

Mixture Models for Single-Cell Assays with Application to Vaccine Studies: Supplementary Materials

GREG FINAK^{1*}, ANDREW MCDAVID¹, PRATIP CHATTOPADHYAY³, MARIA DOMINGUEZ³, STEVE DE ROSA^{1,2}, MARIO ROEDERER³, RAPHAEL GOTTARDO¹

¹ *Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center (FHCRC), Seattle, WA*

² *HIV Vaccine Trials Network, Fred Hutchinson Cancer Research Center (FHCRC), Seattle, WA*

³ *Vaccine Research Center, NIAID, NIH, 40 Convent Drive, Rm 5509, Bethesda, MD 20892*

Appendix A: Constrained beta-binomial model

We can define a model where we constrain the stimulated proportions under the alternative model such that $p^{(s)} > p^{(u)}$. In this case, the only changes required are for the alternative marginal likelihood L_1 defined in the main manuscript by (1). Due to the constraint, the normalizing constant of the prior under the alternative (model \mathcal{M}_1) is not given by $B(\alpha^{(u)}, \beta^{(u)})B(\alpha^{(s)}, \beta^{(s)})$ but requires computing

$$Z(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)}) = \int_0^1 p^{(u)\alpha^{(u)}-1} (1-p^{(u)})^{\beta^{(u)}-1} \int_{p^{(u)}}^1 p^{(s)\alpha^{(s)}-1} (1-p^{(s)})^{\beta^{(s)}-1} dp^{(s)} dp^{(u)}.$$

Using this expression, the constrained alternative marginal likelihood can be written as

$$L_1(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | \mathbf{y}_i) = \binom{N_i^{(u)}}{n_i^{(u)}} \binom{N_i^{(s)}}{n_i^{(s)}} \frac{Z(n_i^{(u)} + \alpha^{(u)}, N_i^{(u)} - n_i^{(u)} + \beta^{(u)}, n_i^{(s)} + \alpha^{(s)}, N_i^{(s)} - n_i^{(s)} + \beta^{(s)})}{Z(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)})}.$$

In general, there is no closed-form expression for $Z(\cdot)$, and a numerical approximation must be used. Let us denote by $\tilde{Z}(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)})$ the approximation. A natural way to estimate \tilde{Z} is to use Monte Carlo integration. Indeed, we can write

$$\tilde{Z}(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)}) = B(\alpha^{(u)}, \beta^{(u)})B(\alpha^{(s)}, \beta^{(s)}) \int_0^1 \frac{p^{(u)\alpha^{(u)}-1} (1-p^{(u)})^{\beta^{(u)}-1}}{B(\alpha^{(u)}, \beta^{(u)})} [1 - F_{\alpha^{(s)}, \beta^{(s)}}(p^{(u)})] dp^{(u)} \quad (1)$$

where $F_{\alpha^{(s)}, \beta^{(s)}}$ is the cumulative distribution function of a beta random variable with parameters $\alpha^{(s)}$ and $\beta^{(s)}$. Using this identity, it can be seen that $\tilde{Z}(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)})$ can be approximated by

$$\tilde{Z}(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)}) \approx B(\alpha^{(u)}, \beta^{(u)})B(\alpha^{(s)}, \beta^{(s)}) \sum_{k=1}^K [1 - F_{\alpha^{(s)}, \beta^{(s)}}(X_k)]$$

where the X_k 's are *iid* beta distributed random variables with parameters $\alpha^{(s)}$, $\beta^{(s)}$ and K is the number of terms used in the Monte Carlo approximation. This approximation works relatively well with our EM implementation and does not significantly increase the computing time. Unfortunately, the number of terms (*i.e.* value of K) required for the approximation to be good might be large and computing such a normalizing

constant at each iteration would significantly slow down our MCMC implementation. As it turns out, a better approximation can be obtained when $\alpha^{(s)}$ and $\beta^{(s)}$ are integers. In this case, the cdf function in (1) can be calculated exactly using integration by parts, as follows,

$$F_{\alpha^{(s)}, \beta^{(s)}}(p^{(u)}) = \sum_{j=\beta^{(s)}}^{\beta^{(s)}+\alpha^{(s)}-1} \frac{(\beta^{(s)} + \alpha^{(s)} - 1)!}{j!(\beta^{(s)} + \alpha^{(s)} - j)!} (1 - p^{(u)})^j p^{(u)\beta^{(s)} + \alpha^{(s)} - j}.$$

