

0.1 Joint model over frequency estimates and regression with a subdivision specific child effect

Below summarises my first attempt at a joint model over both frequency estimation and regression that takes into account structure at the level of the child within and across subdivisions of the data. It is basically the joint model that the extended frequency model plus would-be regression step approximates, but with dependence between haplotypes. That is, each child has a subdivision specific child effect. Ie, it doesn't capture the fact that a child with relatively low immunity in 2008 is likely to have relatively low immunity in 2009.

Due to the dependence between both $\pi_{1r}, \dots, \pi_{Kr}$ and $\pi_{k1}, \dots, \pi_{kR}$, instead of modelling π_{child} as a realisation from a distribution centred about $\pi_{\text{population}} \in \mathbb{S}_R$, π_{child} ought to be modelled jointly for $k = 1, \dots, K$ as a realisation, $\pi_{\text{child}} \in (0, 1)^{K \times R}$, from a matrix variate distribution centred about $\pi_{\text{population}} \in (0, 1)^{K \times R}$, where both π_{child} and $\pi_{\text{population}}$ are a $K \times R$ matrices of dependent haplotype frequencies. I therefore need a distribution over matrices in which the columns sum to one (each row is a realisation from a dirichlet), and the rows are correlated. Unfortunately, I can't model $\pi_{\text{child } k} \in (0, 1)^R \sim \text{Dirichlet}(\pi_{\text{child}} \in (0, 1)^R)$ and $\pi_{\text{child}} \in (0, 1)^R \sim \text{Dirichlet}(\pi_{\text{population}} \in (0, 1)^R)$, since $\pi_{\text{population}}$ needs to be a $K \times R$ in order to do regression with correlation across $K \times R$. I therefore assumed $\pi_{\text{child } k}$ were independent given $\pi_{\text{population } k}$. The correlation between $k = 1, \dots, K$ is accounted for by \mathbf{V}_k (scale matrix among over rows) in the regression. This is basically the joint model that the extended model plus regression using the approach

$$\begin{aligned} \rho(\pi_{\text{population}}, \pi_{\text{child}}, \mathbf{a}, \mathbf{m}, \mathbf{B}, \mathbf{V} \mid \mathbf{y}) &\propto \prod_{\text{child}=1}^{\#\text{children}} \left\{ \prod_{i=1}^{I_{\text{child } k}} \left\{ \prod_{j=1}^J \{\rho(\mathbf{y}_{ijk \text{ child}} \mid \mathbf{a}_{ik \text{ child}})\} \right. \right. \\ &\quad \times \rho(\mathbf{a}_{ik \text{ child}} \mid m_{ik \text{ child}}, \pi_{k \text{ child}}) \rho(m_{ik \text{ child}}) \} \\ &\quad \times \prod_{k=1}^K \left\{ \rho(\pi_{k \text{ child}} \mid \pi_{k \text{ population}}) \right\} \Bigg\} \\ &\quad \times \rho(\theta_{\text{population}} \mid \mathbf{B}, \mathbf{V}_k, \mathbf{V}_r) \left| \det \left(\left\{ \frac{\partial \pi_{\text{population}}}{\partial \theta_{\text{population}}} \right\}^{-1} \right) \right| \\ &\quad \times \rho(\mathbf{B}) \rho(\mathbf{V}_k) \rho(\mathbf{V}_r) \end{aligned} \tag{1}$$

where $\rho(\mathbf{y}_{ijk \text{ child}} \mid \mathbf{a}_{ik \text{ child}})$, $\rho(\mathbf{a}_{ik \text{ child}} \mid m_{ik \text{ child}}, \pi_{k \text{ child}})$, and $\rho(m_{ik \text{ child}})$ are equivalent to previous definitions and

$$\rho(\pi_{\text{child } k} \mid \pi_{\text{population } k}) = \text{Dirichlet}(\pi_{\text{population } k}) \tag{2}$$

$$\rho(\theta_{\text{population}}) = \mathcal{MN}_{KR-1}(\mathbf{X}\mathbf{B}, \mathbf{V}_k, \mathbf{V}_r) \tag{3}$$

$$\rho(\mathbf{B}) = \mathcal{MN}_{PR-1}(\mathbf{M}, \mathbf{U}_p, \mathbf{U}_r), \text{ say, and} \tag{4}$$

$$\rho(\mathbf{V}_k) = \mathcal{W}_K^{-1}(\Psi_{\mathbf{k}}, \nu_k) \tag{5}$$

$$\rho(\mathbf{V}_r) = \mathcal{W}_{R-1}^{-1}(\Psi_{\mathbf{r}}, \nu_r) \tag{6}$$

