# Linear models II: multi locus models and variance decomposition

### Christoph Lippert<sup>1</sup> Oliver Stegle<sup>2</sup>

<sup>1</sup> Microsoft Research, Los Angeles, USA
 <sup>2</sup> Max-Planck-Institutes Tübingen, Germany



Basel 09. September 2012



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

September 2012 1

4 D N 4 B N 4 B N 4 B N

#### Overview

Single marker association model with random effect term



#### Shortcomings

- Weak effects are not captured by single-marker analysis.
- Complex traits are controlled by more than a single SNP.

#### Overview

Single marker association model with random effect term



- Shortcomings
  - Weak effects are not captured by single-marker analysis.
  - Complex traits are controlled by more than a single SNP.



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

#### Overview

Single marker association model with random effect term



- Shortcomings
  - Weak effects are not captured by single-marker analysis.
  - Complex traits are controlled by more than a single SNP.



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

Generalization to multiple genetic factors



Generalization to multiple genetic factors

$$\mathbf{y} = \sum_{\substack{s=1\\\text{genetic effect}}}^{S} \mathbf{x}_{s}\beta_{s} + \underbrace{\mathbf{u}}_{\text{random effect covariates}} + \underbrace{\boldsymbol{\epsilon}}_{\text{noise}}$$

Challenge: N << S: explicit estimation of all β<sub>s</sub> is not feasible.
 Solutions

- Regularize  $\beta_s$  (Ridge regression, LASSO)
- Variance component modeling

Wu et al., 2011

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition

Generalization to multiple genetic factors

$$\mathbf{y} = \sum_{\substack{s=1\\\text{genetic effect}}}^{S} \mathbf{x}_{s}\beta_{s} + \underbrace{\mathbf{u}}_{\text{random effect covariates}} + \underbrace{\boldsymbol{\epsilon}}_{\text{noise}}$$

- ▶ Challenge:  $N \ll S$ : explicit estimation of all  $\beta_s$  is not feasible.
- Solutions
  - Regularize  $\beta_s$  (Ridge regression, LASSO)
  - Variance component modeling

[Wu et al., 2011]

Generalization to multiple genetic factors

$$\mathbf{y} = \sum_{\substack{s=1\\\text{genetic effect}}}^{S} \mathbf{x}_{s}\beta_{s} + \underbrace{\mathbf{u}}_{\text{random effect covariates}} + \underbrace{\boldsymbol{\epsilon}}_{\text{noise}}$$

- ▶ Challenge:  $N \ll S$ : explicit estimation of all  $\beta_s$  is not feasible.
- Solutions
  - Regularize  $\beta_s$  (Ridge regression, LASSO)
  - Variance component modeling

[Wu et al., 2011]

C. Lippert & O. Stegle I

< ロ > < 同 > < 回 > < 回 > < 回 > <

Outline

### Outline

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

▲□▶ ▲□▶ ▲ □▶ ▲ □ ● ● ● ●

#### Outline

#### Variance component models

#### Variance component models for correlated traits

Mixed model Lasso

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition

#### Multi locus models

Random effect models

For now, let's drop the random effect term

$$\mathbf{y} = \sum_{s=1}^{S} \mathbf{x}_s \beta_s + \boldsymbol{\epsilon}.$$

For mathematical convenience, we choose a shared Gaussian prior on the weights and Gaussian noise

$$p(\beta_1, \dots, \beta_S) = \prod_{s=1}^{S} \mathcal{N}\left(\beta_s \mid 0, \sigma_g^2\right) \quad p(\boldsymbol{\epsilon}) = \mathcal{N}\left(\boldsymbol{\epsilon} \mid \boldsymbol{0}, \sigma_e^2 \mathbf{I}\right)$$

• Marginalize out the weights  $\beta_1, \ldots, \beta_S$ 

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\mathbf{y} \left| \sum_{s=1}^{S} \mathbf{x}_{s} \beta_{s}, \sigma_{e}^{2} \mathbf{I} \right) \prod_{s=1}^{S} \mathcal{N}\left(\beta_{s} | 0, \sigma_{g}^{2}\right) d\boldsymbol{\beta} \right.$$
$$= \mathcal{N}\left(\mathbf{y} \left| \mathbf{0}, \sigma_{g}^{2} \sum_{s=1}^{S} \mathbf{x}_{s} \mathbf{x}_{s}^{\mathrm{T}} + \sigma_{e}^{2} \mathbf{I} \right.\right)$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

#### Multi locus models

Random effect models

For now, let's drop the random effect term

$$\mathbf{y} = \sum_{s=1}^{S} \mathbf{x}_s \beta_s + \boldsymbol{\epsilon}.$$

 For mathematical convenience, we choose a shared Gaussian prior on the weights and Gaussian noise

$$p(\beta_1, \dots, \beta_S) = \prod_{s=1}^{S} \mathcal{N}\left(\beta_s \mid 0, \sigma_g^2\right) \quad p(\boldsymbol{\epsilon}) = \mathcal{N}\left(\boldsymbol{\epsilon} \mid \boldsymbol{0}, \sigma_e^2 \mathbf{I}\right)$$

