

Machine Learning and Statistics in Genetics and Genomics

VII: Continuous latent variable models

Christoph Lippert

Microsoft Research
eScience group

Los Angeles , USA

Microsoft
Research

Current topics in computational biology

UCLA

Winter quarter 2014

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

Model

Applications

Modeling unobserved cellular phenotypes in genetic analyses

Model

Applications

A unifying view

Summary

Outline

Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

- Model

- Applications

Modeling unobserved cellular phenotypes in genetic analyses

- Model

- Applications

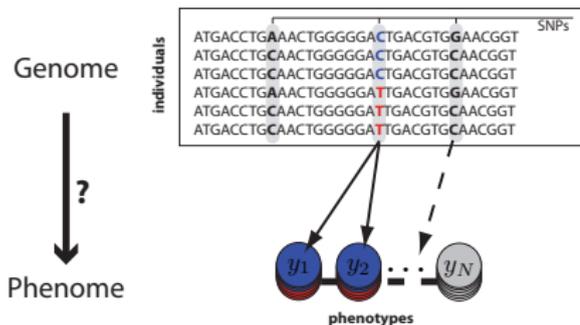
A unifying view

Summary

Why latent variables ?

Causal influences on phenotypes

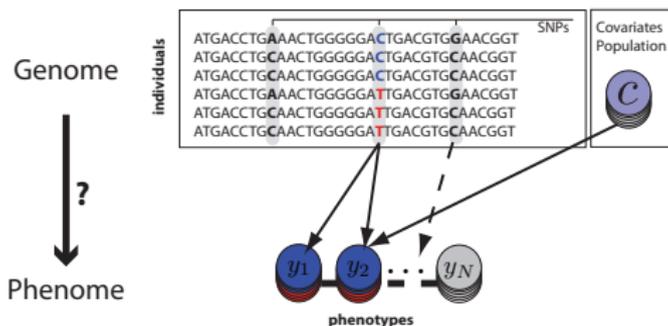
- ▶ Genotype
 - ▶ Primary variable of interest
- ▶ Known confounding factors
 - ▶ Covariates
 - ▶ Population structure
- ▶ Unknown (latent) confounders
 - ▶ Sample handling
 - ▶ Sample history
 - ▶ Subtle environmental perturbations



Why latent variables ?

Causal influences on phenotypes

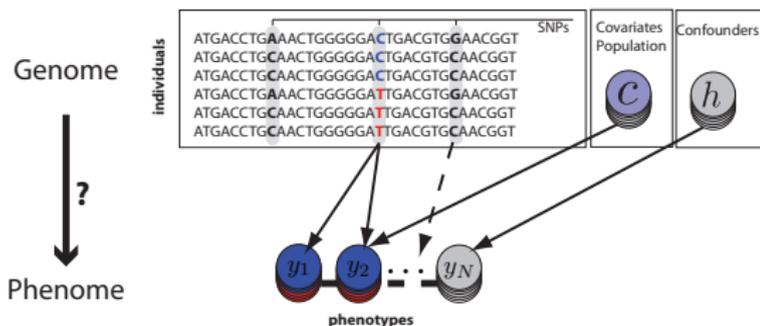
- ▶ Genotype
 - ▶ Primary variable of interest
- ▶ Known confounding factors
 - ▶ Covariates
 - ▶ Population structure
- ▶ Unknown (latent) confounders
 - ▶ Sample handling
 - ▶ Sample history
 - ▶ Subtle environmental perturbations



Why latent variables ?

Causal influences on phenotypes

- ▶ Genotype
 - ▶ Primary variable of interest
- ▶ Known confounding factors
 - ▶ Covariates
 - ▶ Population structure
- ▶ Unknown (latent) confounders
 - ▶ Sample handling
 - ▶ Sample history
 - ▶ Subtle environmental perturbations



Outline

Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

Model

Applications

Modeling unobserved cellular phenotypes in genetic analyses

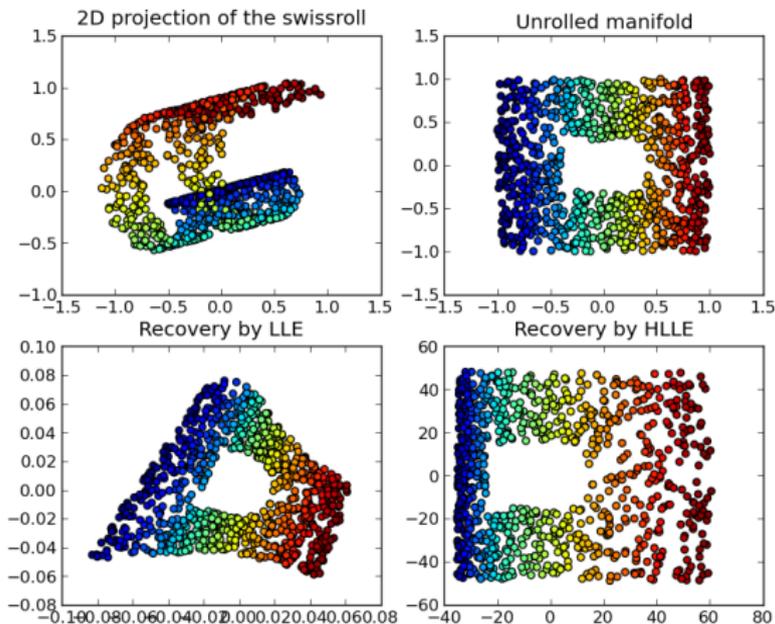
Model

Applications

A unifying view

Summary

Manifolds and dimension reduction



(from Olivier Grisel, Generated using the Modular Data Processing toolkit and matplotlib.)

Linear dimension reduction

- ▶ Map G dimensional data on K dimensional manifold; $K \ll G$

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G} + \underbrace{\boldsymbol{\Psi}}_{N \times G}$$

- ▶ \mathbf{H} : latent factors in low-dimensional space
- ▶ \mathbf{W} : weights for factors on data dimensions
- ▶ $\boldsymbol{\Psi}$: noise, $\psi_{n,g} \sim \mathcal{N}(0, \sigma^2)$.
- ▶ **Challenge:** neither \mathbf{W} nor \mathbf{H} known!
- ▶ Depending on assumptions on \mathbf{W} and \mathbf{H} :
 - ▶ Principle component analysis (PCA)
 - ▶ Independent component analysis (ICA)
 - ▶ ...

Linear dimension reduction

- ▶ Map G dimensional data on K dimensional manifold; $K \ll G$

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G} + \underbrace{\boldsymbol{\Psi}}_{N \times G}$$

- ▶ \mathbf{H} : latent factors in low-dimensional space
- ▶ \mathbf{W} : weights for factors on data dimensions
- ▶ $\boldsymbol{\Psi}$: noise, $\psi_{n,g} \sim \mathcal{N}(0, \sigma^2)$.
- ▶ **Challenge:** neither \mathbf{W} nor \mathbf{H} known!
- ▶ Depending on assumptions on \mathbf{W} and \mathbf{H} :
 - ▶ Principle component analysis (PCA)
 - ▶ Independent component analysis (ICA)
 - ▶ ...

Linear dimension reduction

- ▶ Map G dimensional data on K dimensional manifold; $K \ll G$

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G} + \underbrace{\boldsymbol{\Psi}}_{N \times G}$$

- ▶ \mathbf{H} : latent factors in low-dimensional space
- ▶ \mathbf{W} : weights for factors on data dimensions
- ▶ $\boldsymbol{\Psi}$: noise, $\psi_{n,g} \sim \mathcal{N}(0, \sigma^2)$.
- ▶ **Challenge:** neither \mathbf{W} nor \mathbf{H} known!
- ▶ Depending on assumptions on \mathbf{W} and \mathbf{H} :
 - ▶ Principle component analysis (PCA)
 - ▶ Independent component analysis (ICA)
 - ▶ ...

Linear dimension reduction

PCA

PCA is corresponds to a noise-free version of the model

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G}$$

- ▶ PCA components (\mathbf{H}) correspond to directions of maximum data variance in the original dataset:
 - ▶ Covariance matrix: $\mathbf{C} = \mathbf{Y}\mathbf{Y}^\top$
 - ▶ Eigenvalue/Eigen vectors $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$
 - ▶ Projection matrix $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$
 - ▶ Principle components $\mathbf{H}_n = \mathbf{P} \cdot \mathbf{Y}_n$.

Linear dimension reduction

PCA

PCA is corresponds to a noise-free version of the model

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G}$$

- ▶ PCA components (\mathbf{H}) correspond to directions of maximum data variance in the original dataset:
 - ▶ Covariance matrix: $\mathbf{C} = \mathbf{Y}\mathbf{Y}^\top$
 - ▶ Eigenvalue/Eigen vectors $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$
 - ▶ Projection matrix $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$
 - ▶ Principle components $\mathbf{H}_n = \mathbf{P} \cdot \mathbf{Y}_n$.

Linear dimension reduction

PCA

PCA is corresponds to a noise-free version of the model

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G}$$

- ▶ PCA components (\mathbf{H}) correspond to directions of maximum data variance in the original dataset:
 - ▶ Covariance matrix: $\mathbf{C} = \mathbf{Y}\mathbf{Y}^\top$
 - ▶ Eigenvalue/Eigen vectors $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$
 - ▶ Projection matrix $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$
 - ▶ Principle components $\mathbf{H}_n = \mathbf{P} \cdot \mathbf{Y}_n$.

Linear dimension reduction

PCA

PCA is corresponds to a noise-free version of the model

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G}$$

- ▶ PCA components (\mathbf{H}) correspond to directions of maximum data variance in the original dataset:
 - ▶ Covariance matrix: $\mathbf{C} = \mathbf{Y}\mathbf{Y}^\top$
 - ▶ Eigenvalue/Eigen vectors $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$
 - ▶ Projection matrix $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$
 - ▶ Principle components $\mathbf{H}_n = \mathbf{P} \cdot \mathbf{Y}_n$.

Linear dimension reduction

PCA

PCA is corresponds to a noise-free version of the model

$$\underbrace{\mathbf{Y}}_{N \times G} = \underbrace{\mathbf{H}}_{N \times K} \underbrace{\mathbf{W}}_{K \times G}$$

- ▶ PCA components (\mathbf{H}) correspond to directions of maximum data variance in the original dataset:
 - ▶ Covariance matrix: $\mathbf{C} = \mathbf{Y}\mathbf{Y}^\top$
 - ▶ Eigenvalue/Eigen vectors $\mathbf{C}\mathbf{v}_i = \lambda_i\mathbf{v}_i$
 - ▶ Projection matrix $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$
 - ▶ Principle components $\mathbf{H}_n = \mathbf{P} \cdot \mathbf{Y}_n$.

Linear dimension reduction

Bayesian PCA and GPLVM

Assumption: data dimensions or sample dimension independent **given H and W** .