Where the two inverse wishart prior distributions, $\mathcal{W}^{-1}(\cdot)$, each have two parameters, a scale matrix ($\Psi_{\mathbf{k}} \in \mathbb{R}_0^{+K \times K}$ and $\Psi_{\mathbf{r}} \in \mathbb{R}_0^{+R-1 \times R-1}$), that either encodes the prior knowledge about the correlation across studies ($\Psi_{\mathbf{k}}$), or the correlation

across haplotypes (Ψ_r), and a tuning parameter, (ν_k and ν_r) that specifies strength of the prior belief.

0.2 A note regarding child-level frequencies based on the Ugandan data.

Note that, when there was no child effect, π_{child_r} , were generally well aligned (there were some peculiar results - for example *pfert*-76, AL, 2012 - perhaps due to poor mixing, but I didn't look at the trace for every child in every partition). When there was a child effect, π_{child_r} variation was high. For the most part all children had the same major and minor allele. For one study, however (*pfmdr*-86, AL, 2010) had different children appeared to have different major and minor alleles. This is probably due to the random sample of 10 children used to generate the plot (which differs to the 16 children used for the m_i and a_i plots). For example, for *pfmdr*-86, AL, 2010, only 34 of 200 samples were wild type. If a given child have multiple wild type infections (which is compatible with the hypothesis of inter-individual immunity if the episodes are further apart than the AL half life within 2010), they might not experience much shrinkage, hence have different major/minor alleles to the overall mean. To test this I would need to investigate the major/minor relationship in light of the infections per child (the mutatedness, the number of episodes and the duration between them). I'm not going to do this, as the goal in this study was to concentrate on $\pi_{\text{population}}$, treating π_{child} as a nuisance parameter. In any case, there were too few data per child to report on π_{child} in a meaningful manner (many children suffered only one episode per partition).

0.3 A note regarding spurious precision under the exchangeability assumption.

It seems unlikely that spurious precision will have an effect on the frequency trends (recall that, the frequency trends based on point estimates (Figure X) are practically identical to the trends based on the full posterior densities (Figure X), despite the difference in uncertainty being 100%).

0.4 A note regarding a child effect on the MOI.

Note that, neither a population nor child estimate of the MOI ($m_{\text{population}}$ nor m_{child} , respectively) feature in the child effect model (equation (??)) because we opt to treat $m_{i \text{ child}}$ as a nuisance parameter, included to aid estimation of $\pi_{\text{population}}$. That said, human attributes such as acquired immunity are thought to impact the MOI distribution (REF). We could allow different children to have different MOI propensities (replace $\rho(m_{i \text{ child}})$ with $\rho(m_{i \text{ child}} | m_{\text{child}})$, say, and let $m_{\text{child}} \sim \text{Geometric}(\lambda^{-1}, 0.05)$, say); however, there is a scarcity of information in the data to support MOI estimates (at best, we can discern demonstrably multiclonal infections), so we opt not to include m_{child} .

0.5 Future work

0.5.1 Regression

Need to read about how to incorporate correlations due to time. Also how to implement this model. p.480(Mendeley) p.380 Bayesian Data Analysis onwards. How does this impact the choice of conjugate priors? Need to read into it. $\beta_r | \mathbf{V}_r$ can no-longer be equivalent to $\mathcal{N}_P(\mu_0, \sigma^2 \mathbf{V}_0)$ since \mathbf{V}_r is $K \times K$ and the prior is a P variate normal distribution, but can use an inverse wishart prior on \mathbf{V}_r

