

Jiaqi's Thesis Progress Report (Updated Mar. 29)

Jiaqi Bi^a

^a*Western University,
Schulich School of Medicine & Dentistry,
Department of Epidemiology and Biostatistics*

1. To Do List

1. MCEM - Both univariate and multivariate missing data
2. Add a section that talk about when dealing with missing data without using joint distribution of frailty and missing covariate

2. Notations

List of Notations

i	Individual index
j	Family (Cluster) index
$h(\cdot)$	Hazard fuction
$h_0(\cdot)$	Baseline hazard function
t	Failure/Event or Censored Time
n	Number of individuals

3. Weibull Parametric Approach (Discussions on why use parametric?)

From the beginning of the discussion, I have obtained the model, i.e., the hazard function is

$$h_{ij}(t_{ij}|z_j) = h_0(t_{ij}) \exp(\beta_1 x_{1,ij} + \beta_2 x_{2,ij}) z_j \quad (1)$$

There are total n_j individuals in family j , where $i = 1, \dots, n_j$, and total J families that $j = 1, \dots, J$. $x_{1,ij}$ is the genotype, or say mutation gene status for individual i in family j . $x_{2,ij}$ is the PRS for individual i in family j . The frailty term z_j , has a pdf of $f(z)$, which can be Gamma, log-normal, or other common frailty distributions. The support of $f(z)$ is always non-negative. The Weibull baseline hazard function is defined as

$$h_0(t_{ij}) = \alpha \lambda t_{ij}^{\lambda-1} \quad (2)$$

Email address: jbi23@uwo.ca (Jiaqi Bi)

where λ is the shape parameter and α is the scale parameter. Let $\xi_{ij} = \exp(\beta_1 x_{1,ij} + \beta_2 x_{2,ij})$, the hazard function is

$$h_{ij}(t_{ij}|x_{ij}, g_{ij}, z_j) = \alpha \lambda t_{ij}^{\lambda-1} \xi_{ij} z_j \quad (3)$$

The survival function $S(t)$ can be obtained through cumulative hazard function $H(t)$

$$H(t_{ij}|x_{ij}, g_{ij}, z_j) = \int_0^t h_{ij}(u|x_{ij}, g_{ij}, z_j) du \quad (4)$$

$$= \alpha \xi_{ij} z_j \lambda \int_0^t u^{\lambda-1} du \quad (5)$$

$$= \alpha \xi_{ij} z_j \lambda \cdot \frac{1}{\lambda} t_{ij}^\lambda = \alpha \xi_{ij} z_j t_{ij}^\lambda \quad (6)$$

and the survival function

$$S(t_{ij}|x_{ij}, g_{ij}, z_j) = \exp(-H(t_{ij}|x_{ij}, g_{ij}, z_j)) = \exp(-\alpha \xi_{ij} z_j t_{ij}^\lambda) \quad (7)$$

Let $\boldsymbol{\theta} = \{\beta_1, \beta_2, \alpha, \lambda, \boldsymbol{\phi}\}$, where $\boldsymbol{\phi}$ is the parameter vector for the frailty distribution of the choice. Therefore, the likelihood can be written as

$$L(\boldsymbol{\theta}) = \prod_{j=1}^J \int_0^\infty \prod_{i=1}^{n_j} (\alpha \lambda t_{ij}^{\lambda-1} \xi_{ij} z_j)^{\delta_{ij}} \exp(-\alpha \xi_{ij} z_j t_{ij}^\lambda) f(z) dz \quad (8)$$

$$= \prod_{j=1}^J \int_0^\infty \prod_{i=1}^{n_j} h(t_{ij}|\mathbf{x}_{ij}, z_j)^{\delta_{ij}} \exp(-H(t_{ij}|\mathbf{x}_{ij}, z_j)) f(z) dz \quad (9)$$

So the log-likelihood is

$$\ell(\boldsymbol{\theta}) = \sum_{j=1}^J \log \left[\int_0^\infty \prod_{i=1}^{n_j} h(t_{ij}|\mathbf{x}_{ij}, z_j)^{\delta_{ij}} \exp(-H(t_{ij}|\mathbf{x}_{ij}, z_j)) f(z) dz \right] \quad (10)$$

4. Gamma Frailty

The Laplace transform of the frailty $z \sim \text{Gamma}(k, k)$, for the simplicity of the mathematical expression, the following Laplace transform will ignore the subscript, denote $\mathcal{L}(f(z)) = \phi(s)$ where $s = \sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij})$:

$$\phi(s) = \int_0^\infty e^{-sz} f(z) dz \quad (11)$$

$$= \int_0^\infty e^{-sz} \frac{k^k}{\Gamma(k)} z^{k-1} e^{-kz} dz \quad (12)$$

Using the Gamma property: $\int_0^\infty z^{n-1} e^{-az} dz = \frac{\Gamma(n)}{a^n}$, $\phi(s)$ can be further written as

$$\phi(s) = \frac{k^k}{\Gamma(k)} \int_0^\infty e^{-(s+k)z} z^{k-1} dz = \frac{k^k}{\Gamma(k)} \cdot \frac{\Gamma(k)}{(s+k)^k} = \left(1 + \frac{s}{k}\right)^{-k} \quad (13)$$

- 21 The second derivative is $\frac{d^2\phi(s)}{ds^2} = \int_0^\infty (-z)^2 e^{-sz} f(z) dz$.
 The third derivative is $\frac{d^3\phi(s)}{ds^3} = \int_0^\infty (-z)^3 e^{-sz} f(z) dz, \dots$ Therefore, its d -th derivative, denote $\phi(s)^{(d)}$:

$$\phi(s)^{(d)} = (-1)^d \int_0^\infty z^d e^{-sz} f(z) dz \quad (14)$$

$$= (-1)^d \frac{(k+d-1)!}{(k-1)!(s+k)^d} \left(1 + \frac{s}{k}\right)^{-k} \quad (15)$$

