

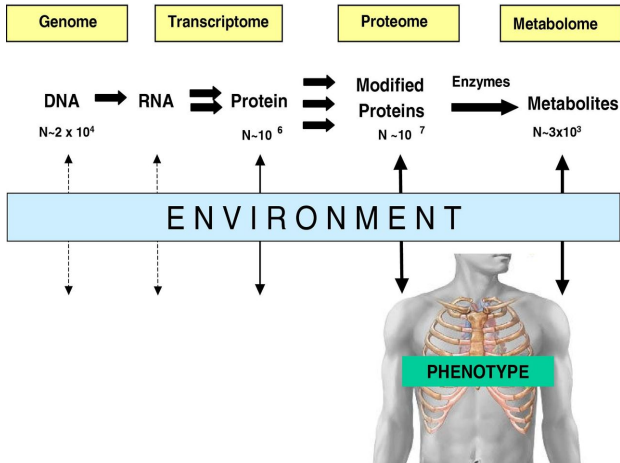
Bayesian Canonical Correlation Analysis for Multiple Groups

Lin Qiu, Lin Lin, Vernon M. Chinchilli

Department of Public Health Sciences and Statistics
Pennsylvania State University

July 29, 2019

Background



(Newgard, 2017)

Omics data

- **GWAS:** Genome-wide association studies (GWAS) have thus far explained only a small fraction of polygenic diseases like, diabetes and cardiovascular diseases.
- **Interest:**
genomics,epigenomics,transcriptomics,proteomics,metabolomics and microbiomics, each of these data types provides a different snapshot of the underlying biological system, and combining multiple data types has shown to be very valuable for the improvement of medical and health care services and for the implementation of preventive and precision medicine (Pan et al., 2010; Chen et al.,2013);
- **Aims:** Integrative statistical and computational analysis tools are needed;

Omics data

Issues to deal with:

- small sample size and large dimension (small n and large p);
- data types are different (continuous vs count);
- data structures are sparse;

Outline for section 1

1 Introduction

- Canonical Correlation Analysis
- Probabilistic CCA
- Problem Statement

2 Model Specification

3 Theoretical Results

4 Empirical Studies

- Simulation Studies
- Real Application

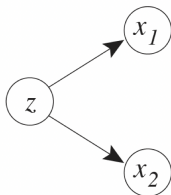
Canonical Correlation Analysis

We study the task by first modeling dependencies between two data sets of co-occurring or paired samples (x,y) . The underlying assumption is that variation within either data set alone is more noisy, or at least less interesting than variation that is in common. This task has been classically solved by Canonical Correlation Analysis (CCA) (Hotelling, 1936).

Suppose we have two random vectors $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)^T$ and $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_q)^T$. CCA aims to find two projection directions, $\mathbf{u} \in \mathbb{R}^p$ and $\mathbf{v} \in \mathbb{R}^q$ so that

- $\rho = \arg \max_{\mathbf{u}, \mathbf{v}} \text{Corr}(\mathbf{u}^T \mathbf{x}, \mathbf{v}^T \mathbf{y}) = \arg \max_{\mathbf{u}, \mathbf{v}} \frac{\mathbf{u}^T \Sigma_{xy} \mathbf{v}}{\sqrt{(\mathbf{u}^T \Sigma_{xx} \mathbf{u})(\mathbf{v}^T \Sigma_{yy} \mathbf{v})}}$, where $\Sigma_{xx}, \Sigma_{yy}, \Sigma_{xy}$ are covariance and cross-covariance matrices.
- This maximization is equivalent to $\max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \Sigma_{xy} \mathbf{v}$ subject to $\mathbf{u}^T \Sigma_{xx} \mathbf{u} = 1$ and $\mathbf{v}^T \Sigma_{yy} \mathbf{v} = 1$.
- The new variables $\eta_1 = \mathbf{u}_1^T \mathbf{x}, \xi_1 = \mathbf{v}_1^T \mathbf{y}$ are called the first pair of canonical variables or latent variables and $\rho_1 = \text{Corr}(\eta_1, \xi_1)$ is the first canonical correlation.