Probabilistic PCA

$$p(\mathbf{Y}|\mathbf{H}, \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n | \mathbf{h}_n \mathbf{W}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{H}) = \prod_{n=1}^N \mathcal{N}(\mathbf{h}_n | 0, \sigma_h^2 \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n | 0, \sigma_h^2 \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I})$$

GPLVM

$$p(\mathbf{Y}|\mathbf{H}, \mathbf{W}) = \prod_{g=1}^G \mathcal{N}(\mathbf{y}_{:,g} | \mathbf{H} \mathbf{w}_g, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{g=1}^G \mathcal{N}(\mathbf{w}_{:,g} | 0, \sigma_h^2 \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{H}) = \prod_{g=1}^G \mathcal{N}\left(\mathbf{y}_{:,g} | 0, \underbrace{\sigma_h^2 \mathbf{H} \mathbf{H}^\top}_{N \times N} + \sigma^2 \mathbf{I}\right)$$

Linear dimension reduction

Bayesian PCA and GPLVM

Assumption: data dimensions or sample dimension independent **given \mathbf{H} and \mathbf{W}** .

Probabilistic PCA

$$p(\mathbf{Y}|\mathbf{H}, \mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n | \mathbf{h}_n \mathbf{W}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{H}) = \prod_{n=1}^N \mathcal{N}(\mathbf{h}_n | 0, \sigma_h^2 \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{n=1}^N \mathcal{N}(\mathbf{y}_n | 0, \sigma_h^2 \mathbf{W} \mathbf{W}^\top + \sigma^2 \mathbf{I})$$

GPLVM

$$p(\mathbf{Y}|\mathbf{H}, \mathbf{W}) = \prod_{g=1}^G \mathcal{N}(\mathbf{y}_{:,g} | \mathbf{H} \mathbf{w}_g, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{g=1}^G \mathcal{N}(\mathbf{w}_{:,g} | 0, \sigma_h^2 \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{H}) = \prod_{g=1}^G \mathcal{N}\left(\mathbf{y}_{:,g} | 0, \underbrace{\sigma_h^2 \mathbf{H} \mathbf{H}^\top}_{N \times N} + \sigma^2 \mathbf{I}\right)$$

GPLVM

- ▶ Marginal likelihood

$$p(\mathbf{Y}|\mathbf{H}) = \prod_{g=1}^G \mathcal{N}\left(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H}\mathbf{H}^\top + \sigma^2 \mathbf{I}\right)$$

- ▶ Inference of most probable “hidden factors” \mathbf{H} :

$$\hat{\mathbf{H}} = \operatorname{argmax}_{\mathbf{H}} \log \prod_{g=1}^G \mathcal{N}\left(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H}\mathbf{H}^\top + \sigma^2 \mathbf{I}\right)$$

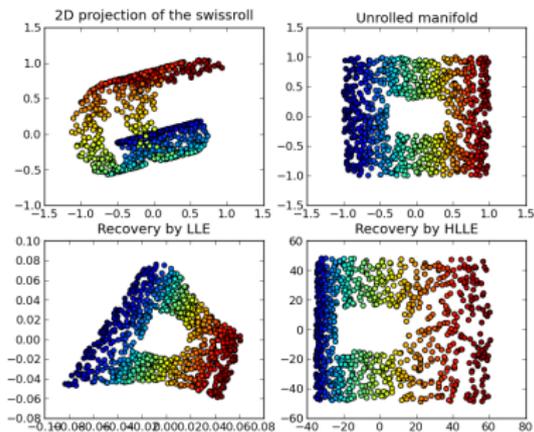
- ▶ Marginal likelihood

$$p(\mathbf{Y}|\mathbf{H}) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H} \mathbf{H}^\top + \sigma^2 \mathbf{I} \right)$$

- ▶ Inference of most probable “hidden factors” \mathbf{H} :

$$\hat{\mathbf{H}} = \underset{\mathbf{H}}{\operatorname{argmax}} \log \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H} \mathbf{H}^\top + \sigma^2 \mathbf{I} \right)$$

Are linear relationships sufficient?

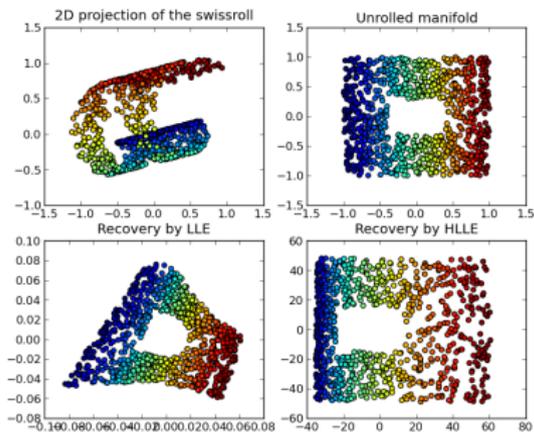


(from Olivier Grisel, Generated using the Modular Data Processing toolkit and matplotlib.)

- ▶ Non-linear generalizations, introducing general kernel function instead of linear covariance

$$\hat{H} = \operatorname{argmax}_H \log \prod_{g=1}^G \mathcal{N}(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H} \mathbf{H}^\top + \sigma^2 \mathbf{I})$$

Are linear relationships sufficient?

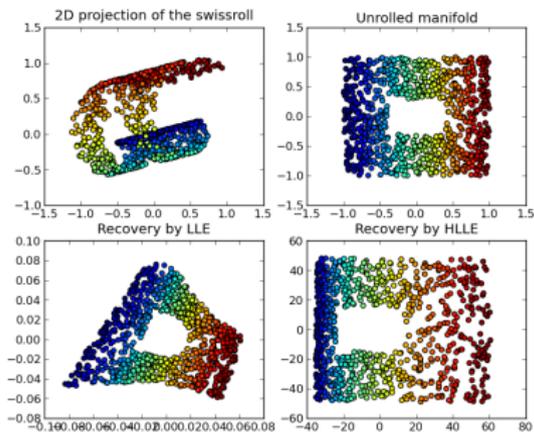


(from Olivier Grisel, Generated using the Modular Data Processing toolkit and matplotlib.)

- ▶ Non-linear generalizations, introducing general kernel function instead of linear covariance

$$\hat{H} = \operatorname{argmax}_H \log \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{H} \mathbf{H}^\top + \sigma^2 \mathbf{I} \right)$$

Are linear relationships sufficient?



(from Olivier Grisel, Generated using the Modular Data Processing toolkit and matplotlib.)

- ▶ Non-linear generalizations, introducing general kernel function instead of linear covariance

$$\hat{H} = \operatorname{argmax}_H \log \prod_{g=1}^G \mathcal{N}(\mathbf{y}_{:,g} \mid \mathbf{0}, \sigma_h^2 \mathbf{K}_{H,H} + \sigma^2 \mathbf{I})$$

Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

Model

Applications

Modeling unobserved cellular phenotypes in genetic analyses

Model

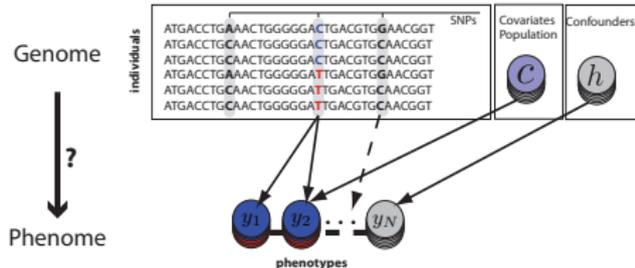
Applications

A unifying view

Summary

Confounders in eQTL studies

- ▶ Confounders in eQTL studies
 - ▶ Experimental procedures
 - ▶ Gene regulation
 - ▶ (Translation)
- ▶ Standard to take *known factors* (gender, population structure) into account.
- ▶ It is key to account for **hidden factors** as well.
 - ▶ Sample preparation
 - ▶ Sample history



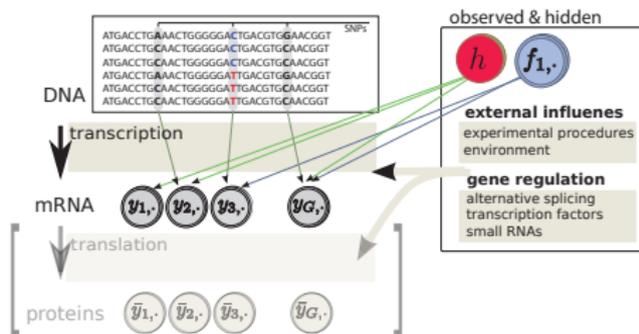
Confounders in eQTL studies

- ▶ Confounders in eQTL studies
 - ▶ **Experimental procedures**
 - ▶ Gene regulation
 - ▶ (Translation)

▶ Standard to take *known factors* (gender, population structure) into account.

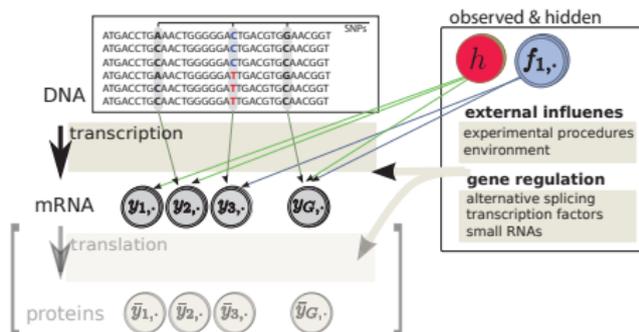
▶ It is key to account for **hidden factors** as well.

