JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

#### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULT

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨー のへで

Conclusion and Future work

# CLUSTERING WITH MULTIPLE DISTANCE METRICS - MULTI-LEVEL MIXTURE MODELS WITH PROFILE

TRANSFORMATIONS.

### Rebecka Jörnsten

### Department of Mathematical Statistics Chalmers and Gothenburg University

jornsten@chalmers.se, http://www.stat.rutgers.edu/~rebecka

September 16, 2009

## 1 CLUSTER ANALYSIS

**2** WITHIN-CLUSTER PARAMETERIZATIONS

BETWEEN-CLUSTER PARAMETERIZATIONS: MULTI-LEVEL MIXTURE MODELING





JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### OUTLINE

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

### GENERAL PROBLEM

- Popular approach for dimension reduction
- Wide range of applications: engineering, geological data, social networks, high-throughput biology

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

Conclusion and Future work

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

# THINKING ABOUT THE PROBLEMS IN TERMS OF MODEL SELECTION

### CLUSTERING

- How many clusters?
- Can we find an efficient <u>within-cluster</u> parameterization? (E.g., are cluster profiles flat, or increasing?)
- Can we find an efficient <u>between-cluster</u> parameterization? (Do cluster profiles share a similar shape, or coincide for a subset of data dimensions?)

#### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

# CLUSTERING GENE EXPRESSION DATA

### TASKS

- To group genes by similarity of expression profiles across experimental conditions
- Assign functions to genes 'guilt by association'

### GOALS

- We need an objective interpretation of cluster profiles gene function inferences
- Avoid "Waste of parameters"

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULT:

### MATHEMATICAL FORMULATION

For each gene g we observe a feature vector  $\mathbf{x}_g$ :

$$\mathbf{x}_{g} \mid R_{g} = k \sim MVN(\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}),$$

where  $R_g$  is the cluster membership indicator.

- We seek efficient/interpretable parameterizations of the cluster profiles;  $\mu_k = W \theta_k$
- A sparse representation of μ<sub>k</sub> is obtained if we set some of θ<sub>k</sub> to 0.
- The design matrix W is chosen to represent a scientific question of interest.

Fitting the model: EM with a *modified M-step* incorporating the  $\theta_k$  constraints (see Jornsten and Keles, Biostatistics, 2008).

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

BETWEEN-CLUSTER PARAMETERIZA-TIONS: MULTI-LEVEL MIXTURE MODELING

RESULT

### MODEL SELECTION

Searching for an optimal sparse within-cluster representation:

- Classification EM (CEM) approach, see Jornsten and Keles, 2008.
- Simultaneous selection via rate-distortion theory, see Jornsten, JCGS, 2009.

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

# Multi-level Mixture models: Between-cluster parameterizations

- Cell-line data
- There seem to be more neuron-specific clusters than glia specific
- $\mathcal{MIX}_{\mathcal{L}}$ : Mixture models for multi-factor experiments

- mRNA expression time course after spinal cord injury
- Clusters seem to share a similar shape, and differ in terms of offset and scale.
- $\mathcal{MIX}_{\mathcal{T}}$ : Mixture models with profile transformations.





Jörnsten. Clustering with Multiple distance metrics

### OUTLINE

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

BETWEEN-CLUSTER PARAMETERIZA-TIONS: MULTI-LEVEL MIXTURE MODELING

RESULT

# $\mathcal{MIX}_{\mathcal{T}}$ - Model formulation

### Multi-level structure

- For each gene g we observe a feature vector  $\mathbf{x}_{g}$ .
- Indicators  $R_g = k$ ,  $U_g = l$ : cluster k and sub-cluster  $l \in \{1, \dots, L_k\}$ .

$$\mathbf{x}_{g} \mid R_{g} = k, U_{g} = l \sim MVN(A_{kl}\boldsymbol{\mu}_{k} + b_{kl}, A_{kl}\boldsymbol{\Sigma}_{k}A_{kl}')$$

- Within-cluster parameterization:  $\mu_k = W \theta_k$
- Between-cluster parameterization: affine transformation  $(A_{kl}, b_{kl}) \rightarrow$  still Gaussian mixture.