0.5.2 Joint model

To capture dependence between haplotypes and across subdivisions of the data at the level of the child the following model (figure ??) was proposed,

$$\begin{aligned}
\rho(\mathbf{B}, \mathbf{V}, \boldsymbol{\mu}_{\text{children}}, \boldsymbol{\theta}_{\text{children}1:K}, \mathbf{a}, \mathbf{m} \mid \mathbf{y}) \propto & \prod_{\text{child}=1}^{291} \left\{ \prod_{k=1}^K \left\{ \prod_{i=1}^{I_{\text{child}_k}} \left\{ \prod_{j=1}^J \{ \rho(\mathbf{y}_{ijk \text{ child}} \mid \mathbf{a}_{ik \text{ child}}) \} \right. \right. \right. \\
& \times \rho(\mathbf{a}_{ik \text{ child}} \mid m_{ik \text{ child}}, f(\boldsymbol{\theta}_{k \text{ child}})) \rho(m_{ik \text{ child}}) \} \\
& \times \rho(\boldsymbol{\theta}_{\text{child}_k} \mid \mathbf{B}, \mathbf{V}, \boldsymbol{\mu}_{\text{child}}) \} \\
& \times \rho(\boldsymbol{\mu}_{\text{child}}) \} \\
& \times \rho(\mathbf{B}) \rho(\mathbf{V})
\end{aligned} \tag{7}$$

- \mathbf{B} is a $P \times R - 1$ matrix of regression coefficients;
- \mathbf{V} is a $R - 1 \times R - 1$ covariance matrix;
- $\boldsymbol{\mu}_{\text{children}}$ denotes the collection of $\boldsymbol{\mu}_{\text{child}}$, where each $\boldsymbol{\mu}_{\text{child}}$ denotes a $1 \times R - 1$ random effect vector associated with a specific child, for $\text{child} = 1, \dots, 291$;
- $\boldsymbol{\theta}_{\text{children}1:K}$ denotes the collection of $\boldsymbol{\theta}_{\text{child}_k}$, where each $\boldsymbol{\theta}_{\text{child}_k}$ denotes a $1 \times R - 1$ vector of haplotype frequencies mapped onto the real line, for $\text{child} = 1, \dots, 291$ and $k = 1, \dots, K$;

and where

- $\rho(\mathbf{y}_{ijk \text{ child}} \mid \mathbf{a}_{ik \text{ child}}), \rho(\mathbf{a}_{ik \text{ child}} \mid m_{ik \text{ child}}, \boldsymbol{\pi}_{k \text{ child}}), \rho(m_{ik \text{ child}})$ are equivalent to the definitions for $\rho(\mathbf{y}_{ij} \mid \mathbf{a}_i), \rho(\mathbf{a}_i \mid m_i, \boldsymbol{\pi}), \rho(m_i)$ (equations X, X and X)
- $\rho(\boldsymbol{\theta}_{\text{child}_k} \mid \mathbf{B}, \mathbf{V}, \boldsymbol{\mu}_{\text{child}}) = \mathcal{N}_{R-1}(\mathbf{x}_k \cdot \mathbf{B} + \boldsymbol{\mu}_{\text{child}}, \mathbf{V})$
- $\rho(\mathbf{B}) = \mathcal{MN}_{PR-1}(\mathbf{M}, \mathbf{U}, (\mathbf{X}^T \mathbf{X})^{-1})$
- $\rho(\mathbf{V}) = \mathcal{W}_{R-1}^{-1}(\boldsymbol{\Psi}, v)$
- $\rho(\boldsymbol{\mu}_{\text{child}}) \sim \mathcal{N}_{R-1}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

This model assumes the child effect is fixed year upon year, or across different durations since last treatment, and that $\boldsymbol{\theta}_{\text{child } k}$ are conditionally independent given \mathbf{B}, \mathbf{V} and $\boldsymbol{\theta}_{\text{child}}$, where the non-diagonal components of \mathbf{V} capture dependence between haplotypes, not subdivisions. Hence, the only dependence between subdivisions is due $\boldsymbol{\theta}_{\text{child}}$. The key parameter of interest is \mathbf{B} (since \mathbf{B} describes the trends over time or duration since last treatment).