

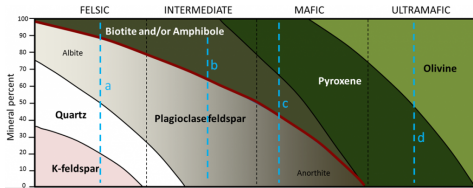
Predictive Modeling of Compositional Data with Supervised Log-Ratios

Jing Ma

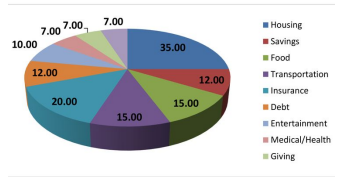
Division of Public Health Sciences
Fred Hutchinson Cancer Center

1 July 2022
ICSA China Conference

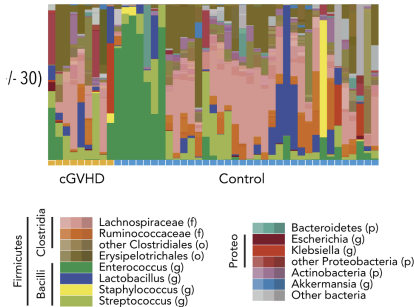
Compositional Data are Everywhere



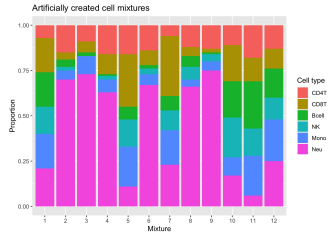
Geology



Sociology



Microbiome: Markey et al., Blood, 20'



Single cell transcriptomics

A vector $X = (X_1, \dots, X_p)$ representing proportions of some whole is subject to the constraint

$$X_1 + \dots + X_p = 1$$

A vector $X = (X_1, \dots, X_p)$ representing proportions of some whole is subject to the constraint

$$X_1 + \dots + X_p = 1$$

Predictive modeling

- ▶ Predictors $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,p})$: compositional
- ▶ Outcome y_i : continuous or binary

A vector $X = (X_1, \dots, X_p)$ representing proportions of some whole is subject to the constraint

$$X_1 + \dots + X_p = 1$$

Predictive modeling

- ▶ Predictors $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,p})$: compositional
- ▶ Outcome y_i : continuous or binary

Scientific question:

- ▶ Define biomarker(s) using a small set of variables that predict disease risk

A vector $X = (X_1, \dots, X_p)$ representing proportions of some whole is subject to the constraint

$$X_1 + \dots + X_p = 1$$

Predictive modeling

- ▶ Predictors $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,p})$: compositional
- ▶ Outcome y_i : continuous or binary

Scientific question:

- ▶ Define biomarker(s) using a small set of variables that predict disease risk

Challenges:

- ▶ The unit-sum constraint makes it difficult to interpret the effect of predictors on the response

Additive log-ratio transform:

$$\text{alr}(X) = \left(\log \frac{X_1}{X_p}, \dots, \log \frac{X_{p-1}}{X_p} \right)$$

Log-contrast regression:

$$\mathbb{E}[y_i | \mathbf{x}_i] = (\boldsymbol{\theta}^{\text{alr}})^\top \text{alr}(\mathbf{x}_i)$$

¹Lin et al., Biometrika, 14'; Shi et al., AOAS, 16'

²Wang and Zhao, AOAS, 17'; Bien et al., Scientific Reports, 21'

Additive log-ratio transform:

$$\text{alr}(X) = \left(\log \frac{X_1}{X_p}, \dots, \log \frac{X_{p-1}}{X_p} \right)$$

Log-contrast regression:

$$\mathbb{E}[y_i | \mathbf{x}_i] = (\boldsymbol{\theta}^{\text{alr}})^\top \text{alr}(\mathbf{x}_i) = \boldsymbol{\beta}^\top \log(\mathbf{x}_i),$$

subject to $\boldsymbol{\beta}^\top \mathbf{1} = 0$.

