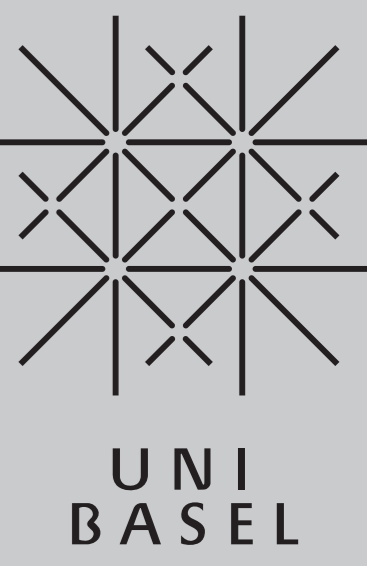


Copula Mixture Model for Dependency-seeking Clustering

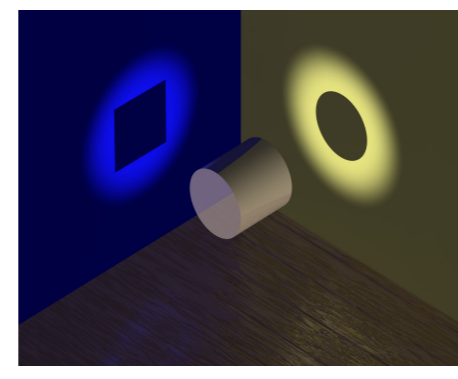
Melanie Rey, Volker Roth

Department of Mathematics and Computer Science, University of Basel, Switzerland

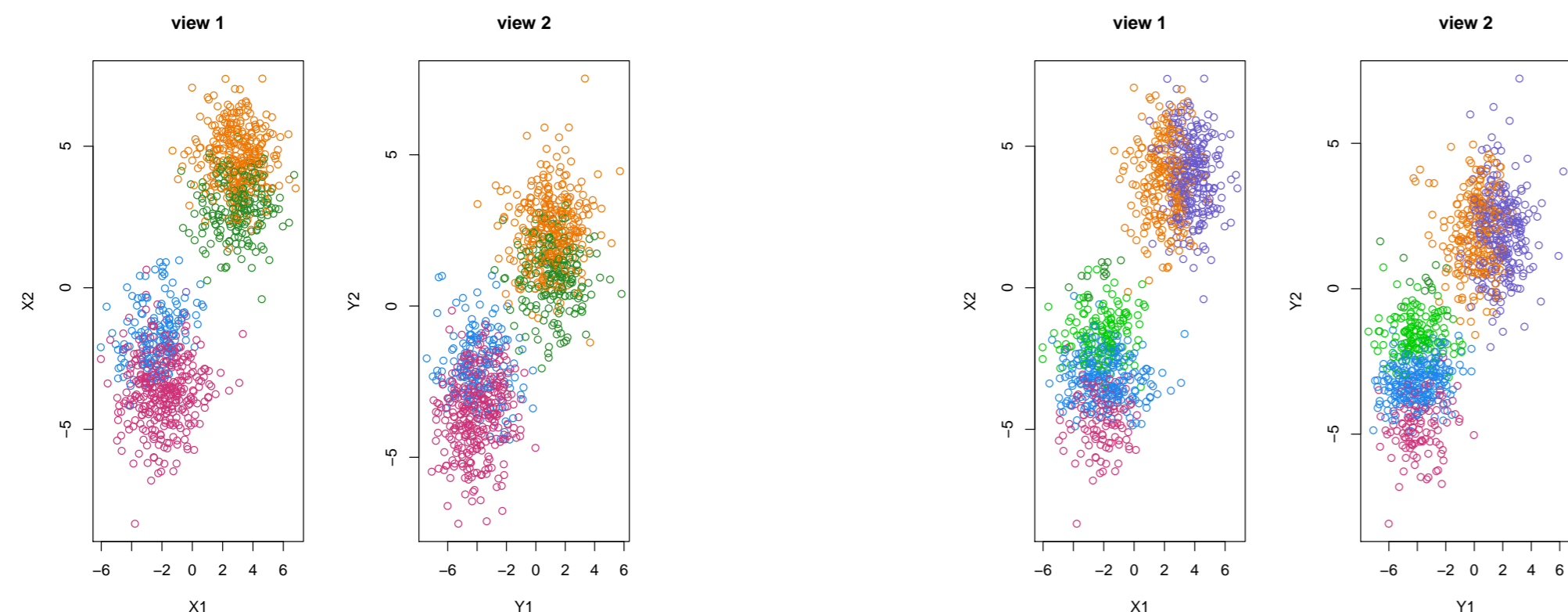


Dependency-seeking Clustering

- ▶ Clustering of co-occurring samples from different data sources called views.



