

# Variational Bayes for Model Averaging for Multivariate models using Compositional predictors

---

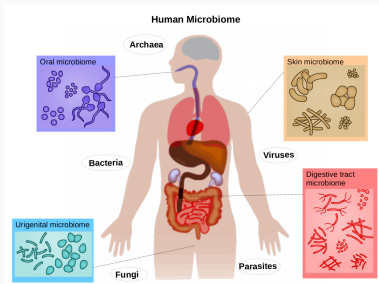
Alex Lewin

(Work with Darren Scott, LSHTM)

December 2020



# Microbiome

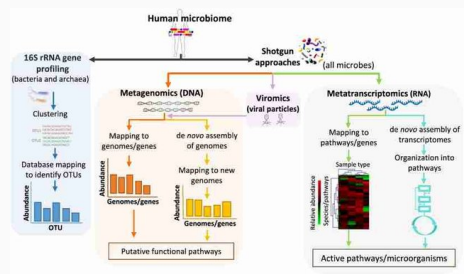


Microbiome: collection of all the microbes inside and on surface of the human body

Microbes interact with immune system, weight regulation etc. etc.

Quantification: next generation sequencing of microbiome species

Obtain relative abundances of different species



## Compositional data as Covariates

Microbiome:  $p$  variables  $X_1, \dots, X_p$  with  $\sum_j X_j = 1$

Usual to use some form of log transform for microbiome variables as predictors.

Eg, with reference category, use **log-ratio** transform:

$$y = \sum_{j=1}^{p-1} \theta_j \log(X_j/X_p) + \epsilon$$

We use an alternative parametrisation:

$$y = \sum_{j=1}^p \theta_j \log(X_j) + \epsilon$$

with **constraint**  $\sum_{j=1}^p \theta_j = 0$

# Variable selection with Compositional data

## Spike-and-slab prior:

Variable selection performed through binary indicators

$$\xi_j = \begin{cases} 1 & \implies \theta_j \neq 0 \\ 0 & \implies \theta_j = 0 \end{cases}$$

Two issues regarding priors:

1. Prior on non-zero  $\theta$  ?  
Need constraint  $\sum_{j=1}^p \theta_j = 0$ .
2. Prior on  $\xi$  ?  
We cannot have  $|\xi| = 1$ .

## Priors for Variable selection with Compositional data

Singular MVN on non-zero  $\theta$ :

$$\theta_{\xi} | \xi \sim N(T_{\xi} \mu_{\xi}, \sigma^2 T_{\xi} T_{\xi}^T)$$

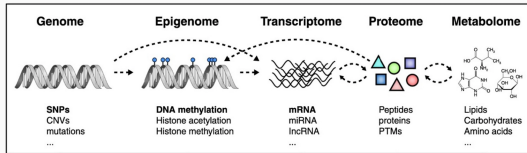
where  $T_{\xi} = \mathbb{I}_{|\xi|} - \frac{1}{|\xi|} \mathbb{J}_{|\xi|}$  ensures the sum to zero constraint.

For binary indicators, use a truncated distribution:

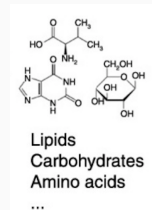
$$\xi \propto \prod_{j=1}^p \kappa_j^{\xi_j} (1 - \kappa_j)^{(1 - \xi_j)} \times \mathbb{I}[|\xi| \neq 1]$$

# High-dimensional outcomes, high-dimensional predictors

Metabolomics: small molecules, products of cellular processes



Our aim: linking metabolome ~ microbiome



# High-dimensional outcomes, high-dimensional predictors

Frame the problem as a multivariate linear regression model:

$$Y|X, Z \sim \mathcal{MN}(X_\gamma B_\gamma + Z_\xi \Theta_\xi, \mathbb{I}, C_\eta)$$

- Sparse variable selection on matrices of associations ( $B_\gamma$  and  $\Theta_\xi$ )
- Sparse covariance selection ( $C_\eta$ )
- Model averaging over all combinations of  $\gamma, \xi, \eta$  provides flexible modelling of  $B, \Theta, C$

Our previous work on sparse variable and covariance selection using MCMC:

