

# Probabilistic Modelling of Expression Variation in Modern eQTL Studies

Max ZwießeBele

Eberhard Karls Universität Tübingen

March 24, 2013

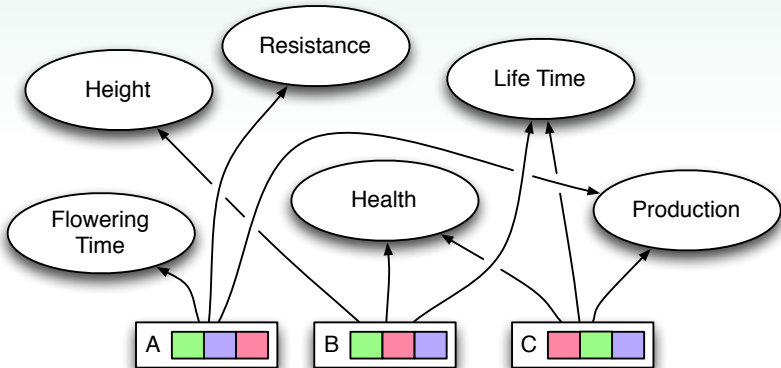


MAX-PLANCK-GESELLSCHAFT

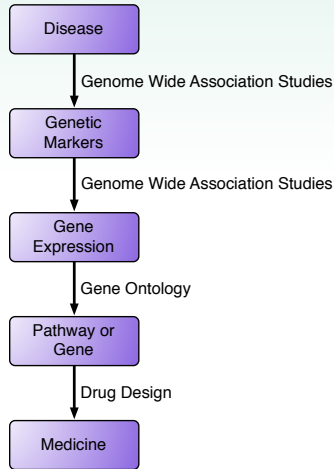


BIOLOGISCHE KYBERNETIK

# Why Genome-wide Association Studies?



# Why Genome-wide Association Studies?



# Outline

- 1 Why Genome-wide Association Studies?
- 2 Genome-wide Association Studies
- 3 Bayesian Modelling of Confounding Variation
- 4 Applications
- 5 Summary and Discussion



# Outline

- 1 Why Genome-wide Association Studies?
- 2 Genome-wide Association Studies
- 3 Bayesian Modelling of Confounding Variation
- 4 Applications
- 5 Summary and Discussion