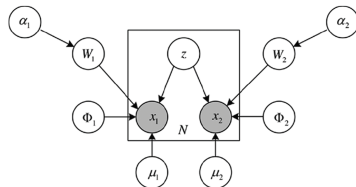
Probabilistic Interpretation of CCA



$$x_1 = W_1 z + \mu_1 + \varepsilon_1$$

$$x_2 = W_2 z + \mu_2 + \varepsilon_2$$

(a) 1a:PCCA



(b) 1b:BCCA

- 1a treats CCA as a generative model (Bach and Jordan, 2005).
- 1b shows the hierarchical Bayesian CCA model (Wang, 2007).

Inter-battery factor analysis(IBFA)

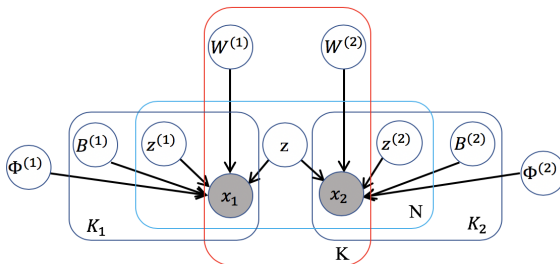


Figure: 1c:IBFA

$$\mathbf{x}^{(m)} \sim N(\mathbf{W}^{(m)}\mathbf{z} + \mathbf{B}^{(m)}\mathbf{z}^{(m)}, \boldsymbol{\Sigma}^{(m)}).$$

- 1c, IBFA model has view-specific latent variable (Browne,1979;Klami and Kaski,2006,2008).

Literature review

Automatic relevance determination(ARD;Neal,1996) prior

ARD is a Normal-Gamma prior for the projection weights.

$$ARD(\mathbf{W}^{(m)}|\alpha_0, \beta_0) = \prod_{k=1}^K p(\mathbf{w}_k^{(m)}|\alpha_k^{(m)})p(\alpha_k^{(m)}|\alpha_0, \beta_0),$$

$$\alpha_k^{(m)} \sim \text{Gamma}(\alpha_0, \beta_0),$$

$$\mathbf{W}_k^{(m)} \sim N(\mathbf{0}, (\alpha_k^{(m)})^{-1} \mathbf{I}).$$

The number of components is automatically selected by pushing $\alpha_k^{(m)}$ of unnecessary towards infinity (Klami and Kaski, 2007;Nagarajan,2008).

Group-wise Sparsity(Klami et al.,2013)

$\mathbf{y} = [\mathbf{z}; \mathbf{z}^{(1)}; \mathbf{z}^{(2)}] \in \mathbb{R}^{K_c \times 1}$, where $K_c = K + K_1 + K_2$,

IBFA model becomes to

$$\mathbf{y} \sim N(\mathbf{0}, \mathbf{I}), \mathbf{x} \sim N(\mathbf{W}\mathbf{y}, \mathbf{\Sigma}), \mathbf{\Sigma} \in \mathbb{R}^{D \times D}, D = D_1 + D_2, \mathbf{W} \in \mathbb{R}^{D \times K_c}.$$

Group Bayesian CCA(GBCCA)

A group may correspond to one view of the same set of objects, one of many data sets tied by co-occurrence, or a set of alternative variables collected from statistics table to measure one property of interest.

- **Aims:** Find components that capture joint variability between data sets instead of individual variables. Given a collection of $\mathbf{X}_1, \dots, \mathbf{X}_M$ of M data sets of dimensionalities D_1, \dots, D_M , the task is to find $K < \sum_{m=1}^M D_m$ components that describe the dependencies between data sets \mathbf{X}_m .