- ▶ Sample preparation
- ▶ Sample history



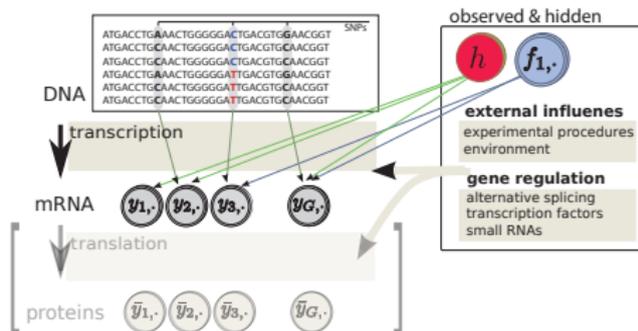
Confounders in eQTL studies

- ▶ Confounders in eQTL studies
 - ▶ **Experimental procedures**
 - ▶ Gene regulation
 - ▶ (Translation)
- ▶ Standard to take *known factors* (gender, population structure) into account.
- ▶ It is key to account for **hidden factors** as well.
 - ▶ Sample preparation
 - ▶ Sample history



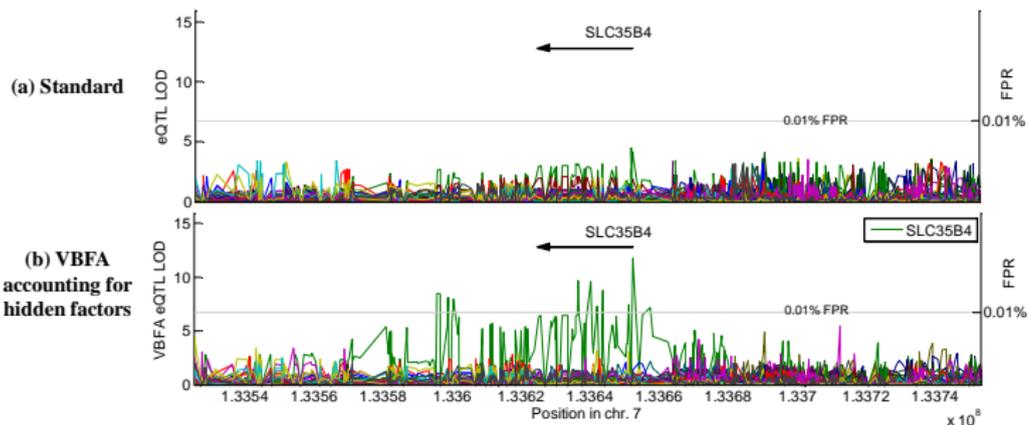
Confounders in eQTL studies

- ▶ Confounders in eQTL studies
 - ▶ **Experimental procedures**
 - ▶ Gene regulation
 - ▶ (Translation)
- ▶ Standard to take *known factors* (gender, population structure) into account.
- ▶ It is key to account for **hidden factors** as well.
 - ▶ Sample preparation
 - ▶ Sample history



Motivating examples

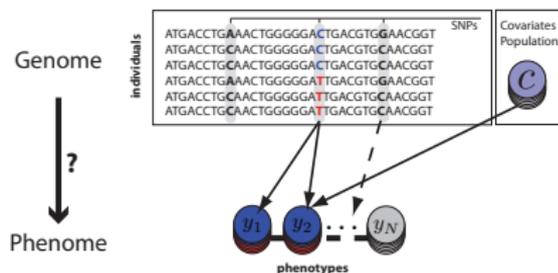
HapMAP II, 3 populations, 90 individuals each, 40K genes, 3 million SNPs
(here: small region in chromosome 7)



Association model

Direct effects of confounders

- ▶ Start with standard association model.
- ▶ Include (K , few) global hidden factors (confounders) in the model.
- ▶ Factors $H = \{h_{:,1}, \dots, h_{:,K}\}$ need to be learned from the expression data.
- ▶ How to control model complexity, i.e. choose the number of factors?

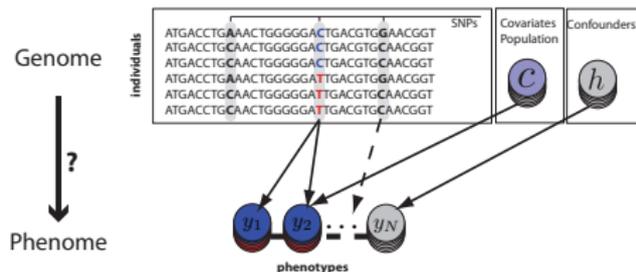


$$\mathbf{y}_{:,g} = \underbrace{\sum_{n=1}^N (\theta_{n,g} \mathbf{s}_n)}_{\text{genetic}} + \underbrace{\psi}_{\text{noise}}$$

Association model

Direct effects of confounders

- ▶ Start with standard association model.
- ▶ Include (K , few) global hidden factors (confounders) in the model.
- ▶ Factors $H = \{h_{:,1}, \dots, h_{:,K}\}$ need to be learned from the expression data.
- ▶ How to control model complexity, i.e. choose the number of factors?

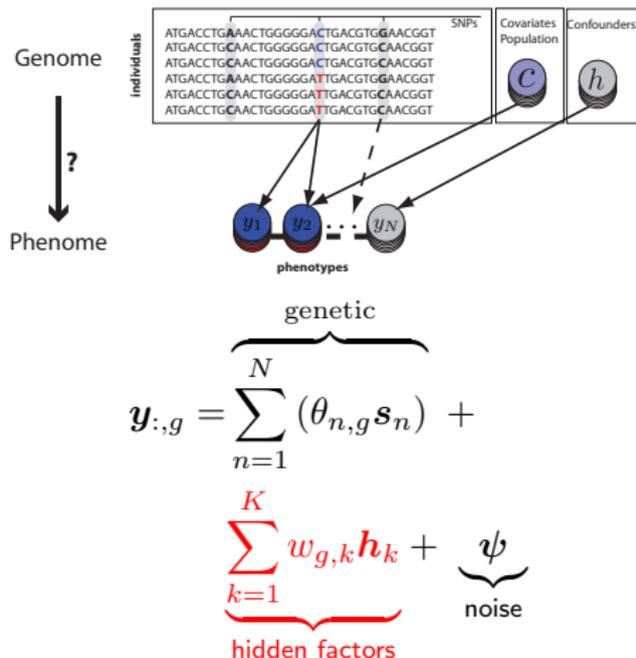


$$\mathbf{y}_{:,g} = \underbrace{\sum_{n=1}^N (\theta_{n,g} \mathbf{s}_n)}_{\text{genetic}} + \underbrace{\sum_{k=1}^K w_{g,k} \mathbf{h}_k}_{\text{hidden factors}} + \underbrace{\psi}_{\text{noise}}$$

Association model

Direct effects of confounders

- ▶ Start with standard association model.
- ▶ Include (K , few) global hidden factors (confounders) in the model.
- ▶ Factors $\mathbf{H} = \{h_{:,1}, \dots, h_{:,K}\}$ need to be learned from the expression data.
- ▶ How to control model complexity, i.e. choose the number of factors?



Association model

Gaussian process formulation

- ▶ Data likelihood of the linear generative model assuming Gaussian noise

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}, \mathbf{W}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s + \sum_{k=1}^K w_{g,k} \mathbf{h}_k, \sigma^2 \mathbf{I} \right)$$

- ▶ Specify Gaussian priors on the factor weight

$$P(w_{g,k}) = \mathcal{N} (w_{g,k} \mid 0, \sigma_h^2)$$

- ▶ Integrating out the weights yields a Gaussian process marginal likelihood, factorizing over genes given \mathbf{H}

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s, \sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma^2 \mathbf{I} \right)$$

Association model

Gaussian process formulation

- ▶ Data likelihood of the linear generative model assuming Gaussian noise

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}, \mathbf{W}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s + \sum_{k=1}^K w_{g,k} \mathbf{h}_k, \sigma^2 \mathbf{I} \right)$$

- ▶ Specify Gaussian priors on the factor weight

$$P(w_{g,k}) = \mathcal{N} (w_{g,k} \mid 0, \sigma_h^2)$$

- ▶ Integrating out the weights yields a Gaussian process marginal likelihood, factorizing over genes given \mathbf{H}

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s, \sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma^2 \mathbf{I} \right)$$

Association model

Gaussian process formulation

- ▶ Data likelihood of the linear generative model assuming Gaussian noise

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}, \mathbf{W}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s + \sum_{k=1}^K w_{g,k} \mathbf{h}_k, \sigma^2 \mathbf{I} \right)$$

- ▶ Specify Gaussian priors on the factor weight

$$P(w_{g,k}) = \mathcal{N} (w_{g,k} \mid 0, \sigma_h^2)$$

- ▶ Integrating out the weights yields a Gaussian process marginal likelihood, factorizing over genes given \mathbf{H}

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s, \sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma^2 \mathbf{I} \right)$$

Association model

Including population structure

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s, \underbrace{\sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma^2 \mathbf{I}}_{N \times N} \right)$$

- ▶ Including additional covariance term for population structure.
- ▶ Association test for single SNP.
- ▶ Challenge: Refitting σ_g^2 , σ_h^2 , \mathbf{H} for every SNP.

Association model

Including population structure

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \sum_{s=1}^S \mathbf{x}_s \theta_s, \underbrace{\sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma_g^2 \mathbf{X} \mathbf{X}^\top + \sigma^2 \mathbf{I}}_{N \times N} \right)$$

- ▶ Including additional covariance term for population structure.
- ▶ Association test for single SNP.
- ▶ Challenge: Refitting σ_g^2 , σ_h^2 , \mathbf{H} for every SNP.

Association model

Including population structure

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \mathbf{x}_i \theta_i, \underbrace{\sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma_g^2 \mathbf{X} \mathbf{X}^\top + \sigma^2 \mathbf{I}}_{N \times N} \right)$$

- ▶ Including additional covariance term for population structure.
- ▶ Association test for single SNP.
- ▶ Challenge: Refitting σ_g^2 , σ_h^2 , \mathbf{H} for every SNP.

Association model

Including population structure

$$P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K) = \prod_{g=1}^G \mathcal{N} \left(\mathbf{y}_{:,g} \mid \mathbf{x}_i \theta_i, \underbrace{\sigma_h^2 \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^\top + \sigma_g^2 \mathbf{X} \mathbf{X}^\top + \sigma^2 \mathbf{I}}_{N \times N} \right)$$

- ▶ Including additional covariance term for population structure.
- ▶ Association test for single SNP.
- ▶ Challenge: Refitting σ_g^2 , σ_h^2 , \mathbf{H} for every SNP.

Association model

- ▶ Fit parameter once on null model

$$\hat{\mathbf{H}}, \hat{\sigma}_g^2, \hat{\sigma}_h^2, \hat{\sigma}^2 = \underset{\hat{\mathbf{H}}, \hat{\sigma}_h^2, \hat{\sigma}_g^2, \hat{\sigma}^2}{\operatorname{argmax}} P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \Theta_K)$$

- ▶ Association testing for fixed hyperparameters

$$p(\mathbf{y}_{:,g} | \mathbf{x}_i, \mathbf{H}, \Theta_K) = \mathcal{N} \left(\mathbf{y}_{:,g} | \mathbf{x}_i \theta_i, \alpha^2 (\hat{\sigma}^2 \hat{\mathbf{H}} \hat{\mathbf{H}}^\top + \hat{\sigma}_g^2 \mathbf{X} \mathbf{X}^\top) + \sigma^2 \mathbf{I} \right)$$

Association model

- ▶ Fit parameter once on null model

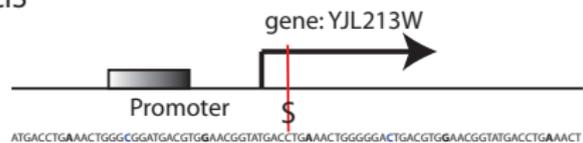
$$\hat{\mathbf{H}}, \hat{\sigma}_g^2, \hat{\sigma}_h^2, \hat{\sigma}^2 = \underset{\hat{\mathbf{H}}, \hat{\sigma}_h^2, \hat{\sigma}_g^2, \hat{\sigma}^2}{\operatorname{argmax}} P(\mathbf{Y} | \mathbf{X}, \mathbf{H}, \boldsymbol{\Theta}_K)$$