# Genome-wide Association Studies

Find Function  $f$ , which maps

- Genotype Variants  $\mathbf{S}$
- Gene Expression  $\mathbf{Y}$

$$\Rightarrow \mathbf{Y} = f(\mathbf{S}) .$$

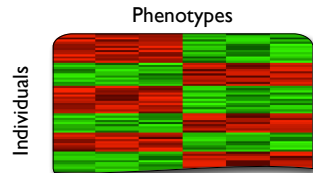


# Genome-wide Association Studies

Find Function  $f$ , which maps

- Genotype Variants  $\mathbf{S}$
- Gene Expression  $\mathbf{Y}$

$$\Rightarrow \mathbf{Y} = f(\mathbf{S}) .$$

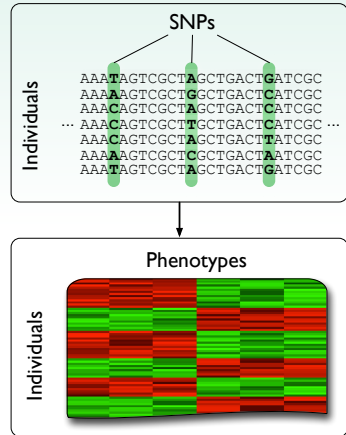


# Genome-wide Association Studies

Find Function  $f$ , which maps

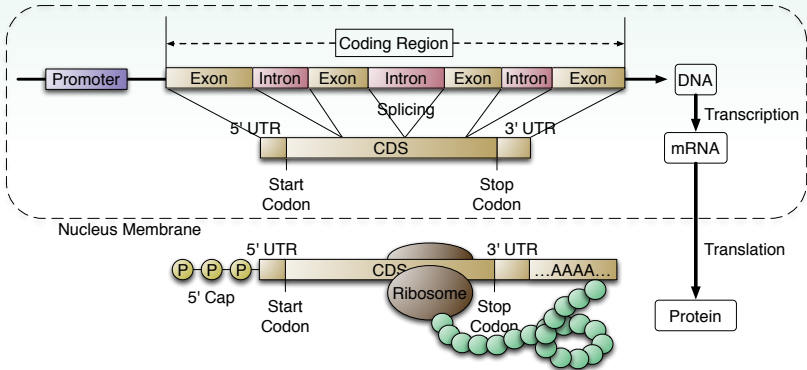
- Genotype Variants  $\mathbf{S}$
- Gene Expression  $\mathbf{Y}$

$$\Rightarrow \mathbf{Y} = f(\mathbf{S}) .$$





# Gene Expression as quantitative Trait



# Confounding Variation (for more see [3])

Confounders:

“Non Genotype Variant Influences on Gene Expression”

Genetically driven:

- Combinatorial Effects of Expressed Genes
- Population Stratification
- Gender

Environmentally driven:

- Treatment Control Experiments
- Exposure to Stress:
  - Physical (pressure, height...)
  - Chemical
  - Human Made (Sea Water quality)



# Confounding Variation (for more see [3])

Confounders:

“Non Genotype Variant Influences on Gene Expression”

Genetically driven:

- Combinatorial Effects of Expressed Genes
- Population Stratification
- Gender

Environmentally driven:

- Treatment Control Experiments
- Exposure to Stress:
  - Physical (pressure, height...)
  - Chemical
  - Human Made (Sea Water quality)



# Confounding Variation (for more see [3])

Confounders:

“Non Genotype Variant Influences on Gene Expression”

Genetically driven:

- Combinatorial Effects of Expressed Genes
- Population Stratification
- Gender

Environmentally driven:

- Treatment Control Experiments
- Exposure to Stress:
  - Physical (pressure, height...)
  - Chemical
  - Human Made (Sea Water quality)



# Confounding Variation (for more see [3])

Confounders:

“Non Genotype Variant Influences on Gene Expression”

Genetically driven:

- Combinatorial Effects of Expressed Genes
- Population Stratification
- Gender

Environmentally driven:

- Treatment Control Experiments
- Exposure to Stress:
  - Physical (pressure, height...)
  - Chemical
  - Human Made (Sea Water quality)



# Confounding Variation (for more see [3])

Confounders:

“Non Genotype Variant Influences on Gene Expression”

Genetically driven:

- Combinatorial Effects of Expressed Genes
- Population Stratification
- Gender

Environmentally driven:

- Treatment Control Experiments
- Exposure to Stress:
  - Physical (pressure, height...)
  - Chemical
  - Human Made (Sea Water quality)



# Confounders and Probabilistic Modelling

Confounding Variation not known:

- Number  $Q$  (Dimension) of Confounders
- Abundance of Variation **X** itself
- Combinatorial Effects uncertain

⇒ Bayesian Modelling of Confounders (Latent Variables)



# Outline

- 1 Why Genome-wide Association Studies?
- 2 Genome-wide Association Studies
- 3 Bayesian Modelling of Confounding Variation**
- 4 Applications
- 5 Summary and Discussion

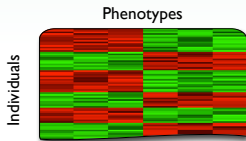




# Notation

Expression Matrix:

$$\mathbf{Y} \in \mathbb{R}^{N \times D}$$



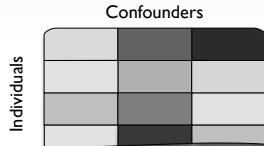
SNP Matrix:

$$\mathbf{S} \in \mathbb{R}^{N \times K}$$



Confounder Matrix:

$$\mathbf{X} \in \mathbb{R}^{N \times Q}$$



$N$  : Samples

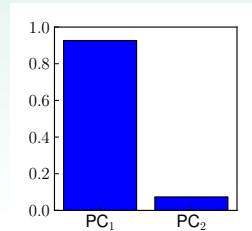
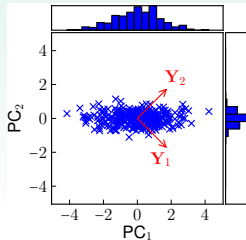
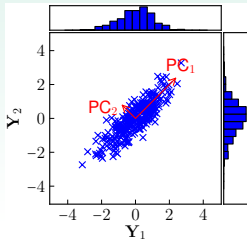
$D$  : Genes