Outline for section 2

1 Introduction

- Canonical Correlation Analysis
- Probabilistic CCA
- Problem Statement

2 Model Specification

3 Theoretical Results

4 Empirical Studies

- Simulation Studies
- Real Application

Model: GBCCA, **DGBCCA**(Dynamic GBCCA)

Assume a collection of observations $y_i \in \mathbb{R}^D$ for $i = 1, \dots, N$, $\mathbf{Y} \in \mathbb{R}^{N \times D}$, and the D variables are split into M groups or subsets. For notational simplicity, assume the first D_1 variables correspond to the first group G_1 , D_M variables to G_M . Let $\mathbf{Y} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}]$, $\mathbf{x}^{(1)} \in \mathbb{R}^{N \times D_1}$ and $\mathbf{x}^{(M)} \in \mathbb{R}^{N \times D_M}$ for the N observations, $\mathbf{W} \in \mathbb{R}^{D \times K}$, $\mathbf{Z} \in \mathbb{R}^{N \times K}$. The model for the m^{th} group of the i^{th} sample is

$$\mathbf{x}_i^{(m)} \sim N(\mathbf{Z}_i \mathbf{W}^{(m)T}, \boldsymbol{\tau}_m^{-1} \mathbf{I}),$$

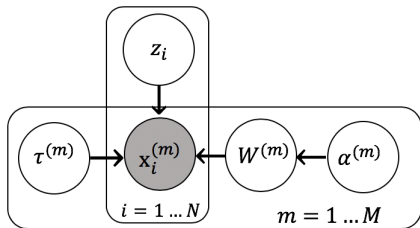
$$\mathbf{W}_{:,k}^{(m)T} \sim N(\mathbf{0}, \boldsymbol{\alpha}_k^{(m)-1} \mathbf{I}),$$

$$\boldsymbol{\alpha}_k^{(m)} \sim \text{Gamma}(\alpha_0, \beta_0),$$

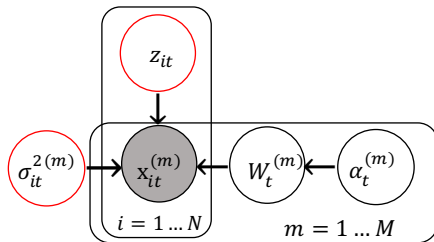
$$\boldsymbol{\tau}_m \sim \text{Gamma}(\alpha_0, \beta_0),$$

$$\mathbf{Z}_i \sim N(\mathbf{0}, \mathbf{I}).$$

GBCCA,DGBCCA model



(a) GBCCA



(b) DGBCCA

Outline for section 3

1 Introduction

- Canonical Correlation Analysis
- Probabilistic CCA
- Problem Statement

2 Model Specification

3 Theoretical Results

4 Empirical Studies

- Simulation Studies
- Real Application

GBCCA model: Inference

Full joint likelihood

$$\begin{aligned} p(\mathbf{X}, \mathbf{W}, \tau, \mathbf{Z}) &= p(\mathbf{W})p(\tau)p(\mathbf{X}|\mathbf{Z}, \mathbf{W}, \tau) \\ &= p(\mathbf{W}) \prod_{m=1}^M \text{Gamma}(\tau_m | \alpha_\tau, \beta_\tau) \prod_{i=1}^N (N(\mathbf{z}_i | 0, \mathbf{I}) N(\mathbf{x}_i^{(m)} | \mathbf{z}_i \mathbf{W}^{(m)}, \tau_m)) \end{aligned}$$

We assume the weight matrix \mathbf{W} is made sparse by a group-wise ARD prior,

$$p(\mathbf{W} | \alpha) = \prod_{m=1}^M \prod_{k=1}^K \prod_{d=1}^{D_m} N(w_{k,d}^{(m)} | 0, \alpha_k^{(m)-1})$$

The ARD makes groups of variables inactive for specific components by driving $\alpha_{m,k}^{-1}$ to zero.

GBCCA model: Variational updates

The variational distribution $Q(\boldsymbol{\Theta})$ is assumed to factorize over the component variables θ_i in $\boldsymbol{\theta}$,

$$Q(\boldsymbol{\theta}) = \prod_{i=1} q_i(\theta_i)$$

In order to minimize the Kullback-Leibler (KL) divergence between the approximating distribution $Q(\boldsymbol{\theta})$ and the true posterior $P(\boldsymbol{\theta}|D)$. The best distribution q_j^* for each of the factors q_j can be expressed as

$$q_j^*(\theta_j|D) = \frac{e^{E_{i \neq j}[\ln p(\boldsymbol{\theta}, D)]}}{\int e^{E_{i \neq j}[\ln p(\boldsymbol{\theta}, D)]} d\theta_j}$$

$$\ln q_j^*(\theta_j|D) = E_{i \neq j}[\ln p(\boldsymbol{\theta}, D)] + \text{constant}$$

