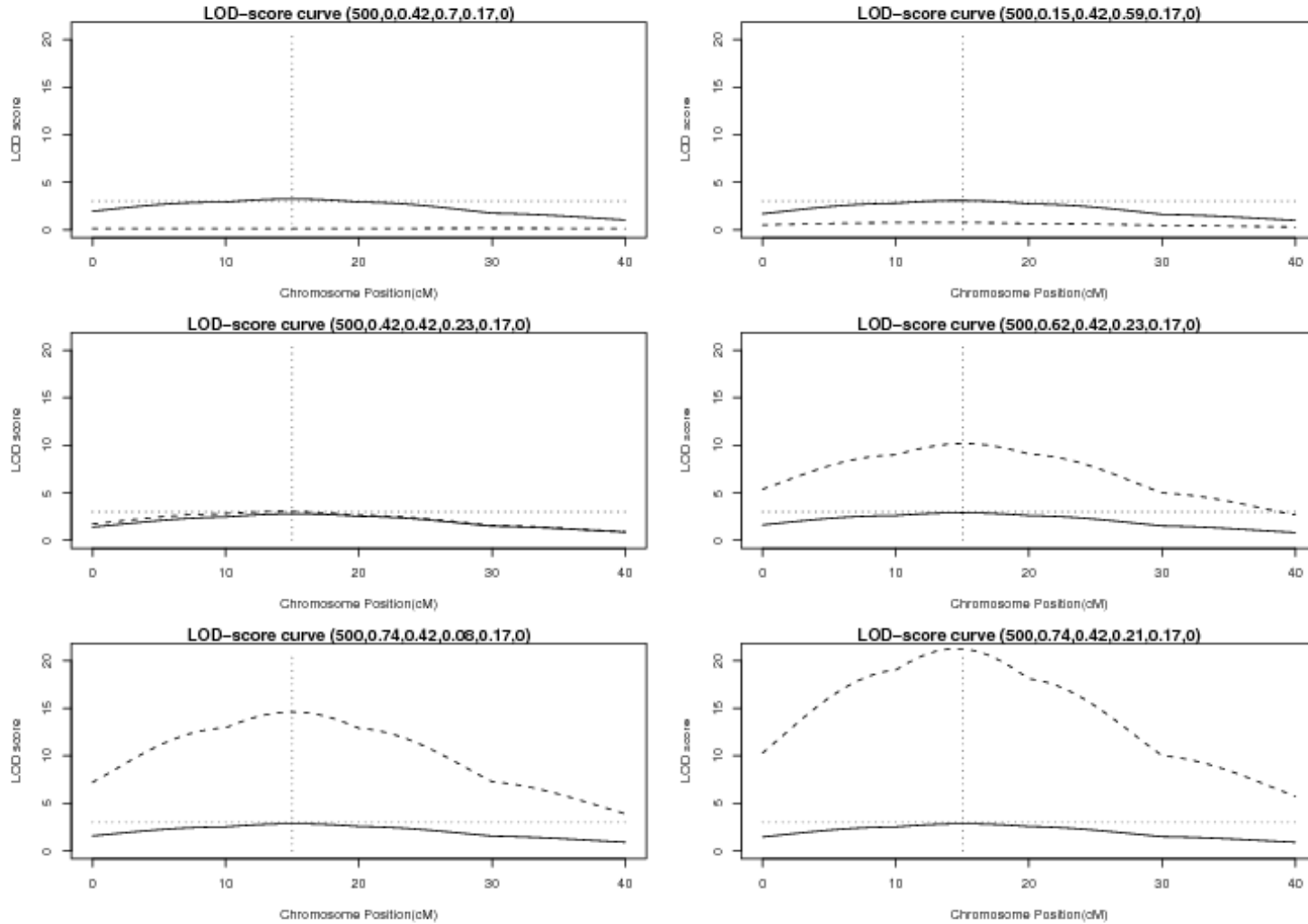


Figure 24: Scenario II mean LOD-score curve with family 500 & $P < E$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G2_P)$ means other heritability of P due to $G2$, and ρ_ϵ means the correlation between non-family deviations of E and P .