Then using this identity, we obtain

$$\begin{aligned} Z(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)}) &= B(\alpha^{(s)}, \beta^{(s)}) \sum_{j=\beta^{(s)}}^{\beta^{(s)}+\alpha^{(s)}-1} \frac{(\beta^{(s)} + \alpha^{(s)} - 1)!}{j!(\beta^{(s)} + \alpha^{(s)} - j)!} \int_0^1 (1 - p^{(u)})^{\beta^{(u)}-1+j} p^{(u)\alpha^{(u)}-1+\beta^{(s)}+\alpha^{(s)}-j} dp^{(u)} \\ &= B(\alpha^{(s)}, \beta^{(s)}) \sum_{j=\beta^{(s)}}^{(\beta^{(s)}+\alpha^{(s)}-1)} \frac{(\beta^{(s)} + \alpha^{(s)} - 1)!}{j!(\beta^{(s)} + \alpha^{(s)} - j)!} B(\beta^{(u)} + j) B(\alpha^{(u)} + \beta^{(s)} + \alpha^{(s)} - j). \end{aligned}$$

Typically, in ICS data $\alpha^{(u)}$ is relatively small leading to relatively few terms in the sum. However, the use of this exact identity in our MCMC algorithm requires the use of discrete priors on $\alpha^{(s)}$ and $\beta^{(s)}$, which can be restrictive in terms of fit (*e.g.*, if the true $\alpha^{(s)}$ is less than one) and can render mixing in the MCMC more difficult. In addition, even though the computation is exact and much faster for small values of $\alpha^{(s)}$, which is typically the case with ICS data, it is still more demanding than the unconstrained model. In our case, we have decided to use the unconstrained model and simply fix the z_i to zero if the empirical proportion for the un-stimulated sample, $p^{(u)}$, is greater than that of the stimulation sample, $p^{(s)}$. Indeed, in the one-sided case, if $p^{(u)} > p^{(s)}$ the associated individual should be a non-responder and thus $z_i = 0$. In our experience, this computational shortcut performs just as well as the true one-sided implementation while being computationally much less demanding.

Appendix B: Computational details for the beta-binomial model

Marginal likelihood derivations

For a given subject i , the null marginal likelihood is obtained after integrating out the prior from the likelihood for model \mathcal{M}_0 , as follows,

$$\begin{aligned} L_0(\alpha^{(u)}, \beta^{(u)} | n_i^{(s)}, n_i^{(u)}) &= \int_0^1 \Pr(n_i^{(s)}, n_i^{(u)} | p^{(u)}) \pi(p^{(u)} | \alpha^{(u)}, \beta^{(u)}) dp^{(u)} \\ &= \int_0^1 \Pr(n^{(s)} | p^{(u)}) \Pr(n^{(u)} | p^{(u)}) \pi(p^{(u)} | \alpha^{(u)}, \beta^{(u)}) dp^{(u)} \\ &= \int_0^1 \binom{N^{(s)}}{n^{(s)}} p^{(s)n^{(s)}} (1 - p^{(u)})^{N^{(s)}-n^{(s)}} \binom{N^{(u)}}{n^{(u)}} p^{(u)n^{(u)}} (1 - p^{(u)})^{N^{(u)}-n^{(u)}} \\ &\quad \frac{1}{B(\alpha^{(u)}, \beta^{(u)})} p^{(u)\alpha^{(u)}-1} (1 - p^{(u)})^{\beta^{(u)}-1} dp^{(u)} \\ &= \binom{N^{(s)}}{n^{(s)}} \binom{N^{(u)}}{n^{(u)}} \frac{1}{B(\alpha^{(u)}, \beta^{(u)})} \int_0^1 p^{(u)n^{(s)}+n^{(u)}+\alpha^{(u)}-1} (1 - p^{(u)})^{N^{(s)}+N^{(u)}-n^{(s)}-n^{(u)}+\beta^{(u)}-1} dp^{(u)} \\ &= \binom{N^{(s)}}{n^{(s)}} \binom{N^{(u)}}{n^{(u)}} \frac{B(n^{(s)} + n^{(u)} + \alpha^{(u)}, N^{(s)} + N^{(u)} - n^{(s)} - n^{(u)} + \beta^{(u)})}{B(\alpha^{(u)}, \beta^{(u)})}. \end{aligned}$$

