

# Web-based Supplementary Materials for Combinatorial Mixture Models for Single-Cell Assays with Application to Vaccine Studies

by Greg Finak, Steve De Rosa, Mario Roederer and Raphael Gottardo

## Web Appendix A: HVTN065 vaccine trial ICS data description

HVTN065 is a phase 1 (safety and immunogenicity) trial of GeoVax HIV/AIDS DNA and MVA vaccine in 120 individuals (100 vaccinees, 20 placebo recipients, parts A and B). CD4 and CD8 T-cell epitope specific immune responses were measured via the ICS assay. Cytokines measured in the ICS assay included IFN $\gamma$ , TNF $\alpha$ , IL2, and IL4, and antigens included three Env, three Gag, and three Pol peptide pools. Results of the trial have been published in Nilu et al. (2006).

## Web Appendix B: Computational details for the beta-binomial model

### Marginal likelihood derivations

#### 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 cycle 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 tune for each parameter using the approach of (Greg add the citation).
2. Update  $w$  by Gibbs sampling using the full conditional,

$$(w|\dots) \sim$$

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

$$(z_i|\dots) \sim$$

For each updated parameter, step 1 above involves the calculation of the following acceptance ratio

*add formula.*

where  $\pi$  is the prior distribution of the corresponding parameter. In our case each parameter has the same exponential prior with mean  $10^{31}$ .

---

<sup>1</sup>Check that this is correct

## Web Appendix C: 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

$$L_1(\alpha_u, \beta_u, \alpha_s, \beta_s | \mathbf{y}) = \prod_{i=1}^P \binom{N_{ui}}{n_{ui}} \binom{N_{si}}{n_{si}} \frac{Z(n_{ui} + \alpha_u, N_{ui} - n_{ui} + \beta_u, n_{si} + \alpha_s, N_{si} - n_{si} + \beta_s)}{Z(\alpha_u, \beta_u, \alpha_s, \beta_s)}. \quad (1)$$

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

$$\tilde{Z}(\alpha_u, \beta_u, \alpha_s, \beta_s) = B(\alpha_u, \beta_u)B(\alpha_s, \beta_s) \int_0^1 \frac{p_u^{\alpha_u-1} (1-p_u)^{\beta_u-1}}{B(\alpha_u, \beta_u)} (1 - F_{\alpha_s, \beta_s}(p_u)) dp_u \quad (2)$$

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

$$\tilde{Z}(\alpha_u, \beta_u, \alpha_s, \beta_s) \approx B(\alpha_u, \beta_u)B(\alpha_s, \beta_s) \sum_{k=1}^K (1 - F_{\alpha_s, \beta_s}(X_k))$$

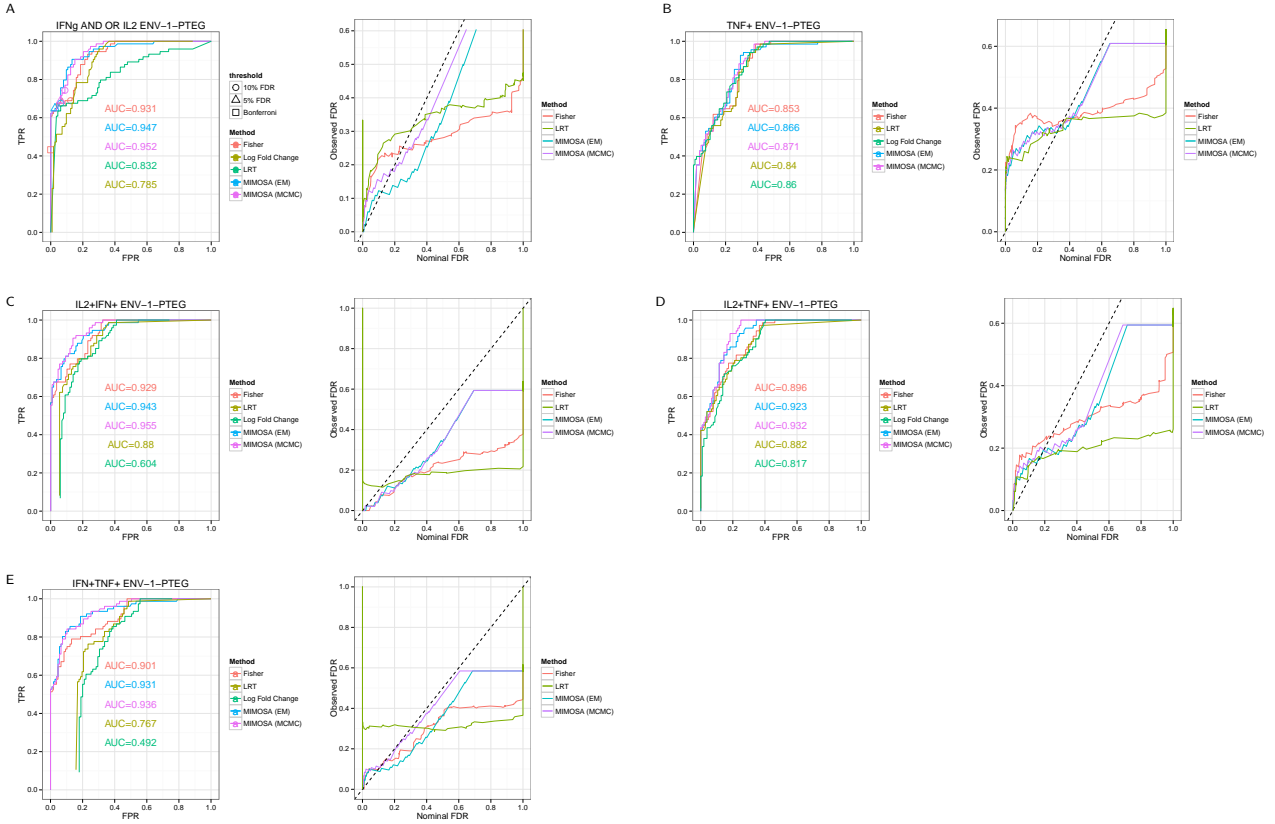
where the  $X_k$ 's are *iid* beta distributed random variables with parameters  $\alpha_s$ ,  $\beta_s$  and  $K$  is the number of terms used in the Monte Carlo approximation. This approximation works relatively well with our EM implementation and does not significantly increase the computing time. Unfortunately, the number of terms (*i.e.* value of  $I$ ) 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 (2) 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  $\alpha_s$  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  $\alpha_s$ ,



