Linear models III: Improved linear mixed models

Christoph Lippert¹ Oliver Stegle²

 1 Microsoft Research, Los Angeles, USA 2 Max-Planck-Institutes Tübingen, Germany



Basel 09. September 2012



C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

September 2012 1

イロト イポト イヨト イヨト

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 models the probability
 - that two individuals share causal SNPs.



[Kang et al., 2008, 2010, Lippert et al., 2011]

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 models the probability
 - that two individuals share causal SNPs.



[Kang et al., 2008, 2010, Lippert et al., 2011]

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data models the probability
 - that two individuals share causal SNPs.



[Kang et al., 2008, 2010, Lippert et al., 2011]

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 models the probability
 - causal SNPs.



[Kang et al., 2008, 2010, Lippert et al., 2011]

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- ▶ Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





[Kang et al., 2008, 2010, Lippert et al., 2011]

(日) (周) (三) (三)

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





[Kang et al., 2008, 2010, Lippert et al., 2011]

(日) (周) (三) (三)

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





(日) (周) (三) (三)

 $\mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\beta};\sigma_{*}^{2}\mathbf{I})$

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





(日) (周) (三) (三)

 $\mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\beta};\sigma_{\epsilon}^{2}\mathbf{I})$

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





(日) (周) (三) (三)

 $\mathcal{N}\left(\mathbf{y}|\mathbf{X}\boldsymbol{\beta};\sigma_{e}^{2}\mathbf{I}\right)$

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





[Kang et al., 2008, 2010, Lippert et al., 2011]

(日) (周) (三) (三)

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- Confounding structure leads to false positives.
 - Population structure
 - Family structure
 - Cryptic relatedness
- Genetic similarity matrix K
 - Estimated from SNP data
 - models the probability that two individuals share causal SNPs.





[Kang et al., 2008, 2010, Lippert et al., 2011]

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

Goal: genotype-phenotype relatioships

- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



Goal: genotype-phenotype relatioships

- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



A = A = A = A = A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

- E - N

- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



A = A = A = A = A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

- E - N

- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



• • • • • • • • • • • •

- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



(日) (同) (三) (三)

- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



(日) (同) (三) (三)

A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- Goal: genotype-phenotype relatioships
- Associations via linkage disequislibrium
- SNPs confounded by population structure
- Spurious associations if not taken into account
- Alternatively condition on causal SNPs



Image: A matrix

A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.

C. Lippert & O. Stegle

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ▶ For linear similarities, a LMM is equivalent to a *linear regression*.
- ► All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

 $\mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\mathbf{K}+\sigma_{e}^{2}\mathbf{I}\right).$

[Listgarten et al., 2012]

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ののの

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ► For linear similarities, a LMM is equivalent to a *linear regression*.
- ► All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

 $\mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathrm{T}}+\sigma_{e}^{2}\mathbf{I}\right).$

[Listgarten et al., 2012]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ► For linear similarities, a LMM is equivalent to a *linear regression*.
- ► All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

$$\mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}\right).$$

$$\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_{e}^{2}\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}
ight)\mathsf{d}m{ heta}.$$

[Listgarten et al., 2012]

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ► For linear similarities, a LMM is equivalent to a *linear regression*.
- All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

$$\mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}\right).$$

$$\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_{e}^{2}\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}
ight)\mathsf{d}m{ heta}.$$

[Listgarten et al., 2012]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ► For linear similarities, a LMM is equivalent to a *linear regression*.
- ► All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

$$\mathcal{N}\left(\mathbf{y}|\mathbf{x}eta;\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}
ight).$$

 $\propto\int\mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+\tilde{\mathbf{X}} heta;\sigma_{e}^{2}\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(heta|\mathbf{0};\sigma_{g}^{2}\mathbf{I}
ight)\mathsf{d} heta.$

Can we do better than using all SNPs for correction?

[Listgarten et al., 2012]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙

Bayesian linear regression view

- A LMM accounts for model misspecification when testing a univariate model when in reality the trait is multi-factorial.
- ► For linear similarities, a LMM is equivalent to a *linear regression*.
- ► All SNPs are used as regression covariates.
- Uncertainty about identity of regulative SNPs expressed by using a Bayesian prior distribution on the regression weights.

$$\mathcal{N}\left(\mathbf{y}|\mathbf{x}\beta;\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}
ight).$$

$$\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_e^2\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_g^2\mathbf{I}
ight)\mathsf{d}m{ heta}.$$

Can we do better than using all SNPs for correction? \rightarrow Select SNPs by their association to the phenotype.