- ▶ Association testing for fixed hyperparameters

$$p(\mathbf{y}_{:,g} | \mathbf{x}_i, \mathbf{H}, \boldsymbol{\Theta}_K) = \mathcal{N} \left(\mathbf{y}_{:,g} | \mathbf{x}_i \theta_i, \alpha^2 (\hat{\sigma}^2 \hat{\mathbf{H}} \hat{\mathbf{H}}^\top + \hat{\sigma}_g^2 \mathbf{X} \mathbf{X}^\top) + \sigma^2 \mathbf{I} \right)$$

Evaluation on real data

CIS

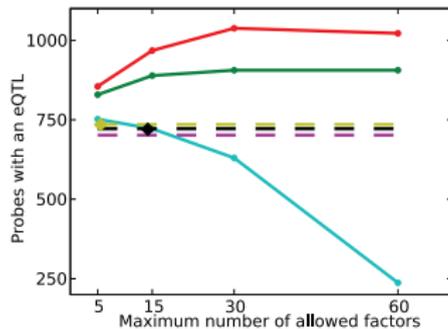
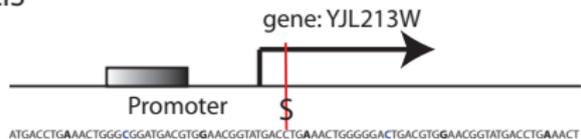


TRANS



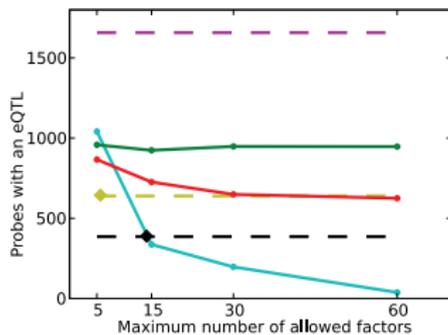
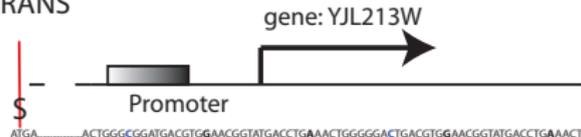
Evaluation on real data

CIS



(a) Yeast *cis* eQTLs

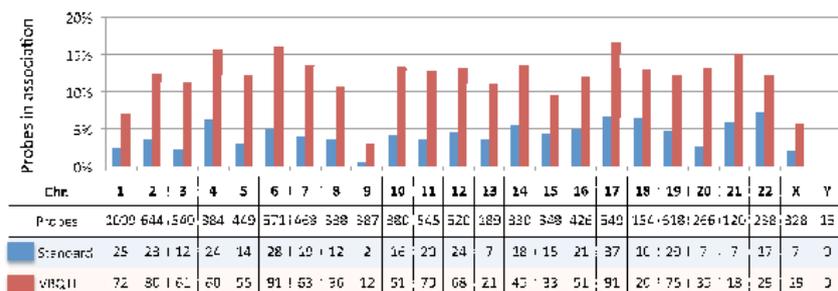
TRANS



(b) Yeast *trans* eQTLs

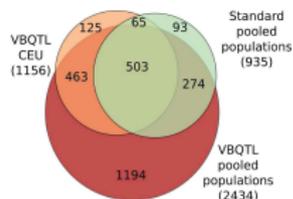
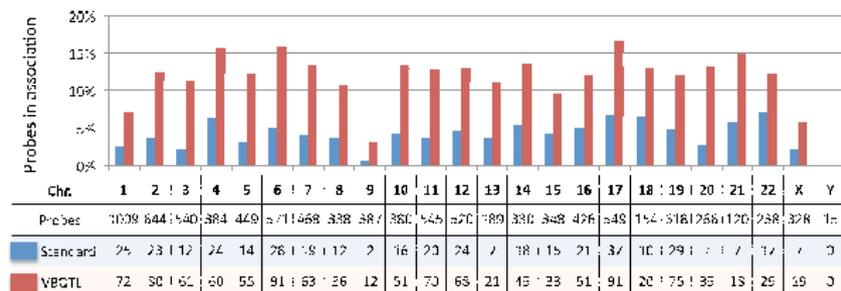
Evaluation on real data

Application to HapMap II expression analysis

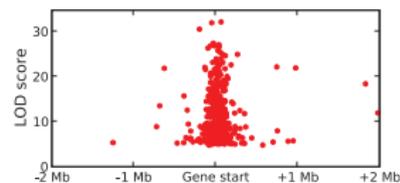


Evaluation on real data

Application to HapMap II expression analysis



(a) Probes with a VB:QIL in pooled population



(f) *cis* VB:QIL location and strength relative to gene start

Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

Model

Applications

Modeling unobserved cellular phenotypes in genetic analyses

Model

Applications

A unifying view

Summary

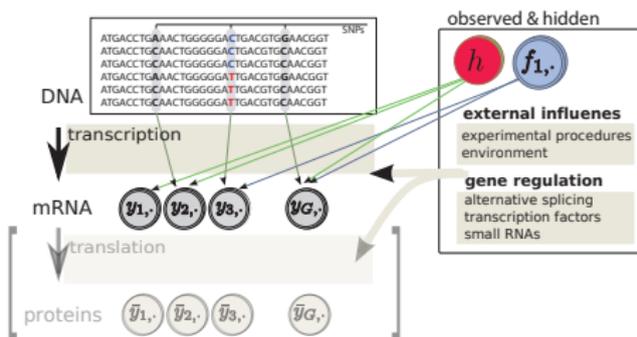
Association studies

Regulatory and external factors

- ▶ Confounding factors
 - ▶ Accounting for hidden confounders

$$y_{g,j} = \underbrace{b_{n,g}}_{\text{genetic}} (\underbrace{\theta_{n,g}}_{\text{genetic}} \underbrace{s_{n,j}}_{\text{genetic}}) + \underbrace{v_g f_j}_{\text{known factors}} + \underbrace{w_g x_j}_{\text{hidden factors}} + \underbrace{\psi_{n,g}}_{\text{noise}}$$

- ▶ Accounting for regulatory factors?



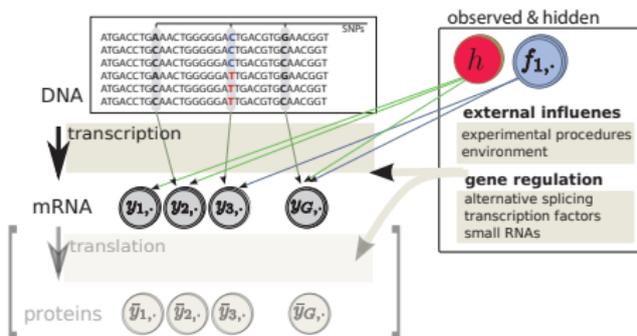
Association studies

Regulatory and external factors

- ▶ Confounding factors
 - ▶ Accounting for hidden confounders

$$y_{g,j} = \underbrace{b_{n,g}}_{\text{genetic}} (\underbrace{\theta_{n,g}}_{\text{genetic}} \underbrace{s_{n,j}}_{\text{genetic}}) + \underbrace{v_g f_j}_{\text{known factors}} + \underbrace{w_g x_j}_{\text{hidden factors}} + \underbrace{\psi_{n,g}}_{\text{noise}}$$

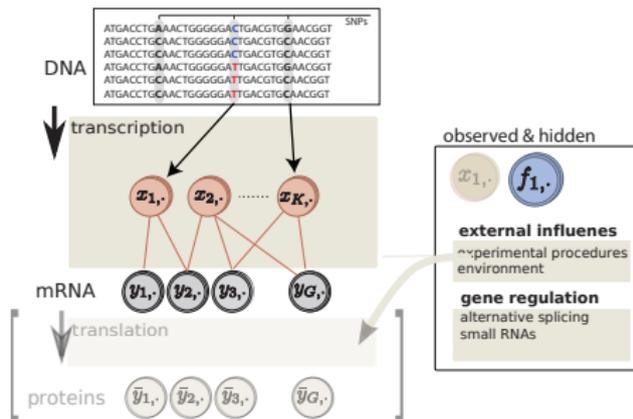
- ▶ Accounting for regulatory factors?



Association studies

Regulatory factors

- ▶ Account for regulatory factors:
 - ▶ Transcription factors
 - ▶ Pathway components
- ▶ Hypothesis: intermediate cellular factors mediate the observed association signals of target genes.
 - ▶ Measuring X ?
 - ▶ Difficult and expensive.
 - ▶ Learn the unobserved factors X .

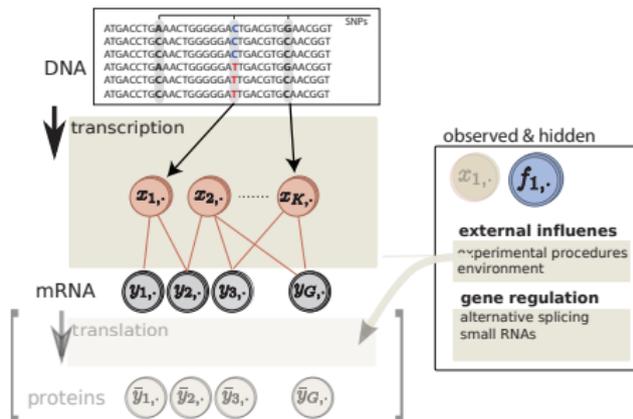


[Parts et al. 2011]

Association studies

Regulatory factors

- ▶ Account for regulatory factors:
 - ▶ Transcription factors
 - ▶ Pathway components
- ▶ Hypothesis: intermediate cellular factors mediate the observed association signals of target genes.
- ▶ Measuring X ?
 - ▶ Difficult and expensive.
- ▶ Learn the unobserved factors X .

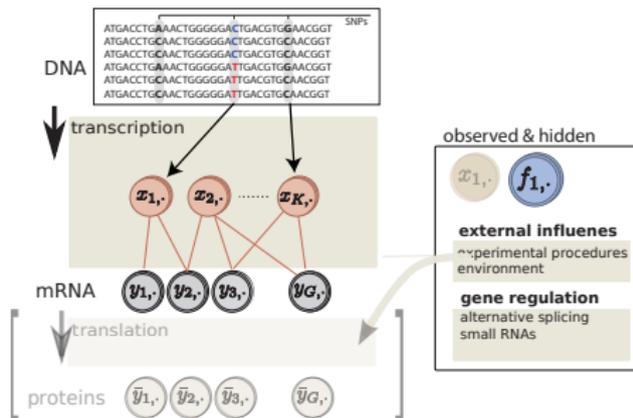


[Parts et al. 2011]

Association studies

Regulatory factors

- ▶ Account for regulatory factors:
 - ▶ Transcription factors
 - ▶ Pathway components
- ▶ Hypothesis: intermediate cellular factors mediate the observed association signals of target genes.
- ▶ Measuring X ?
 - ▶ Difficult and expensive.
- ▶ **Learn** the unobserved factors X .



[Parts et al. 2011]

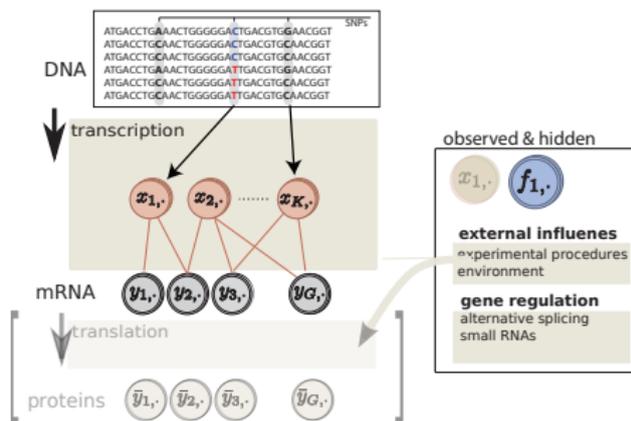
Association studies

Transcription factor effects

- ▶ Inference of regulatory factors using a bilinear model

$$\underbrace{\mathbf{Y}}_{\text{Expr.}} = \underbrace{\mathbf{W}}_{\text{Weights}} \cdot \underbrace{\mathbf{X}}_{\text{Factors}} + \underbrace{\boldsymbol{\Psi}}_{\text{Noise}}$$

- ▶ \mathbf{W} is sparse; each factor regulates only a subset of all genes.
- ▶ Incorporate biological prior knowledge to that factor are interpretable:
 - ▶ Transcription factor binding affinities.
- ▶ Inferred factors summarize co-expression clusters.



Association studies

Transcription factor effects

- ▶ Inference of regulatory factors using a bilinear model

$$\underbrace{\mathbf{Y}}_{\text{Expr.}} = \underbrace{\mathbf{W}}_{\text{Weights}} \cdot \underbrace{\mathbf{X}}_{\text{Factors}} + \underbrace{\Psi}_{\text{Noise}}.$$

- ▶ \mathbf{W} is sparse; each factor regulates only a subset of all genes.
- ▶ Incorporate biological prior knowledge to that factor are interpretable:
 - ▶ Transcription factor binding affinities.
- ▶ Inferred factors summarize co-expression clusters.



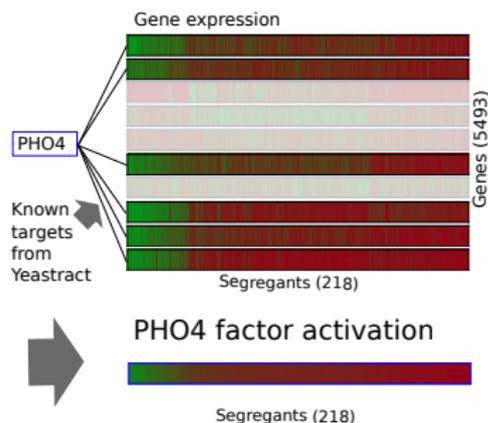
Association studies

Transcription factor effects

- ▶ Inference of regulatory factors using a bilinear model

$$\underbrace{\mathbf{Y}}_{\text{Expr.}} = \underbrace{\mathbf{W}}_{\text{Weights}} \cdot \underbrace{\mathbf{X}}_{\text{Factors}} + \underbrace{\boldsymbol{\Psi}}_{\text{Noise}}$$

- ▶ \mathbf{W} is sparse; each factor regulates only a subset of all genes.
- ▶ Incorporate biological prior knowledge to that factor are interpretable:
 - ▶ Transcription factor binding affinities.
- ▶ Inferred factors summarize co-expression clusters.



Sparse factor analysis

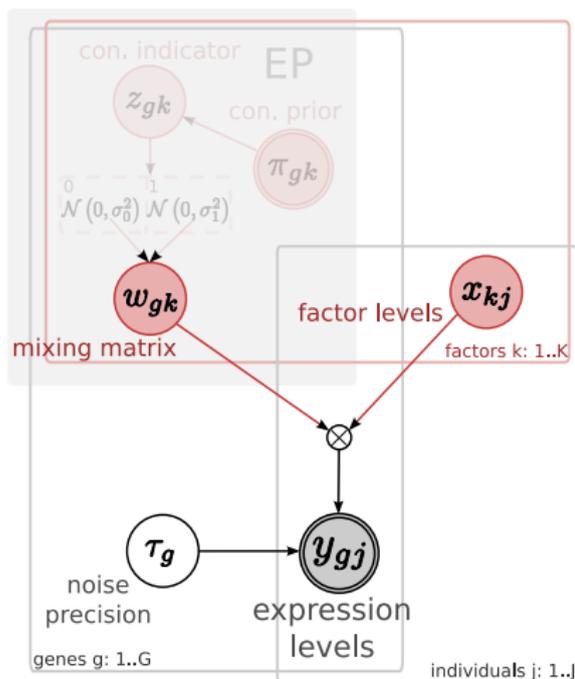
Probabilistic model

- ▶ Graphical model for $\mathbf{Y} = \mathbf{W} \cdot \mathbf{X} + \Psi$.
- ▶ Indicators $z_{g,k}$ determine the sparsity pattern:

$$P(w_{g,k} | z_{g,k} = 0) = \mathcal{N}(w_{g,k} | 0, \sigma_0^2)$$

$$P(w_{g,k} | z_{g,k} = 1) = \mathcal{N}(w_{g,k} | 0, \sigma_1^2).$$

- ▶ Prior knowledge is encoded in $\pi_{g,k} = P(z_{g,k} = 1)$.
- ▶ Standard conjugate priors for the remaining random variables.



Sparse factor analysis

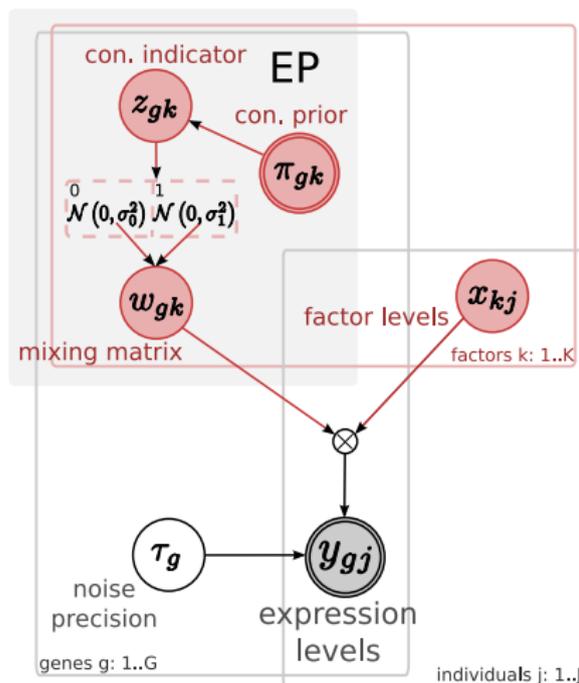
Probabilistic model

- ▶ Graphical model for $\mathbf{Y} = \mathbf{W} \cdot \mathbf{X} + \Psi$.
- ▶ Indicators $z_{g,k}$ determine the sparsity pattern:

$$P(w_{g,k} | z_{g,k} = 0) = \mathcal{N}(w_{g,k} | 0, \sigma_0^2)$$

$$P(w_{g,k} | z_{g,k} = 1) = \mathcal{N}(w_{g,k} | 0, \sigma_1^2).$$

- ▶ Prior knowledge is encoded in $\pi_{g,k} = P(z_{g,k} = 1)$.
- ▶ Standard conjugate priors for the remaining random variables.



Sparse factor analysis

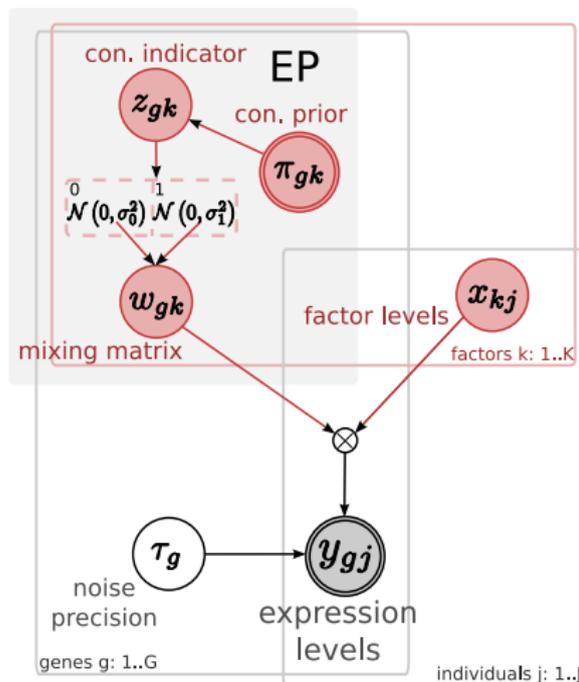
Probabilistic model

- ▶ Graphical model for $\mathbf{Y} = \mathbf{W} \cdot \mathbf{X} + \Psi$.
- ▶ Indicators $z_{g,k}$ determine the sparsity pattern:

$$P(w_{g,k} | z_{g,k} = 0) = \mathcal{N}(w_{g,k} | 0, \sigma_0^2)$$

$$P(w_{g,k} | z_{g,k} = 1) = \mathcal{N}(w_{g,k} | 0, \sigma_1^2).$$

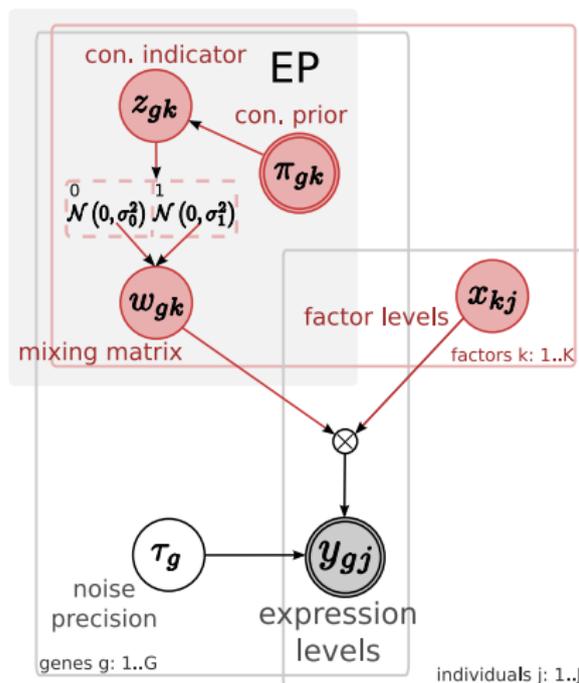
- ▶ Prior knowledge is encoded in $\pi_{g,k} = P(z_{g,k} = 1)$.
- ▶ Standard conjugate priors for the remaining random variables.



Sparse factor analysis

Inference

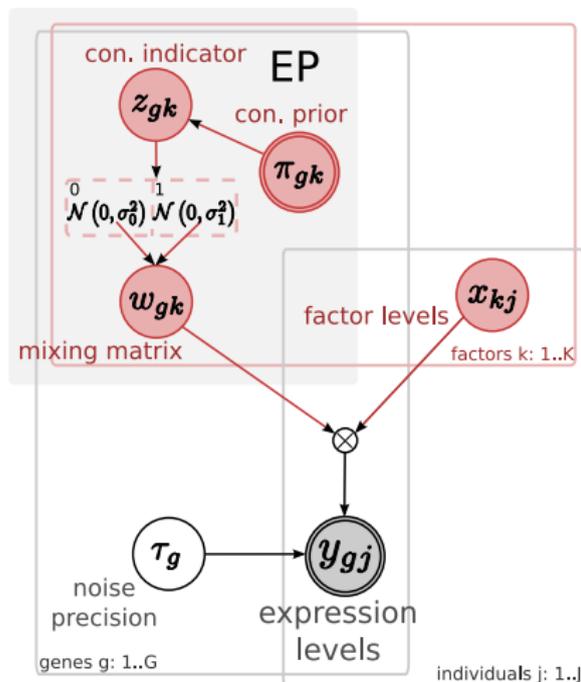
- ▶ Scale is challenging: thousands of genes, hundreds of TFs.
- ▶ Efficient deterministic approximate inference:
 1. Variational Bayesian inference for the core model.
 2. Expectation Propagation for the sparsity submodel.



Sparse factor analysis

Inference

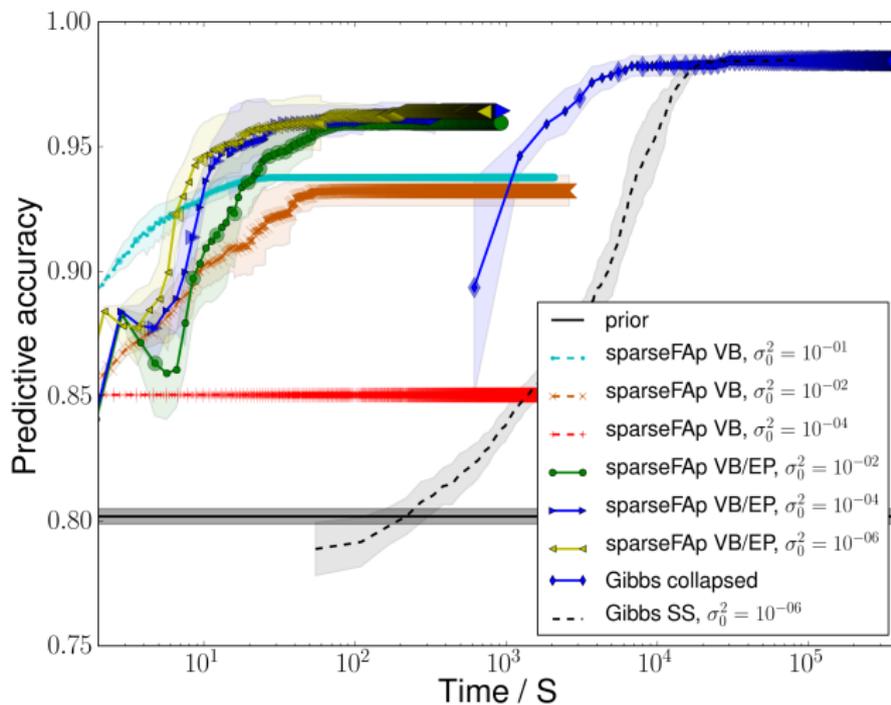
- ▶ Scale is challenging: thousands of genes, hundreds of TFs.
- ▶ Efficient deterministic approximate inference:
 1. Variational Bayesian inference for the core model.
 2. Expectation Propagation for the sparsity submodel.



Sparse factor analysis

Comparison to MCMC on (small!) simulated dataset

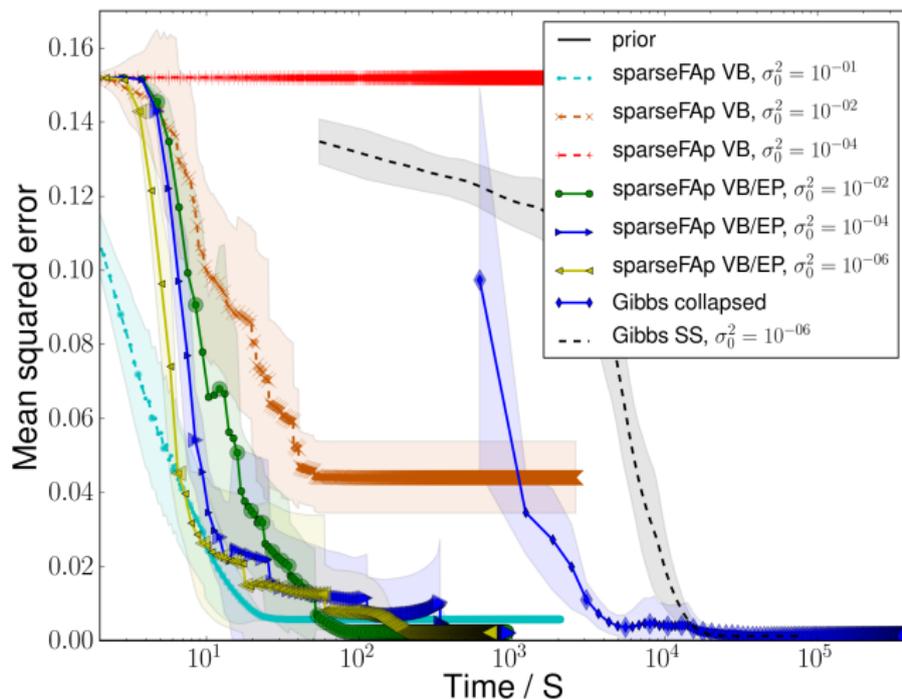
Recovery of the true simulated **network structure**.



Sparse factor analysis

Comparison to MCMC on (small!) simulated dataset

Recovery of the true simulated **factor activations**.



Application to yeast

Dataset

- ▶ Applied the approach to 108 yeast strains (crosses), genotyped and expression profiled in 2 conditions (ethanol and glucose).
- ▶ Alternative types of prior information available
 - ▶ TF binding affinities (Yeasttract).
 - ▶ Pathway information (KEGG).

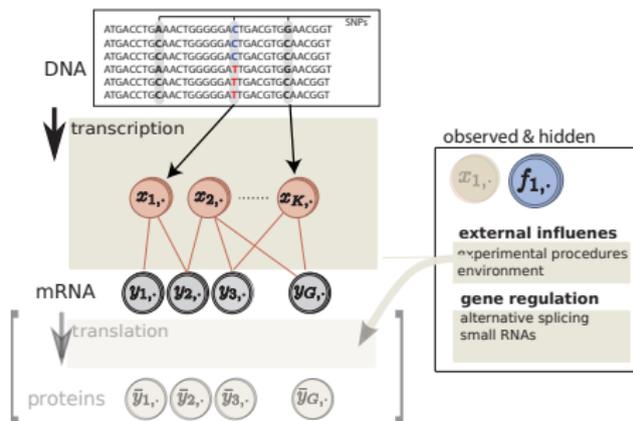
[Smith et al. 2008]

Application to yeast

Factor associations

Biological hypotheses

- ▶ Genetic variation (SNPs) may regulate factor activations.
- ▶ Similarly for the environment condition (glucose/ethanol).
- ▶ Interaction effects between SNPs, factors and genes.

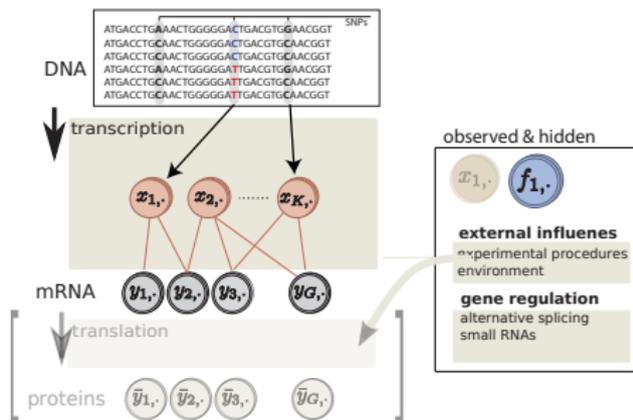


Application to yeast

Factor associations

Biological hypotheses

- ▶ Genetic variation (SNPs) may regulate factor activations.
- ▶ Similarly for the environment condition (glucose/ethanol).
- ▶ Interaction effects between SNPs, factors and genes.

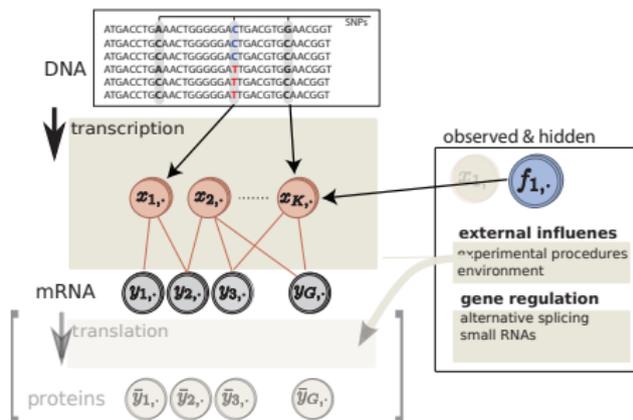


Application to yeast

Factor associations

Biological hypotheses

- ▶ Genetic variation (SNPs) may regulate factor activations.
- ▶ Similarly for the environment condition (glucose/ethanol).
- ▶ Interaction effects between SNPs, factors and genes.