Let $\boldsymbol{\theta} = (\beta_1, \beta_2, \alpha, \lambda, k)$ for Gamma frailty model, the log-likelihood is then written as

$$\ell(\boldsymbol{\theta}) = \sum_{j=1}^k \log \left[\int_0^\infty \prod_{i=1}^{n_j} (h(t_{ij}|\mathbf{x}_{ij}, z_j))^{\delta_{ij}} \exp(-H(t_{ij}|\mathbf{x}_{ij}, z_j)) f(z_j) dz_j \right] \quad (16)$$

$$= \sum_{j=1}^J \log \left[\int_0^\infty \prod_{i=1}^{n_j} (z_j h(t_{ij}|\mathbf{x}_{ij}))^{\delta_{ij}} \exp(-z_j H(t_{ij}|\mathbf{x}_{ij})) f(z_j) dz_j \right] \quad (17)$$

$$= \sum_{j=1}^J \log \left[\prod_{i=1}^{n_j} (h(t_{ij}|\mathbf{x}_{ij}))^{\delta_{ij}} \int_0^\infty z_j^{d_j} \exp(-z_j \sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij})) f(z_j) dz_j \right] \quad (18)$$

$$= \sum_{j=1}^J \log \left[\prod_{i=1}^{n_j} (h(t_{ij}|\mathbf{x}_{ij}))^{\delta_{ij}} \frac{(k+d_j-1)!}{(k-1)!(\sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij}) + k)^{d_j}} \left(1 + \frac{\sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij})}{k}\right)^{-k} \right] \quad (19)$$

$$= \sum_{j=1}^J \log \left[\prod_{i=1}^{n_j} ((h(t_{ij}|\mathbf{x}_{ij}))^{\delta_{ij}}) \frac{(k+d_j-1)!}{k!k^{d_j-1}} \left(1 + \frac{\sum_{i=1}^{n_j} (H(t_{ij}|\mathbf{x}_{ij}))}{k}\right)^{-k-d_j} \right] \quad (20)$$

$$= \sum_{j=1}^J \log \left[h(t_{ij}|\mathbf{x}_{ij})^{\delta_{ij}} \frac{(k+d_j-1)!}{k!k^{d_j-1}} \left(1 + \frac{\sum_{i=1}^{n_j} (H(t_{ij}|\mathbf{x}_{ij}))}{k}\right)^{-k-d_j} \right] \quad (21)$$

$$= \sum_{j=1}^J \left[\sum_{i=1}^{n_j} (\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij})) + \log \left(\frac{(k+d_j-1)!}{k!k^{d_j-1}} \left(1 + \frac{\sum_{i=1}^{n_j} (H(t_{ij}|\mathbf{x}_{ij}))}{k}\right)^{-k-d_j} \right) \right] \quad (22)$$

- 22 For each family j , the ascertainment A_j is defined to be the probability of the proband p
 23 being ascertained by the age a_{j_p} at examination. Applying the ascertainment correction for
 24 the log-likelihood in family j :

$$\tilde{\ell}_j(\boldsymbol{\theta}) = \ell_j(\boldsymbol{\theta}) - \log A_j(\boldsymbol{\theta}) \quad (23)$$

where $\tilde{\ell}$ is the log-likelihood with ascertainment correction, and ℓ is the crude log-likelihood. Define \mathbf{x}_{j_p} the covariate matrix for proband in family j . Note we can still apply Laplace

transform here, such that

$$A_j(\boldsymbol{\theta}) = 1 - S_{j_p}(a_{j_p}|\mathbf{x}_{j_p}) \quad (24)$$

$$= 1 - \int_0^\infty S_{j_p}(a_{j_p}|\mathbf{x}_{j_p}, z_j) f(z_j) dz_j \quad (25)$$

$$= 1 - \int_0^\infty \exp(-z_j \cdot H_{j_p}(a_{j_p}|\mathbf{x}_{j_p})) f(z_j) dz_j \quad (26)$$

$$= 1 - \left(1 + \frac{H_{j_p}(a_{j_p}|\mathbf{x}_{j_p})}{k}\right)^{-k} \quad (27)$$

25 5. Log-Normal Frailty

The log-normal frailty is not the power-variance-function (PVF) family, so there is no closed form for Laplace transform or expressions for survivors. But we are able to estimate the Laplace transform using Gauss Hermite Quadrature. We typically standardize the log-normal frailty Z as

$$E(\log Z) = 0 \quad (28)$$

$$\text{Var}(\log Z) = \sigma^2 \quad (29)$$

26 That is, $z \sim \text{log-Normal}(0, \sigma^2)$. The probability density function $f(z)$ is then

$$f(z) = \frac{1}{\sqrt{2\pi}\sigma} z^{-1} \exp\left(-\frac{\log(z)^2}{2\sigma^2}\right) \quad (30)$$

27 The Laplace transform is then

$$\phi(s) = \mathcal{L}(f_Z)(s) = \int_0^\infty \exp(-sz) \cdot f(z) dz \quad (31)$$

Using variable transformation, let $y = \frac{\log(z)}{\sqrt{2}\sigma}$, then $z = \exp(\sqrt{2}\sigma y)$, and $dz = \sqrt{2}\sigma \exp(\sqrt{2}\sigma y) dy$. Therefore, for d -th derivative:

$$\phi(s)^d = \int_{-\infty}^\infty z^d \exp(-sz) \cdot \frac{1}{\exp(\sqrt{2}\sigma y) \sigma \sqrt{2\pi}} \cdot \exp(-y^2) \cdot \sqrt{2}\sigma \exp(\sqrt{2}\sigma y) dy \quad (32)$$

$$= \int_{-\infty}^\infty \exp(\sqrt{2}\sigma y)^d \exp(-s \exp(\sqrt{2}\sigma y)) \cdot \frac{1}{\sqrt{\pi}} \exp(-y^2) dy \quad (33)$$