¹Lin et al., *Biometrika*, 14'; Shi et al., *AOAS*, 16'

²Wang and Zhao, *AOAS*, 17'; Bien et al., *Scientific Reports*, 21'

Additive log-ratio transform:

$$\text{alr}(X) = \left(\log \frac{X_1}{X_p}, \dots, \log \frac{X_{p-1}}{X_p} \right)$$

Log-contrast regression:

$$\mathbb{E}[y_i | \mathbf{x}_i] = (\boldsymbol{\theta}^{\text{alr}})^\top \text{alr}(\mathbf{x}_i) = \boldsymbol{\beta}^\top \log(\mathbf{x}_i),$$

subject to $\boldsymbol{\beta}^\top \mathbf{1} = 0$.

High-dimensional extensions: compositional Lasso¹, tree-guided compositional Lasso²

¹Lin et al., *Biometrika*, 14'; Shi et al., *AOAS*, 16'

²Wang and Zhao, *AOAS*, 17'; Bien et al., *Scientific Reports*, 21'

Additive log-ratio transform:

$$\text{alr}(X) = \left(\log \frac{X_1}{X_p}, \dots, \log \frac{X_{p-1}}{X_p} \right)$$

Log-contrast regression:

$$\mathbb{E}[y_i | \mathbf{x}_i] = (\boldsymbol{\theta}^{\text{alr}})^\top \text{alr}(\mathbf{x}_i) = \boldsymbol{\beta}^\top \log(\mathbf{x}_i),$$

subject to $\boldsymbol{\beta}^\top \mathbf{1} = 0$.

High-dimensional extensions: compositional Lasso¹, tree-guided compositional Lasso²

Limitation: alr coefficients need to be interpreted w.r.t. a reference variable, while constrained regression suffers from prediction accuracy.

¹Lin et al., *Biometrika*, 14'; Shi et al., *AOAS*, 16'

²Wang and Zhao, *AOAS*, 17'; Bien et al., *Scientific Reports*, 21'

Pairwise log-ratios³

$$\mathbb{E}[y_i | \mathbf{x}_i] = \sum_{1 \leq j < k \leq p} \theta_{j,k}^{\text{plr}} \log \frac{x_{i,j}}{x_{i,k}}$$

The log-contrast coefficient $\beta = C^T \theta^{\text{plr}}$ where for $p = 4$

$$C^T = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 1 & 0 \\ 0 & -1 & 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 0 & -1 & -1 \end{pmatrix}.$$

³Bates and Tibshirani, Biometrics, 19'

Pairwise log-ratios³

$$\mathbb{E}[y_i | \mathbf{x}_i] = \sum_{1 \leq j < k \leq p} \theta_{j,k}^{\text{plr}} \log \frac{X_{i,j}}{X_{i,k}}$$

The log-contrast coefficient $\beta = C^T \theta^{\text{plr}}$ where for $p = 4$

$$C^T = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 1 & 0 \\ 0 & -1 & 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 0 & -1 & -1 \end{pmatrix}.$$

Limitation: this model is not identifiable due to co-linearity of predictors, e.g.

$$\log \frac{X_1}{X_2}, \quad \log \frac{X_1}{X_3}, \quad \log \frac{X_2}{X_3}$$

³Bates and Tibshirani, Biometrics, 19'

Balance is the log-ratio between two geometric means

$$B(\mathbf{X}; I_+, I_-) = \log \frac{g(\mathbf{X}_{I_+})}{g(\mathbf{X}_{I_-})} = \frac{\sum_{j \in I_+} \log X_j}{|I_+|} - \frac{\sum_{j \in I_-} \log X_j}{|I_-|}$$

Balance regression searches for the best subsets I_+ and I_- :

$$\mathbb{E}[y_i | \mathbf{x}_i] = \theta_0 + \theta_1 B(\mathbf{x}_i; I_+, I_-)$$

Balance is the log-ratio between two geometric means

$$B(\mathbf{X}; I_+, I_-) = \log \frac{g(\mathbf{X}_{I_+})}{g(\mathbf{X}_{I_-})} = \frac{\sum_{j \in I_+} \log X_j}{|I_+|} - \frac{\sum_{j \in I_-} \log X_j}{|I_-|}$$

Balance regression searches for the best subsets I_+ and I_- :

$$\mathbb{E}[y_i | \mathbf{x}_i] = \theta_0 + \theta_1 B(\mathbf{x}_i; I_+, I_-)$$

selbal⁴ uses greedy search to find the best subsets by adding one variable at a time from the best k -part balances for $k \geq 2$.

