

Scan variance QTL for GxE interaction

Xiaoran, Dmitri

- 1 Two types of variance loci (VLC)
- 2 identify variance loci
- 3 Simulations
- 4 Proof
- 5 Pit falls

GxE induced variance loci

Let \mathbf{y} be affected by $\mathbf{g} \in \{0, 1, 2\}$, unobserved \mathbf{u} , and interaction \mathbf{gu}

$$\mathbf{y} = m_0 + a\mathbf{g} + b\mathbf{gu} + c\mathbf{u} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, v_\epsilon), \quad \mathbf{g} \perp \mathbf{u} \quad (1)$$

Fitting the mean of \mathbf{y} given \mathbf{g} , but since \mathbf{u} is unseen, we have

$$\mathbf{y} = m_0 + a\mathbf{g} + \mathbf{e} \quad \text{or} \quad E(\mathbf{y}|\mathbf{g}) = m_0 + a\mathbf{g} \quad (2)$$

Unexplained $(c + b\mathbf{g})\mathbf{u} + \epsilon$ flood into residual \mathbf{e} , as a result,

$$V(\mathbf{y}|\mathbf{g}) = E(\mathbf{e}^2|\mathbf{g}) = (v_\epsilon + c^2v_u) + 2bcv_u\mathbf{g} + b^2v_u\mathbf{g}^2 \quad (3)$$

The phenotype variance depends on the genotype, thus \mathbf{g} is a "variance loci"; when $2bcv_u \neq 0$ or $b^2v_u > 0$, we know $b \neq 0$, and the presence of GxE, $b\mathbf{gu}$, even if we do not observe \mathbf{u} .

Suppose \mathbf{g} directly determine the mean and variance of \mathbf{y} ,

$$\mathbf{y} \sim \mathcal{N}(m_0 + a\mathbf{g}, \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2)) \quad (4)$$

again, explain the mean of \mathbf{y} with \mathbf{g} gets

$$\mathbf{y} = m_0 + a\mathbf{g} + \mathbf{e}, \quad \text{or} \quad E(\mathbf{y}|\mathbf{g}) = m_0 + a\mathbf{g} \quad (5)$$

The squared leftover \mathbf{e}^2 is Gamma distributed,

$$\mathbf{e}^2 \sim \text{Gamma}(\text{shape} = 2, \text{scale} = \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2)) \quad (6)$$

$$V(\mathbf{y}|\mathbf{g}) = E(\mathbf{e}^2|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2) \quad (7)$$

When $\alpha \neq 0$ or $\beta \neq 0$, \mathbf{g} is a variance loci by definition, without an origin.

Connection between two types of variance loci

If $\mathbf{y} = m_0 + a\mathbf{g} + b\mathbf{g}\mathbf{u} + c\mathbf{u} + \epsilon$, \mathbf{g} is a GxE induced variance loci, we see

$$V(\mathbf{y}|\mathbf{g}) = (v_e + c^2 v_u) + 2bcv_u \mathbf{g} + b^2 v_u \mathbf{g}^2$$

If we force \mathbf{g} to be a direct variance loci such that

$$V(\mathbf{y}|\mathbf{g}) = \exp(\tilde{\mu} + \tilde{\alpha}\mathbf{g} + \tilde{\beta}\mathbf{g}^2).$$

As a result, the following approximation holds

$$\tilde{\mu} \approx \log(v_e), \quad \tilde{\alpha} \approx 2bcv_u/v_e, \quad \tilde{\beta} \approx b^2 v_u/v_e - 2b^2 c^2 v_u^2/v_e^2 \quad (8)$$

where $v_e = v_e + c^2 v_u$ is part of the variance unaffected by \mathbf{g} . When white noise $v_e = 0$, the above simplifies to

$$\tilde{\mu} \approx 2 \log c + \log v_u, \quad \tilde{\alpha} \approx 2b/c, \quad \tilde{\beta} \approx -b^2/c^2 \quad (9)$$

Connection between two types of variance loci

If $\mathbf{y} \sim \mathcal{N}(m_0 + \mathbf{a}\mathbf{g}, \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2))$, \mathbf{g} is a direct variance locus,

$$V(\mathbf{y}|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2)$$

Perceive \mathbf{g} as a GxE induced variance loci with latent environment $\tilde{\mathbf{u}}$,

$$V(\mathbf{y}|\mathbf{g}) = (v_{\tilde{\epsilon}} + \tilde{c}^2 v_{\tilde{\mathbf{u}}}) + 2\tilde{b}\tilde{c}v_{\tilde{\mathbf{u}}}\mathbf{g} + \tilde{b}^2 v_{\tilde{\mathbf{u}}}\mathbf{g}^2,$$

As a result, the following approximation holds

$$(v_{\tilde{\epsilon}} + \tilde{c}^2 v_{\tilde{\mathbf{u}}}) \approx e^{\mu}, \quad 2\tilde{b}\tilde{c}v_{\tilde{\mathbf{u}}} \approx \alpha e^{\mu}, \quad \tilde{b}^2 v_{\tilde{\mathbf{u}}} \approx (\alpha^2 + 2\beta)e^{\mu}/2 \quad (10)$$

If we let $v_{\tilde{\epsilon}} = 0$, naturally, $\tilde{c}^2 v_{\tilde{\mathbf{u}}} \approx e^{\mu}$, and consequently

$$\mu \approx 2 \log \tilde{c} + \log v_{\tilde{\mathbf{u}}}, \quad \alpha \approx 2\tilde{b}/\tilde{c}, \quad \beta \approx -\tilde{b}^2/\tilde{c}^2, \quad (11)$$

Are they equivalent?

Direct variance loci are too generic

- given v_u , b , c uniquely determine $\tilde{\mu}$, $\tilde{\alpha}$, and $\tilde{\beta}$;
- given v_u , μ , α , and β can result in unreal \tilde{b} , since $\tilde{b}^2 \propto (\alpha^2 + 2\beta)e^\mu$, and $\alpha^2 + 2\beta \in (-\infty, \infty)$

Direct variance loci do not necessarily indicate GxE.

Direct variance loci are too rigid, if restricting $\beta = 0$

- restricting $V(\mathbf{y}|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g})$ is a mainstream practice.
- no risk of \tilde{b} being unreal, but imposes $\tilde{b}/\tilde{c} = \alpha$.

Direct variance loci with $\beta = 0$ do not capture all GxE.

Are they equivalent?

How about restricted GxE vQTL ($b = 0$)

- restricted $V(\mathbf{y}|\mathbf{g}) = a_1 + a_2\mathbf{g}$ is popular;
 - $a_2\mathbf{g} \approx 2bcv_u\mathbf{g}$ also capture GxE effect b ;
 - when MAF is low, absorb $a_2\mathbf{g}$ part of $b^2v_u\mathbf{g}^2$, since $\mathbf{g} \not\perp \mathbf{g}^2$,
- lower power when environmental main effect is low ($c \approx 0$), or when MAF is moderate, causing \mathbf{g} and \mathbf{g}^2 nearly independent.
- less specific when environmental main effect is big ($|c| > 0$), since a significant $|a_2|$ does not necessarily imply GxE effect $|b|$ is large.
- deviated from the nature result of GxE, since $b^2v_u\mathbf{g}^2$ is deduced from the premises that $b\mathbf{g}u$ exists.

Are they equivalent?

`img/sim/vQTL_space.pdf`

The essence of GxE caused vQTL

Fitting the squared residual e^2 with squared SNP g^2

$$e^2 = (v_\epsilon + c^2 v_u) + 2bcv_u g + b^2 v_u g^2 \quad (12)$$

implies

$$\frac{\partial e^2}{\partial g} = 2bcv_u + 2b^2 v_u g \quad (13)$$

as a result, rejecting $h_0 : b^2 v_u \leq 0$ concludes that the derivative of e^2 wrt. g is positive, or, the curve of e^2 over g must bend upwards.

The goal is to screen for GxE candidate

- favor the nature model deduced from GxE, over Gaussian vQTL.
 - $V(\mathbf{y}|\mathbf{g}) = (v_\epsilon + c^2 v_u) + 2bcv_u \mathbf{g} + b^2 v_u \mathbf{g}^2$
- consider the quadratic term \mathbf{g}^2 ,
 - model $V(\mathbf{y}|\mathbf{g}) = a_1 + a_2 \mathbf{g} + a_3 \mathbf{g}^2$, test $h_0 : a_2 = a_3 = 0$.
- exploit the mutual absorption between \mathbf{g} and \mathbf{g}^2
 - model use $V(\mathbf{y}|\mathbf{g}) = a_1 + a_3 \mathbf{g}^2$, test $h_0 : a_3 \leq 0$
 - $a_3 \mathbf{g}^2 \approx b^2 v_u \mathbf{g}^2$ is GxE specific;
 - absorb $a_2 \approx 2bcv_u \mathbf{g}$ when \mathbf{g} and \mathbf{g}^2 is near parallel;
 - more powerful one tail t-test, while expecting $a_3 \geq 0$.

Simulated Truth:

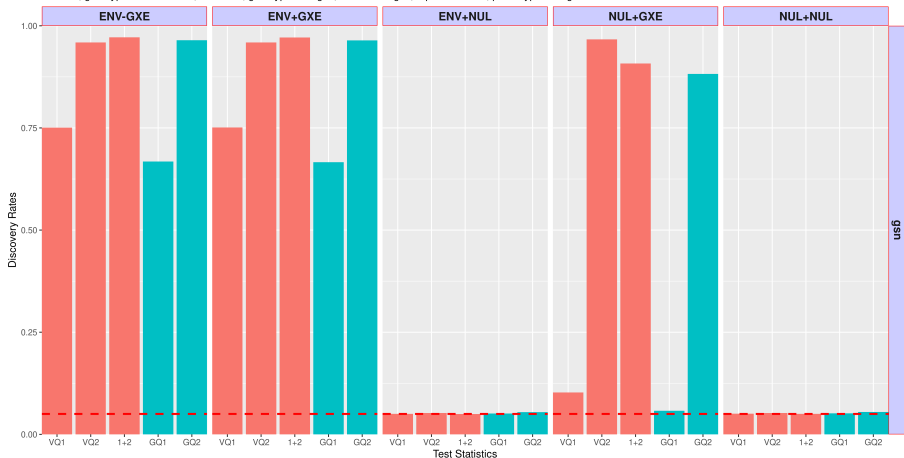
- for power check:
 - ENV \pm GXE: $\mathbf{y} = \mathbf{c}\mathbf{u} \pm \mathbf{b}\mathbf{g}\mathbf{u}$, environmental main and GxE;
 - NUL + GXE: $\mathbf{y} = \mathbf{0}\mathbf{u} + \mathbf{b}\mathbf{g}\mathbf{u}$, GxE only;
- for type 1 error check
 - ENV + NUL: $\mathbf{y} = \mathbf{c}\mathbf{u} + \mathbf{0}\mathbf{g}\mathbf{u}$, environmental main only;
 - NUL + NUL: $\mathbf{y} = \mathbf{0}\mathbf{u} + \mathbf{0}\mathbf{g}\mathbf{u}$, total null;

Competing Test Statistics

- VQ1: $V(\mathbf{y}|\mathbf{g}) = \mathbf{a}_1 + \mathbf{a}_2\mathbf{g}$, mainstream linear vQTL 2-tail test
- VQ2: $V(\mathbf{y}|\mathbf{g}) = \mathbf{a}_1 + \mathbf{a}_3\mathbf{g}^2$, proposed 1-tail test ($\mathbf{a}_3 \geq 0$)
- 1+2: $V(\mathbf{y}|\mathbf{g}) = \mathbf{a}_1 + \mathbf{a}_2\mathbf{g} + \mathbf{a}_3\mathbf{g}^2$, full model
- GQ1: $V(\mathbf{y}|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g})$, mainstream Gaussian vQTL 2-tail test
- GQ2: $V(\mathbf{y}|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2)$, full Gaussian vQTL 2-tail test

Genotype normal, unrealistically eliminate absorption between g and g^2

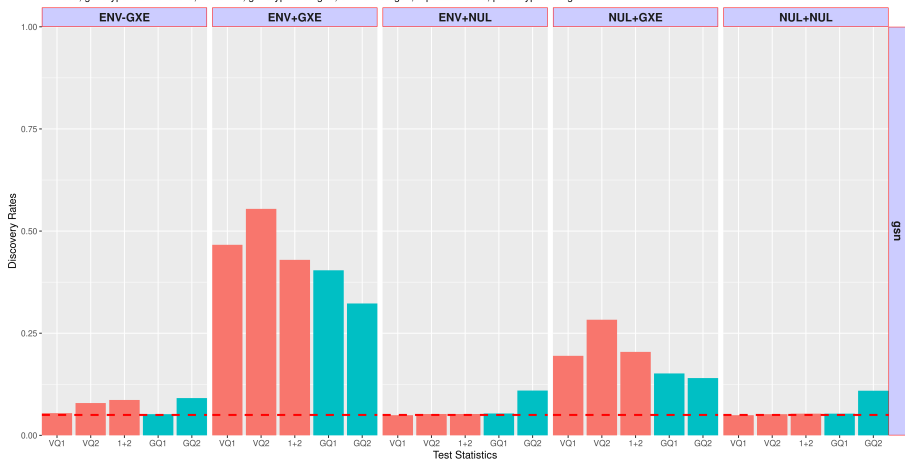
N=2000, genotype main effect=0, noise=10, genotype dist=gs0, enviro dist=gsn, repeats=2e+05, phenotype dist=gsn



Gaussian Environment, Binomial Genotype

Genotype is Binomial(0.05, 2), $g \neq g^2$, heavy absorption.

N=2000, genotype main effect=0, noise=10, genotype dist=g05, enviro dist=gsn, repeats=2e+05, phenotype dist=gsn



Gaussian Environment, 1000G project data

Genotype is drawn from 1000G, $MAF > 0.05$, realistic absorption.

N=1500, genotype main effect=0, noise=10, genotype dist=k05, enviro dist=gsn, repeats=2e+05, phenotype dist=gsn



Pseudo ground truth of GxE

- ENV: $env = alc + smk + edu + inc$
- GWA: $bmi = snp + age * sex + env$
- GXE: $bmi = snp + age * sex + env + snp : env$
- test GWA = GXE, rejection counts snp as a Pseudo true GxE SNP

Q: how many pseudo true GxE SNP passes various vQTL screening tests – the positive predictive value (PPV), relative to no screening at all?

Address issues with binary outcome

I1: T2D case, I2: T2D non-case

Pseudo Real Data analysis for BMI

img/bmi.png

Proof: polynomial Gaussian vQTL \rightarrow interaction

If \mathbf{g} is a polynomial Gaussian vQTL that

$$V(\mathbf{y}|\mathbf{g}) = \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2) \quad (14)$$

$$\rightarrow \frac{V(\mathbf{y}|\mathbf{g})}{e^\mu} = 1 + \alpha\mathbf{g} + \frac{(\alpha^2 + 2\beta)}{2}\mathbf{g}^2 + \sum_{n=3}^{\infty} \frac{s^{(n)}}{n!}\mathbf{g}^n \quad (15)$$

$$s^{(n)} = \begin{cases} s^{(n-1)}\alpha + 2s^{(n-2)}\beta, & n > 0 \\ 1, & n = 0 \\ 0, & n < 0 \end{cases} \quad (16)$$

Here we Taylor expand $f(\mathbf{g}) = e^{\alpha\mathbf{g} + \beta\mathbf{g}^2}$ at $\mathbf{0}$.

Proof: polynomial Gaussian vQTL \rightarrow interaction

When $\mathbf{g} \sim N(0, v_g)$, $v_g < 1$, $|\alpha| < 1$ and $|\beta| < 1$, approximation

$$V(\mathbf{y}|\mathbf{g}) \approx e^\mu + \alpha e^\mu \mathbf{g} + \frac{(\alpha^2 + 2\beta)e^\mu}{2} \mathbf{g}^2 \quad (17)$$

is accurate. If we force the vQTL to be interaction induced, that is,

$$V(\mathbf{y}|\mathbf{g}) = c^2 v_u + 2bcv_u \mathbf{g} + b^2 v_u \mathbf{g}^2,$$

Naturally, $e^\mu = c^2 v_u$, and

$$\mu = 2 \log c + \log v_u, \quad \alpha = \frac{2b}{c}, \quad \beta = -\frac{b^2}{c^2}, \quad (18)$$

Proof: interaction \rightarrow polynomial Gaussian vQTL

If \mathbf{g} is interaction induced vQTL, $\mathbf{y} = \mathbf{a}\mathbf{g} + \mathbf{b}\mathbf{g}\mathbf{u} + \mathbf{c}\mathbf{u} + \epsilon$, we know

$$\begin{aligned} V(\mathbf{y}|\mathbf{g}) &= v_\epsilon + c^2 v_u + 2bcv_u \mathbf{g} + b^2 v_u \mathbf{g}^2 \\ \rightarrow \log V(\mathbf{y}|\mathbf{g}) &= \log(v_\epsilon + c^2 v_u) + \frac{2bcv_u}{v_\epsilon + c^2 v_u} \mathbf{g} + \left[\frac{b^2 v_u}{v_\epsilon + c^2 v_u} - \frac{2b^2 c^2 v_u^2}{(v_\epsilon + c^2 v_u)^2} \right] \mathbf{g}^2 \\ &+ \sum_{n=3}^{\infty} \left[\frac{(-1)^n (2b^2 v_u)^{n-2}}{n(v_\epsilon + c^2 v_u)^{n-1}} + \frac{(-1)^{n-1} (2bcv_u)^n}{n(v_\epsilon + c^2 v_u)^n} \right] \mathbf{g}^n \\ &\approx \log(v_\epsilon + c^2 v_u) + \frac{2bcv_u}{v_\epsilon + c^2 v_u} \mathbf{g} + \left[\frac{b^2 v_u}{v_\epsilon + c^2 v_u} - \frac{2b^2 c^2 v_u^2}{(v_\epsilon + c^2 v_u)^2} \right] \mathbf{g}^2 \end{aligned}$$

Here we Taylor expand $f(\mathbf{g}) = \log V(\mathbf{y}|\mathbf{g})$ at $\mathbf{0}$.

The approximation up to the second order

$$\log \phi \approx \log v_u + 2 \log c + \frac{2b}{c} \mathbf{g} - \frac{b^2}{c^2} \mathbf{g}^2 = \tilde{\mu} + \tilde{\alpha} \mathbf{g} + \tilde{\beta} \mathbf{g}^2 \quad (19)$$

is accurate when $\mathbf{g} \sim N(0, v_g)$, $v_g < 1$, and $|c| > |b|$. Thus,

$\mathbf{y} = a\mathbf{g} + b\mathbf{g}\mathbf{u} + c\mathbf{u}$, give rise to $V(\mathbf{y}|\mathbf{g}) = e^{\tilde{\mu} + \tilde{\alpha}\mathbf{g} + \tilde{\beta}\mathbf{g}^2}$; $\tilde{\mu}$ is the effect of \mathbf{u} on \mathbf{y} ; $\tilde{\alpha}$ and $\tilde{\beta}$ relate to the effect size of interaction $\mathbf{g}\mathbf{u}$ on \mathbf{y} relative to the effect size of \mathbf{u} on \mathbf{y} .

A polynomial Gaussian vQTL induce interaction between genotype \mathbf{g} and latent variable $\tilde{\mathbf{u}}$, when $\beta = -0.25\alpha^2$.

Non-centered interaction

When \mathbf{u} is the difference from non-zero mean u_0

$$\mathbf{y} = a\mathbf{g} + b\mathbf{g}(\mathbf{u} + u_0) + c(\mathbf{u} + u_0), \quad \mathbf{g} \perp \mathbf{u} \quad (20)$$

Explain the mean of \mathbf{y} with $\mathbf{u} + u_0$ unseen,

$$\begin{aligned} \mathbf{y} &= cu_0 + (a + bu_0)\mathbf{g} + \mathbf{e}, \quad \text{or} \\ E(\mathbf{y}|\mathbf{g}) &= cu_0 + (a + bu_0)\mathbf{g} \end{aligned} \quad (21)$$

Notice spurious genetic effect bu_0 . Still, $(c + b\mathbf{g})\mathbf{u}$ flow into \mathbf{e} , therefore vQTL is unaffected.

Suppose any one of the phenotypes follows a multivariate normal. Under the basic null, we say

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \sigma_0^2 \mathbf{I} + \boldsymbol{\Sigma}) \quad (22)$$

$$\boldsymbol{\mu} = \mathbf{X}_c \boldsymbol{\alpha}, \quad (23)$$

$$\boldsymbol{\Sigma} = \sigma_t^2 \mathbf{K}_t + \sigma_m^2 \mathbf{K}_m + \sigma_s^2 \mathbf{K}_s + \sigma_u^2 \mathbf{K}_u \quad (24)$$

where \mathbf{X}_c is the matrix of confounders; and \mathbf{K}_* are kernels of transcriptom, methylation, expression, seurm metabolom, and urine metabolome. This null model left out exposome.