**Bottolo, Banterle, et al. doi: [10.1101/467019](https://doi.org/10.1101/467019) on bioRxiv**

# Estimating the posterior: Variational Bayes

## Full model

$$\begin{aligned} p(y, \theta) = & \left\{ \prod_t p(y_t | \beta_t, \sigma_t^2, \rho_t) \right\} \times \left\{ \prod_t \prod_s p(\beta_{ts} | w_t, \gamma_{ts}) \right\} \times \left\{ \prod_t p(\theta_t | \Sigma_t(w_t, T), \xi_t) \right\} \times \\ & \left\{ \prod_t \prod_s p(\gamma_{ts} | \omega_s) \right\} \times \left\{ \prod_t \prod_j p(\xi_{tj} | \kappa_j) \right\} \times \left\{ \prod_s p(\omega_s) \right\} \times \left\{ \prod_j p(\kappa_j) \right\} \times \\ & \left\{ \prod_t p(\sigma_t^2 | \tau, \nu) \prod_{k < t} p(\rho_{tk} | \sigma_t^2, \tau, \eta_{tk}) \right\} \times \\ & \left\{ \prod_t \prod_{k < t} p(\eta_{tk} | \lambda) \right\} \times \left\{ \prod_t p(w_t | a_w, b_w) \right\} \times p(\lambda) p(b_w) p(\tau) \end{aligned}$$

## Mean field approximation

$$\begin{aligned} q(\mathbf{z}) = & \left\{ \prod_t \prod_s q(\beta_{ts}, \gamma_{ts}) \right\} \times \left\{ \prod_t q(\theta_t, \xi_t) \right\} \times \left\{ \prod_s q(\omega_s) \right\} \times \left\{ \prod_j q(\kappa_j) \right\} \\ & \left\{ \prod_t q(\sigma_t^2) \prod_{k < t} q(\rho_{tk}, \eta_{tk} | \sigma_t^2) \right\} \times \\ & \left\{ \prod_t q(w_t | a_w, b_w) \right\} \times q(\lambda) q(b_w) q(\tau) \end{aligned}$$



## Mean field approximation for variable selection

We use CAVI (co-ordinate ascent variational inference).

Unconstrained	Constrained
Need joint distribution for $\beta, \gamma$	Need joint distribution for $\theta, \xi$
Product of Normal and Bernoulli	Product of Singular Normal and Truncated Bernoulli
Treat variables ( $j$ ) independently	Variables ( $j$ ) are dependent
Mean field: $\prod_{j=1}^p q(\beta_j   \gamma_j) q(\gamma_j)$	Mean field: $q(\theta   \xi) q(\xi)$
VB (CAVI) obtains $E_q(\beta_j   \gamma_j)$ , $Var_q(\beta_j   \gamma_j)$ , $E_q(\gamma_j)$	VB (CAVI) obtains $E_q(\theta   \xi)$ , $Var_q(\theta   \xi)$
	MCMC needed to get $E_q(\xi)$

---

**Algorithm 1** Pseudo-Algorithm for Hybrid VB-MCMC

---

**for**  $i = 1, \dots, n_{VB}$  **do**

CAVI update obtains  $E_q(\beta)$

CAVI update obtains  $E_q(\gamma)$

...

**for**  $j = 1, \dots, n_{MCMC}$  **do**

Reversible Jump Birth/Death move propose update to  $\xi$

Accept/reject using Metropolis-Hastings step

**end for**

Calculate  $\hat{E}_q(\xi)$  Monte Carlo average

**end for**

---

# Bayesian Model Averaging

Bayesian framework allows for model averaging over all explored models.

Posterior means of binary indicators  $\gamma, \xi, \eta$  are marginal posterior probabilities of inclusion (MPPI):

*Using mean field approximation,  $E_q(\gamma)$  are estimates of posterior means.*

*Monte Carlo average for compositional indicators:*

$$\hat{E}_q(\boldsymbol{\xi}) = \frac{1}{N_{iter}} \sum_{t=1}^{N_{iter}} \boldsymbol{\xi}^{(t)}$$

Also obtain shrunk estimates of regression coefficients (include uncertainty from model selection):

*Mean field approximation:  $E_q(\boldsymbol{\beta}), E_q(\boldsymbol{\theta})$*

# Simulation Study

Predictors (real data from  $n=514$  adult cohort):

- 40 microbiome species
- sex, alcohol intake, diet intake
- 60 noise variables

Simulate 10 response variables from multi-variate Normal

Associations:

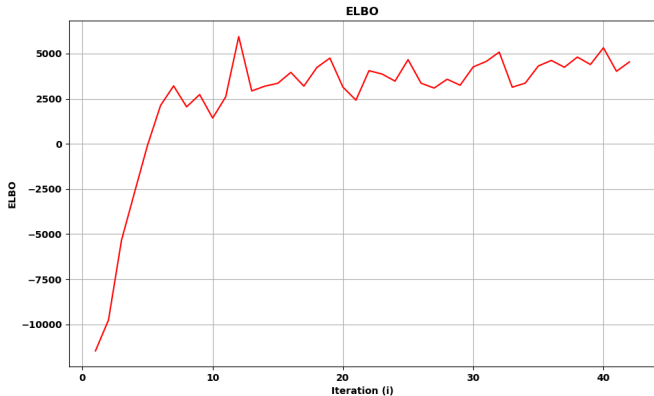
- 8 microbiome species associated with 5 responses each
- sex, alcohol, diet associated with 10, 8, 4 responses