Similarly, the alternative marginal likelihood for a given subject is defined as,

$$\begin{aligned}
L_1(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | n_i^{(s)}, n_i^{(u)}) &= \int_0^1 \int_0^1 \Pr(n_i^{(s)}, n_i^{(u)} | p^{(u)}, p^{(s)}) \pi(p^{(u)}, p^{(s)} | \alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)}) dp^{(s)} dp^{(u)} \\
&= \int_0^1 \int_0^1 \binom{N^{(s)}}{n^{(s)}} p^{(s)n^{(s)}} (1-p^{(s)})^{N^{(s)}-n^{(s)}} \binom{N^{(u)}}{n^{(u)}} p^{(u)n^{(u)}} (1-p^{(u)})^{N^{(u)}-n^{(u)}} \\
&\quad \frac{1}{B(\alpha^{(u)}, \beta^{(u)})} p^{(u)\alpha^{(u)}-1} (1-p^{(u)})^{\beta^{(u)}-1} \frac{1}{B(\alpha^{(s)}, \beta^{(s)})} p^{(s)\alpha^{(s)}-1} (1-p^{(s)})^{\beta^{(s)}-1} dp^{(s)} dp^{(u)} \\
&= \binom{N^{(s)}}{n^{(s)}} \binom{N^{(u)}}{n^{(u)}} \frac{1}{B(\alpha^{(u)}, \beta^{(u)})} \frac{1}{B(\alpha^{(s)}, \beta^{(s)})} \int_0^1 p^{(u)n^{(u)}+\alpha^{(u)}-1} (1-p^{(u)})^{N^{(u)}-n^{(u)}+\beta^{(u)}-1} \\
&\quad \int_0^1 p^{(s)n^{(s)}+\alpha^{(s)}-1} (1-p^{(s)})^{N^{(s)}-n^{(s)}+\beta^{(s)}-1} dp^{(s)} dp^{(u)} \\
&= \binom{N^{(s)}}{n^{(s)}} \binom{N^{(u)}}{n^{(u)}} \frac{B(n^{(u)} + \alpha^{(u)}, N^{(u)} - n^{(u)} + \beta^{(u)}) B(n^{(s)} + \alpha^{(s)}, N^{(s)} - n^{(s)} + \beta^{(s)})}{B(\alpha^{(u)}, \beta^{(u)}) B(\alpha^{(s)}, \beta^{(s)})}.
\end{aligned}$$

MCMC algorithm

In what follows, we use $(x|y)$ to denote the conditional distribution of x given y . In particular, we use $(x|\dots)$ to denote the distribution of x conditional on everything else in the model. Our MCMC algorithms cycles through the following steps:

1. Update each $\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}$ and $\beta^{(s)}$ by Metropolis-Hastings using a Gaussian symmetric proposal where the variance of the proposal is tuned for each parameter using the approach of Gelman *and others* (2004).
2. Update w by Gibbs sampling using the full conditional,

$$(w|\dots) \sim \text{Beta}\left(\sum_i z_i, \sum_i (1 - z_i)\right).$$

3. for each i , update z_i by Gibbs sampling using the following full conditional,

$$(z_i|\dots) \sim B(1, p_i).$$

where,

$$p_i = \frac{w \cdot L_1(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | n_i^{(u)}, n_i^{(s)})}{w \cdot L_1(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | n_i^{(u)}, n_i^{(s)}) + (1 - w) \cdot L_0(\alpha^{(u)}, \beta^{(u)} | n_i^{(u)}, n_i^{(s)})}.$$

For each updated parameter, step 1 above involves the calculation of the following acceptance ratio, (e.g. $\alpha^{(u)}$, below),

$$\frac{L(\alpha^{(u)\text{new}}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | \dots) \pi(\alpha^{(u)\text{new}})}{L(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | \dots) \pi(\alpha^{(u)})}$$

where L is the complete marginal likelihood conditional on \mathbf{z} defined as,

$$L(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | \mathbf{n}^{(u)}, \mathbf{n}^{(s)}, \mathbf{z}) = \prod_i L_{z_i}(\alpha^{(u)}, \beta^{(u)}, \alpha^{(s)}, \beta^{(s)} | \mathbf{n}^{(u)}, \mathbf{n}^{(s)})$$

and π is the prior distribution of $\alpha^{(u)}$. The obvious changes in the above expression are made for the acceptance ratios of $\alpha^{(s)}, \beta^{(s)}, \beta^{(u)}$. In our case each parameter has the same exponential prior with mean 1,000

Finally, \tilde{z}_i is calculated as $\frac{\sum_{b+1}^T z_i^t}{T}$, where T is the number of MCMC iterations and b is the number of burn-in iterations.