⁴Rivera-Pinto, mSystems, 18'

Balance is the log-ratio between two geometric means

$$B(X; I_+, I_-) = \log \frac{g(X_{I_+})}{g(X_{I_-})} = \frac{\sum_{j \in I_+} \log X_j}{|I_+|} - \frac{\sum_{j \in I_-} \log X_j}{|I_-|}$$

Balance regression searches for the best subsets I_+ and I_- :

$$\mathbb{E}[y_i | \mathbf{x}_i] = \theta_0 + \theta_1 B(\mathbf{x}_i; I_+, I_-)$$

selbal⁴ uses greedy search to find the best subsets by adding one variable at a time from the best k -part balances for $k \geq 2$.

selbal prioritizes sparse models, but **exhaustive search is time consuming.**

⁴Rivera-Pinto, mSystems, 18'

CoDaCoRe⁵ uses continuous relaxation to find the best subsets. For a vector of assignment weights \mathbf{w} , let

$$\tilde{\mathbf{w}} = \frac{2}{1 + \exp(-\mathbf{w})} - 1.$$

Let $\tilde{\mathbf{w}}^+ = \text{ReLU}(\tilde{\mathbf{w}})$ and $\tilde{\mathbf{w}}^- = \text{ReLU}(-\tilde{\mathbf{w}})$. The relaxed balance is

$$\tilde{B}(X; \mathbf{w}) = \frac{\sum_j \tilde{w}_j^+ \log X_j}{\sum_j \tilde{w}_j^+} - \frac{\sum_j \tilde{w}_j^- \log X_j}{\sum_j \tilde{w}_j^-}$$

⁵Gordon-Rodriguez et al., Bioinformatics, 22'

CoDaCoRe⁵ uses continuous relaxation to find the best subsets. For a vector of assignment weights \mathbf{w} , let

$$\tilde{\mathbf{w}} = \frac{2}{1 + \exp(-\mathbf{w})} - 1.$$

Let $\tilde{\mathbf{w}}^+ = \text{ReLU}(\tilde{\mathbf{w}})$ and $\tilde{\mathbf{w}}^- = \text{ReLU}(-\tilde{\mathbf{w}})$. The relaxed balance is

$$\tilde{B}(X; \mathbf{w}) = \frac{\sum_j \tilde{w}_j^+ \log X_j}{\sum_j \tilde{w}_j^+} - \frac{\sum_j \tilde{w}_j^- \log X_j}{\sum_j \tilde{w}_j^-}$$

Hard thresholding

$$\hat{I}_+ = \{j : \tilde{w}_j^+ > \tau\}, \quad \hat{I}_- = \{j : \tilde{w}_j^- < -\tau\}$$

⁵Gordon-Rodriguez et al., Bioinformatics, 22'

CoDaCoRe⁵ uses continuous relaxation to find the best subsets. For a vector of assignment weights \mathbf{w} , let

$$\tilde{\mathbf{w}} = \frac{2}{1 + \exp(-\mathbf{w})} - 1.$$

Let $\tilde{\mathbf{w}}^+ = \text{ReLU}(\tilde{\mathbf{w}})$ and $\tilde{\mathbf{w}}^- = \text{ReLU}(-\tilde{\mathbf{w}})$. The relaxed balance is

$$\tilde{B}(\mathbf{X}; \mathbf{w}) = \frac{\sum_j \tilde{w}_j^+ \log X_j}{\sum_j \tilde{w}_j^+} - \frac{\sum_j \tilde{w}_j^- \log X_j}{\sum_j \tilde{w}_j^-}$$