### CLUSTERING WITH MULTIPLE DISTANCE METRICS

- Here, consider sign-flip and/or scale transformations:  $\mu_{kl} = T_{kl}(\mu_k) = \beta_{kl}(\mu_k + \alpha_{kl}), \quad \mathbf{\Sigma}_{kl} = \beta_{kl}^2 \mathbf{\Sigma}_k$
- Equivalent to clustering with three distance metrics simultaneously: 1 |correlation|, 1 correlation, and euclidean distance!

うしん 同一人用 人用 人用 人口 マ

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

# Model fitting

- Fitting the model: EM with a *profile M-step* incorporating the constraints.
- Separate the tasks into
  - **(**) The *top-level* estimation problem:  $(\mu_k = W_k \theta_k, \boldsymbol{\Sigma}_k)$
  - One transformation parameter estimation problem: (α<sub>kl</sub>, β<sub>kl</sub>) (where β<sub>k1</sub> = 1, α<sub>k1</sub> = 0, ∀k).
- How can we do this? Well, given (α<sub>kl</sub>, β<sub>kl</sub>), maximizing the conditional likelihood translates to an application of the inverse profile transformation, and

$$\left(\frac{\mathbf{x}_g}{\beta_{kl}} - \alpha_{kl}\right) | R_g = k, U_g = l \sim MVN(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Solve the aggregate problem first.

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

## The aggregate problem

For each cluster k, and its sub-clusters  $l = 1, \cdots, L_k$ , compute

$$T_{kl}^{-1}(\mathbf{x}_{g}) = \frac{\mathbf{x}_{g}}{\beta_{kl}} - \alpha_{kl}$$

Update

$$\boldsymbol{\Sigma}_{k} = \frac{\sum_{g,l} \eta_{gkl} (T_{kl}^{-1}(\mathbf{x}_{g}) - \boldsymbol{\mu}_{k}) (T_{kl}^{-1}(\mathbf{x}_{g}) - \boldsymbol{\mu}_{k})'}{\sum_{g,l} \eta_{gkl}}$$

**2** Condition on  $\Sigma_k$ , solve for the constrained  $\theta_k$ :

$$\boldsymbol{\theta}_{k} = \frac{\sum_{g,l} \eta_{gkl} (W_{k}^{\prime} \boldsymbol{\Sigma}_{k}^{-1} W_{k})^{-1} W_{k}^{\prime} \boldsymbol{\Sigma}_{k}^{-1} (T_{kl}^{-1}(\mathbf{x}_{g}))}{\sum_{g,l} \eta_{gkl}}$$

Note that  $\mathbf{x}_g$  contributes to the top-level parameters in multiple ways, moderated by the sub-cluster membership probabilities.

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

JÖRNSTEN. Clustering with Multiple Distance metrics

#### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

### ESTIMATING THE TRANSFORM PARAMETERS

Condition on  $(\mu_k, \mathbf{\Sigma}_k)$ , and consider

$$\begin{split} \max_{\alpha_{kl},\beta_{kl}} \sum_{g} \eta_{gkl} - \frac{1}{2} \frac{(\mathbf{x}_{g} - \beta_{kl}\boldsymbol{\mu}_{k} - \alpha'_{kl}\mathbf{1}_{J})'\boldsymbol{\Sigma}_{k}^{-1}(\mathbf{x}_{g} - \beta_{kl}\boldsymbol{\mu}_{k} - \alpha'_{kl}\mathbf{1}_{J})}{\beta_{kl}^{2}} + \\ - \frac{J}{2}\log(\beta_{kl}^{2}) + c, \text{ where } \alpha'_{kl} = \beta_{kl}\alpha_{kl}, \ c \text{ is a constant.} \end{split}$$