Implementation Details for MCMC Algorithm

We used the method of Raftery and Lewis (1992) and Raftery (Gilks *and others*, 1996) to determine the number of iterations, based on a short pilot run of the sampler. For each dataset presented in the manuscript, we calculated that no more than about 100,000 iterations with 50,000 burn-in iterations was sufficient to estimate standard posterior quantities. To leave some margin, we used 200,000 iterations after 50,000 burn-in iterations for each dataset explored here.

Appendix C: Computational details for the Dirichlet-multinomial model

Marginal likelihood derivations

Because our Dirichlet-multinomial is a direct extension of the beta-binomial model, the marginal likelihoods are obtained in the exact same fashion. For a given subject, the null marginal likelihood is defined as

$$\begin{aligned}
L_0(\boldsymbol{\alpha}^{(u)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)}) &= \int \cdots \int \text{Pr}(\mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)} | \mathbf{p}^{(u)}) \pi(\mathbf{p}^{(u)} | \boldsymbol{\alpha}^{(u)}) d\mathbf{p}^{(u)} \\
&= \int \cdots \int \frac{\mathbf{N}^{(s)}!}{\prod_k n_{ik}^{(s)}!} \prod_k p_k^{(u)n_{ik}^{(s)}} \frac{\mathbf{N}^{(u)}!}{\prod_k n_{ik}^{(u)}!} \prod_k p_k^{(u)n_{ik}^{(u)}} \frac{1}{B(\boldsymbol{\alpha}^{(u)})} \prod_k p_k^{(u)\alpha_k^{(u)}-1} d\mathbf{p}^{(u)} \\
&= \frac{\mathbf{N}^{(s)}!}{\prod_k n_{ik}^{(s)}!} \frac{\mathbf{N}^{(u)}!}{\prod_k n_{ik}^{(u)}!} \frac{1}{B(\boldsymbol{\alpha}^{(u)})} \int_0^1 \prod_k p_k^{(u)n_{ik}^{(s)}+n_{ik}^{(u)}+\alpha_k^{(u)}-1} d\mathbf{p}^{(u)} \\
&= \frac{\mathbf{N}^{(s)}!}{\prod_k n_{ik}^{(s)}!} \frac{\mathbf{N}^{(u)}!}{\prod_k n_{ik}^{(u)}!} \frac{B(\mathbf{n}_i^{(s)} + \mathbf{n}_i^{(u)} + \boldsymbol{\alpha}^{(u)})}{B(\boldsymbol{\alpha}^{(u)})}
\end{aligned}$$

Similarly, the alternative marginal likelihood for a given subject is:

$$\begin{aligned}
L_1(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \mathbf{n}_i^{(u)}, \mathbf{n}_i^{(s)}) &= \frac{\mathbf{N}^{(s)}!}{\prod_k n_{ik}^{(s)}!} \frac{\mathbf{N}^{(u)}!}{\prod_k n_{ik}^{(u)}!} \frac{1}{B(\boldsymbol{\alpha}^{(u)})B(\boldsymbol{\alpha}^{(s)})} \int \cdots \int \prod_k p_k^{(u)n_{ik}^{(u)}+\alpha_k^{(u)}-1} \int \cdots \int \prod_k p_k^{(s)n_{ik}^{(s)}+\alpha_k^{(s)}-1} d\mathbf{p}^{(u)} d\mathbf{p}^{(s)} \\
&= \frac{\mathbf{N}^{(s)}!}{\prod_k n_{ik}^{(s)}!} \frac{\mathbf{N}^{(u)}!}{\prod_k n_{ik}^{(u)}!} \frac{B(\mathbf{n}_i^{(u)} + \boldsymbol{\alpha}^{(u)})B(\mathbf{n}_i^{(s)} + \boldsymbol{\alpha}^{(s)})}{B(\boldsymbol{\alpha}^{(u)})B(\boldsymbol{\alpha}^{(s)})}.
\end{aligned}$$