Hard thresholding

$$\hat{I}_+ = \{j : \tilde{w}_j^+ > \tau\}, \quad \hat{I}_- = \{j : \tilde{w}_j^- < -\tau\}$$

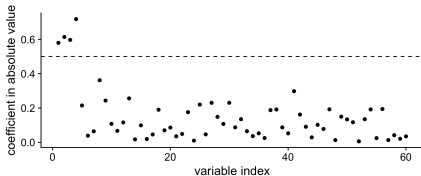
CoDaCoRe is efficient, **but tends to select too many variables.**

⁵Gordon-Rodriguez et al., Bioinformatics, 22'

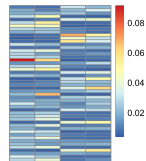
Our Framework: Supervised Log-Ratios

Input: (\mathbf{x}_i, y_i) for $i = 1, \dots, n$.

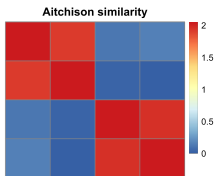
Step 1: Screen



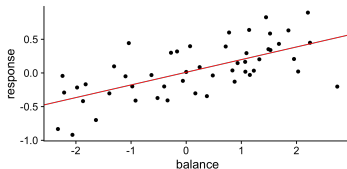
Step 2: Reduce



Step 3: Cluster



Step 4: Predict



Output: two subsets I_+ , I_- of variables for defining the balance.

Step 1: Screen

Let \mathbf{z}_i denote the clr-transformed version of \mathbf{x}_i , where

$$\mathbf{Z} = \left(\log \frac{X_1}{g(\mathbf{X})}, \dots, \log \frac{X_p}{g(\mathbf{X})} \right)^\top$$

Let $\mathbf{z}^{(j)}$ denote the vector of observations from the j -th variable. Variables are screened by thresholding their univariate effect on \mathbf{y} :

$$\left| \frac{(\mathbf{y} - \bar{\mathbf{y}})^\top (\mathbf{z}^{(j)} - \bar{\mathbf{z}}^{(j)})}{\|\mathbf{z}^{(j)} - \bar{\mathbf{z}}^{(j)}\|^2} \right| > \tau$$

The threshold τ is chosen by cross-validation.

Step 3: Cluster

Let C_τ be the collection of indices containing selected variables.

The Aitchison variation on the reduced data matrix is defined as

$$\hat{A}(\tau)_{j,k} = \frac{1}{n} \sum_{i=1}^n \left(\log \frac{x_{i,j}}{x_{i,k}} - \frac{1}{n} \sum_{i'=1}^n \log \frac{x_{i',j}}{x_{i',k}} \right)^2, \quad j, k \in C_\tau.$$

The Aitchison similarity is

$$\hat{S}(\tau)_{j,k} = \max_{j',k'} \left\{ \hat{A}(\tau)_{j',k'} \right\} - \hat{A}(\tau)_{j,k}, \quad j, k \in C_\tau.$$

Clustering returns two subsets of variables for defining the balance.

$$\log \frac{X_j}{X_p} = \alpha_{0,j} + \alpha_{1,j}U + \epsilon_j, \quad j = \{1, \dots, p\} \setminus \{p\} \quad (1)$$

$$y = \beta_0 + \beta_1 U + \epsilon, \quad (2)$$

where for $c_1, c_2 > 0$ the coefficients $\alpha_{1,j}$ satisfy

$$\alpha_{1,j} = 0, \quad j \notin I_+ \cup I_-,$$

$$\alpha_{1,j} = c_1, \quad j \in I_+,$$

$$\alpha_{1,j} = -c_2, \quad j \in I_-,$$

$$\sum_{j=1}^p \alpha_{1,j} = 0.$$

Here p is an inactive variable that belongs to $I_0 = \{1, \dots, p\} \setminus \{I_+ \cup I_-\}$.