$K$  : SNPs

$Q$  : Confounders



# Principal Component Analysis



$$\mathbf{Y} = \mathbf{XV} \quad , \quad \text{where}$$

$$\mathbf{V}^{-1} = \text{eig}(\mathbf{Y}^T \mathbf{Y})$$

$$\mathbf{X} = \mathbf{YV}_Q^{-1} \quad , \quad \text{where}$$

$$\mathbf{V}_Q^{-1} = \{\mathbf{v}_q : 1 \leq q \leq Q \text{ largest eigenvalues}\}$$



# Gaussian Process Latent Variable Model [4, 5]

Consider Linear Generative Model:

$$\mathbf{Y} = \mathbf{XV} \quad (1)$$

Prior over  $\mathbf{V}$ , instead of  $\mathbf{X}$ :

$$p(\mathbf{V}) = \prod_{d=1}^D \mathcal{N}(\mathbf{v}_d | \mathbf{0}, \mathbf{I}_Q) \quad (2)$$

Integrate out  $\mathbf{V}$  to get:

$$\begin{aligned} p(\mathbf{Y}|\mathbf{X}) &= \int p(\mathbf{Y}|\mathbf{X}, \mathbf{V}) p(\mathbf{V}|\mathbf{X}) d\mathbf{V} \\ &= \prod_{d=1}^D \mathcal{N}(\mathbf{y}_d | \mathbf{0}, \mathbf{K}^{NN} + \beta^{-1} \mathbf{I}_N) \quad , \text{ where} \end{aligned} \quad (3)$$

$\mathbf{K}_{ij}^{NN} = k(\mathbf{x}_i, \mathbf{x}_j)$  Covariance Function.



# Gaussian Process Latent Variable Model [4, 5]

Learn Latent Variables  $\mathbf{X}$  through Maximum A Posteriori (MAP):

$$\{\hat{\mathbf{X}}, \hat{\boldsymbol{\theta}}\}_{\text{MAP}} = \arg \max_{\mathbf{X}, \boldsymbol{\theta}} \ln p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) \quad , \quad \text{where} \quad (5)$$

$$\ln p(\mathbf{Y}|\mathbf{X}) = -\frac{DN}{2} \ln(2\pi) - \frac{D}{2} \ln |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) \quad , \quad \text{with} \quad (6)$$

$$\mathbf{K} = \mathbf{K}^{NN} + \beta^{-1} \mathbf{I}_N \quad (7)$$

But, we want Bayesian estimate

$$p(\mathbf{Y}) = \int p(\mathbf{Y}|\mathbf{X}) p(\mathbf{X}) d\mathbf{X} \quad (8)$$



## Variational Bayesian GPLVM [1]

$$p(\mathbf{Y}) = \int p(\mathbf{Y}|\mathbf{X})p(\mathbf{X}) d\mathbf{X} \quad (9)$$

Intractable in  $\mathbf{X} \Rightarrow$  Variational Approximation  $q(\mathbf{X})$ :

$$\ln p(\mathbf{Y}) \geq \mathcal{F}(q) = \int q(\mathbf{X}) \ln \frac{p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})}{q(\mathbf{X})} d\mathbf{X} \quad (10)$$

$$= \int q(\mathbf{X}) \ln p(\mathbf{Y}|\mathbf{X}) d\mathbf{X} - \int q(\mathbf{X}) \ln \frac{p(\mathbf{X})}{q(\mathbf{X})} d\mathbf{X} \quad (11)$$

$$= \int q(\mathbf{X}) \ln p(\mathbf{Y}|\mathbf{X}) d\mathbf{X} - \text{KL}(q||p) \quad (12)$$

$$= \tilde{\mathcal{F}}(q) - \text{KL}(q||p)$$

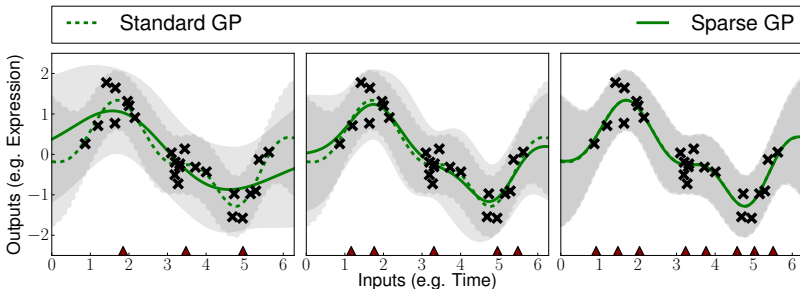


# Variational Bayesian GPLVM [1]

$$\ln p(\mathbf{Y}) \geq \mathcal{F}(q) = \tilde{\mathcal{F}}(q) - \text{KL}(q||p) \quad (14)$$

$$\tilde{\mathcal{F}}(q) = \int q(\mathbf{X}) \ln p(\mathbf{Y}|\mathbf{X}) d\mathbf{X} \quad (15)$$

Sparse Approximation for  $p(\mathbf{Y}|\mathbf{X})$ :

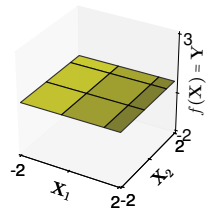
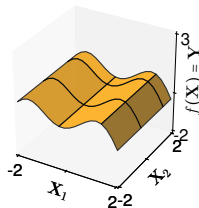
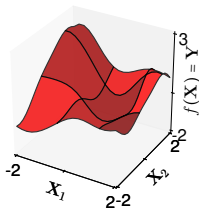


# Automatic Relevance Determination [6]

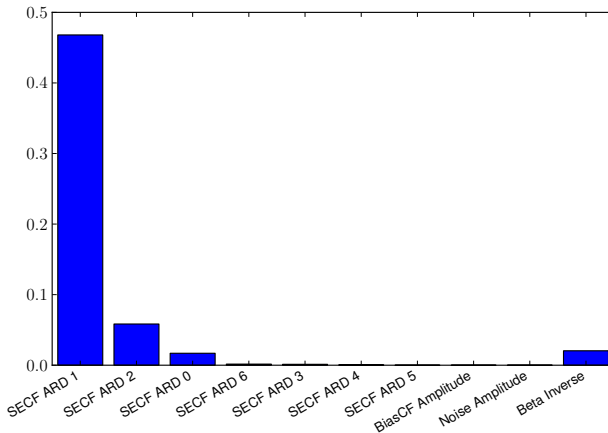
Consider Squared Exponential Covariance Function

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left\{ \frac{1}{2} \boldsymbol{\alpha}^\top \|\mathbf{x} - \mathbf{x}'\|^2 \right\} \quad (16)$$

$\boldsymbol{\alpha} \in \mathbb{R}^Q$  contains Relevance for each Dimension  $1 \leq q \leq Q$  of  $\mathbf{X}$   
 $\mathbf{X} \in \mathbb{R}^{N \times Q}$  with rows  $\mathbf{x} \in \mathbb{R}^Q$

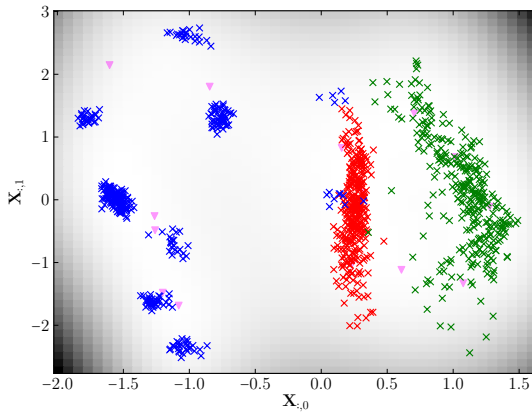


# Intuition

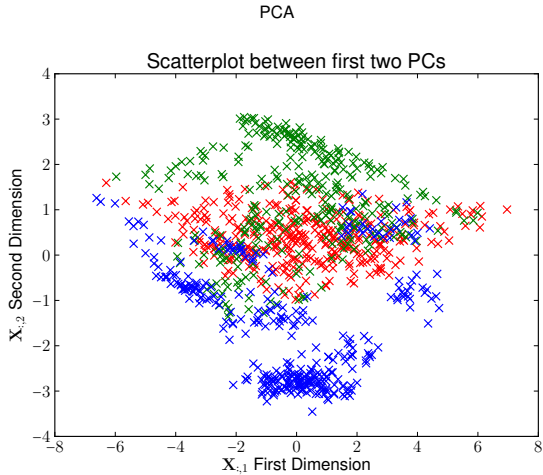




# Intuition

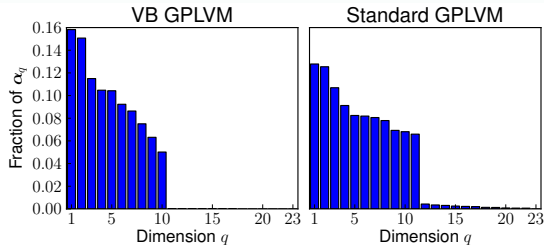


# Intuition



# Learning the right Dimensionality

Simulated Dataset (from Linear Mixed Model) with  $Q = 10$   
Confounders:



# Outline

- 1 Why Genome-wide Association Studies?
- 2 Genome-wide Association Studies
- 3 Bayesian Modelling of Confounding Variation
- 4 Applications**
- 5 Summary and Discussion



# PANAMA

Linear Mixed Model Assumed:

$$\underbrace{\mathbf{Y}}^{N \times D} = \underbrace{\mathbf{S}}^{N \times K} \mathbf{W} + \underbrace{\mathbf{X}}^{N \times Q} \mathbf{V} \quad (17)$$

$N$  : Samples       $D$  : Genes       $K$  : SNPs       $Q$  : Confounders

Iterative Approach through GPLVM:

- ① Learn Confounders by GPLVM
- ② Add Confounders correlated with Genotype
- ③ Fit the new Model
- ④ Continue with 1. until no Confounders are added



# Variational Bayesian GPLVM for Confounder Detection

Same Linear Mixed Model Assumed:

$$\mathbf{Y} = \mathbf{S}\mathbf{W} + \mathbf{X}\mathbf{V} \quad (18)$$

No Iterative Approach necessary

⇒ ARD + Bayesian Approximation

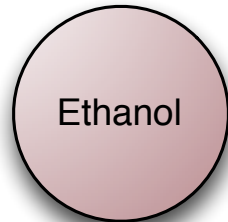
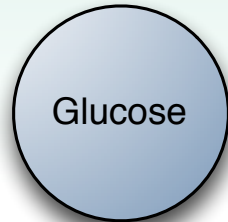
- Learn Confounders  $\mathbf{X}$  through Variational Bayesian GPLVM
- Fit Linear Mixed Model with learned Confounders.



# Yeast Dataset [7]

## *Saccharomyces Cerevisiae*:

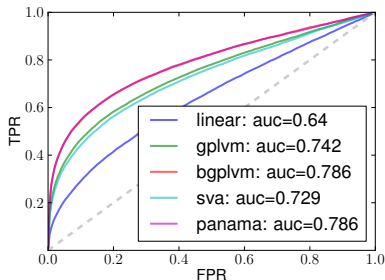
- Expression Profiled  $N = 109$   
Samples over  $D = 5493$   
Genes
- $K = 2956$  SNPs Genotyped
- 2 Environments:  
Glucose and Ethanol
- We only use Glucose Part.