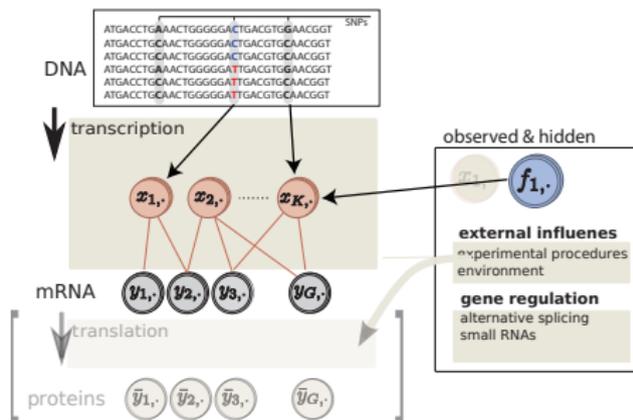


Application to yeast

Factor associations

Biological hypotheses

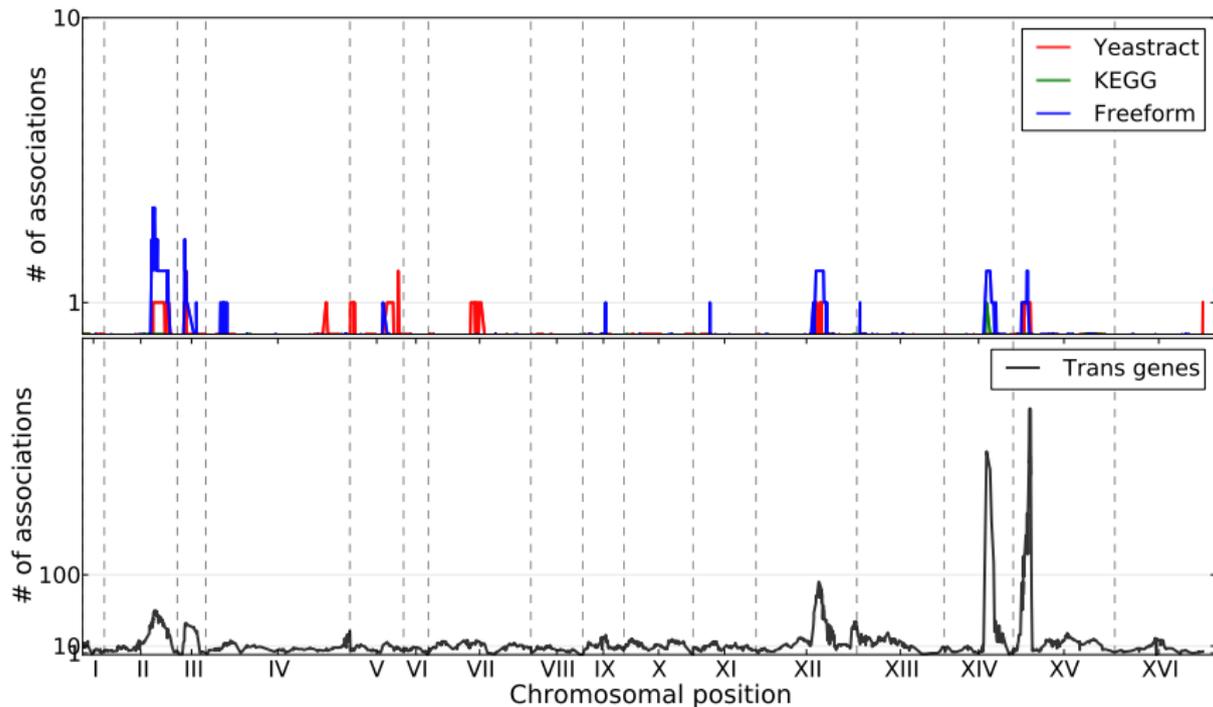
- ▶ Genetic variation (SNPs) may regulate factor activations.
- ▶ Similarly for the environment condition (glucose/ethanol).
- ▶ Interaction effects between SNPs, factors and genes.



Application to yeast

Factor associations

- ▶ Genome-wide association density.

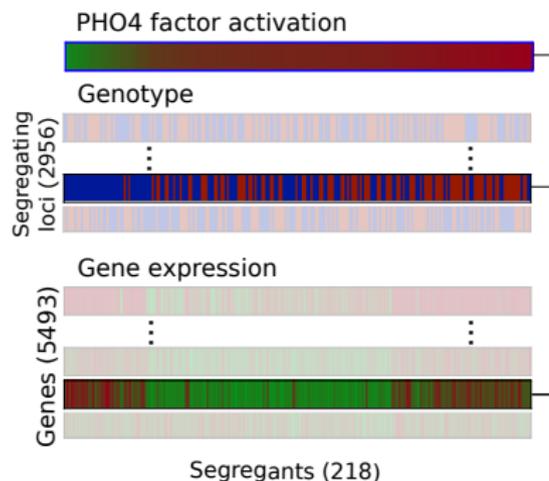


Application to yeast

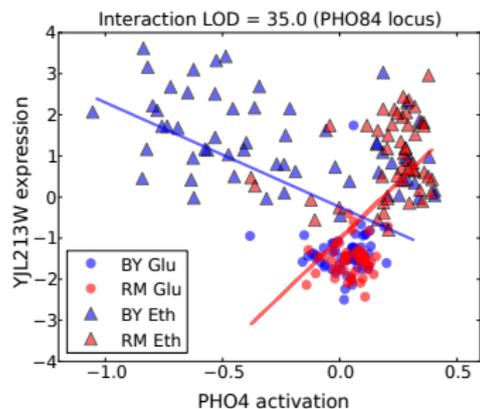
Factor interactions

- ▶ Interaction tests between all gene/SNP/factor triplets.

(c)



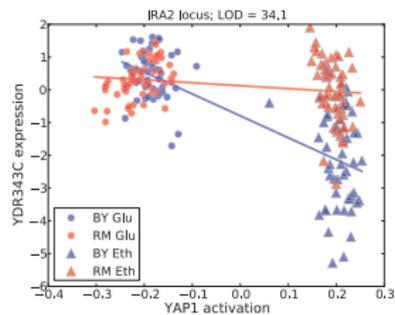
Genotype-factor interaction



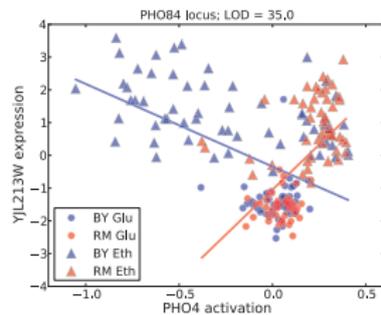
Application to yeast

Factor interactions

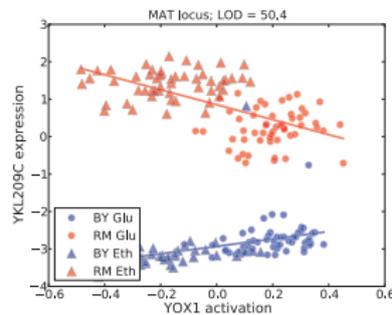
► Example interactions.



(a) YAP1-IRA2 interaction



(b) PHO4-PHO84 interaction

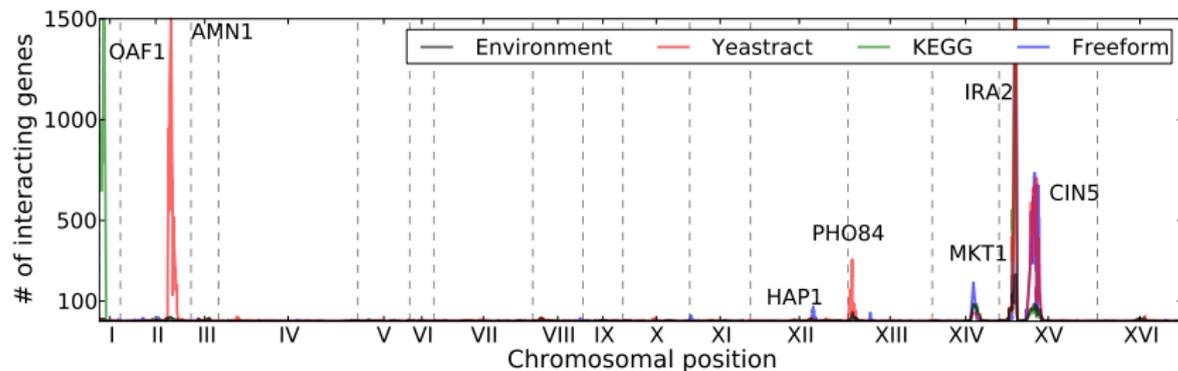


(c) MAT-YOX1 interaction

Application to yeast

Factor interactions

- ▶ Genome-wide interaction density.



Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

- Model

- Applications

Modeling unobserved cellular phenotypes in genetic analyses

- Model

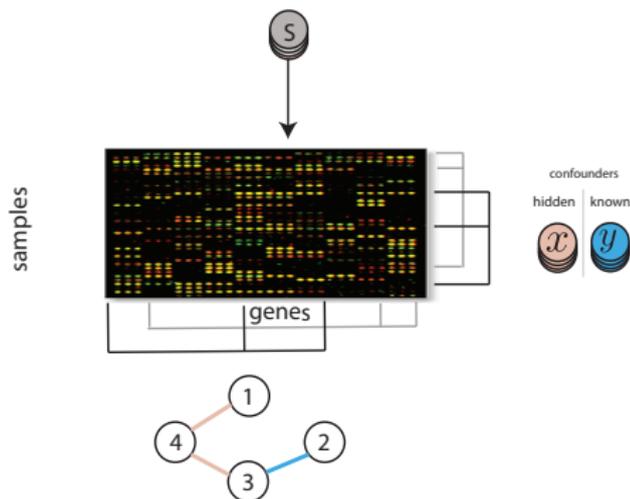
- Applications

A unifying view

Summary

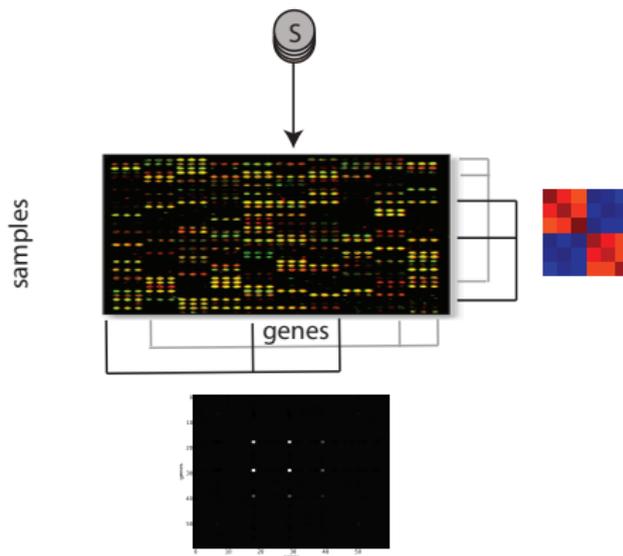
Matrix variate normal distributions

- ▶ Confounders induces structure between samples.
- ▶ Gene regulation induces structure between genes.
- ▶ Matrix variate models for joint correction.
- ▶



Matrix variate normal distributions

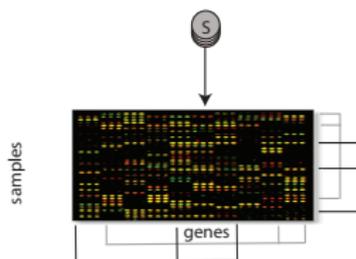
- ▶ Confounders induces structure between samples.
- ▶ Gene regulation induces structure between genes.
- ▶ Matrix variate models for joint correction.