• Marginalize out the weights  $\beta_1, \ldots, \beta_S$ 

$$p(\mathbf{y} | \mathbf{X}, \sigma_{\mathbf{g}}^{2}, \sigma_{\mathbf{e}}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\mathbf{y} \left| \sum_{s=1}^{S} \mathbf{x}_{s} \beta_{s}, \sigma_{\mathbf{e}}^{2} \mathbf{I} \right) \prod_{s=1}^{S} \mathcal{N}\left(\beta_{s} | 0, \sigma_{\mathbf{g}}^{2}\right) \mathrm{d}\boldsymbol{\beta} \right.$$
$$= \mathcal{N}\left(\mathbf{y} \left| \mathbf{0}, \sigma_{\mathbf{g}}^{2} \sum_{s=1}^{S} \mathbf{x}_{s} \mathbf{x}_{s}^{\mathrm{T}} + \sigma_{\mathbf{e}}^{2} \mathbf{I} \right.\right)$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

#### Multi locus models

Random effect models

For now, let's drop the random effect term

$$\mathbf{y} = \sum_{s=1}^{S} \mathbf{x}_s \beta_s + \boldsymbol{\epsilon}.$$

 For mathematical convenience, we choose a shared Gaussian prior on the weights and Gaussian noise

$$p(\beta_1, \dots, \beta_S) = \prod_{s=1}^{S} \mathcal{N}\left(\beta_s \mid 0, \sigma_g^2\right) \quad p(\boldsymbol{\epsilon}) = \mathcal{N}\left(\boldsymbol{\epsilon} \mid \boldsymbol{0}, \sigma_e^2 \mathbf{I}\right)$$

• Marginalize out the weights  $\beta_1, \ldots, \beta_S$ 

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\mathbf{y} \left| \sum_{s=1}^{S} \mathbf{x}_{s} \beta_{s}, \sigma_{e}^{2} \mathbf{I} \right) \prod_{s=1}^{S} \mathcal{N}\left(\beta_{s} | 0, \sigma_{g}^{2}\right) d\boldsymbol{\beta} \right.$$
$$= \mathcal{N}\left(\mathbf{y} \left| \mathbf{0}, \sigma_{g}^{2} \sum_{s=1}^{S} \mathbf{x}_{s} \mathbf{x}_{s}^{T} + \sigma_{e}^{2} \mathbf{I} \right)$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

#### Multi locus models Remarks

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \mathcal{N}\left(\mathbf{y} | \mathbf{0}, \sigma_{g}^{2} \sum_{\substack{s=1\\\mathbf{K}_{g}}}^{S} \mathbf{s}_{s} \mathbf{x}_{s}^{\mathrm{T}} + \sigma_{e}^{2} \mathbf{I}\right)$$
(1)



- Closely related to Kinship explaining population structure.
- Inference can be done my maximum likelihood.
- The ratio of σ<sup>2</sup><sub>g</sub> and σ<sup>2</sup><sub>e</sub> defines the narrow sense heritability of the trait



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

< ロ > < 同 > < 回 > < 回 > < 回 > <

#### Multi locus models Remarks

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_{g}^{2} \underbrace{\sum_{s=1}^{S} \mathbf{s}_{s} \mathbf{x}_{s}^{T}}_{\mathbf{K}_{g}} + \sigma_{e}^{2} \mathbf{I})$$
(1)



- Closely related to Kinship explaining population structure.
- Inference can be done my maximum likelihood.
- The ratio of σ<sup>2</sup><sub>g</sub> and σ<sup>2</sup><sub>e</sub> defines the narrow sense heritability of the trait



Image: A match a ma

#### Multi locus models Remarks

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_{g}^{2} \sum_{\substack{s=1\\\mathbf{K}_{g}}}^{S} \mathbf{s}_{s} \mathbf{x}_{s}^{\mathrm{T}} + \sigma_{e}^{2} \mathbf{I})$$
(1)



- Closely related to Kinship explaining population structure.
- Inference can be done my maximum likelihood.
- The ratio of σ<sup>2</sup><sub>g</sub> and σ<sup>2</sup><sub>e</sub> defines the narrow sense heritability of the trait



Image: A match a ma

#### Multi locus models Remarks

$$p(\mathbf{y} | \mathbf{X}, \sigma_{g}^{2}, \sigma_{e}^{2}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_{g}^{2} \sum_{\substack{s=1\\\mathbf{K}_{g}}}^{S} \mathbf{s}_{s} \mathbf{x}_{s}^{\mathrm{T}} + \sigma_{e}^{2} \mathbf{I})$$
(1)