$$B(X; I_+, I_-) = \tilde{\alpha}_0 + (c_1 + c_2)U + \tilde{\epsilon},$$

where

$$\tilde{\alpha}_0 = \frac{1}{|I_+|} \sum_{j \in I_+} \alpha_{0,j} - \frac{1}{|I_-|} \sum_{j \in I_-} \alpha_{0,j}, \quad \tilde{\epsilon} = \frac{1}{|I_+|} \sum_{j \in I_+} \epsilon_j - \frac{1}{|I_-|} \sum_{j \in I_-} \epsilon_j.$$

The response y is also linear in $B(X; I_+, I_-)$

$$y = \beta_0 - \tilde{\alpha}_0 \frac{\beta_1}{c_1 + c_2} + \frac{\beta_1}{c_1 + c_2} B(X; I_+, I_-) + \epsilon - \frac{\beta_1}{c_1 + c_2} \tilde{\epsilon}.$$

Let $Z_j = \log(X_j) - \log g(X)$ denote the clr-transformed data. Then

$$Z_j - \mathbb{E}[Z_j] = \alpha_{1,j}U + \frac{1}{p} \sum_{k=1}^p (\epsilon_j - \epsilon_k)$$

⇒ univariate regression can distinguish active from inactive variables

Let $Z_j = \log(X_j) - \log g(X)$ denote the clr-transformed data. Then

$$Z_j - \mathbb{E}[Z_j] = \alpha_{1,j}U + \frac{1}{p} \sum_{k=1}^p (\epsilon_j - \epsilon_k)$$

⇒ univariate regression can distinguish active from inactive variables

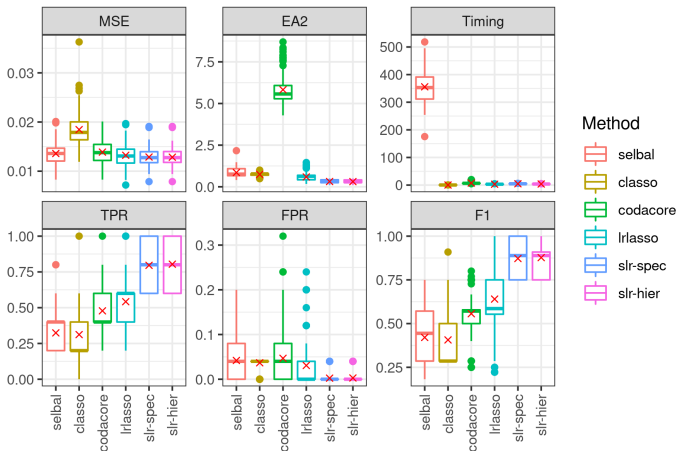
Aitchison Variation

$$\text{Var}\left(\log \frac{X_j}{X_k}\right) = \begin{cases} 2\sigma_\epsilon^2 & j \in I_+, k \in I_+ \\ (c_1 + c_2)^2 \sigma_U^2 + 2\sigma_\epsilon^2 & j \in I_+, k \in I_- \\ 2\sigma_\epsilon^2 & j \in I_-, k \in I_- \end{cases}$$