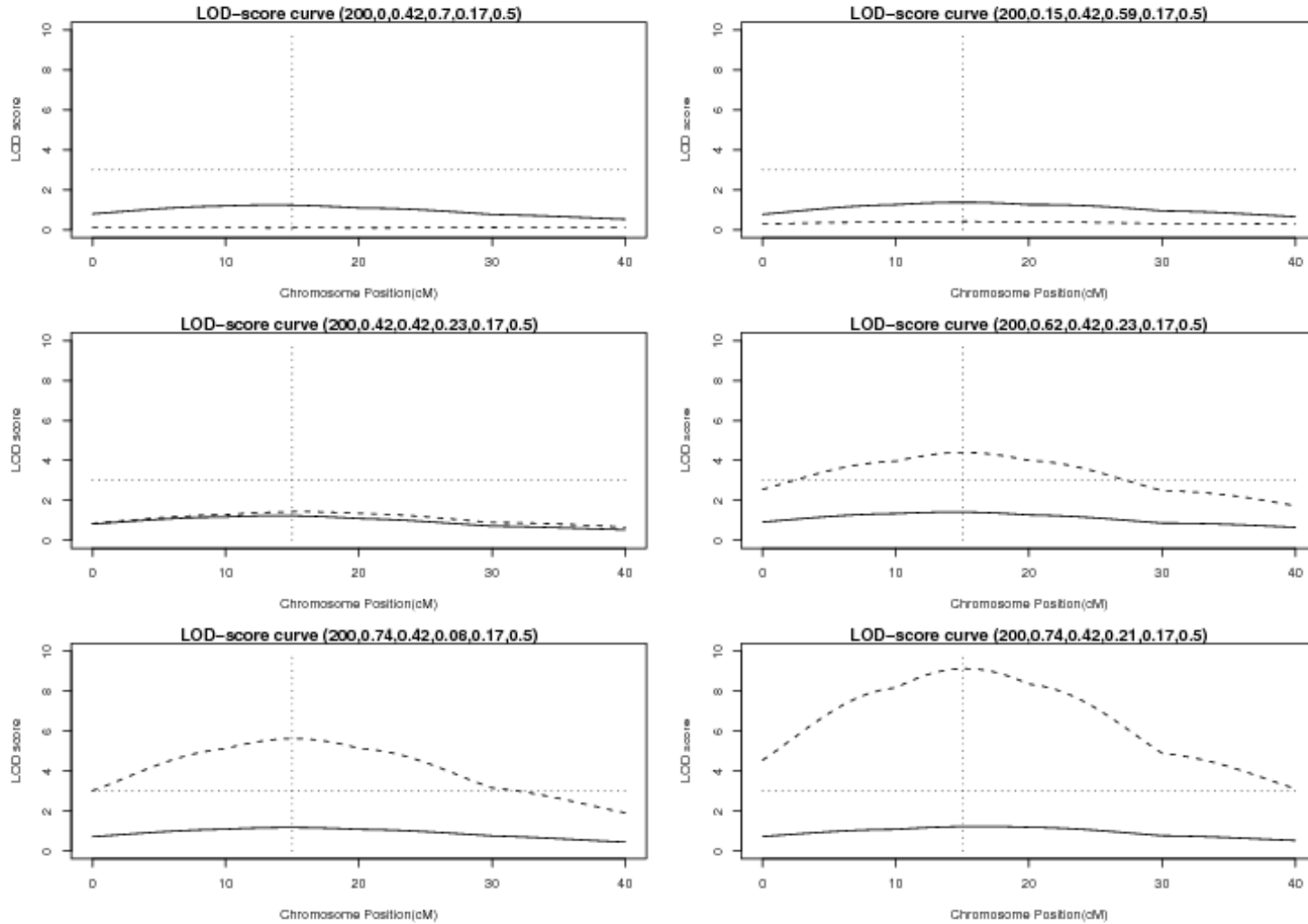


Figure 25: Scenario II mean LOD-score curve with family 200 & $P < E$ & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .

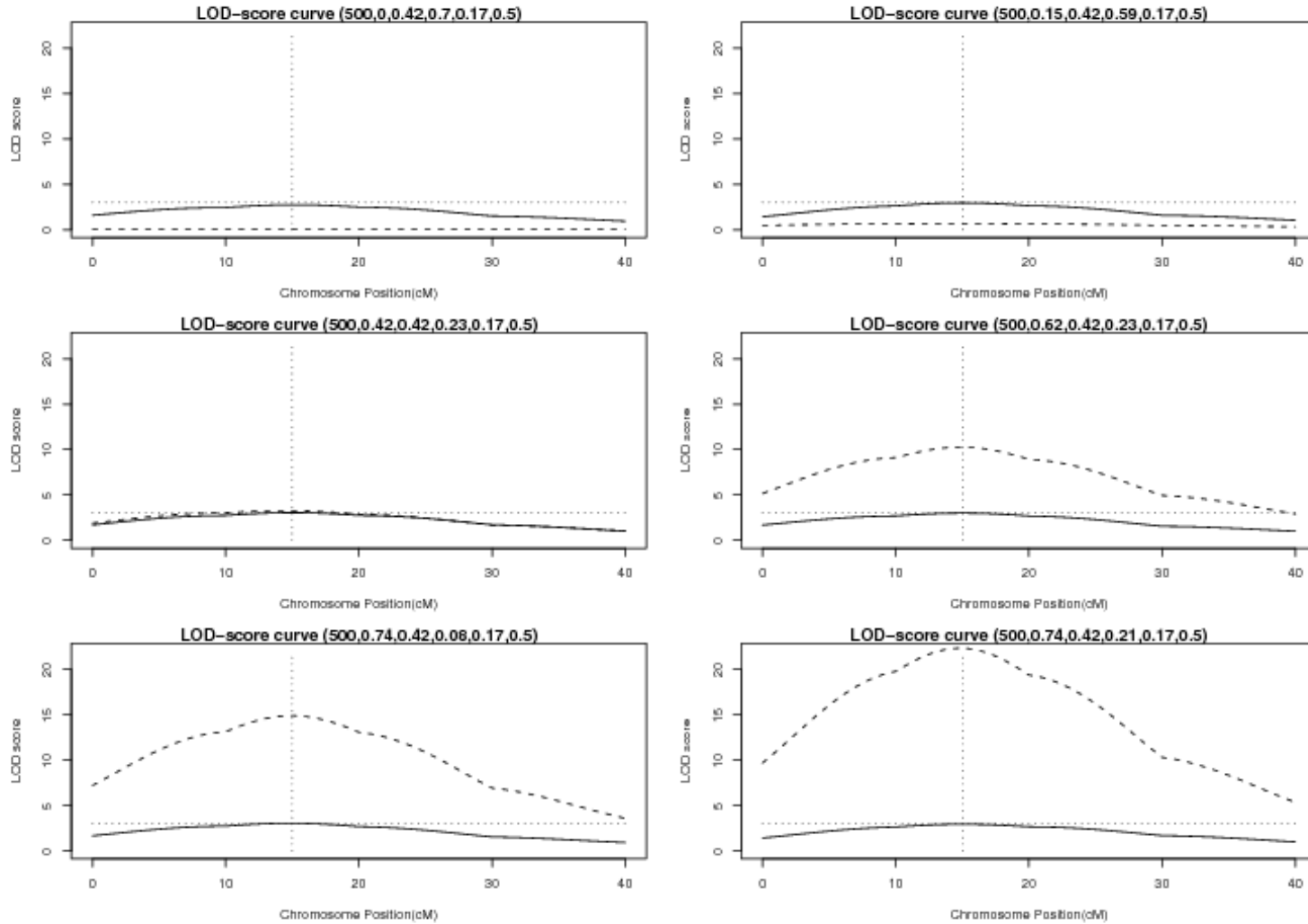


Figure 26: Scenario II mean LOD-score curve with family 500 & $P < E$ & $\rho_\epsilon = 0.5$, where solid line is phenotype, dashed line is endophenotype, horizontal dotted line is LOD-score=3, and vertical dotted line is the position of disease gene. The title in each figure, (fam, $h(G1_E)$, $h(G1_P)$, $h(G2_E)$, $h(G3_P)$, ρ_ϵ), to express these parameters in each situation, where fam means the numbers of family members, $h(G1_E)$ means the heritability of E due to $G1$, $h(G1_P)$ means the heritability of P due to $G1$, $h(G2_E)$ means other heritability of E due to $G2$, $h(G3_P)$ means other heritability of P due to $G3$, and ρ_ϵ means the correlation between non-family deviations of E and P .