Solutions:

$$\alpha_{kl}' = \frac{\sum_{g} \eta_{gkl} \mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{g}}{\sum_{g} \eta_{gkl} \mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{1}_{J}} - \beta_{kl} \frac{\mathbf{1}_{J}' \boldsymbol{\Sigma}_{k} \boldsymbol{\mu}_{k}}{\mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{1}_{J}}.$$

$$\boldsymbol{\Theta} \quad \beta_{kl}^{2} - A\beta_{kl} - B = 0, \text{ where}$$

$$A = \frac{1}{J \sum_{g} \eta_{gkl}} \left[ \frac{\mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\mu}_{k} \sum_{g} \eta_{gkl} \mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{g}}{\mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{1}_{J}} - \sum_{g} \eta_{gkl} \mathbf{x}_{g}' \boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\mu}_{k} \right]$$

$$B = \frac{1}{J \sum_{g} \eta_{gkl}} \left[ \sum_{g} \eta_{gkl} \mathbf{x}_{g}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{g} - \frac{(\sum_{g} \eta_{gkl} \mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{g})^{2}}{\sum_{g} \eta_{gkl} \mathbf{1}_{J}' \boldsymbol{\Sigma}_{k}^{-1} \mathbf{1}_{J}} \right]$$

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

#### Outlini

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

BETWEEN-CLUSTER PARAMETERIZA-TIONS: MULTI-LEVEL MIXTURE MODELING

RESULT

Conclusion and Future work

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

# Making it work!

Setting up the model is the easy bit...

- Starting values: The search space  $(K, \{L_k, k = 1, \cdots, K\})$  is huge ...
- We parameterize this as (M, K), where  $M = \sum_{k} L_{k}$  is the total number of clusters.

Initialize with M kmeans clusters

- Use PAM (robust kmeans) with 1 correlation or 1 – |correlation| to group the *M* centroids into *K* top-level clusters.
- 3 Add an additional random element by initializing  $\Sigma_k = |z| \times Cov(\tilde{\mathbf{x}}|R_g = k)$ , where  $z \sim N(r, s)$ .
- Search forward for *M*, backwards for *K* ∈ {*M*, · · · , 1} to minimize BIC.
- More details in the papers (regularization, parameter value initialization,...).

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

BIC as a function of the number of clusters, for the standard method (K), the sign-flip (F), and sign-flip+scale (B) models.



JÖRNSTEN. Clustering with Multiple Distance metrics

#### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULTS

Conclusion and Future work

We select 9 clusters with an efficient *between-cluster* parameterization, compared to 5 using the standard mixture model.

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・



The  $MIX_T$  with M = 5 has K = 2 profile clusters, with sub-clusters (1,2,3), where cluster 3 is a sign-flip cluster, and sub-clusters (4,5), both with positive  $\beta_{kl}$ .

JÖRNSTEN. Clustering with Multiple Distance metrics

#### Outlini

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULTS



JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULTS

Conclusion and Future work

Clusters 7,2,5,6: increased level of response - distinct functional gene classes in each cluster. Clusters 1,3,4: wave of early immune response Clusters 8,9: under-expression of neuron-specific genes, marker of neuron death after injury.

- $\mathcal{MIX}_{\mathcal{T}}$  allows for clustering via multiple distance metrics simultaneously
  - Sub-clusters: Mahalanobis distance
  - Top-level: 1 cor or 1 |cor|.