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

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 less 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.

## Web Appendix D: Computational details for the Dirichlet-multinomial model

### Marginal likelihood derivations

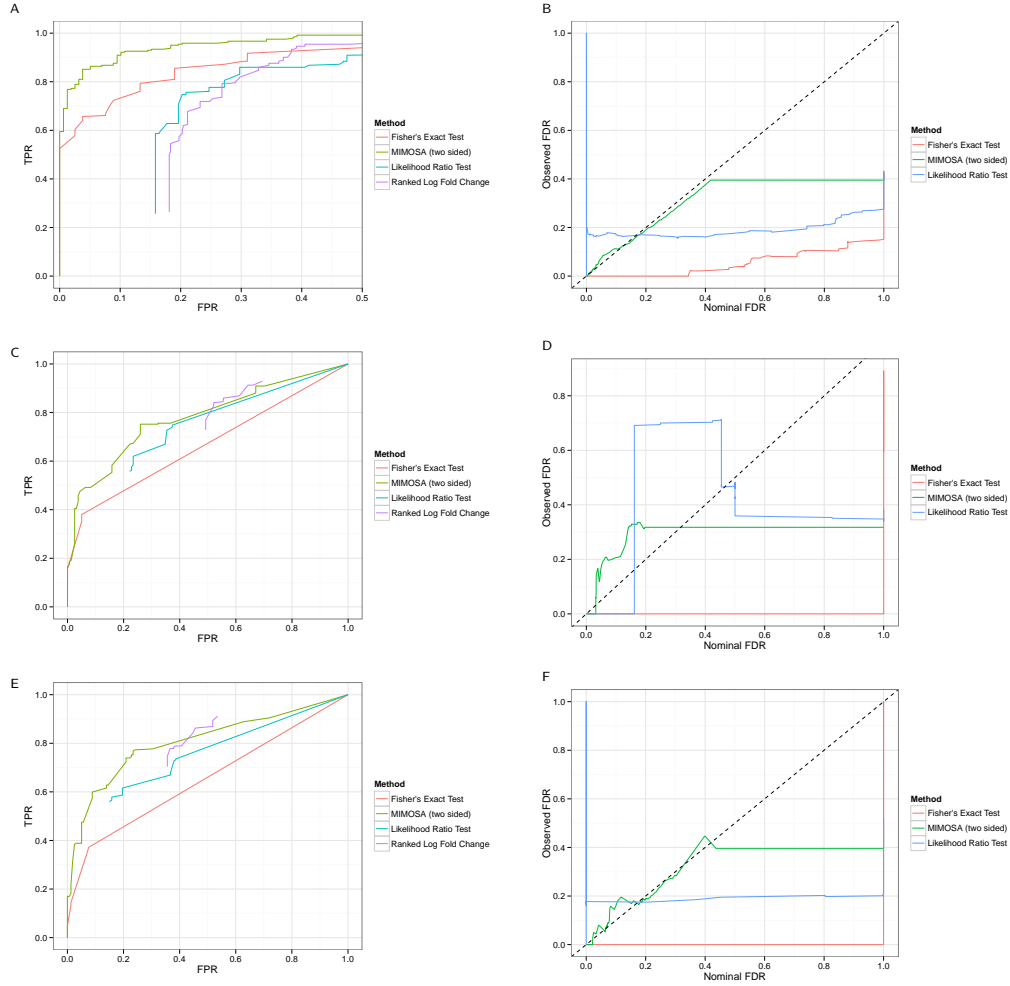
Because our Dirichlet-multinomial is a direct extension of the beta-binomial model, the marginal likelihoods are obtained in the exact same fashion. The derivations is described below,

### MCMC algorithm

GREG: Look at the what I did above, and do something similar here.

## References

Nilu G. et al., 2006, J Virol, 80, 4717



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