MCMC algorithm

The MCMC algorithm for the Dirichlet-multinomial model is analogous to the beta-binomial, above. The parameter vectors $\boldsymbol{\alpha}^{(s)}, \boldsymbol{\alpha}^{(u)}$ are updated component-wise:

1. Update each $\alpha_k^{(s)}, \alpha_k^{(u)}$ using a Gaussian symmetric proposal distribution with the variance of each proposal tuned using the approach of Gelman *and others* (2004).
2. Update w by Gibbs sampling using the full conditional,

$$(w | \cdots) \sim \text{Beta}\left(\sum_i z_i, \sum_i (1 - z_i)\right)$$

3. For each i , update z_i using the full conditional,

$$(z_i | \cdots) \sim B(1, p_i)$$

where

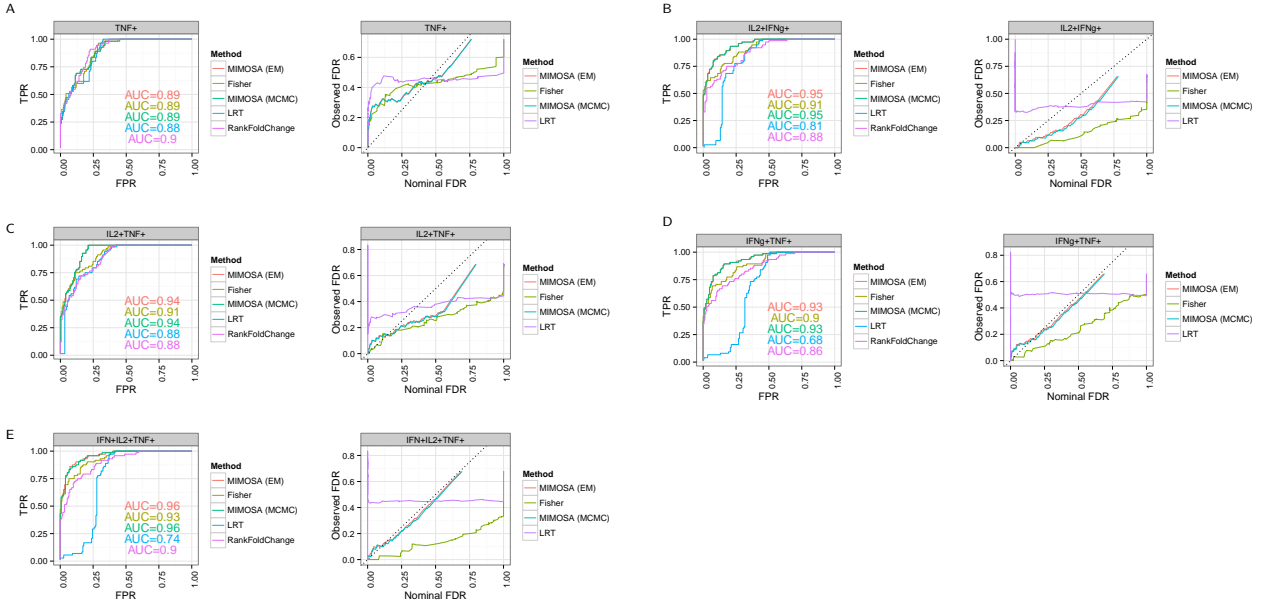
$$p_i = \frac{w \cdot L_1(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)})}{w \cdot L_1(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)}) + (1 - w) \cdot L_0(\boldsymbol{\alpha}^{(u)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)})}$$

For each parameter component updated in step 1 above, compute the acceptance ratio (e.g. $\alpha_k^{(u)}$, below):