- Closely related to Kinship explaining population structure.
- Inference can be done my maximum likelihood.
- The ratio of σ<sup>2</sup><sub>g</sub> and σ<sup>2</sup><sub>e</sub> defines the narrow sense heritability of the trait

$$h = \frac{\sigma_{\rm g}^2}{\sigma_{\rm g}^2 + \sigma_{\rm e}^2}.$$



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

Heritability Heritability estimated on 107 *A. thaliana* phenotypes

Global genetic heritability



Linear models II: multi-locus models and variance decomposition

September 2012 8

Heritability Heritability estimated on 107 A. thaliana phenotypes

Estimate can be restricted to a genomic region such as a single chromosome, etc.

$$\mathcal{N}\!\left(\mathbf{y} \,|\, \mathbf{0}, \sigma_{\mathrm{g}}^2 \sum_{s \in \mathsf{Chrom}} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} \!+\! \sigma_{\mathrm{e}}^2 \mathbf{I}\right)$$



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

### Window-based composite variance analysis Region-based testing

- Just fitting a particular region ignores the genome-wide context
- Variance dissection with region-based separation



- Explained variance components can be read off subject to suitable normalization of the covariances K<sub>w</sub> and K<sub>q</sub>.
- "Local" heritability

$$h(W) = \frac{\sigma_w^2}{\sigma_w^2 + \sigma_g^2 + \sigma_e^2}$$

C. Lippert & O. Stegle

### Window-based composite variance analysis Region-based testing

- Just fitting a particular region ignores the genome-wide context
- Variance dissection with region-based separation

$$p(\mathbf{y} \mid W) = \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \sigma_w^2 \underbrace{\sum_{s \in W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}}}_{\mathbf{K}_w} + \sigma_g^2 \underbrace{\sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}}}_{\mathbf{K}_g} + \sigma_e^2 \mathbf{I})$$

- Explained variance components can be read off subject to suitable normalization of the covariances K<sub>w</sub> and K<sub>g</sub>.
- "Local" heritability

$$h(W) = \frac{\sigma_w^2}{\sigma_w^2 + \sigma_g^2 + \sigma_e^2}$$

C. Lippert & O. Stegle

#### Window-based composite variance analysis



### Window-based composite variance analysis Significance testing

- Analogously to fixed effect testing, the significance of a specific window can be tested.
- Likelihood-ratio statistics to score the relevance of a particular genomic region W

$$\mathsf{LOD}(W) = \frac{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \sigma_w^2 \sum_{s \in W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \qquad \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}$$

 P-values can be obtained from permutation statistics or analytical approximation (variants of score tests or likelihood ratio tests).

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 12

(日) (周) (三) (三)

#### Window-based composite variance analysis Significance testing

- Analogously to fixed effect testing, the significance of a specific window can be tested.
- $\blacktriangleright$  Likelihood-ratio statistics to score the relevance of a particular genomic region W

$$\mathsf{LOD}(W) = \frac{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \sigma_w^2 \sum_{s \in W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \qquad \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}$$

 P-values can be obtained from permutation statistics or analytical approximation (variants of score tests or likelihood ratio tests).

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 12

### Window-based composite variance analysis Significance testing

- Analogously to fixed effect testing, the significance of a specific window can be tested.
- Likelihood-ratio statistics to score the relevance of a particular genomic region W

$$\mathsf{LOD}(W) = \frac{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \sigma_w^2 \sum_{s \in W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}{\mathcal{N}\left(\mathbf{y} \mid \mathbf{0}, \qquad \sigma_g^2 \sum_{s \notin W} \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}} + \sigma_{\mathrm{e}}^2 \mathbf{I}\right)}$$

 P-values can be obtained from permutation statistics or analytical approximation (variants of score tests or likelihood ratio tests).

Phenotype prediction

Best linear unbiased prediction

- Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.
- ►  $P(\mathbf{y}^* | \mathbf{y}) = \mathcal{N}\left(\mathbf{y}^* | \boldsymbol{\mu}^*, \sigma_{g}^2 \mathbf{V}_{g}^* + \sigma_{e}^2 \mathbf{I}\right)$ ► Predicitve mean:  $\boldsymbol{\mu}^* = \underbrace{\sigma_{g}^2 \mathbf{K}_{g}^{\star,:} \left(\sigma_{g}^2 \mathbf{K}_{g} + \sigma_{e}^2 \mathbf{I}\right)^{-1} \mathbf{y}}_{\text{BLUP}}$ ► Predictive Variance:  $\mathbf{V}^* = \mathbf{K}^{\star,*} - \sigma^2 \mathbf{K}^{\star,:} \left(\sigma^2 \mathbf{K}_{e} + \sigma^2 \mathbf{I}\right)^{-1} \mathbf{K}^{:,*}$

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

- 4 回 ト 4 三 ト 4 三 ト

Phenotype prediction

Best linear unbiased prediction

Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

$$\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N}\left(\mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}_{g}^{\star} + \sigma_{e}^{2} \mathbf{I}\right)$$

► Predicitve mean:  $\mu^* = \underbrace{\sigma_g^2 \mathbf{K}_g^{\star,:} \left(\sigma_g^2 \mathbf{K}_g + \sigma_e^2 \mathbf{I}\right)^{-1} \mathbf{y}}_{\mathbf{y}}$ 