28 **Definition 1** (Gauss-Hermite Quadrature). *The integrand part can be solved using Gauss-*
29 *Hermite Quadrature. In numerical analysis, the method can be applied in the following form:*

$$\int_{-\infty}^\infty \exp(-x^2) f(x) dx \approx \sum_{i=1}^n \omega_i f(x_i) \quad (34)$$

30 where n is number of sample points used, and x_i is the roots of Hermite polynomial $H_n(x)$

such that $i = 1, \dots, n$, and the weights ω_i is

$$\omega_i = \frac{2^{n-1} n! \sqrt{n}}{n^2 [H_{n-1}(x_i)]^2} \quad (35)$$

Applying Definition 1, the integral of the Laplace transform is then

$$\phi(s)^d = \frac{1}{\sqrt{\pi}} \sum_{q=1}^{N_q} \omega_q \exp(-s \exp(\sqrt{2}\sigma y_q)) \exp(\sqrt{2}\sigma y_q)^d \quad (36)$$

where q denotes the q -th element of Gauss Hermite Quadrature, i.e., ω_q denotes the q -th weight, y_q denotes the q -th node, and N_q denotes the total number of quadratures. Thus, substituting into the log-likelihood:

$$\ell_j(\boldsymbol{\theta}) = \sum_{i=1}^{n_j} \delta_{ij} \log(h(t_{ij}|\mathbf{x}_{ij})) + \log \left(\frac{1}{\sqrt{\pi}} \sum_{q=1}^{N_q} \left[\omega_q \exp(\sqrt{2}\sigma y_q)^{d_j} \exp \left(- \sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij}) \exp(\sqrt{2}\sigma y_q) \right) \right] \right) \quad (37)$$

Similarly, the ascertainment correction in the log-normal frailty can be written as

$$A_j(\boldsymbol{\theta}) = 1 - \int_{-\infty}^{\infty} \exp(-z H(a_{j_p}|\mathbf{x}_{j_p})) f(z) dz \quad (38)$$

$$= 1 - \sum_{q=1}^{N_q} \omega_q \exp \left(- \left(\sum_{i=1}^{n_j} H(a_{j_p}|\mathbf{x}_{j_p}) \exp(\sqrt{2}\sigma y_{q_p}) \right) \right) \quad (39)$$

6. Log-Likelihood with Missing Data

6.1. Reviews on Missing Data

In this subsection, the notations are **distinct** to all other sections or subsections. The missing data problem was firstly brought by Rubin [1], and further targetted as a major statistical problem which many methodologists have developed different statistical tools to handle the missing data. Such as the practical book written by Rubin [2], and some comprehensive reviews on current missing data problems by Baraldi and Enders [3]. The missing data mechanism was introduced by Little and Rubin [4]. There are three missing data mechanisms, which are Missing Completely at Random (MCAR), Missing at Random (MAR), and Missing Not at Random (MNAR). There are some reviews on the missing data which rigorously present the statistical concept of three types of the missing mechanism [5].

Definition 2. (MCAR) Denote Y as the complete data matrix, and M as the missing data indicator matrix. Define y_{ij} and m_{ij} as i -th row (observation) and j -th column (variable) for the matrix Y and M . The conditional distribution of the missingness is said to be

$$f(m_i|y_i, \phi) = f(m_i|\phi) \quad (40)$$

That is, for the parameters of this distribution, m_i does not depend on any observed or missing data.

Example 1. (MCAR Example) There is a blind box with 500 indexed balls (No. 1 to 500) and their weights are unknown. We randomly draw 100 balls and measure their weights and record them in the Excel file. The Excel file contains two columns called Index and Weight, only those randomly selected balls will have Weights being filled. Those weights of unselected balls are called MCAR.

Definition 3. (MAR) Denote $y_{i,obs}$ as the observed y , and $y_{i,mis}$ as the missing y . Note that $y_i = (y_{i,obs}, y_{i,mis})$. The missing component is defined to be MAR if m only depends on $y_{i,obs}$. That is,

$$f(m_i|y_i, \phi) = f(m_i|y_{i,obs}, \phi) \quad (41)$$

Example 2. (MAR Example) In a psychological study, participants are asked to complete a survey so the scientist can profile their personalities. One question that asks participants to report their Mood status being good or bad. Male participants are typically too shy to answer this question, which yields some responses being missing. This is called the MAR, that the missingness on Mood status depends on the participant's gender, but not on the missing Mood itself.

Definition 4. (MNAR) In the MNAR, the missingness depends on the missing data itself, which is

$$f(m_i|y_i, \phi) = f(m_i|y_{i,mis}, y_{i,obs}, \phi) \quad (42)$$

In this case, the analysis needs to be conducted with caution. The missingness should be included in the likelihood construction.

Example 3. (MNAR Example) There is a study on participants' incomes. Person A makes \$200,000 per year, so they decide to report this amount without hesitancies. Person B makes \$10,000 per year, so they are not willing to provide this information, which this response is left as blank. This type of missing depends on the missing data itself, that Person B refuses to provide the response due to the response being comparatively low.

6.2. Log-Likelihood and Missing at Random

6.3. Log-Likelihood and Missing Not at Random

Recall the completed parametric survival log-likelihood with ascertainment correction:

$$\ell_C(\boldsymbol{\theta}) = \sum_{j=1}^J \sum_{i=1}^{n_j} \delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}, z_j) - H(t_{ij}|\mathbf{x}_{ij}, z_j) \quad (43)$$

$$- \sum_{j=1}^J \log(1 - S_{j_p}(a_{j_p}|\mathbf{x}_{j_p}, z_j)) \quad (44)$$

$$= \sum_{j=1}^J \sum_{i=1}^{n_j} \delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}) z_j - H(t_{ij}|\mathbf{x}_{ij}) z_j \quad (45)$$

$$- \sum_{j=1}^{n_j} \log(1 - \exp(z_j H_{j_p}(a_{j_p}|\mathbf{x}_{j_p}))) \quad (46)$$

When accounting for the frailty and missing data, the E-step is:

$$E(\ell_C(\boldsymbol{\theta})|\boldsymbol{\theta}^{(r)}) = \sum_{j=1}^J \sum_{i=1}^{n_j} \int_{\mathbf{x}_{mis}, z_j} \left(\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}, z_j) - H(t_{ij}|\mathbf{x}_{ij}, z_j) \right) \quad (47)$$

$$\times f(\mathbf{x}_{ij,mis}, z_j | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)}) d\mathbf{x}_{ij,mis} dz_j \quad (48)$$

$$- \sum_{j=1}^J \int_{\mathbf{x}_{mis}, z_j} \log(1 - \exp(z_j H_{j_p}(a_{j_p}|\mathbf{x}_{j_p}))) \quad (49)$$

$$\times f(\mathbf{x}_{ij,mis}, z_j | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)}) d\mathbf{x}_{ij,mis} dz_j \quad (50)$$

such that we need to integrate out the joint density of the frailty term and the missing covariate from $f(\mathbf{x}_{ij,mis}, z_j | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)})$. There are selections of the frailty density, such as Gamma distribution, log-normal distribution, and etc. In general, let's write $f(z_j|\nu)$ for the frailty distribution may be chosen with some parameters ν . Proposed by Herring et al. [6], the joint distribution of the frailty and the missing data can be adapted as:

$$f(z_j, \mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)}) = f(\mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, z_j, \boldsymbol{\theta}^{(r)}) f(z_j | \boldsymbol{\theta}^{(r)}) \quad (51)$$

while this has a form of efficient sampling, so we can write

$$f(z_j, \mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)}) \propto f(t_{ij}, \delta_{ij} | \mathbf{x}_{mis,ij}, \mathbf{x}_{obs,ij}, z_j, a_{j_p}, \boldsymbol{\beta}^{(r)}) \quad (52)$$

$$\times f(\mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, \psi^{(r)}) f(z_j | \nu^{(r)}) \quad (53)$$

Clearly, we know $f(t_{ij}, \delta_{ij} | \mathbf{x}_{mis,ij}, \mathbf{x}_{obs,ij}, z_j, \boldsymbol{\beta}^{(r)})$ is the likelihood of one single observation j in family j , also we know the distribution of $f(\mathbf{x}_{ij}|\psi)$, as well as the frailty distribution $f(z_j|\nu)$. Therefore, in our case, we can write

$$f(z_j, \mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, \boldsymbol{\theta}^{(r)}) \propto f(z_j | \nu^{(r)}) \left[\prod_{i=1}^{n_j} f(\mathbf{x}_{mis,ij} | \mathbf{x}_{obs,ij}, \psi^{(r)}) \right] \quad (54)$$

$$\times h^{(r)}(t_{ij} | \mathbf{x}_{ij}, z_j)^{\delta_{ij}} \exp(-H^{(r)}(t_{ij} | \mathbf{x}_{ij}, z_j)) \quad (55)$$

In general, without the specification of the frailty distribution, the E-step in MCEM can be written as

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) = \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \sum_{i=1}^{n_j} \left(\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}^{(m)}, z_j^{(m)}) - H(t_{ij}|\mathbf{x}_{ij}^{(m)}, z_j^{(m)}) \right) \quad (56)$$

$$+ \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \log(1 - \exp(z_j H_{j_p}(a_{j_p}|\mathbf{x}_{j_p}))) \quad (57)$$

$$+ \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \sum_{i=1}^{n_j} \log f(\mathbf{x}_{ij,mis}^{(m)}, z_j^{(m)}|\mathbf{x}_{obs,ij}, \boldsymbol{\theta}) \quad (58)$$

$$= \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \sum_{i=1}^{n_j} \left(\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}^{(m)}, z_j^{(m)}) - H(t_{ij}|\mathbf{x}_{ij}^{(m)}, z_j^{(m)}) \right) \quad (59)$$

$$+ \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \log(1 - \exp(z_j H_{j_p}(a_{j_p}|\mathbf{x}_{j_p}))) \quad (60)$$

$$+ \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \sum_{i=1}^{n_j} \log f(\mathbf{x}_{mis,ij}^{(m)}|\mathbf{x}_{obs,ij}, \psi) + \sum_{j=1}^J \frac{1}{M_j} \sum_{m=1}^{M_j} \sum_{i=1}^{n_j} \log f(z_j^{(m)}|v) \quad (61)$$

where for example, when $x_{mis,ij,1}$ is missing and $x_{obs,ij,2}$ is observed, we take

$$f(x_{mis,ij,1}^{(m)}|x_{obs,ij,2}, \psi) = \frac{1}{\tilde{\psi}\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x_{mis,ij,1}^{(m)} - \mu}{\tilde{\psi}}\right)^2\right) \quad (62)$$

such that $\hat{\mu} = \psi_0 + \psi_1 x_{obs,ij,2}$, and $\tilde{\psi}$ is the standard error of the conditional normal distribution. We can apply the kinship matrix for a family study, that we will deal with a multivariate normal distributed $\mathbf{x}_{mis,j,1}^{(m)}$ for family j .

$$f(\mathbf{x}_{mis,j,1}^{(m)}|\mathbf{x}_{obs,j,2}, \psi, K) = (2\pi)^{(-\frac{n_j}{2})} \det(\Sigma)^{(-\frac{1}{2})} \exp\left(-\frac{1}{2}(\mathbf{x}_{mis,j,1}^{(m)} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}_{mis,j,1}^{(m)} - \boldsymbol{\mu})\right) \quad (63)$$

where $\Sigma = \tilde{\psi}_g^2 K + \tilde{\psi}_e^2$, and K is the kinship matrix with diagonals of 1 and is known. Also, $\hat{\boldsymbol{\mu}} = \boldsymbol{\psi}_0 + \boldsymbol{\psi}_1 \mathbf{x}_{obs,j,2}$. Therefore, we may sample from $\mathbf{x}_{mis,j,1} \sim MVN(\boldsymbol{\mu}, \tilde{\psi}_g^2 K + \tilde{\psi}_e^2)$. Moreover, we can give the specific distribution of the frailty in the E-step by sampling the frailty term z_j , for example, when z_j is assumed to be Gamma distribution, then

$$f(z_j^{(m)}|k) = \frac{k^k}{\Gamma(k)} (z_j^{(m)})^{k-1} \exp(-k z_j^{(m)}) \quad (64)$$

84 When z_j is log-normal distributed with the mean 0, we have

$$f(z_j^{(m)}|v) = \frac{1}{z_j^{(m)}v\sqrt{2\pi}} \exp\left(-\frac{(\log z_j^{(m)})^2}{2v^2}\right) \quad (65)$$

85 Thus, the $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)})$ function will yield a closed form that can be optimized via M-step, with
 86 the iterations that update all parameters using Gibb's sampler.

87 7. Missing PRS and MCEM

88 7.1. Not Considering the Family Correlations

We propose an MCEM framework in terms of estimating the distribution of the PRS. The PRS was calculated to infer the relationship between a phenotype and multiple genetic loci, while these information were not gained if one was not involved in the original GWAS. Thus, we propose to sample the PRS using the information that we have already obtained through the study. Denote $\mathbf{x}_{j,1}$ as the PRS scores vector in family j , and $\mathbf{x}_{j,2}$ the mutation status vector in family j . Take \mathbf{p}_j as the proband indicator vector in family j , \mathbf{c}_j is the current age for patients in family j . So we now have a design matrix when modelling the missing PRS, call it $\mathbf{W} = (\log(\mathbf{t}_j), \boldsymbol{\delta}_j, \log(\mathbf{t}_j) \odot \boldsymbol{\delta}_j, \mathbf{p}_j, \mathbf{c}_j, \mathbf{x}_{j,2})$. We can make the assumption on the conditional distribution of the PRS, take $\mathbf{X}_{j,1}|\mathbf{W} \sim MVN(\mathbf{W}\boldsymbol{\psi} + \mathbf{u}, \sigma^2\mathbf{I})$. We are interested in modelling the PRS while accounting for the between family variance, so $\mathbf{u} \sim MVN(0, \sigma_u^2\mathbf{I})$. Thus, the E-step for Gamma frailty model with ascertainment correction is then

$$E_{\mathbf{x}_{j,1,mis}}(\ell(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)})) = \sum_{j=1}^J \left[\sum_{i=1}^{n_j} \int_{\mathbf{x}_{j,1,mis}} (\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij})) + \right. \quad (66)$$

$$+ \log\left(\frac{(k+d_j-1)!}{k!k^{d_j-1}} \left(1 + \frac{\sum_{i=1}^{n_j} (H(t_{ij}|\mathbf{x}_{ij}))}{k}\right)^{-k-d_j}\right) - \quad (67)$$

$$\left. - \log(A_j(\boldsymbol{\theta})) + \log f(x_{ij,1,mis}|w_{ij}, \boldsymbol{\psi}) dx_{ij,1,mis} \right] \quad (68)$$

Taking a sample of size M when we sample $f(x_{ij,1,mis}|w_{ij}, \boldsymbol{\psi})$ for each subject i in family j , $(x_{ij,i,mis}^{(1)}, \dots, x_{ij,i,mis}^{(M)})$. This leads to an E-step:

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) = \sum_{j=1}^J \left[\sum_{i=1}^{n_j} \sum_{m=1}^M (\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}^{(m)})) + \right. \quad (69)$$

$$+ \log\left(\frac{(k+d_j-1)!}{k!k^{d_j-1}} \left(1 + \frac{\sum_{i=1}^{n_j} (H(t_{ij}|\mathbf{x}_{ij}^{(m)}))}{k}\right)^{-k-d_j}\right) - \quad (70)$$

$$\left. - \log(A_j(\boldsymbol{\theta})) + \log f(x_{ij,1,mis}^{(m)}|w_{ij}, \boldsymbol{\psi}) \right] \quad (71)$$

Similarly, the expectation with respect to the missing PRS in log-normal frailty can be written as

$$E_{\mathbf{X}_{j,1,mis}}(\ell(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)})) = \sum_{j=1}^J \left[\sum_{i=1}^{n_j} \int_{\mathbf{X}_{j,1,mis}} (\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij})) + \right. \quad (72)$$

$$+ \log \left(\frac{1}{\sqrt{\pi}} \sum_{q=1}^{N_q} \left[\omega_q \exp(\sqrt{2}\sigma y_q)^{d_j} \exp\left(-\sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij}) \exp(\sqrt{2}\sigma y_q)\right) \right] \right) - \quad (73)$$

$$\left. - \log(A_j(\boldsymbol{\theta})) + \log f(x_{ij,1,mis}|w_{ij}, \boldsymbol{\psi}) dx_{ij,1,mis} \right] \quad (74)$$

and the E-step:

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) = \sum_{j=1}^J \left[\sum_{i=1}^{n_j} \sum_{m=1}^M (\delta_{ij} \log h(t_{ij}|\mathbf{x}_{ij}^{(m)})) + \right. \quad (75)$$

$$+ \log \left(\frac{1}{\sqrt{\pi}} \sum_{q=1}^{N_q} \left[\omega_q \exp(\sqrt{2}\sigma y_q)^{d_j} \exp\left(-\sum_{i=1}^{n_j} H(t_{ij}|\mathbf{x}_{ij}^{(m)}) \exp(\sqrt{2}\sigma y_q)\right) \right] \right) - \quad (76)$$

$$\left. - \log(A_j(\boldsymbol{\theta})) + \log f(x_{ij,1,mis}^{(m)}|w_{ij}, \boldsymbol{\psi}) \right] \quad (77)$$