⇒ clustering can distinguish variables in I_+ from those in I_-

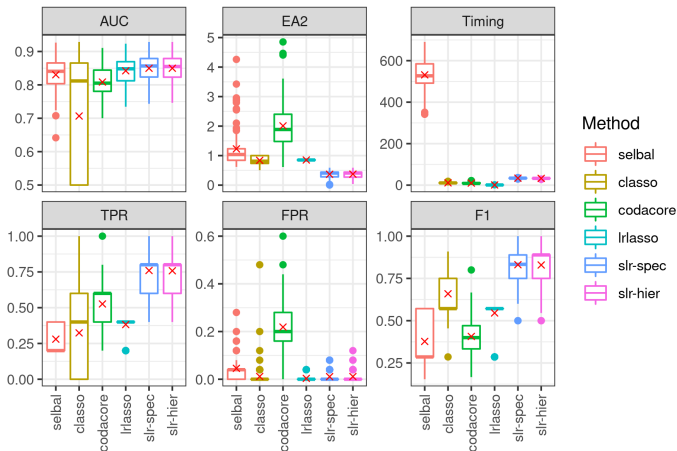
Simulation with Continuous Response

$n = 100, p = 30; I_+ = \{1, 2, 3, 4\}, I_- = \{5\}$



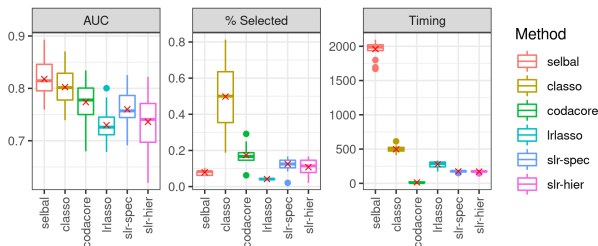
Simulation with Binary Response

$n = 100, p = 30; I_+ = \{1, 2, 3, 4\}, I_- = \{5\}$



Classification of Crohn's Disease

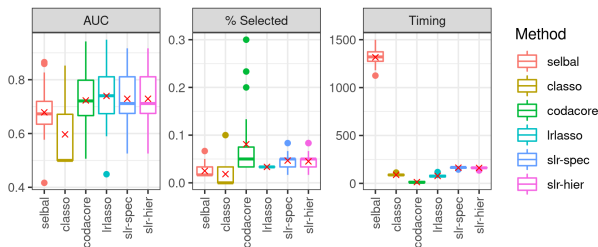
$n = 975$; $p = 48$ genera; y is binary with 662 cases



- ▶ Selbal is the most accurate and also the most time consuming.
- ▶ classo does well in AUC, but returns a non-sparse model.
- ▶ SLR with spectral clustering and CoDaCoRe are comparable.

Classification of HIV Status

$n = 155$; $p = 60$ genera; y is binary with 128 cases



- ▶ SLR selects a sparser model than CoDaCoRe.
- ▶ selbal is the most time consuming.
- ▶ classo do not perform well. l1lasso is the most sparse.

Microbiome and sCD14 Inflammation

$n = 151$; $p = 60$ genera; y is continuous

Taxa	selbal	codacore-1	llasso-1	llasso-2	slr-spec	slr-hier
"g_Faecalibacterium"					+	+
"f_Ruminococcaceae_g_unclassified"					+	+
"g_Subdoligranulum"	+	+	+		+	+
"g_Thalassospira"	+	+			+	+
"f_Defluviitaleaceae_g_Incertae_Sedis"		+			+	+
"f_Lachnospiraceae_g_Incertae_Sedis"	+			+		+
"g_Dorea"	+					
"g_Dialister"		+				
"f_Lachnospiraceae_g_unclassified"	-	-	-		+	+
"g_Catenibacterium"		-			-	-
"g_Mitsuokella"		-			-	-
"g_Bifidobacterium"	-	-			-	-
"g_Collinsella"	-	-		-	-	-
"g_Lachnospira"	-	-				
"k_Bacteria_g_unclassified"		-				
"g_Ruminococcus"		-				
"g_Megasphaera"		-				
"g_Sutterella"		-				
"o_Clostridiales_g_unclassified"		-				

- ▶ Supervised log-ratio can efficiently predict health outcomes from compositional data.
- ▶ SLR leads to interpretable biomarker selection.
- ▶ SLR can be extended to semi-supervised settings.

- ▶ Supervised log-ratio can efficiently predict health outcomes from compositional data.
- ▶ SLR leads to interpretable biomarker selection.
- ▶ SLR can be extended to semi-supervised settings.
- ▶ SLR requires proper zero handling.

- ▶ Supervised log-ratio can efficiently predict health outcomes from compositional data.
- ▶ SLR leads to interpretable biomarker selection.
- ▶ SLR can be extended to semi-supervised settings.
- ▶ SLR requires proper zero handling.
- ▶ Selection of more than one balance?

Acknowledgment



Kristyn Pantoja @TAMU



David Jones @Google

Thank You!

<https://drjingma.com>