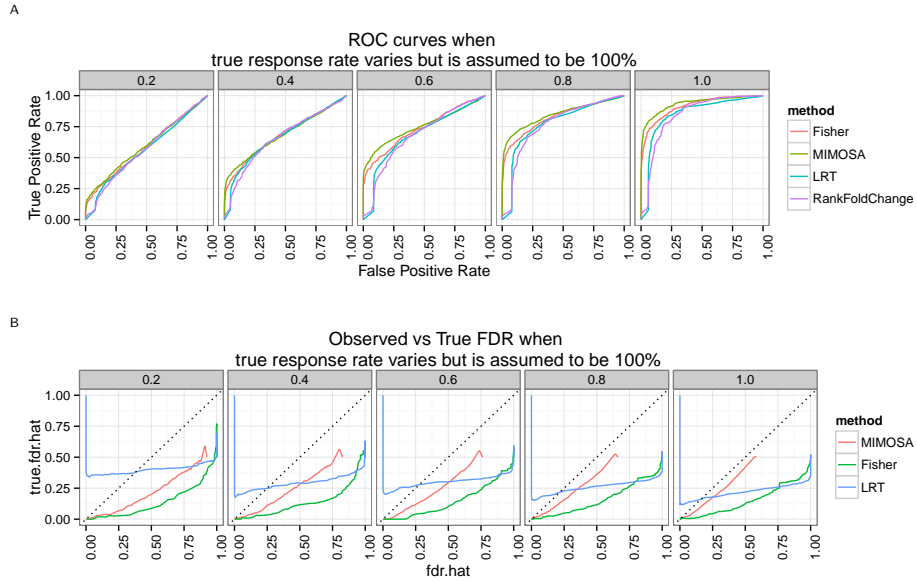
$$\frac{L(\alpha_k^{(u)\text{new}}, \boldsymbol{\alpha}_{\{-k\}}^{(u)}, \boldsymbol{\alpha}^{(s)} | \cdots) \pi(\alpha_k^{(u)\text{new}})}{L(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \cdots) \pi(\alpha_k^{(u)})}$$

where π is the prior distribution of the parameter, and $\boldsymbol{\alpha}_{\{-k\}}^{(u)} = \{\alpha_j^{(u)} : j \neq k\}$. Again, we have used the same exponential prior with mean 1,000 for each parameter. L is the complete marginal likelihood conditional on \mathbf{z} , defined as

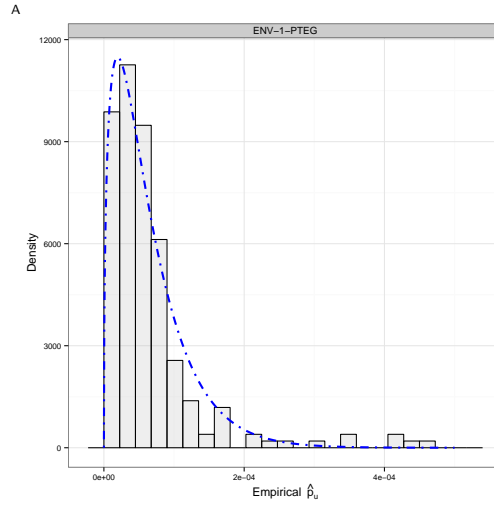
$$L(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)}, \mathbf{z}) = \prod_i L_{z_i}(\boldsymbol{\alpha}^{(u)}, \boldsymbol{\alpha}^{(s)} | \mathbf{n}_i^{(s)}, \mathbf{n}_i^{(u)}).$$



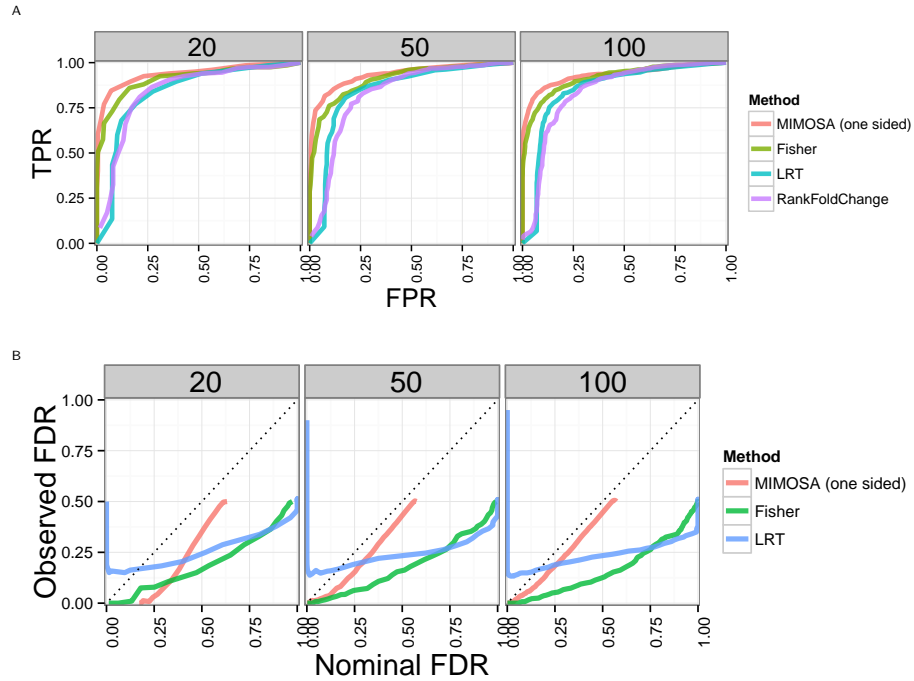
Supplementary Figure 1: Comparison of MIMOSA on other cytokines and cytokine combinations for ENV-1-PTEG stimulated CD4+ T-cells from the HVTN065 trial.