7.2. Considering the Family Correlations

Given that family j has some subjects containing the missing PRS due to the sampling cost (maybe, need to confirm the previous paper), that not all subjects are being sampled for the PRS calculation. Since subjects within one family are correlated in some genetic associations, we intend to sample the missing PRS using a multivariate normal distribution while accounting for the kinship matrix. Using the same setting for the modelling of the PRS, the only difference is the variance-covariance matrix. Here, $\mathbf{X}_{j,1}|\mathbf{W}, \boldsymbol{\psi} \sim MVN(\mathbf{W}\boldsymbol{\psi} + \mathbf{u}, \sigma^2\mathbf{I})$. But, $\mathbf{u} \sim MVN(0, 2\sigma_u^2\mathbf{K})$ where \mathbf{K} is a $\mathbb{R}^{n \times n}$ kinship matrix with a diagonal of 0.5. The variance can be computed via $\text{Var}(X_{j,1}) = 2\sigma_u^2\mathbf{K} + \sigma^2\mathbf{I}$ when both the genetic variation and the residual errors contribute to the variance-covariance structure. Therefore, the marginal distribution of the PRS after integrating out the random effects will be $\mathbf{X}_{j,1} \sim MVN(\mathbf{W}\boldsymbol{\psi}, 2\sigma_u^2\mathbf{K} + \sigma^2\mathbf{I})$, and will be the distribution used for the Monte Carlo sampling.

M Step, Nelder-Mead, why? Gradient free? Variance?

8. Multiple Imputation with Monte Carlo Sampling

8.1. Missing Mechanisms

...A discussion on MCAR, MAR, and MNAR...

8.2. Multiple Imputations *(Be sure to report a conceptual statistical definition in the thesis)*

In the general framework, denote $\mathbf{X}_{j,1,mis}$ as the element of $\mathbf{X}_{j,1}$ when it's missing. Denote $\mathbf{X}_{j,1,obs}$ on the observed element. Adapting the idea from the MCEM, we can draw $\hat{\mathbf{X}}_{j,1,mis}$ from the same assumption of the distribution discussed in the section 7 for M times. But, the missing indicator \mathbf{R} should be addressed to attest the missing mechanism. The conditional distribution can be written as $f(\mathbf{X}_{j,1,mis}|\mathbf{X}_{j,1,obs}, \mathbf{W}, \mathbf{R}_j)$. Under the assumption of the MAR *(Confirm the paper)*,

$$f(\mathbf{X}_{j,1,mis}|\mathbf{X}_{j,1,obs}, \mathbf{W}, \mathbf{R}_j) = f(\mathbf{X}_{j,1,mis}|\mathbf{X}_{j,1,obs}, \mathbf{W}) \quad (78)$$

This results a series of M complete datasets containing all values observed with the fill of the Monte Carlo samples. Then with these M datasets, we run the analysis individually. Thus, the estimate of $\boldsymbol{\theta}$ can be calculated via

$$\hat{\boldsymbol{\theta}}_{MI} = \frac{1}{M} \sum_{m=1}^M \hat{\boldsymbol{\theta}}_m \quad (79)$$

Denote V the variance for m -th complete data inference variance, we first calculate the average variance of these analyses to be the within-imputation variance,

$$\hat{\mathbf{W}} = \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{V}}_m \quad (80)$$

Under the unbiased variance estimator, the between-imputation variance is

$$\hat{\mathbf{B}} = \frac{1}{M-1} \sum_{m=1}^M (\hat{\boldsymbol{\theta}}_m - \hat{\boldsymbol{\theta}}_{MI})(\hat{\boldsymbol{\theta}}_m - \hat{\boldsymbol{\theta}}_{MI})^\top \quad (81)$$

Combining these *(How? Are there any statistical/mathematical proof of this)*, we obtain the estimate of the variance of $\hat{\boldsymbol{\theta}}_{MI}$,

$$\hat{\mathbf{V}}_{MI} = \hat{\mathbf{W}} + (1 + \frac{1}{M})\hat{\mathbf{B}} \quad (82)$$

Discussions on why MI over MCEM or vise versa after the analysis...

9. Monte Carlo EM *(This section needs to be re-written to a definitional MCEM framework, maybe consult papers)*

The complete data log-likelihood for family j is $\ell_j(\boldsymbol{\theta}; h_{ij})$ where $\boldsymbol{\theta}$ consists all baseline parameters, and model coefficients β 's, as well as the frailty parameter ϕ . The E-step for complete data is:

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(r)}) = \int \ell(\boldsymbol{\theta}; h_{ij}) \cdot f(x_{mis,i}|x_{obs,i}, z, \boldsymbol{\theta}^{(r)}, t_{ij}, \delta_{ij}, p_j) dx_{mis,ij} \quad (83)$$

We sample the size m_i for each i -th observation, $x_{i1}^*, \dots, x_{im_i}^*$ from the distribution $f(x_{mis,ij}|\cdot)$, and take $M = 1, \dots, m_i$, such that each X_{iM}^* depends on the iteration number for $r + 1$ iterations. In general:

$$\hat{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(r)}) = \frac{1}{m_i} \sum_{M=1}^{m_i} \ell(x_{iM}^*, x_{obs,ij}, t_{ij}, \boldsymbol{\theta}, z_j) \quad (84)$$

More specifically,

1. We first initialize $m, \theta^{(0)}$, and start the burn-in.
2. Also, we set importance weights $w_t = 1$ for all $t = 1, \dots, m$.
3. At the burn-in iteration s , we generate $x_{miss,1}, \dots, x_{miss,m} \sim N(\mu_X | X_{obs}, \theta^{(s)}, z)$ using MCMC sample.
4. In the E-step, we estimate $Q(\theta | \theta^{(s)})$ by using the importance weights:

$$Q_m(\theta | \hat{\theta}^{(s)}) = \frac{\sum_{t=1}^m w_t \log f(X_{obs}, X_{miss,t} | \theta)}{\sum_{t=1}^m w_t} \quad (85)$$

5. Note the numerator is actually a weighted log-likelihood. In the M-step, we maximize $Q_m(\theta | \hat{\theta}^{(s)})$ to obtain $\hat{\theta}^{(s+1)}$.
6. Repeat (3.) - (5.) for s burn-in iterations.
7. Then re-initialize $\hat{\theta}^{(0)} = \hat{\theta}^{(s)}$
8. We generate $x_{miss,1}, \dots, x_{miss,m} \sim N(\mu_X | X_{obs}, \hat{\theta}^{(0)}, z)$ using MCMC sampler. At iteration $r + 1$
9. Compute the importance weights from the ratio of likelihood

$$w_t = \frac{L(\hat{\theta}^{(r)} | X_{miss,t}, X_{obs})}{L(\hat{\theta}^{(0)} | X_{miss,t}, X_{obs})} \quad (86)$$

10. Thus, the E-step can be written as

$$Q_m(\theta | \hat{\theta}^{(r)}) = \frac{\sum_{t=1}^m w_t \log f(X_{miss,t}, X_{obs} | \theta)}{\sum_{t=1}^m w_t} \quad (87)$$

11. Then M-step: we maximize $Q_m(\theta | \hat{\theta}^{(r)})$ to obtain $\hat{\theta}^{(r+1)}$.

This automated MCEM firstly optimizes the importance weights at burn-ins, then performs the actual EM to find $\hat{\theta}$. This importance weight ensures the imputation step of the missing data actually yields to the real distribution.

10. Correlated Frailty using Kinship Matrix

Family members are correlated within one family, that we denote K as the kinship correlation matrix among all observations. This matrix ensures those individuals not from the same family automatically have a correlation of 0. The likelihood construction needs multivariate form. For $\mathbf{Z} \sim \text{MVN}(0, \sigma^2 K)$, that K has the diagonal of 1. The likelihood is

$$L(\cdot) = \int_{\mathbb{R}^n} \prod_{i=1}^n (h(t|\mathbf{x}_i, \mathbf{z}_i))^{\delta_i} \exp(-H(t|\mathbf{x}_i, \mathbf{z}_i)) f(\mathbf{z}) d\mathbf{z} \quad (88)$$

$$= \int_{\mathbb{R}^n} \prod_{i=1}^n (h(t|\mathbf{x}_i))^{\delta_i} \exp(\mathbf{z}_i)^{\delta_i} \exp(-H(t|\mathbf{x}_i) \exp(\mathbf{z}_i)) f(\mathbf{z}) d\mathbf{z} \quad (89)$$

$$= \prod_{i=1}^n (h(t|\mathbf{x}_i))^{\delta_i} \int_{\mathbb{R}^n} \exp(\delta_i \mathbf{z}_i - H(t|\mathbf{x}_i) \exp(\mathbf{z}_i)) f(\mathbf{z}) d\mathbf{z} \quad (90)$$

146 Applying the Laplace approximation, and taking the log for the likelihood, we obtain

$$\ell(\cdot) = \sum_{i=1}^n \left[\delta_i \log h(t|\mathbf{x}_i) \right] + \sum_{i=1}^n \left[\delta_i \hat{\mathbf{z}} - H(t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}}) \right] - \frac{1}{2} \hat{\mathbf{z}}^\top \Sigma^{-1} \hat{\mathbf{z}} \quad (91)$$

147 such that $\Sigma = \sigma^2 K$. Also, we treat the random effect \mathbf{z} as a vector of parameters, and use
 148 outer-loop to search for the σ , and use inner-loop to search for other parameters (baseline
 149 parameters, and β) including \mathbf{z} . The process can be achieved via Newton-Raphson algorithm.
 150 For computational efficiency, we can set $\Sigma^{-1} = L^\top L$ through Cholesky Decomposition. In
 151 this way, $\mathbf{z}L \sim MVN(0, \sigma^2 I)$. In order to apply NR-algorithm, the gradient and the hessian
 152 are required. The gradient for parameters is:

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^n \delta_i \mathbf{x}_i + \sum_{i=1}^n -H(t_i|\mathbf{x}_i) \mathbf{x}_i \exp(\mathbf{z}) \quad (92)$$

$$\frac{\partial \ell}{\partial \mathbf{z}} = \sum_{i=1}^n \delta_i - (t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}}) - \Sigma^{-1} \hat{\mathbf{z}} \quad (93)$$

$$\frac{\partial \ell}{\partial \alpha} = \sum_{i=1}^n \frac{\delta_i}{\alpha} + \sum_{i=1}^n -\frac{H(t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}})}{\alpha} \quad (94)$$

$$\frac{\partial \ell}{\partial \lambda} = \sum_{i=1}^n \delta_i \left(\frac{1}{\lambda} + \log(t_i) \right) + \sum_{i=1}^n -H(t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}}) \log(t_i) \quad (95)$$

153 The hessian matrix element, i.e., second partial derivative is

$$\frac{\partial^2 \ell}{\partial \beta^\top \partial \beta} = \sum_{i=1}^n -H(t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}}) x_{ij} x_{ik} \quad (96)$$

$$\frac{\partial^2 \ell}{\partial \mathbf{z}^\top \partial \mathbf{z}} = \sum_{i=1}^n -H(t_i|\mathbf{x}_i) \exp(\hat{\mathbf{z}}) - \Sigma^{-1} \quad (97)$$

$$\frac{\partial^2 \ell}{\partial \alpha^2} = \sum_{i=1}^n -\frac{\delta_i}{\alpha^2} \quad (98)$$

$$\frac{\partial^2 \ell}{\partial \lambda^2} = \sum_{i=1}^n -\frac{\delta_i}{\lambda^2} - H(t_i | \mathbf{x}_i) \exp(\hat{\mathbf{z}}) \log(t_i)^2 \quad (99)$$

10.1. Proof of $\Sigma = LL^\top$

Every symmetric positive definite matrix Σ can be decomposed into $\Sigma = LL^\top$, where L is a lower triangular matrix with real and positive diagonal entries.