- ▶ The aim is to cluster the points according to their between-views dependence structure.



- ▶ The probabilistic interpretation of CCA given by [Bach, 2005]:

$$\mathbf{Z} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I}_d),$$

$$(\mathbf{X}, \mathbf{Y}) | \mathbf{Z} \sim \mathcal{N}_{p+q}(\mathbf{WZ} + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

where $\boldsymbol{\mu} = (\boldsymbol{\mu}_x, \boldsymbol{\mu}_y) \in \mathbb{R}^{p+q}$, $\mathbf{W} = \begin{pmatrix} \mathbf{W}_x \\ \mathbf{W}_y \end{pmatrix} \in \mathbb{R}^{(p+q) \times d}$,

$\mathbf{1} \leq d \leq \min(p, q)$ and the covariance matrix $\boldsymbol{\Psi}$ has a block diagonal form:

$$\boldsymbol{\Psi} = \begin{pmatrix} \boldsymbol{\Psi}_x & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Psi}_y \end{pmatrix}.$$

- ▶ Probabilistic dependency-seeking clustering model derived in [Klami, 2006]:

$$\mathbf{Z} \sim \text{Mult}(\boldsymbol{\theta}),$$

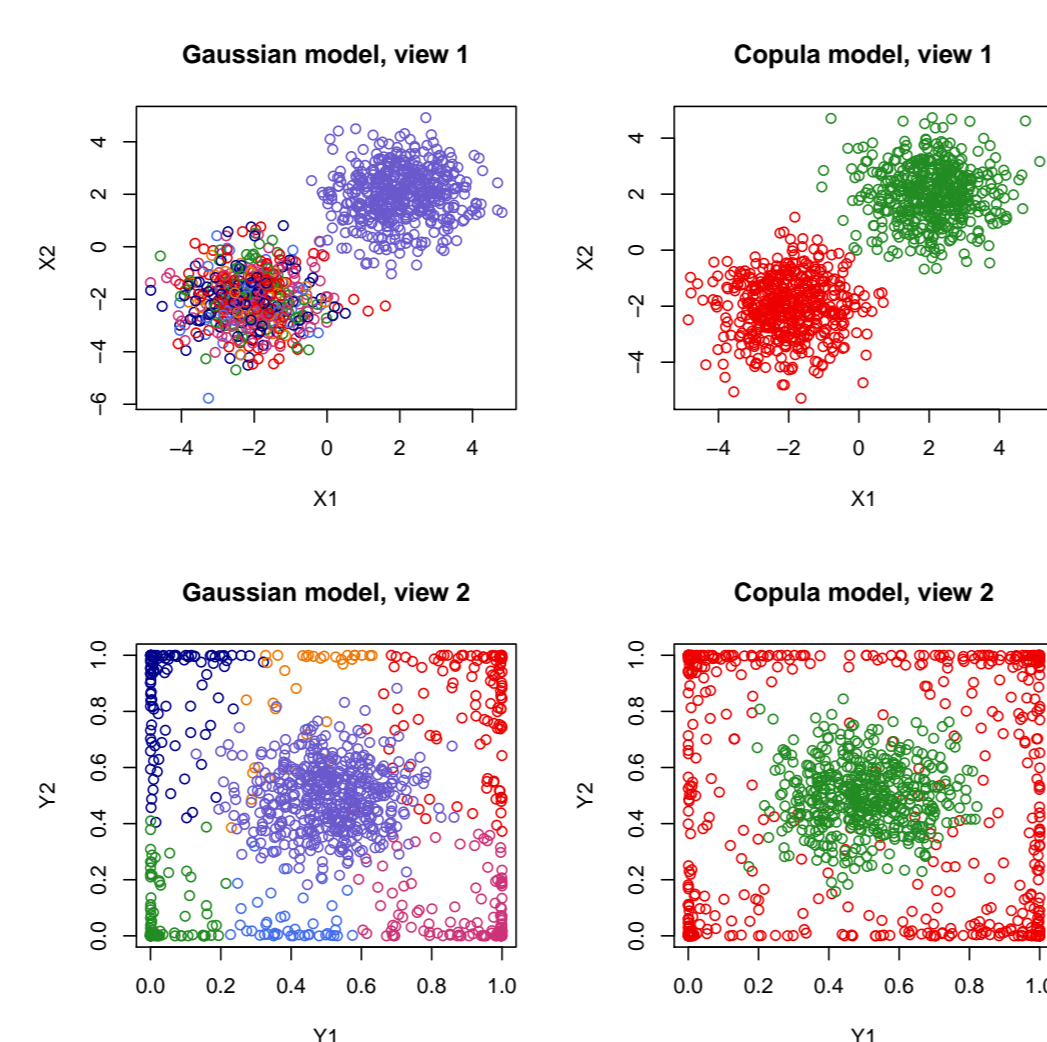
$$(\mathbf{X}, \mathbf{Y}) | \mathbf{Z} \sim \mathcal{N}_{p+q}(\boldsymbol{\mu}_z, \boldsymbol{\Psi}).$$