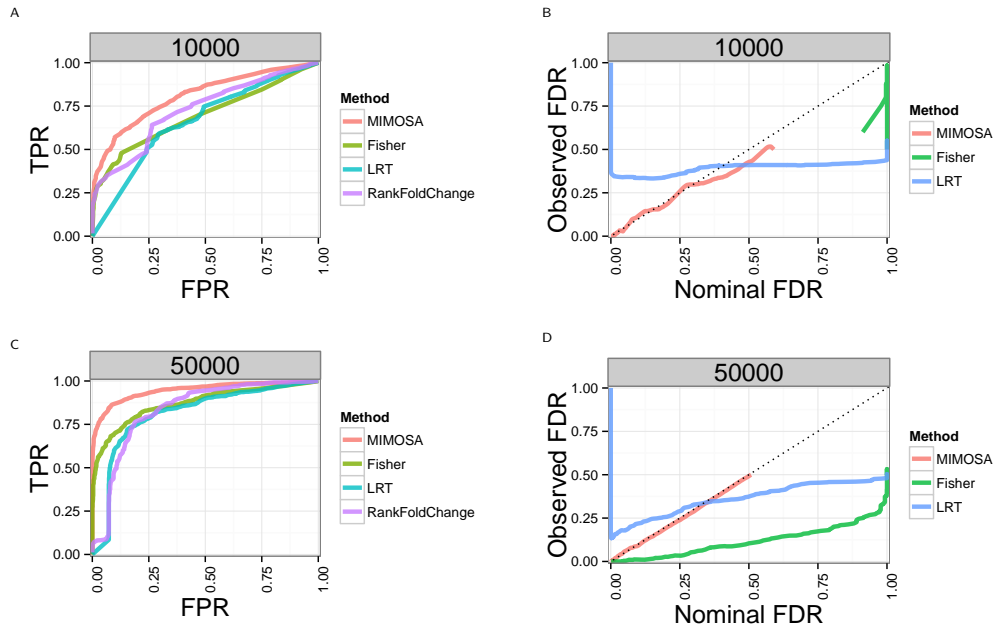
Supplementary Figure 2: Effect of deviations of the assumptions of 100% true responders at the post-vaccine time point on A) ROC and B) FDR curves. The true response rate is shown in the gray box above each plot.



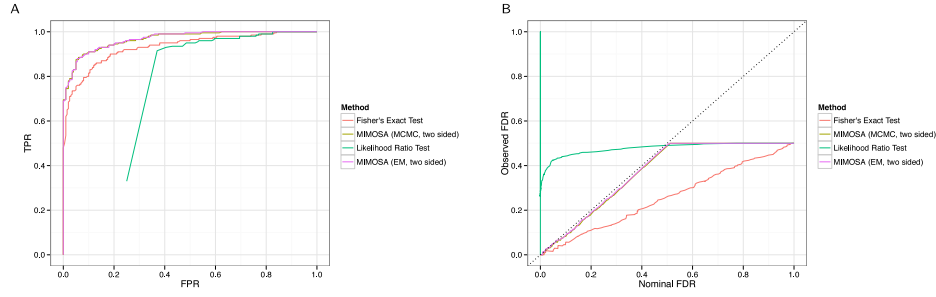
Supplementary Figure 3: Histogram of the empirical proportions of unstimulated cells and overlaid posterior densities of the beta distribution with $\alpha^{(u)}$ and $\beta^{(u)}$ estimated from the data for ENV-1-PTEG stimulated, IFN- γ +, CD4+ T-cells, demonstrating that the assumption of a common distribution for p_{iu} across subjects is reasonable.



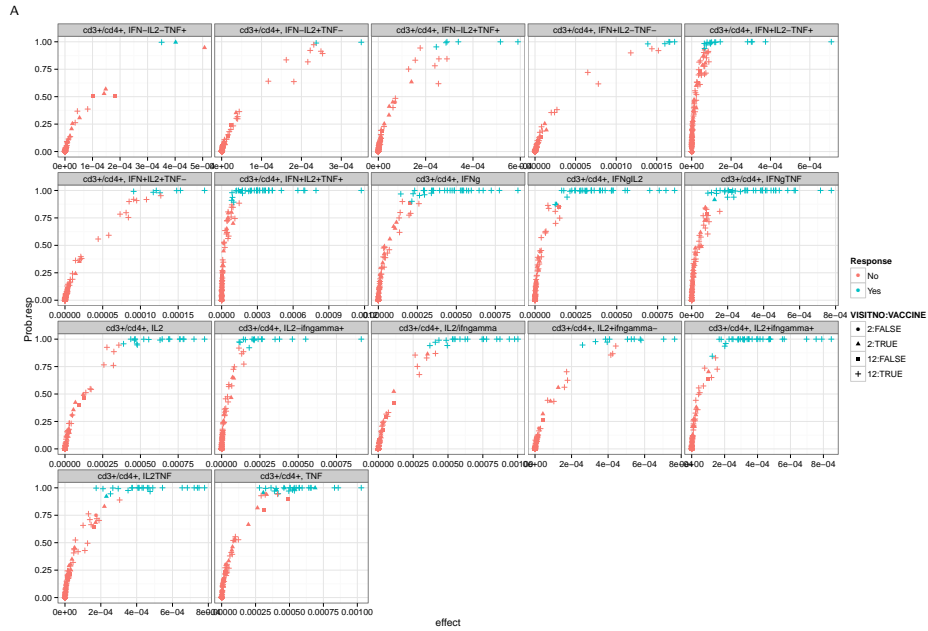
Supplementary Figure 4: One-sided MIMOSA model fit to simulated data with 50K cells and varying values of I (number of observations). A) Average ROC curves from 10 simulations for 20, 50 and 100 observations. B) Average observed and nominal FDR from 10 simulations for 20, 50, and 100 observations.



Supplementary Figure 5: Unconstrained MIMOSA model fit to two-sided data with small counts and to data from a model violating model assumptions. For two-sided data A) the average ROC from 10 simulation with $N=10,000$ cells. B) the average observed and nominal FDR from 10 simulations with $N=10,000$ cells. Data was simulated from a model where proportions were sampled from a truncated normal distribution over $[0, 1]$ rather than a Beta distribution. C) Average ROC for $N=50,000$ cells D) Average observed and nominal FDR for $N=50,000$ cells.



Supplementary Figure 6: Multivariate simulations from a two-sided model. Ten, eight-dimensional data sets were simulated from a two-sided model with an effect sizes of 2.5×10^{-3} and -2.5×10^{-3} in two of the eight dimensions ($N=1,500$). Multivariate MIMOSA was compared against Fisher's exact test, and the likelihood ratio test. A) Average ROC curves for the competing methods over 10 simulations. B) Average observed and nominal false discovery rate for each method over 10 simulations. This figure appears in color in the electronic version of this article.



Supplementary Figure 7: Volcano plots of effect size vs posterior probability of response for each model fit to each cytokine in the ICS data set. Green points indicate subjects called responders for that cytokine and stimulation at a 1% FDR threshold (adjusted across subjects within cytokine subset). Red points indicate non-responders. Visit code 2 is day 0, and 12 is day 182. VACCINE=FALSE indicates a placebo recipient, while TRUE indicates a vaccinee.

References

- GELMAN, ANDREW, CARLIN, JOHN B, STERN, HAL S AND RUBIN, DONALD B. (2004). *Bayesian Data Analysis*. Chapman & Hall/CRC.
- GILKS, WALTER R, RICHARDSON, SYLVIA AND SPIEGELHALTER, D. J. (1996). *Markov Chain Monte Carlo in Practice*, First CRC Press Reprint edition. 2000 N.W. Corporate Blvd. Boca Raton, Florida, 33431: Chapman Hall / CRC Press.
- RAFTERY, ADRIAN E AND LEWIS, STEVEN M. (1992, November). [Practical Markov Chain Monte Carlo]: Comment: One Long Run with Diagnostics: Implementation Strategies for Markov Chain Monte Carlo. *STATISTICAL SCIENCE* **7**(4), 493–497.