Proof. Set-ups:

1. Covariance matrix Σ is by definition symmetric and positive definite, e.g.

$$\Sigma = \begin{pmatrix} \sigma_{X_1}^2 & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \sigma_{X_2}^2 \end{pmatrix} \quad (100)$$

such that $\mathbf{X}\Sigma\mathbf{X}^\top > 0$ always, and this matrix is symmetric.

2. Suppose \mathbf{X} has n observations, then Σ is $n \times n$, the first element is $\sigma_{11} > 0$ by definition (For simplicity, we use σ_{11} rather than it's square to denote the variance). Define $l_{11} = \sqrt{\sigma_{11}}$, to be the first element of L . For the first column of L , let $l_{j1} = \frac{\sigma_{j1}}{l_{11}}$ for $j = 2, \dots$

Induction step: Assume we have first $k-1$ columns of L , consider k -th column

- For the diagonal element $l_{kk} = \sqrt{\sigma_{kk} - \sum_{j=1}^{k-1} l_{kj}^2}$
- For off-diagonals,

$$l_{ik} = \frac{\sigma_{ik} - \sum_{j=1}^{k-1} l_{ij}l_{kj}}{l_{kk}} \quad (101)$$

for $i = k+1, \dots, n$.

with the repetition for each column $k = 2, \dots, n$, the top-left $k \times k$ submatrix of LL^\top matches that of Σ . For example, when $k = 3$,

$$\Sigma = \begin{pmatrix} \sigma_{11} & & \\ & \sigma_{22} & \\ & & \sigma_{33} \end{pmatrix} \quad (102)$$

and

$$L = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \quad (103)$$

then

$$LL^\top = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & l_{31} \\ 0 & l_{22} & l_{32} \\ 0 & 0 & l_{33} \end{pmatrix} = \begin{pmatrix} l_{11}^2 & l_{11}l_{21} & l_{11}l_{31} \\ l_{21}l_{11} & l_{21}^2 + l_{22}^2 & l_{21}l_{31} + l_{22}l_{32} \\ l_{31}l_{11} & l_{31}l_{21} + l_{32}l_{22} & l_{31}^2 + l_{32}^2 + l_{33}^2 \end{pmatrix} \quad (104)$$

172 Take

$$\Sigma = \begin{pmatrix} 4 & 2 & 2 \\ 2 & 3 & 1 \\ 2 & 1 & 3 \end{pmatrix} \quad (105)$$

173 Then by definition of Cholesky Decomposition, we can calculate $l_{11}^2 = \sigma_{11} \implies l_{11} = \sqrt{4} = 2$,
 174 and $l_{21} = \frac{\sigma_{21}}{l_{11}} = 2/2 = 1$, and $l_{31} = 1$. Similarly for l_{22}, l_{32}, l_{33} . Therefore,

$$L = \begin{pmatrix} 2 & 0 & 0 \\ 1 & \sqrt{2} & 0 \\ 1 & 0 & \sqrt{2} \end{pmatrix} \quad (106)$$

175 which implies

$$LL^\top = \begin{pmatrix} 2 & 0 & 0 \\ 1 & \sqrt{2} & 0 \\ 1 & 0 & \sqrt{2} \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 \\ 0 & \sqrt{2} & 0 \\ 0 & 0 & \sqrt{2} \end{pmatrix} = \begin{pmatrix} 4 & 2 & 2 \\ 2 & 3 & 1 \\ 2 & 1 & 3 \end{pmatrix} = \Sigma \quad (107)$$

176

□

177 Essentially, the Cholesky Decomposition transforms the multivariate normal to a stan-
 178 dard multivariate normal. When $\mathbf{Z} \sim \mathcal{N}(0, \Sigma)$, let $\Sigma = \mathbf{L}\mathbf{L}^\top$, then $\mathbf{Y} = \mathbf{L}^{-1}\mathbf{Z} \sim \mathcal{N}(0, \mathbf{I})$
 179 that \mathbf{I} is the identity matrix, since $\mathbf{L}^{-1}\Sigma(\mathbf{L}^{-1})^\top = \mathbf{L}^{-1}\mathbf{L}\mathbf{L}^\top(\mathbf{L}^{-1})^\top = \mathbf{I}$. This will simplify
 180 the computational process.

References

- [1] Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- [2] Donald B Rubin. *Multiple imputation for nonresponse in surveys*. John Wiley & Sons, New York, 1987. doi: 10.1002/9780470316696.
- [3] Amanda N Baraldi and Craig K Enders. An introduction to modern missing data analyses. *Journal of school psychology*, 48(1):5–37, 2010.
- [4] Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John Wiley & Sons, 2019.
- [5] Miriam Seoane Santos, Ricardo Cardoso Pereira, Adriana Fonseca Costa, Jastin Pompeu Soares, João Santos, and Pedro Henriques Abreu. Generating synthetic missing data: A review by missing mechanism. *IEEE Access*, 7:11651–11667, 2019.
- [6] Amy H Herring, Joseph G Ibrahim, and Stuart R Lipsitz. Frailty models with missing covariates. *Biometrics*, 58(1):98–109, 2002.