Matrix variate normal distributions

- ▶ Confounders induces structure between samples.
- ▶ Gene regulation induces structure between genes.
- ▶ Matrix variate models for joint correction.
- ▶

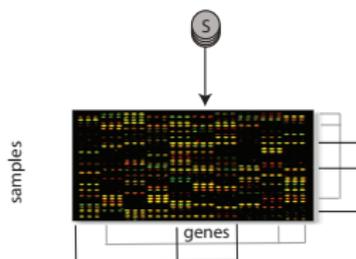
Matrix variate normal distributions



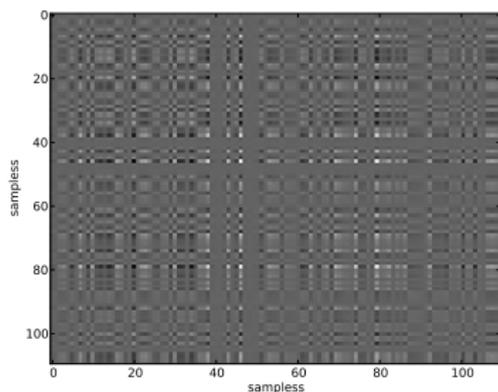
Row covariances (samples) K

Column covariances (genes)
 Λ^{-1}

Matrix variate normal distributions



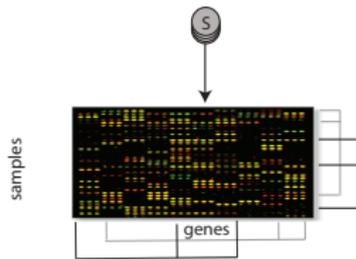
Row covariances (samples) K



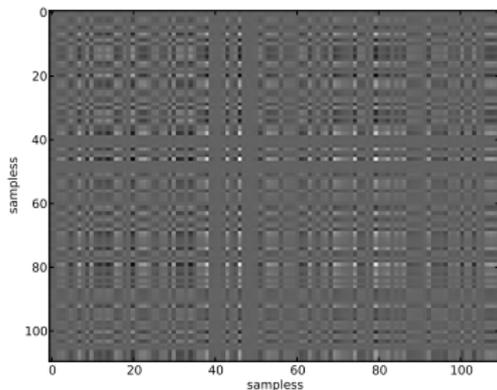
Column covariances (genes)

$$\Lambda^{-1}$$

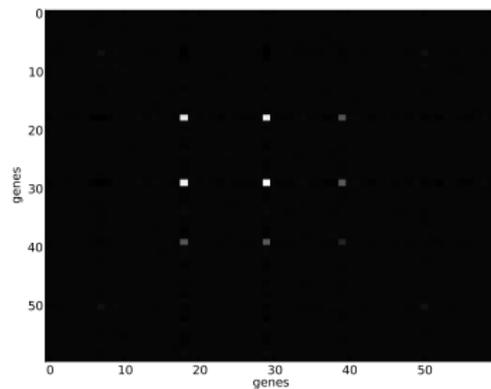
Matrix variate normal distributions



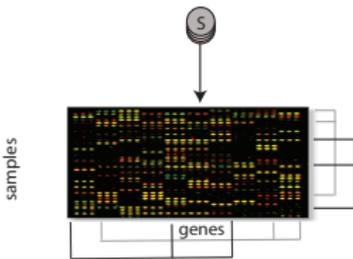
Row covariances (samples) K



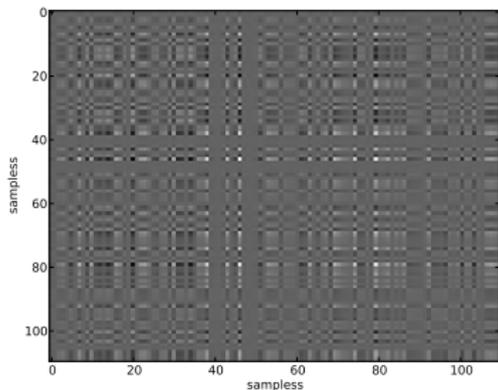
Column covariances (genes) Λ^{-1}



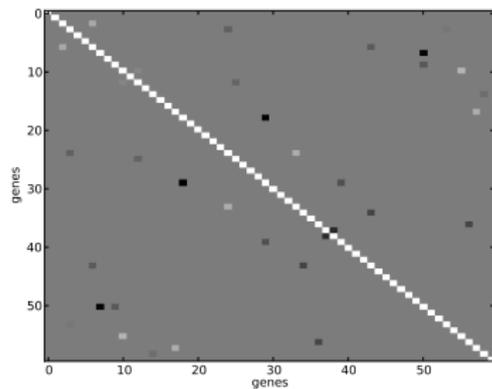
Matrix variate normal distributions



Row covariances (samples) K



Column covariances (genes) Λ^{-1}



Matrix variate distributions

- ▶ Matrix variate distribution with column covariance $\mathbf{\Lambda}^{-1}$ and row covariance \mathbf{K}

$$p(\mathbf{Y} \mid \boldsymbol{\mu} = \mathbf{0}, \mathbf{\Lambda}^{-1}, \mathbf{K}) \propto \exp\left(-0.5\text{tr}\left[\mathbf{\Lambda}\mathbf{Y}^{\top}\mathbf{K}^{-1}\mathbf{Y}\right]\right)$$

- ▶ Equivalent to Kronecker covariance structure on vectorized data

$$p(\text{vec}\mathbf{Y} \mid \boldsymbol{\mu} = \mathbf{0}, \mathbf{\Lambda}^{-1}, \mathbf{K}) = \mathcal{N}\left(\text{vec}\mathbf{Y} \mid \mathbf{0}, \mathbf{\Lambda}^{-1} \otimes \mathbf{K}\right)$$

Matrix variate distributions

- ▶ Matrix variate distribution with column covariance $\mathbf{\Lambda}^{-1}$ and row covariance \mathbf{K}

$$p(\mathbf{Y} \mid \boldsymbol{\mu} = \mathbf{0}, \mathbf{\Lambda}^{-1}, \mathbf{K}) \propto \exp \left(-0.5 \text{tr} \left[\mathbf{\Lambda} \mathbf{Y}^\top \mathbf{K}^{-1} \mathbf{Y} \right] \right)$$

- ▶ Equivalent to Kronecker covariance structure on vectorized data

$$p(\text{vec} \mathbf{Y} \mid \boldsymbol{\mu} = \mathbf{0}, \mathbf{\Lambda}^{-1}, \mathbf{K}) = \mathcal{N} \left(\text{vec} \mathbf{Y} \mid \mathbf{0}, \mathbf{\Lambda}^{-1} \otimes \mathbf{K} \right)$$

Drawing samples from a MN

Kronecker matrix product

1. sample

$$Y \sim \prod_{r=1}^R \prod_{g=1}^G \mathcal{N}(0, 1)$$

2. $Y = \text{chol}(\mathbf{K}) \cdot Y$

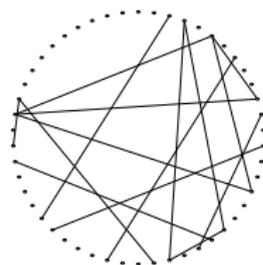
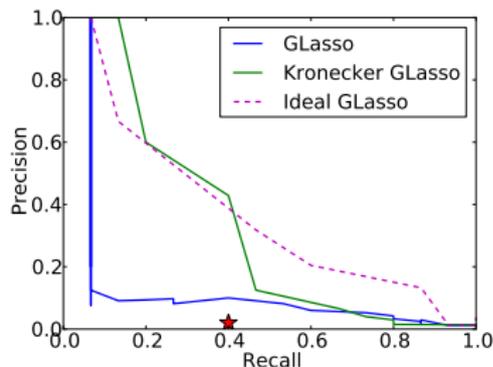
3. $Y = Y \cdot \text{chol}(\mathbf{\Lambda}^{-1})^\top$

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11} \mathbf{B} & a_{12} \mathbf{B} & \dots & a_{1H} \mathbf{B} \\ a_{21} \mathbf{B} & a_{22} \mathbf{B} & \dots & a_{2H} \mathbf{B} \\ \dots & \dots & \dots & \vdots \\ a_{G1} \mathbf{B} & a_{G2} \mathbf{B} & \dots & a_{GH} \mathbf{B} \end{pmatrix}$$

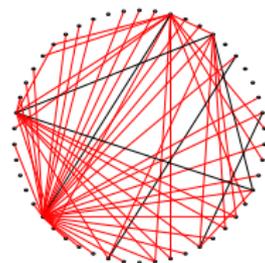
Simulated example

better recovery of the true graph

- Recovery of simulated network. Weak confounding variation (20% variance in row covariance).



ground truth

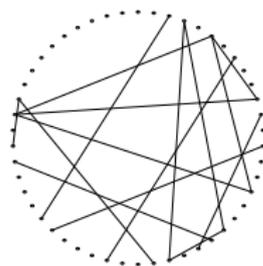
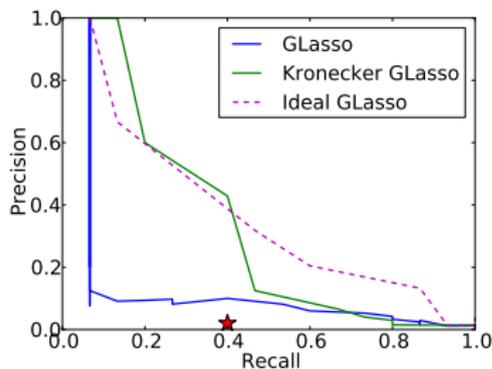


GLasso

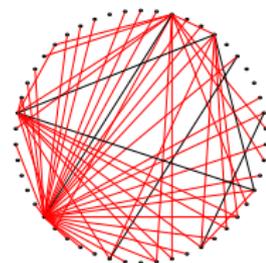
Simulated example

better recovery of the true graph

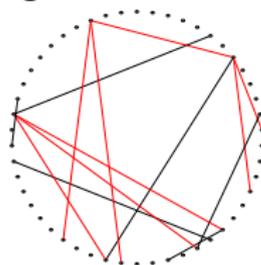
- ▶ Recovery of simulated network. Weak confounding variation (20% variance in row covariance).



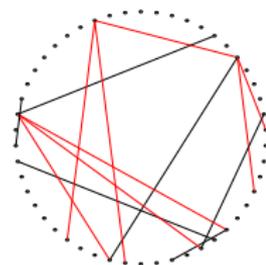
ground truth



GLasso



Kronecker GLasso



Ideal GLasso

Outline

Motivation

Dimension reduction and the Gaussian Process Latent Variable Model (GPLVM)

Modeling hidden confounders in GWAS

- Model

- Applications

Modeling unobserved cellular phenotypes in genetic analyses

- Model

- Applications

A unifying view

Summary

Summary

- ▶ **Latent variables** can have a dramatic effect in GWAs.
- ▶ In expression studies, **confounders can be estimated from data**.
- ▶ Cellular features can be learnt from the expression profiles.
- ▶ Duality of regulatory relationships and confounders in **matrix variate normal models**.