Simulation Study Parameters:

- $\text{SNR} = 2$
- residual correlation = 0.3
- residual dependence block diagonal

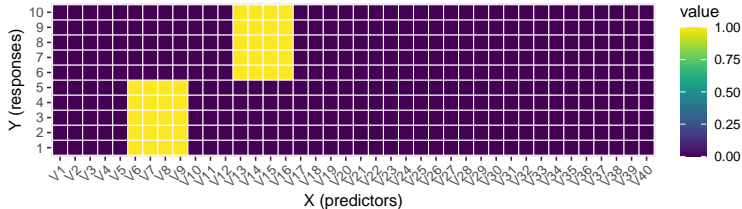
# Results

Monitoring convergence of Variational updates:

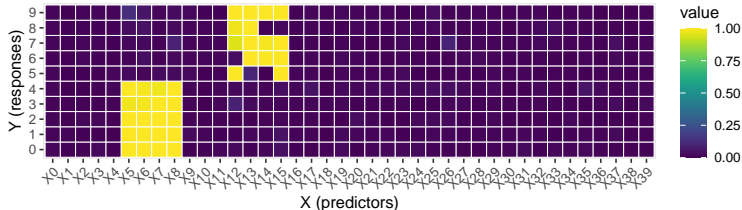


# Results

## True associations with compositional predictors

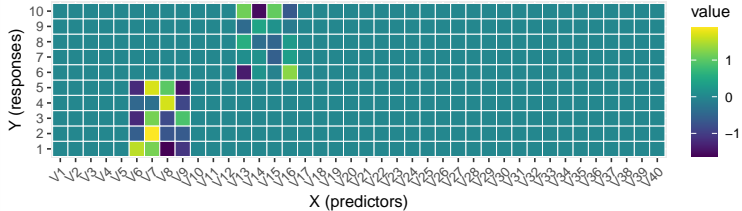


## Estimated MPPI for compositional predictors

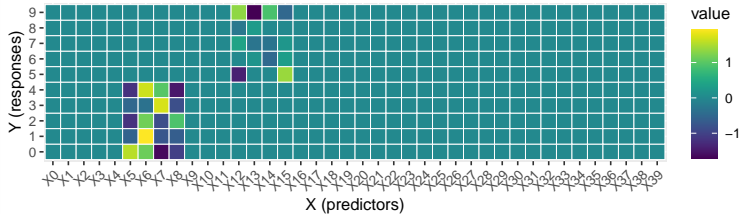


# Results

### True regression coefficients for compositional predictors



### Posterior mean regression coefficients



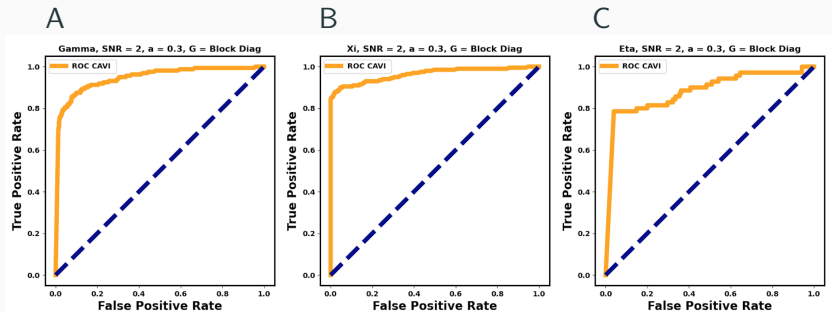
# Results

ROC curves averaged over 5 simulated data sets.

A: Sparse variable selection on non-compositional predictors

B: Sparse variable selection on compositional predictors

C: Sparse covariance selection





# Summary

- Bayesian modelling of multivariate responses: sparse variable selection and covariance selection
- Bayesian model averaging provides flexible modelling of associations
- Variational Bayes to speed up computation for high-dimensional data
- Hybrid VB-MCMC to deal with compositional predictors

# Acknowledgements

## Thank you!

- Darren Scott
- Leonardo Bottolo
- Marco Banterle
- Sylvia Richardson

Funded by MRC LID Doctoral training grant, MRC Methodology Grant, DynaHealth Healthy Aging EU Horizon 2020 grant.

Variable and covariance selection:

**Bottolo, Banterle, et al. doi: 10.1101/467019 on bioRxiv**

Microbiome data transformation:

**Lin et al. Biometrika 2013.**

Hybrid VB-MCMC:

**Ye et al. Statistics and Computing 2020.**