# Simulated Dataset by [2]

- Simulated Dataset founded on real Yeast [7] eQTL Experiment

⇒ Ground Truth for Confounding Variation and Associations:

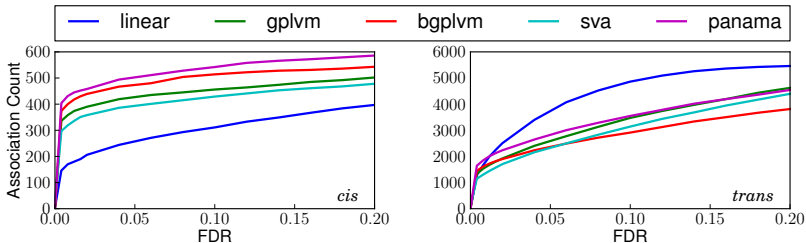




# Simulated Dataset by [2]

● *cis* and *trans* Associations:

⇒ Less False Positives:



# Manifold Relevance Determination for eQTL Studies

Combine Variational Bayesian GPLVMs to determine Private/Shared Information:

$$\ln p(\mathbf{Y}^{\kappa}) \geq \ln \sum_{\ell \in \kappa} \int q(\mathbf{X}) \ln p(\mathbf{Y}^{\ell} | \mathbf{X}) d\mathbf{X} - \text{KL}(q \| p) \quad , \text{ where} \quad (19)$$

$$\mathbf{Y}^{\kappa} = \{\mathbf{Y}^{\ell} : \ell \in \kappa\} \quad (20)$$

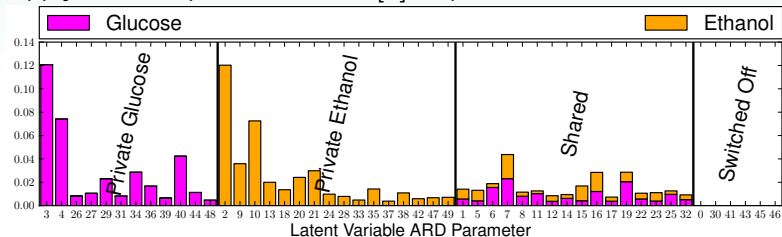
Thus,

- One Confounder Matrix  $\mathbf{X}$
- ARD Parameters  $\alpha^{\ell}$  per Experiment  $\mathbf{Y}^{\ell}$



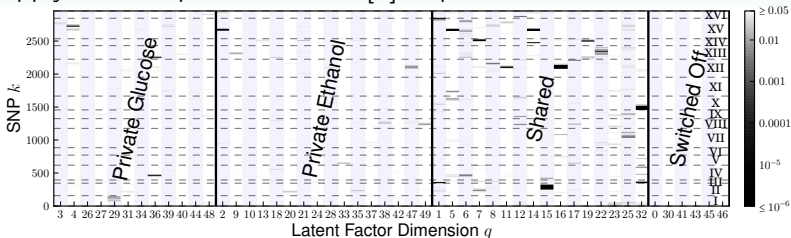
# Manifold Relevance Determination for eQTL Studies

Apply to both parts of Smith [7] Experiment:



# Manifold Relevance Determination for eQTL Studies

Apply to both parts of Smith [7] Experiment:



# Manifold Relevance Determination for eQTL Studies

Private Glucose	#Hits	GO
cellular macromolecule catabolic process	7	0044265
macromolecule catabolic process	7	0009057
ncRNA processing	6	0034470
ncRNA metabolic process	6	0034660
tRNA metabolic process	4	0006399
Private Ethanol	#Hits	GO
response to abiotic stimulus	6	0009628
regulation of cell cycle	5	0051726
pseudohyphal growth	4	0007124
cell growth	4	0016049
growth of unicellular organism as a thread of attached cells	4	0070783
regulation of cell cycle process	4	0010564
Shared	#Hits	GO
cell cycle	42	0007049
cell cycle process	34	0022402
tRNA metabolic process	14	0006399
cytokinesis	13	0000910
DNA replication	13	0000260
coenzyme metabolic process	13	0006732

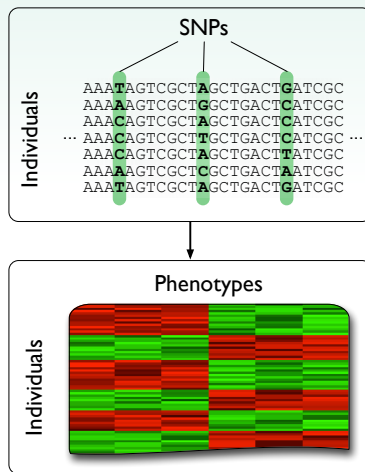


# Outline

- 1 Why Genome-wide Association Studies?
- 2 Genome-wide Association Studies
- 3 Bayesian Modelling of Confounding Variation
- 4 Applications
- 5 Summary and Discussion