- ▶ The latent variable \mathbf{Z} plays now the role of the cluster assignment.
- ▶ $\boldsymbol{\Psi}$ still has a block structure.
- ▶ This special form implies that given the cluster assignment the two views are independent and thereby enforces the cluster structure to capture all the dependencies.

Motivation for Copula Mixture

- ▶ The above model assumes that \mathbf{X} and \mathbf{Y} are conditionally multivariate Gaussian.
- ▶ When applied to non-normally distributed data these models have to increase the number of clusters to achieve a reasonable fit.

- ▶ The components of these mixtures will not only be used to reflect differences in dependence structures but will also be used to approximate a non-Gaussian distribution.



Copula Mixture Model

- ▶ We assume that the joint density is a Dirichlet process prior mixture

$$f_{(\mathbf{X}, \mathbf{Y})}(\mathbf{x}, \mathbf{y}) = \int \int f_{(\mathbf{X}, \mathbf{Y})|\theta, \mathbf{P}}(\mathbf{x}, \mathbf{y}) d\mu_{\theta, \mathbf{P}} d\mu_{\mathbf{G}}(\boldsymbol{\lambda}, \mathbf{G}_0).$$

- ▶ We specify the margins and the dependence separately:

- ▶ The margins can be arbitrary continuous cdfs:

$$\mathbf{X}^j | \theta = \mathbf{X}^j | \theta^j \sim F_{\mathbf{X}^j | \theta^j}^j, \quad j = 1, \dots, p,$$

$$\mathbf{Y}^j | \theta = \mathbf{Y}^j | \theta^j \sim F_{\mathbf{Y}^j | \theta^j}^j, \quad j = 1, \dots, q.$$

- ▶ The dependence structure is then specified by a Gaussian copula \mathbf{C}_P^G with block diagonal correlation matrix \mathbf{P} .

- ▶ Finally the conditional cdf is: $F_{(\mathbf{X}, \mathbf{Y})|\theta, \mathbf{P}}(\mathbf{x}, \mathbf{y}) = \mathbf{C}_P^G(F_{\mathbf{X}^1 | \theta^1}^1(x^1), \dots, F_{\mathbf{X}^p | \theta^p}^p(x^p), F_{\mathbf{Y}^1 | \theta^1}^1(y^1), \dots, F_{\mathbf{Y}^q | \theta^q}^q(y^q)).$

Gaussian Copula and Copula densities

- ▶ Consider a multivariate Gaussian variable $\mathbf{X} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then by Sklar's Theorem there exists a unique copula \mathbf{C} such that:

$$\Phi_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{x}^1, \dots, \mathbf{x}^d) = \mathbf{C}(\Phi_{\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}}(\mathbf{x}^1), \dots, \Phi_{\boldsymbol{\mu}_d, \boldsymbol{\Sigma}_{dd}}(\mathbf{x}^d)).$$

- ▶ The copula of $\mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the same as the copula of $\mathcal{N}_d(\mathbf{0}, \mathbf{P})$, where $\mathbf{P} = \mathcal{P}(\boldsymbol{\Sigma})$.

- ▶ The Gaussian copula is given by:

$$\mathbf{C}_P^G(\mathbf{u}) = \Phi_P(\Phi^{-1}(u^1), \dots, \Phi^{-1}(u^d)).$$

- ▶ Consider a multivariate cdf \mathbf{F} with copula \mathbf{C} and marginal cdfs F^1, \dots, F^d :

$$\mathbf{F}(\mathbf{x}^1, \dots, \mathbf{x}^d) = \mathbf{C}(F^1(x^1), \dots, F^d(x^d)).$$

If \mathbf{F} has a density then it can be expressed as:

$$f(\mathbf{x}^1, \dots, \mathbf{x}^d) = c(F^1(x^1), \dots, F^d(x^d)) \prod_{j=1}^d f^j(x^j),$$

where $c(\mathbf{u}^1, \dots, \mathbf{u}^d) = \frac{\partial \mathbf{C}(\mathbf{u}^1, \dots, \mathbf{u}^d)}{\partial u^1 \dots \partial u^d}$ is the copula density of \mathbf{C} .