► Predictive Variance:  $\mathbf{V}_{g}^{\star} = \mathbf{K}_{g}^{\star,\star} - \sigma_{g}^{2}\mathbf{K}_{g}^{\star,:} \left(\sigma_{g}^{2}\mathbf{K}_{g} + \sigma_{e}^{2}\mathbf{I}\right)^{-1}\mathbf{K}_{g}^{:,\star}$ 

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

### Phenotype prediction

Best linear unbiased prediction

Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

$$\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N}\left(\mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}_{g}^{\star} + \sigma_{e}^{2} \mathbf{I}\right)$$

► Predicitve mean: 
$$\mu^{\star} = \underbrace{\sigma_{g}^{2} \mathbf{K}_{g}^{\star,:} \left(\sigma_{g}^{2} \mathbf{K}_{g} + \sigma_{e}^{2} \mathbf{I}\right)^{-1} \mathbf{y}}_{\mathsf{BLUP}}$$

► Predictive Variance:  $\mathbf{V}_{g}^{\star} = \mathbf{K}_{g}^{\star,\star} - \sigma_{g}^{2}\mathbf{K}_{g}^{\star,:} \left(\sigma_{g}^{2}\mathbf{K}_{g} + \sigma_{e}^{2}\mathbf{I}\right)^{-1}\mathbf{K}_{g}^{:,\star}$ 

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

### Phenotype prediction

Best linear unbiased prediction

Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

► 
$$P(\mathbf{y}^* | \mathbf{y}) = \mathcal{N} \left( \mathbf{y}^* | \boldsymbol{\mu}^*, \sigma_g^2 \mathbf{V}_g^* + \sigma_e^2 \mathbf{I} \right)$$
  
► Predicitve mean:  $\boldsymbol{\mu}^* = \underbrace{\sigma_g^2 \mathbf{K}_g^{\star,:} \left( \sigma_g^2 \mathbf{K}_g + \sigma_e^2 \mathbf{I} \right)^{-1} \mathbf{y}}_{\text{BLUP}}$   
► Predictive Variance:  $\mathbf{V}_g^* = \mathbf{K}_g^{\star,*} - \sigma_g^2 \mathbf{K}_g^{\star,:} \left( \sigma_g^2 \mathbf{K}_g + \sigma_e^2 \mathbf{I} \right)^{-1} \mathbf{K}_g^{:,*}$ 

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition Sept

#### Outline

Variance component models

#### Variance component models for correlated traits

Mixed model Lasso

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 14

イロト 不得 トイヨト イヨト

э

#### Overview

- Frequently, we are interested in the genetic architecture of related traits.
- Example: flowering time in 10C and 16C

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

イロト イポト イヨト イヨト

3

#### Overview

- Frequently, we are interested in the genetic architecture of related traits.
- Example: flowering time in 10C and 16C



#### Multi-trait mixed models

• Extend variance component models to pairs of traits for environments  $e \in \{0, 1\}$ 



Again, prior is shared across SNPs but environment specific

$$p(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_S) = \prod_{s=1}^{S} \mathcal{N}\left(\boldsymbol{\beta}_s \left| \mathbf{0}, \begin{bmatrix} \sigma_{g,0}^2 & \rho_{0,1} \\ \rho_{0,1} & \sigma_{g,1}^2 \end{bmatrix} \right)$$

genetic variance and correlation.

One noise level per environment

$$\boldsymbol{\epsilon}_0 \sim \mathcal{N}(0, \sigma_{e,0}^2 \mathbf{I}) \ \boldsymbol{\epsilon}_1 \sim \mathcal{N}(0, \sigma_{e,1}^2 \mathbf{I}).$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

#### Multi-trait mixed models

• Extend variance component models to pairs of traits for environments  $e \in \{0, 1\}$ 



Again, prior is shared across SNPs but environment specific

$$p(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_S) = \prod_{s=1}^{S} \mathcal{N}\left(\boldsymbol{\beta}_s \left| \mathbf{0}, \begin{bmatrix} \sigma_{g,0}^2 & \rho_{0,1} \\ \rho_{0,1} & \sigma_{g,1}^2 \end{bmatrix} \right)$$

genetic variance and correlation.

One noise level per environment

$$\boldsymbol{\epsilon}_0 \sim \mathcal{N}(0, \sigma_{e,0}^2 \mathbf{I}) \ \boldsymbol{\epsilon}_1 \sim \mathcal{N}(0, \sigma_{e,1}^2 \mathbf{I}).$$

C. Lippert & O. Stegle I

Linear models II: multi-locus models and variance decomposition

September 2012 16

#### Multi-trait mixed models

• Extend variance component models to pairs of traits for environments  $e \in \{0, 1\}$ 



Again, prior is shared across SNPs but environment specific

$$p(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_S) = \prod_{s=1}^{S} \mathcal{N}\left(\boldsymbol{\beta}_s \left| \mathbf{0}, \begin{bmatrix} \sigma_{g,0}^2 & \rho_{0,1} \\ \rho_{0,1} & \sigma_{g,1}^2 \end{bmatrix} \right)$$

genetic variance and correlation.

One noise level per environment

$$\boldsymbol{\epsilon}_0 \sim \mathcal{N}(0, \sigma_{e,0}^2 \mathbf{I}) \ \boldsymbol{\epsilon}_1 \sim \mathcal{N}(0, \sigma_{e,1}^2 \mathbf{I}).$$

C. Lippert & O. Stegle

Multi-trait mixed models

Marginalized multi trait variance component model

 Because of the Gaussian assumption, again noise and weights can be marginalized out analytically

$$p\left(\mathbf{y}_{0}, \mathbf{y}_{1} \mid \mathbf{X}, \sigma_{g,0}^{2}, \sigma_{g,1}^{2}, \dots\right) = \\ \mathcal{N}\left(\begin{bmatrix}\mathbf{y}_{0}\\\mathbf{y}_{1}\end{bmatrix} \middle| \begin{array}{c} \mu_{0}\mathbf{I}\\ \mu_{1}\mathbf{I}' \begin{bmatrix} \sigma_{g,0}^{2}\mathbf{K}_{g} & \rho_{0,1}\mathbf{K}_{g}\\ \rho_{0,1}\mathbf{K}_{g} & \sigma_{g,1}^{2}\mathbf{K}_{g} \end{bmatrix} + \begin{bmatrix} \sigma_{e,0}^{2}\mathbf{I} & \mathbf{0}\\ \mathbf{0} & \sigma_{e,1}^{2}\mathbf{I} \end{bmatrix}\right)$$

As before,  $\mathbf{K}_g = \sum_{s=1}^S \mathbf{x}_s \mathbf{x}_s^{\mathrm{T}}$  denotes the genotype covariance.

C. Lippert & O. Stegle

イロト 不得下 イヨト イヨト 二日

Multi-trait mixed models Illustration on two related *A. thaliana* traits

Independent single-marker GWAs on flowering time (10C and 16C)



Multi-trait mixed models Illustration on two related *A. thaliana* traits

Mixed model inference

Genotype covariance

$$\begin{pmatrix} \sigma_{g,0}^2 = 0.2 & \sigma_{0,1} = 0.33 \\ \sigma_{0,1} = 0.33 & \sigma_{g,1}^2 = 0.92 \end{pmatrix}$$

Noise covariance

$$\begin{pmatrix} \sigma_{e,0}^2 = 0.03 & 0 \\ 0 = 0.33 & \sigma_{e,1}^2 = 0.02 \end{pmatrix}$$

 Marginal heritabilities from joint analysis

$$h_0 = 0.88 \ h_1 = 0.95$$



Image: A matrix

I ∃ ►

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

- N

Multi-trait mixed models Illustration on two related *A. thaliana* traits

Mixed model inference

Genotype covariance

$$\begin{pmatrix} \sigma_{g,0}^2 = 0.2 & \sigma_{0,1} = 0.33 \\ \sigma_{0,1} = 0.33 & \sigma_{g,1}^2 = 0.92 \end{pmatrix}$$

Noise covariance

$$\begin{pmatrix} \sigma_{e,0}^2 = 0.03 & 0 \\ 0 = 0.33 & \sigma_{e,1}^2 = 0.02 \end{pmatrix}$$

 Marginal heritabilities from joint analysis

$$h_0 = 0.88 \ h_1 = 0.95$$



Image: A match a ma

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

- - E - N

Multi-trait mixed models Illustration on two related *A. thaliana* traits

Mixed model inference

Genotype covariance

$$\begin{pmatrix} \sigma_{g,0}^2 = 0.2 & \sigma_{0,1} = 0.33 \\ \sigma_{0,1} = 0.33 & \sigma_{g,1}^2 = 0.92 \end{pmatrix}$$

Noise covariance

$$\begin{pmatrix} \sigma_{e,0}^2 = 0.03 & 0 \\ 0 = 0.33 & \sigma_{e,1}^2 = 0.02 \end{pmatrix}$$

 Marginal heritabilities from joint analysis

$$h_0 = 0.88 \ h_1 = 0.95$$



C. Lippert & O. Stegle

#### Multi-trait mixed models

- Multi-trait models can be used to control for population structure in a multi-environment setting
- Common association test affecting both traits

$$\mathcal{N}\bigg(\begin{bmatrix}\mathbf{y}_0\\\mathbf{y}_1\end{bmatrix} \mid \underbrace{\begin{bmatrix}\mu_0\mathbf{I}\\\mu_1\mathbf{I}\end{bmatrix}}_{\text{env effect}} + \underbrace{\begin{bmatrix}\mathbf{x}_s\\\mathbf{x}_s\end{bmatrix}}_{\text{SNP effect}}\beta_s, \mathbf{K}_{GxE}\bigg).$$

Interaction test, specific for one environment

$$\mathcal{N}\bigg(\begin{bmatrix}\mathbf{y}_{0}\\\mathbf{y}_{1}\end{bmatrix} \mid \underbrace{\begin{bmatrix}\mu_{0}\mathbf{I}\\\mu_{1}\mathbf{I}\end{bmatrix}}_{\text{env effect}} + \underbrace{\begin{bmatrix}\mathbf{x}_{s}\\\mathbf{x}_{s}\end{bmatrix}}_{\text{SNP effect}} \frac{\beta_{s}}{\beta_{s}} + \underbrace{\begin{bmatrix}\mathbf{x}_{s}\\\mathbf{0}\end{bmatrix}}_{\text{GxE effect}} \beta_{s}^{I}, \mathbf{K}_{GxE}\bigg).$$
  
where  $\mathbf{K}_{GxE} = \begin{bmatrix}\sigma_{g,0}^{2}\mathbf{K}_{g} & \rho_{0,1}\mathbf{K}_{g}\\\rho_{0,1}\mathbf{K}_{g} & \sigma_{g,1}^{2}\mathbf{K}_{g}\end{bmatrix} + \begin{bmatrix}\sigma_{e,0}^{2}\mathbf{I} & \mathbf{0}\\\mathbf{0} & \sigma_{e,1}^{2}\mathbf{I}\end{bmatrix}$ 

[Korte et al., 2012]

イロト 不得下 イヨト イヨト 二日

### Common effect mapping on two traits



C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 21

#### Outline

Variance component models

Variance component models for correlated traits

Mixed model Lasso

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

イロト 不得下 イヨト イヨト 二日

### Simultaneous Analysis of all SNPs: LassoTibshirani [1996]



$$\begin{array}{lll} \mathbf{y} &=& \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \\ \boldsymbol{\beta} &\sim& p(\boldsymbol{\beta}|\boldsymbol{\lambda}), \quad p(\boldsymbol{\beta}|\boldsymbol{\lambda}) \propto \prod e^{-\boldsymbol{\lambda}|\boldsymbol{\beta}_i|} \\ \boldsymbol{\epsilon} &\sim& \mathcal{N}(0,\sigma_{\mathrm{e}}^2\mathbf{I}) \end{array}$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

September 2012 23

э

### GWAS for Flowering Time in Arabidopsis thaliana

► Linear Model



Lasso



#### Linear Mixed Models

$$\begin{split} \mathbf{y} &= \mathbf{x}\boldsymbol{\beta} + \mathbf{u} + \boldsymbol{\epsilon}, \\ \mathbf{u} &\sim \mathcal{N}(0, \sigma_{\mathrm{g}}^{2}\mathbf{K}) \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(0, \sigma_{\mathrm{e}}^{2}\mathbf{I}) \end{split}$$

 Kinship Matrix K: measures genetic similarity between pairs of samples



25

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012

Linear Mixed Models

$$\mathbf{y} = \mathbf{x}\beta + \mathbf{u} + \boldsymbol{\epsilon},$$

Integrating over the random effects yields:

 $= \mathcal{N} \left( \mathbf{y} | \mathbf{x}_s \beta_s; \sigma_g^2 \mathbf{K} + \sigma_e^2 \mathbf{I} \right)$ 



(日) (同) (三) (三)

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 26

#### LMM-Lasso



$$\begin{split} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{u} + \boldsymbol{\epsilon}, \\ \boldsymbol{\beta} &\sim p(\boldsymbol{\beta}|\boldsymbol{\lambda}), \quad p(\boldsymbol{\beta}|\boldsymbol{\lambda}) \propto \prod e^{-\boldsymbol{\lambda}|\boldsymbol{\beta}_i|} \\ \mathbf{u} &\sim \mathcal{N}(0, \sigma_{\mathrm{g}}^2 \mathbf{K}) \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(0, \sigma_{\mathrm{e}}^2 \mathbf{I}) \end{split}$$

[Rakitsch et al., 2012]

(日)

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

September 2012 27

#### Inference

Minimize the objective:

$$-\log \mathcal{N}\left(\mathbf{y} | \mathbf{X} \boldsymbol{\beta}; \sigma_{g}^{2} \mathbf{K} + \sigma_{e}^{2} \mathbf{I}\right) + \lambda \|\boldsymbol{\beta}\|_{1}$$

Algorithm:

1. Fix 
$$\delta=\sigma_{\rm e}^2/\sigma_{\rm g}^2$$
 on the null model:

$$\begin{aligned} -\log \mathcal{N}\left(\mathbf{y}|\mathbf{0}; \sigma_{g}^{2}(\mathbf{K} + \delta \mathbf{I})\right) &= -\log \mathcal{N}\left(\mathbf{y}|\mathbf{0}; \sigma_{g}^{2}(\mathbf{U}\mathbf{S}\mathbf{U}^{T} + \delta \mathbf{I})\right) \\ &= -\log \mathcal{N}\left(\mathbf{U}^{T}\mathbf{y}|\mathbf{0}; \sigma_{g}^{2}(\mathbf{S} + \delta \mathbf{I})\right) \\ &= \frac{1}{\sigma_{g}^{2}}(\mathbf{U}^{T}\mathbf{y})^{T}(\mathbf{S} + \delta \mathbf{I})^{-1}(\mathbf{U}^{T}\mathbf{y}) \end{aligned}$$



C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

September 2012 28

### Inference (continued)

Minimize the objective:

$$-\log \mathcal{N}\left(\mathbf{y} | \mathbf{X} \boldsymbol{\beta}; \sigma_{g}^{2} (\mathbf{K} + \delta \mathbf{I})\right) + \lambda \| \boldsymbol{\beta} \|_{1}$$

Algorithm:

2. Train  $\beta$ :

$$\begin{split} & \operatorname{argmin}_{\boldsymbol{\beta}} \quad (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\sigma_{\mathrm{g}}^2 (\mathbf{K} + \delta \mathbf{I}))^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \operatorname{argmin}_{\boldsymbol{\beta}} \quad \frac{1}{\sigma_{\mathrm{g}}^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{U}\mathbf{S}\mathbf{U}^T + \delta \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \operatorname{argmin}_{\boldsymbol{\beta}} \quad \frac{1}{\sigma_{\mathrm{g}}^2} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta})^T (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1, \end{split}$$

where

$$\begin{split} \tilde{\mathbf{X}} &= (\mathbf{S} + \delta \mathbf{I})^{-\frac{1}{2}} \mathbf{U}^{\mathrm{T}} \mathbf{X} \\ \tilde{\mathbf{y}} &= (\mathbf{S} + \delta \mathbf{I})^{-\frac{1}{2}} \mathbf{U}^{\mathrm{T}} \mathbf{y} \end{split}$$

C. Lippert & O. Stegle

Linear models II: multi-locus models and variance decomposition

3

### GWAS for Flowering Time in Arabidopsis thaliana

Linear Mixed Model



LMM - Lasso



C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 30

#### Power comparison





- LASSO models perform better than univariate testing.
- Combining mixed models with LASSO outperforms other models.



(a) Precision/Recall

(b) ROC

#### Power comparison



- Power comparison on semi-empirical data
  - LASSO models perform better than univariate testing.
  - Combining mixed models with LASSO outperforms other models.



(a) Precision/Recall

(b) ROC

#### Power comparison



- Power comparison on semi-empirical data
  - LASSO models perform better than univariate testing.
  - Combining mixed models with LASSO outperforms other models.



(日) (同) (三) (三)

(a) Precision/Recall

(b) ROC

### Detection of multiple linked causal variants.pdf



Stegle Linear models II: multi-locus models and variance decomposition

September 2012 32

C. Lippert & O. Stegle

### Detection of multiple linked causal variants.pdf



< □ > < ---->

→ Ξ →

- E - N

### Phenotype prediction with the LMM-Lasso

- Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.
- $\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N} \left( \mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}^{\star} + \sigma_{e}^{2} \mathbf{I} \right)$
- ► Predicitve mean:  $\mu^* = \underbrace{\mathbf{X}^*\beta}_{\mathbf{X}^*} + \underbrace{\mathbf{K}^{\star,:} (\mathbf{K} + \delta \mathbf{I})^{-1} (\mathbf{y} \mathbf{X}\beta)}_{\mathbf{X}^*}$

► Predictive Variance:  $\mathbf{V}^{\star} = \mathbf{K}^{\star,\star} - \mathbf{K}^{\star,:} (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^{:,\star}$ 

### Phenotype prediction with the LMM-Lasso

 Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

$$\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N} \left( \mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}^{\star} + \sigma_{e}^{2} \mathbf{I} \right)$$

► Predicitve mean:  $\mu^* = \underbrace{\mathbf{X}^* \beta}_{\mathbf{X}^*} + \underbrace{\mathbf{K}^{*,:} (\mathbf{K} + \delta \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X} \beta)}_{\mathbf{X}^*}$ 

• Predictive Variance:  $\mathbf{V}^{\star} = \mathbf{K}^{\star,\star} - \mathbf{K}^{\star,:} (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^{:,\star}$ 

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 33

Phenotype prediction with the LMM-Lasso

Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

$$\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N}\left(\mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}^{\star} + \sigma_{e}^{2} \mathbf{I}\right)$$

- ► Predicitve mean:  $\mu^{\star} = \underbrace{\mathbf{X}^{\star}\beta}_{\text{LASSO component}} + \underbrace{\mathbf{K}^{\star,:} (\mathbf{K} + \delta \mathbf{I})^{-1} (\mathbf{y} \mathbf{X}\beta)}_{\text{BLUP component}}$
- Predictive Variance:  $\mathbf{V}^* = \mathbf{K}^{*,*} \mathbf{K}^{*,:} (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^{:,*}$