 With MIX<sub>T</sub> we can choose to investigate clusters at either level of the modeling hierarchy. JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

 $\mathbf{x}_g$  and  $\mathbf{y}_g$  the time-course expression of mild ( $\mathbf{x}$ ) and moderate ( $\mathbf{y}$ ) injury levels. We assume that

$$(\mathbf{x}_{g}, \mathbf{y}_{g})|R_{g} = k, U_{g} = l \sim$$
$$\sim N\left((\beta_{kl}(\boldsymbol{\mu}_{k}^{X} + \alpha_{kl}\mathbf{1}_{J}), \delta_{kl}(\boldsymbol{\mu}_{k}^{Y} + \gamma_{kl}\mathbf{1}_{J})), \boldsymbol{\Sigma}_{kl}\right),$$

where

$$\boldsymbol{\Sigma}_{kl} = \begin{bmatrix} \beta_{kl}^2 \boldsymbol{\Sigma}_k^X & \beta_{kl} \delta_{kl} \boldsymbol{\Sigma}_k^{XY} \\ \beta_{kl} \delta_{kl} \boldsymbol{\Sigma}_k^{YX} & \delta_{kl}^2 \boldsymbol{\Sigma}_k^{Y} \end{bmatrix}.$$

 $\mu_k^Y = \mu_k^X + \Delta_k = W_k^X \theta_k^X + \Delta_k$ , where  $\Delta_k$  is a vector of cluster contrast parameters between the two levels of injury. This parameterization provides a sparse representation of  $\mu_k$  if the cluster profiles  $\mu_k^X \simeq \mu_k^Y$ .

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

#### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level Mixture Modeling

RESULTS

- The joint maximization with respect to scale transformation parameters (α<sub>kl</sub>, β<sub>kl</sub>, γ<sub>kl</sub>, δ<sub>kl</sub>), given the profile shape parameters (μ<sub>k</sub>, Σ<sub>k</sub>), is complex.
- We approximate this optimization by de-coupling the problem in terms of the components (x, y), essentially assuming that Σ<sub>k</sub> is block-diagonal.

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULTS

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Cluster mean profiles. The left sub-panel: mild injury data. The right figure sub-panel: the moderate injury data.



JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

#### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULTS

Conclusion and Future work

The  $MIX_T$  model with M = 7 clusters has K = 4 profile clusters. No sign-flip clusters are detected.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ = 臣 = のへで

We detect several scale-transformation clusters ((1,2), (3,4), and (5,6)).



JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Betweencluster parameterizations: Multi-level mixture modeling

RESULTS

Conclusion and Future work

The sub-cluster separation for the mild injury level data,  $|\beta_{k1} - \beta_{k2}|$ , is smaller than the sub-cluster separation for the moderate injury level data,  $|\delta_{k1} - \delta_{k2}|$ .

### MULTI-LEVEL MIXTURE MODELS

- Between-cluster parameterizations provide an objective inference tool for comparing clusters...
- ... and constitutes a substantial savings in terms of the number of model parameters.
- We detect clusters that a standard model-based clustering method may miss.
- Multi-level mixture models allow for clustering using multiple distance metrics simultaneously.

 Papers are available at http://www.stat.rutgers.edu/~rebecka JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### OUTLINE

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

# FUTURE WORK

- More complex experiments need to automate subset selection as much as possible (e.g. consider different parameterizations simultaneously).
- Time course experiments: no multi-level structure is known a-priori. Define clusters by a unique path through a set of nodes. This allows clusters to coincide for a subset of data dimensions.



JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

- Prof. Ron Hart and Loyal Goff, W.M. Keck center for Collaborative Neuroscience, Rutgers
- Sunduz Keles, U. Wisconsin, Madison
- NSF, EPA

JÖRNSTEN. CLUSTERING WITH MULTIPLE DISTANCE METRICS

### Outline

CLUSTER ANALYSIS

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

Conclusion and Future work

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

JÖRNSTEN. Clustering with multiple distance metrics

#### Outline

Cluster Analysis

WITHIN-CLUSTER PARAMETERIZA-TIONS

Between-Cluster Parameterizations: Multi-level Mixture Modeling

RESULT

Conclusion and Future work

### ▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲□ ● ● ● ●