

A Hierarchical Model for Simulating a Quantitative Phenotype Mediated by Gene Expression

Introduction

This document details an advanced simulation framework for a quantitative phenotype. It extends the foundational model by incorporating a biologically plausible hierarchical structure that reflects the causal chain from genetic variants to an intermediate molecular trait (gene expression) and finally to a complex organismal phenotype.

This two-stage model simulates a phenotype whose genetic architecture is mediated by the expression levels of a predefined set of genes. The expression of each gene is, in turn, regulated by its own set of genetic loci, known as expression quantitative trait loci (eQTLs). This approach enables a more nuanced examination of how genetic regulation influences complex trait variation.

The simulation will utilize real genotype data for individuals, while the gene expression levels and the final phenotype will be simulated based on the parameters specified in the model.

The Simulation Model

Stage 1: Simulating Gene Expression Levels

The expression level of gene j , GE_j , is the sum of its genetic value ($G_{GE,j}$) and an environmental/residual component ($E_{GE,j}$).

$$GE_j = G_{GE,j} + E_{GE,j}$$

The genetic value is determined by its k_j eQTLs:

$$G_{GE,j} = \sum_{i=1}^{k_j} X_{ji} \alpha_{ji}$$

The environmental component is drawn from a normal distribution, $E_{GE,j} \sim \mathcal{N}(0, \sigma_{E_{GE,j}}^2)$, with its variance determined by the gene expression heritability, h_{GE}^2 :

$$\sigma_{E_{GE,j}}^2 = \text{Var}(G_{GE,j}) \left(\frac{1 - h_{GE}^2}{h_{GE}^2} \right)$$

Stage 2: Simulating the Final Quantitative Phenotype

To accurately model the heritability, we decompose the phenotype P into three independent components:

$$P = G_{P,eQTL} + C_{GE,env} + E_{P,resid}$$

- $G_{P,eQTL}$ (**Pure Genetic Component**): The portion of phenotypic variance strictly mediated by the genetic value of gene expression.

$$G_{P,eQTL} = \sum_{j=1}^m b_j G_{GE,j}$$

- $C_{GE,env}$ (**Downstream Non-Genetic Component**): The variance in the phenotype arising from the non-heritable component of gene expression.

$$C_{GE,env} = \sum_{j=1}^m b_j E_{GE,j}$$

- $E_{P,resid}$ (**Residual Component**): All other sources of variance independent of the m gene expression pathways.

$$E_{P,resid} \sim \mathcal{N}(0, \sigma_{E_{P,resid}}^2)$$

Deriving Residual Variance and a Key Constraint

The heritability parameter $h_{P,eQTL}^2$ (your ‘ h^2 ’) is defined as the proportion of total phenotypic variance explained by the pure genetic component:

$$h_{P,eQTL}^2 = \frac{\text{Var}(G_{P,eQTL})}{\text{Var}(P)}$$

The total non-genetic variance is therefore:

$$\text{Var}(C_{GE,env}) + \text{Var}(E_{P,resid}) = \text{Var}(G_{P,eQTL}) \left(\frac{1 - h_{P,eQTL}^2}{h_{P,eQTL}^2} \right)$$

From this, we solve for the variance of the final residual term:

$$\sigma_{E_{P,resid}}^2 = \left[\text{Var}(G_{P,eQTL}) \left(\frac{1 - h_{P,eQTL}^2}{h_{P,eQTL}^2} \right) \right] - \text{Var}(C_{GE,env})$$

Condition for a Valid Simulation

For $\sigma_{E_{P,resid}}^2$ to be non-negative, the total non-genetic variance allowed by $h_{P,eQTL}^2$ must be at least as large as the non-genetic variance already introduced from the gene expression layer. This leads to a simple and critical constraint on the input parameters:

$$h^2_{P,eQTL} \leq h^2_{GE}$$

In words, the heritability of the final phenotype that is mediated by eQTLs cannot exceed the heritability of the intermediate gene expression traits themselves. This check must be performed prior to simulation to ensure the parameter set is mathematically and biologically plausible.