(Bishop, 1999)

GBCCA model: Variational updates

latent variables \mathbf{Z}

$$q(\mathbf{Z}) = \prod_{i=1}^N N(z_i | \mu_i^{(z)}, \Sigma^{(z)}),$$

$$\mu_i^{(z)} = \sum_{m=1}^M \Sigma^{(z)} \langle W^{(m)} \rangle \langle \tau_{(m)} \rangle x_i^{(m)}, \Sigma^{(z)} = (I_k + \sum_{m=1}^M \langle \tau_m \rangle \langle W^{(m)} W^{(m)T} \rangle)^{-1}.$$

projection matrices \mathbf{W}

$$q(\mathbf{W}) = \prod_{m=1}^M \prod_{j=1}^{D_m} N(w_{:,j}^{(m)} | \mu_{m,j}^{(w)}, \Sigma_m^{(w)}),$$

$$\mu_{m,j}^{(w)} = \Sigma_m^{(w)} \langle \tau_m \rangle (\sum_{i=1}^N x_{ij}^{(m)}) \langle \mathbf{Z}_i \rangle,$$

$$\Sigma_m^{(w)} = (\langle \tau_m \rangle \sum_{i=1}^N \langle \mathbf{Z}_i \mathbf{Z}_i^T \rangle + \langle \alpha_m \rangle^{-1})^{-1}.$$

GBCCA model: Variational updates

ARD parameters

$$q(\alpha) = \prod_{m=1}^M \prod_{k=1}^K \text{Gamma}(\alpha_m^\alpha, \beta_{mk}^\alpha),$$

$$\alpha_m^\alpha = \alpha^\alpha + D_m/2, \beta_{mk}^\alpha = \beta^\beta + w_k^{(m)T} w_k^{(m)} / 2.$$

noise precision parameters

$$q(\tau) = \prod_{m=1}^M \text{Gamma}(\tau_m | \alpha_m^\tau, \beta_m^\tau)$$

$$\alpha_m^\tau = \alpha^\tau + D_m N / 2, \beta_m^\tau = \beta^\tau + 1/2 \langle (x_i^{(m)} - Z_i \mathbf{W}^{(m)T})^2 \rangle.$$

GBCCA model: predictive inference

$$p(\mathbf{X}^{(m)} | \mathbf{Y}^{-(m)})$$

The task is to predict unobserved groups based on observed data. Given $\mathbf{Y}^{-(m)}$, $q(\mathbf{W})$ and $q(\tau)$, we can approximate the posterior distribution for the latent variables $q(\mathbf{Z})$ corresponding to $\mathbf{Y}^{-(m)}$ and approximate the mean of the predictive distribution as

$$\begin{aligned} \langle \mathbf{X}^{(m)} | \mathbf{Y}^{-(m)} \rangle &= \langle \mathbf{Z} \mathbf{W}^{(m)} \rangle_{q(\mathbf{W}^{(m)}), q(\mathbf{Z})} \\ &= \mathbf{Y}^{-(m)} \boldsymbol{\tau} \langle \mathbf{W}^{(-m)T} \rangle \boldsymbol{\Sigma}^{-1} \langle \mathbf{W}^{(m)} \rangle \end{aligned}$$

$$\boldsymbol{\tau} = \text{diag}(\langle \tau_j \rangle \mathbf{I}_{D_{j \neq m}}), \boldsymbol{\Sigma} = \mathbf{I}_K + \sum_{j \neq m} \langle \tau_j \rangle \langle \mathbf{W}^{(j)} \mathbf{W}^{(j)T} \rangle.$$

Outline for section 4

1 Introduction

- Canonical Correlation Analysis
- Probabilistic CCA
- Problem Statement

2 Model Specification

3 Theoretical Results

4 Empirical Studies

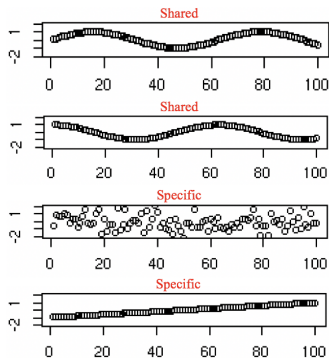
- Simulation Studies
- Real Application

Strategy: $M=2$, GBCCA \rightarrow BIBFA

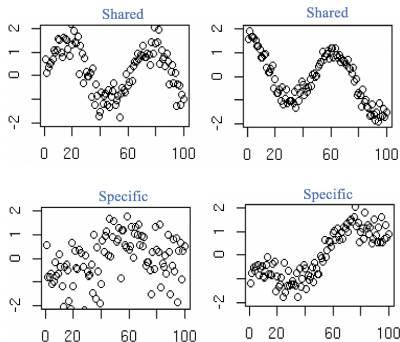
We conducted simulations on the BIBFA model to show that the inference procedure converges to the correct solution.

- We generated two collections of data sets, $\{\mathbf{Y}^1, \mathbf{Y}^2\}$, with $N = 100$ samples and dimensions $(50, 40)$.
- The parameters are set to contain all types of components (view-specific and shared components). $K = 4$ true components, $\alpha = 1, \tau_1 = 3$ and $\tau_2 = 6$.
- We fit the BIBFA model with $K = 6$, other parameters of this model are: $z_n \sim N(\mathbf{0}, \mathbf{1})$; initial $\tau = 1000$; $\alpha_k^{(m)} \sim \text{Gamma}(10^{-14}, 10^{-14})$; $\tau_m \sim \text{Gamma}(10^{-14}, 10^{-14})$.
- The maximum number of iterations is 10000 to train the model parameters. We use the convergence as a relative change of $L(Q)$ falling below 10^{-6} .

M=2, GBCCA \rightarrow BIBFA



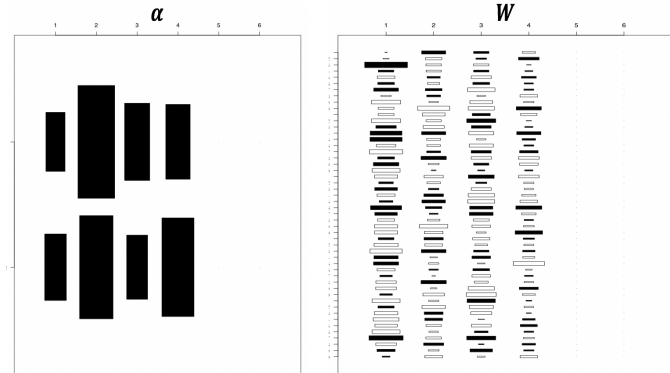
(a) 1a: True latent components z



(b) 1b: Learned latent components z

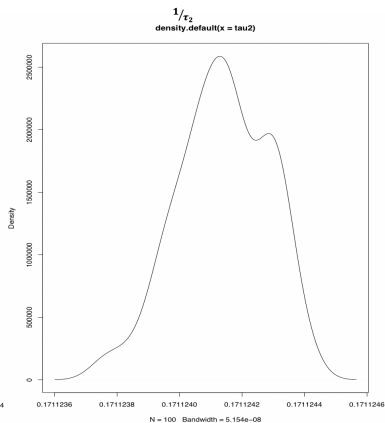
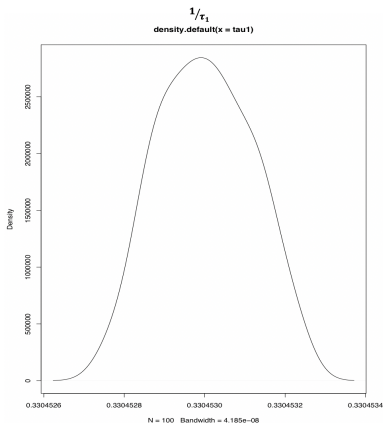
BIBFA learns the correct latent components and discards the access ones

- α is the variance of components, W denotes the elements



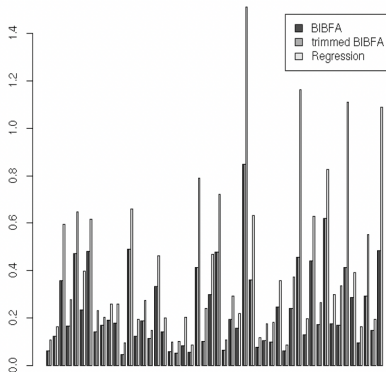
BIBFA finds the correct posterior distributions for model parameters $q(\tau_m)$

- The true $1/\tau_1 = 0.33$ and $1/\tau_2 = 0.17$

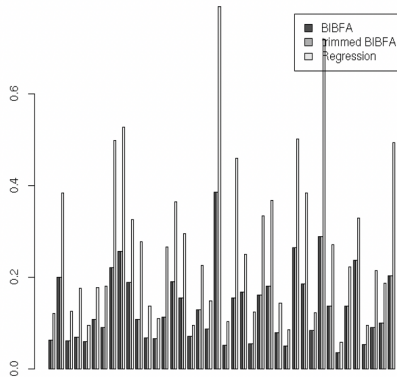


BIBFA prediction

Prediction errors for the features of data set 1



Prediction errors for the features of data set 2

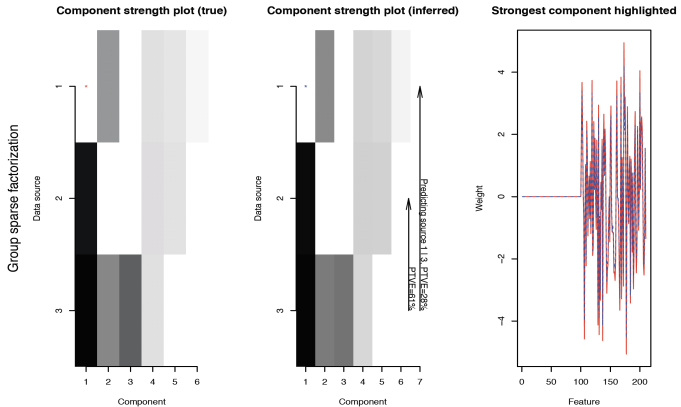


$$rmse = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2}{\frac{1}{n} \sum_{i=1}^n y_i^2}$$

Strategy: $M > 2$, GBCCA

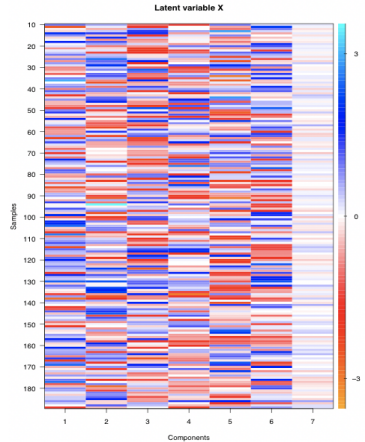
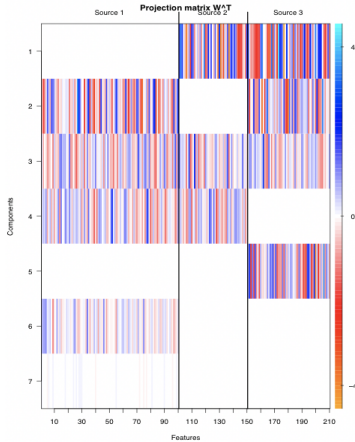
- We generated three collections of datasets, $\{\mathbf{Y}^1, \mathbf{Y}^2, \mathbf{Y}^3\}$, with $N = 200$ samples and dimensions $(100, 50, 60)$.
- $K = 6$, Z and W were drawn from $N(0, 1)$, $\tau = 1$.
- the first 20 samples as the test data, the model was trained under an initial $\tau = 1000$, $\alpha_\tau = 10$, $\beta_\tau = 10$, $\alpha_0 = 10$, $\beta_0 = 1$.
- The total number of Gibbs Sampling steps is 5000, burn-in samples are 2500, and $K = 15$.

GBCCA infers the correct component affiliations

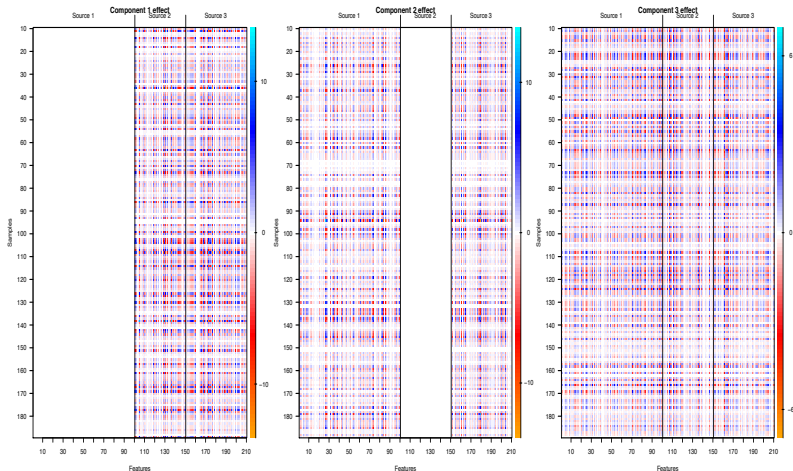


$$ptve = \frac{1 - \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2}{\frac{1}{n} \sum_{i=1}^n y_i^2} * 100\%$$

GBCCA: projection matrix and latent variable effect



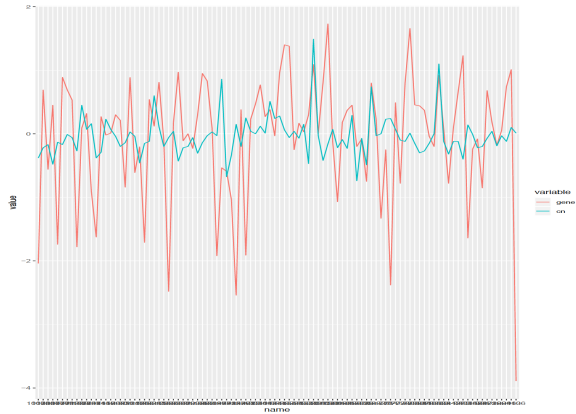
GBCCA: component effects for each feature $Y=ZW$



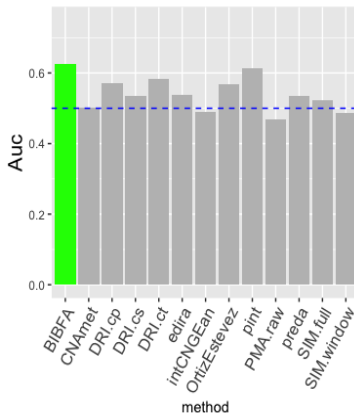
BIBFA:Pollack data set(Lahti et al.,2012)

- This data set contains 41 breast cancer samples, each sample has 4296 genes' expression value and their corresponding copy number value. The data dimension is $Y^{(m)} = 41 \times 4296$.
- There are 38 known cancer genes.
- Gene expression and copy number data can be integrated to identify DNA copy number alternations which induce the changes of expression levels of associated genes (Huang et al.,2011)
- BIBFA model was run by different K values ranging from 4 to 60 and we chose the best K according to the variational lower bound,resulting in K=41.
- Measure the $s_g = \sum_{k=1}^K |\langle W_{g,k}^{(1)} \rangle \langle W_{g,k}^{(2)} \rangle|$ for each gene to rank the genes (Klami et al.,2013), and we repeated the experimental setup to obtain results comparable with Lahti et al(2012).

Pollack data set



BIBFA(M=2,GBCCA)



Summary

- The GBCCA model can handle high dimension problems, while the previous BCCA (Wang, 2007) can only deal with dimension less than 10.
- The GBCCA model can accurately capture the true components for multiple groups, inferring the group dependencies.
- Under $M = 2$, we showed GBCCA outperformed by supervised regression models in predictive tasks.

Thank you! **Email:**luq7@psu.edu **Wechat:**lqconnor