- ▶ The Gaussian copula density has a simple form and the multivariate conditional density \mathbf{f} becomes:

$$f_{(\mathbf{X}, \mathbf{Y})|\theta, \mathbf{P}}(\mathbf{x}, \mathbf{y}) = |\mathbf{P}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \tilde{\mathbf{x}}^T (\mathbf{P}^{-1} - \mathbf{I}) \tilde{\mathbf{x}}\right\} \prod_{j=1}^{p+q} f^j(x^j),$$

where $\tilde{\mathbf{x}}^j = \Phi^{-1}(F^j(x^j))$.

Bayesian inference

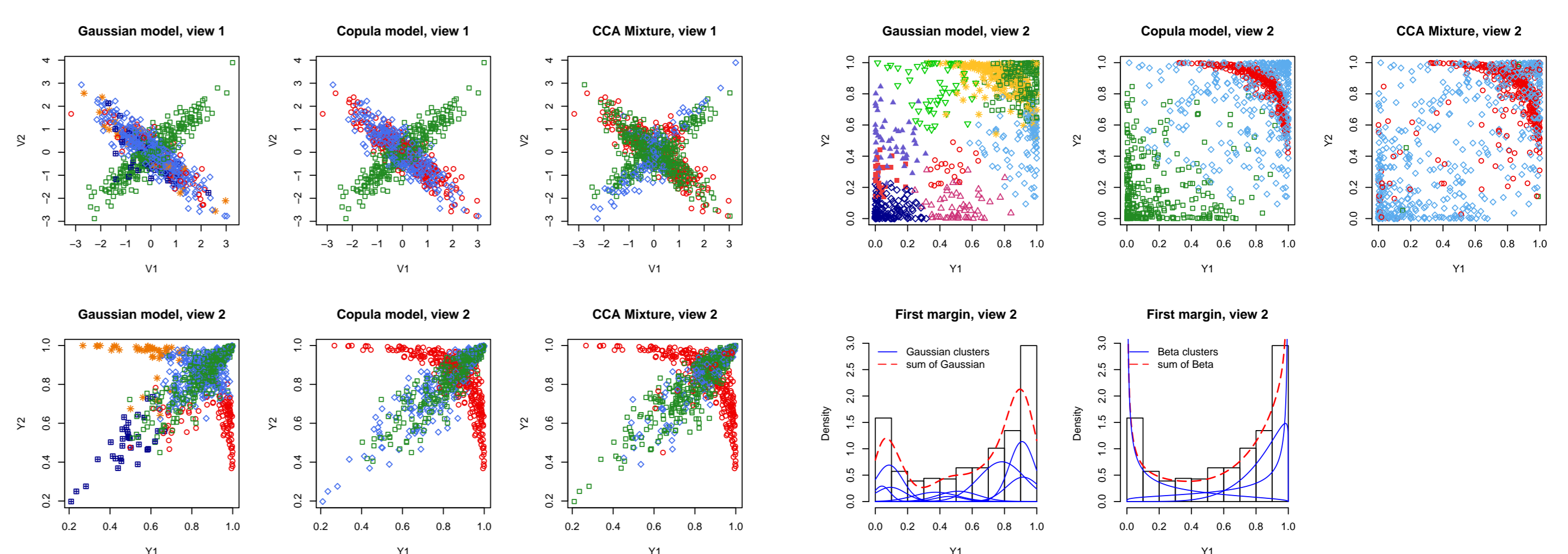
- ▶ Assume *a priori* independence of $\boldsymbol{\theta}$ and \mathbf{P} , specify the priors separately.
- ▶ Specify prior distributions for the blocks \mathbf{P}_x and \mathbf{P}_y , where $\mathbf{P} = \begin{pmatrix} \mathbf{P}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_y \end{pmatrix}$ and assume that \mathbf{P}_x and \mathbf{P}_y are *a priori* independent.
- ▶ For \mathbf{P}_x and \mathbf{P}_y we choose the marginally uniform prior [Barnard, 2000]:

$$\pi(\mathbf{R}, d+1) \propto |\mathbf{R}|^{\frac{d(d-1)}{2}-1} \left(\prod_{i=1}^d |\mathbf{R}_{ii}| \right)^{-\frac{(d+1)}{2}}.$$

- ▶ Inference for $\boldsymbol{\theta}$ and \mathbf{P} is performed using MCMC.

Results

Simulations



- ▶ Real data. Two data sets containing information about the regulation of heat shock in yeast. First view: series of gene expressions for yeast under heat shock measured at 4 time points. Second view: probability scores of binding interactions between the promoter region of different genes and several DNA-binding transcriptional regulators under heat shock.