We include exposome, one at a time, say \mathbf{z}

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \sigma_0^2 \mathbf{I} + \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_{v_i}) \quad (25)$$

$$\boldsymbol{\mu} = \mathbf{X}_c \boldsymbol{\alpha} + \mathbf{z} \boldsymbol{\beta}, \quad (26)$$

$$\boldsymbol{\Sigma} = \sigma_t^2 \mathbf{K}_t + \sigma_m^2 \mathbf{K}_m + \sigma_s^2 \mathbf{K}_s + \sigma_u^2 \mathbf{K}_u \quad (27)$$

$$+ \tau_t^2 \mathbf{V} \mathbf{K}_t + \tau_m^2 \mathbf{V} \mathbf{K}_m + \tau_s^2 \mathbf{V} \mathbf{K}_s + \tau_u^2 \mathbf{V} \mathbf{K}_u \quad (28)$$

$$\mathbf{V} = \mathbf{z}' \mathbf{z} \quad (29)$$

where \mathbf{z} is an exposome variant; \mathbf{V} is the kernel of that variants. This model the alternative, with exposome.

Testing $H_0: \tau_*^2 = 0$ amounts to test, exposome is not interacting with the omics.

Let \mathbf{g} be a vQTL by directly determine the variance of \mathbf{y} ,

$$\mu_{\mathbf{y}} = h^{-1}(m_0 + a\mathbf{g}), \quad \mathbf{v}_{\mathbf{y}} = \theta(\mu_{\mathbf{y}})\phi^{-1}(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2), \quad (30)$$

If \mathbf{y} is **Gaussian**, then $h^{-1}(\mathbf{x}) = \mathbf{x}$ and $\theta(\mu_{\mathbf{y}}) = 1$, and

$$\mathbf{y} = m_0 + a\mathbf{g} + \mathbf{e}, \quad \text{or} \quad E(\mathbf{y}|\mathbf{g}) = \mu_{\mathbf{y}} = m_0 + a\mathbf{g} \quad (31)$$

The squared residual \mathbf{e}^2 is Gamma distributed

$$\mathbf{e}^2 \sim \text{Gamma}(1, \phi^{-1}(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2)), \quad (32)$$

$$\text{or} \quad V(\mathbf{y}|\mathbf{g}) = E(\mathbf{e}^2|\mathbf{g}) = \phi^{-1}(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2) \quad (33)$$

When $\alpha \neq 0$ or $\beta \neq 0$, \mathbf{g} is a vQTL with no specif origin (i.e., due to GxE).

Connection between two types vQTL

If $\mathbf{y} \sim \mathcal{N}(a\mathbf{g}, \exp(\mu + \alpha\mathbf{g} + \beta\mathbf{g}^2))$ and $\beta = -0.25\alpha^2$, \mathbf{g} is a constrained polynomial Gaussian vQTL, then \mathbf{g} interacts with latent $\tilde{\mathbf{u}}$ such that

$$\mathbf{y} = a\mathbf{g} + \tilde{b}\mathbf{g}\tilde{\mathbf{u}} + \tilde{c}\tilde{\mathbf{u}} \quad (34)$$

$$\tilde{b} \approx \frac{\tilde{c}}{2}\alpha, \quad \tilde{c} \approx \sqrt{\frac{e^\mu}{v_{\tilde{\mathbf{u}}}}}, \quad \tilde{\mathbf{u}} \perp \mathbf{g}, \quad v_{\tilde{\mathbf{u}}} = \text{Var}(\mathbf{u}) \quad (35)$$

The goal is to detect Gene by Environment interaction

- with measured environment (i.e., u become visible),
- with variants pre-selected by vQTL scan across genome.

We prefer interaction induced vQTL,

- prefer linear model over log-Gamma regression, the later may detect vQTL not specifically due to interaction.
- include quadratic term g^2 to enhance the detection.

Our analysis of Type 2 diabetes:

- scan vQTL on UK Biobank blood assays for GxE SNPs;
- pairing the SNPs with measured environment in PEGS;
- analyze the pairs for GxE interaction affecting T2D.