### Phenotype prediction with the LMM-Lasso

Given the phenotype values of a set of individuals and the genetic relatedness, we can predict the genetic component of the phenotype of a new individual.

$$\blacktriangleright P(\mathbf{y}^{\star} | \mathbf{y}) = \mathcal{N}\left(\mathbf{y}^{\star} | \boldsymbol{\mu}^{\star}, \sigma_{g}^{2} \mathbf{V}^{\star} + \sigma_{e}^{2} \mathbf{I}\right)$$

- ► Predicitive mean:  $\mu^* = \underbrace{\mathbf{X}^* \boldsymbol{\beta}}_{\text{LASSO component}} + \underbrace{\mathbf{K}^{*,:} (\mathbf{K} + \delta \mathbf{I})^{-1} (\mathbf{y} \mathbf{X} \boldsymbol{\beta})}_{\text{BLUP component}}$ ► Predictive Variance:  $\mathbf{V}^* = \mathbf{K}^{*,*} - \mathbf{K}^{*,:} (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^{:,*}$
- Predictive Variance:  $\mathbf{V}^{\star} = \mathbf{K}^{\star,\star} \mathbf{K}^{\star,:} (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^{:,\star}$

イロト 不得下 イヨト イヨト 二日

### Phenotype prediction with the LMM-Lasso



 Out of sample prediction on FT in A. thaliana

### Phenotype prediction with the LMM-Lasso Summary on a compendium of phenotypes

- Variance explained by LMMHasso Variance explained by LMM-Jasso 0.6 0.! 0.4 0.3 0.2 0.1 0.8.c 0.3 0.4 0.5 0.6 0.8 0.2 0.4 Variance explained by lasso Variance explained by lasso (a) Arabidopsis test variance (b) Mouse test variance Number of Active SNPs by LMMHasso 0 07 09 08 00 17 09 109 Number of Active SNPs by LMM-lasso 200 150 100 50 • 0b 20 40 60 80 100120140160 100 150 200 250 Number of Active SNPs by lasso Number of Active SNPs by lasso (c) Arabidopsis number of SNPs (d) Mouse number of SNPs
- Out of sample predictions on A. thaliana and mouse.

Linear models II: multi-locus models and variance decomposition

### Summary

### Joint modeling of multiple SNPs is compromised by sample size.

- Solutions: shrinkage (LASSO) or marginalization.
- Variance component models allow for estimating the proportion of genetic variance.
- Genetic co-regulation causes phenotype correlation which can be incorporated into variance component models.
- Shrinkage-based LASSO models and variance component modeling can be combined, significantly improving phenotype prediction.

#### Summary

- ► Joint modeling of multiple SNPs is compromised by sample size.
- ► Solutions: shrinkage (LASSO) or marginalization.
- Variance component models allow for estimating the proportion of genetic variance.
- Genetic co-regulation causes phenotype correlation which can be incorporated into variance component models.
- Shrinkage-based LASSO models and variance component modeling can be combined, significantly improving phenotype prediction.

#### Summary

- Joint modeling of multiple SNPs is compromised by sample size.
- Solutions: shrinkage (LASSO) or marginalization.
- Variance component models allow for estimating the proportion of genetic variance.
- Genetic co-regulation causes phenotype correlation which can be incorporated into variance component models.
- Shrinkage-based LASSO models and variance component modeling can be combined, significantly improving phenotype prediction.

### Summary

- Joint modeling of multiple SNPs is compromised by sample size.
- Solutions: shrinkage (LASSO) or marginalization.
- Variance component models allow for estimating the proportion of genetic variance.
- Genetic co-regulation causes phenotype correlation which can be incorporated into variance component models.
- Shrinkage-based LASSO models and variance component modeling can be combined, significantly improving phenotype prediction.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### Summary

- Joint modeling of multiple SNPs is compromised by sample size.
- Solutions: shrinkage (LASSO) or marginalization.
- Variance component models allow for estimating the proportion of genetic variance.
- Genetic co-regulation causes phenotype correlation which can be incorporated into variance component models.
- Shrinkage-based LASSO models and variance component modeling can be combined, significantly improving phenotype prediction.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

#### Acknowledgements

## Mixed model Lasso Barbara Rakitsch, Karsten Borgwardt

C. Lippert & O. Stegle Linear models II: multi-locus models and variance decomposition September 2012 37

イロト イポト イヨト イヨト

#### References I

- A. Korte, B. Vilhjálmsson, V. Segura, A. Platt, Q. Long, and M. Nordborg. A mixed-model approach for genome-wide association studies of correlated traits in structured populations. *Nature Genetics*, 44(9):1066–1071, 2012.
- B. Rakitsch, C. Lippert, O. Stegle, and K. Borgwardt. Lmm-lasso: A lasso multi-marker mixed model for association mapping with population structure correction. *Arxiv preprint* arXiv:1205.6986, 2012.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- M. Wu, S. Lee, T. Cai, Y. Li, M. Boehnke, and X. Lin. Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics*, 2011.

イロト 不得下 イヨト イヨト 二日