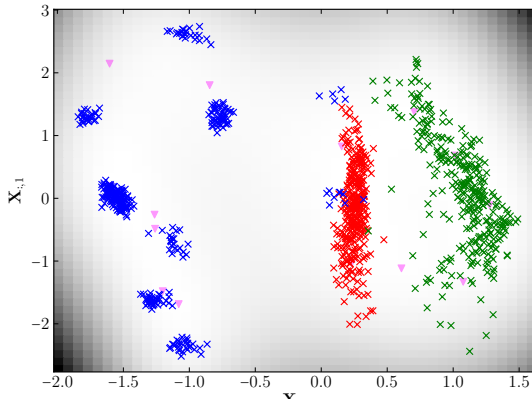


# Summary



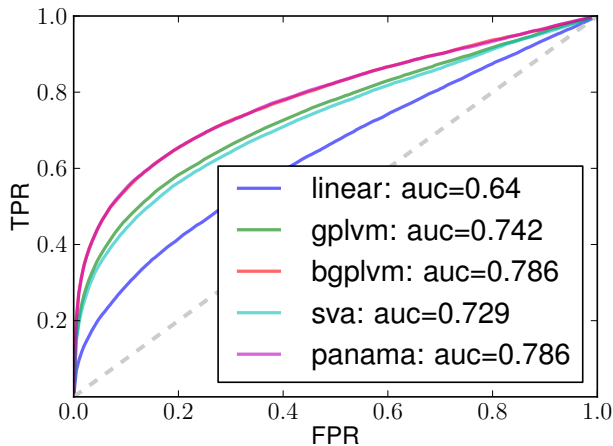
# Summary

$$p(\mathbf{Y}) = \int p(\mathbf{Y}|\mathbf{X})p(\mathbf{X}) d\mathbf{X} \quad (21)$$

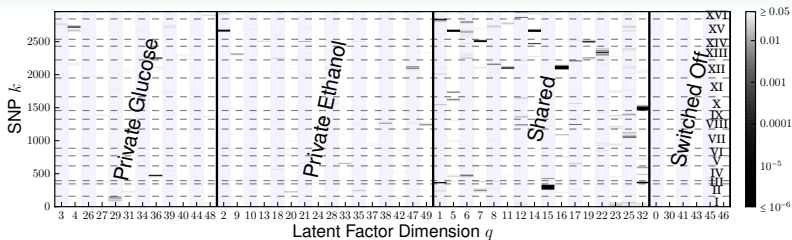




# Summary



# Summary



# Discussion



# References I



A.C. Damianou, M.K. Titsias, and N.D. Lawrence.  
Variational Gaussian Process Dynamical Systems.  
*Neural Information Processing System (NIPS)*, 2011.



Nicolo Fusi, Oliver Stegle, and Neil D. Lawrence.  
Accurate modeling of confounding variation in eQTL studies leads to a great increase in power to detect trans-regulatory effects.  
*Nature Precedings*, 2011.



Joel N Hirschhorn and Mark J Daly.  
Genome-wide association studies for common diseases and complex traits.  
*Nat Rev Genet*, 6(2):95–108, Feb 2005.



N.D. Lawrence.  
Gaussian process latent variable models for visualisation of high dimensional data.  
*Advances in neural information processing systems*, 16:329–336, 2004.



Neil Lawrence.  
Probabilistic non-linear principal component analysis with Gaussian process latent variable models.  
*Journal of Machine Learning Research*, November 2005.



## References II



Carl Edward Rasmussen.  
*Gaussian Processes for Machine Learning.*  
The MIT Press, 2006.



Erin N Smith and Leonid Kruglyak.  
Gene-environment interaction in yeast gene expression.  
*PLoS Biol*, 6(4):e83, Apr 2008.



# From DNA to Protein