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

September 2012 4

[Listgarten et al., 2012]

- Compute similarity matrix based on highly associated SNPs
- Equivalent to linear regression conditioned on these SNPs
- ▶ Conditioning on SNPs in LD blocks association.
 → False negative results!
- Remove SNPs in LD from computation of similarity matrix



[Listgarten et al., 2012]

- Compute similarity matrix based on highly associated SNPs
- Equivalent to linear regression conditioned on these SNPs
- Conditioning on SNPs in LD blocks association.
 → False negative results!
- Remove SNPs in LD from computation of similarity matrix



[Listgarten et al., 2012]

- Compute similarity matrix based on highly associated SNPs
- Equivalent to linear regression conditioned on these SNPs
- Conditioning on SNPs in LD blocks association.
 → False negative results!
- Remove SNPs in LD from computation of similarity matrix



(日) (同) (三) (三)

[Listgarten et al., 2012]

- Compute similarity matrix based on highly associated SNPs
- Equivalent to linear regression conditioned on these SNPs
- Conditioning on SNPs in LD blocks association.
 → False negative results!
- Remove SNPs in LD from computation of similarity matrix



(日) (同) (三) (三)

[Listgarten et al., 2012]

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- SNPs that are in close proximity on the chromosome are highly correlated.
- ▶ Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.

$$egin{aligned} &\mathcal{N}\left(\mathbf{y}|\mathbf{x}eta;\sigma_g^2 ilde{\mathbf{X}} ilde{\mathbf{X}}^{ extsf{T}}+\sigma_e^2\mathbf{I}
ight). \ &\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_e^2\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_g^2\mathbf{I}
ight)\mathsf{d}m{ heta}. \end{aligned}$$

[Lippert et al., 2011, Listgarten et al., 2012]

イロト イポト イヨト イヨト

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

▶ SNPs that are in close proximity on the chromosome are highly correlated.

- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.

$$egin{aligned} &\mathcal{N}\left(\mathbf{y}|\mathbf{x}eta;\sigma_g^2 ilde{\mathbf{X}} ilde{\mathbf{X}}^{ extsf{T}}+\sigma_e^2\mathbf{I}
ight). \ &\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_e^2\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_g^2\mathbf{I}
ight)\mathsf{d}m{ heta}. \end{aligned}$$

[Lippert et al., 2011, Listgarten et al., 2012]

< ロ > < 同 > < 回 > < 回 > < 回 > <

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.

$$egin{aligned} &\mathcal{N}\left(\mathbf{y}|\mathbf{x}eta;\sigma_g^2 ilde{\mathbf{X}} ilde{\mathbf{X}}^{ extsf{T}}+\sigma_e^2\mathbf{I}
ight). \ &\propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}eta+ ilde{\mathbf{X}}m{ heta};\sigma_e^2\delta\mathbf{I}
ight)\cdot\mathcal{N}\left(m{ heta}|\mathbf{0};\sigma_g^2\mathbf{I}
ight)\mathsf{d}m{ heta}. \end{aligned}$$

[Lippert et al., 2011, Listgarten et al., 2012]

소리가 소문가 소문가 소문가 ...

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.
- Proximal contamination on IBD phenotype [WTCCC, 2007]
- Test SNPs are equally distant sampled along the chromosome
- Compute 6 similarity matrices containing same number of SNPs at growing distance to test SNPs.



[Lippert et al., 2011, Listgarten et al., 2012]

イロト イポト イヨト イヨト

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.
- Proximal contamination on IBD phenotype [WTCCC, 2007]
- Test SNPs are equally distant sampled along the chromosome
- Compute 6 similarity matrices containing same number of SNPs at growing distance to test SNPs.



[Lippert et al., 2011, Listgarten et al., 2012]

イロト イポト イヨト イヨト

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.
- Proximal contamination on IBD phenotype [WTCCC, 2007]
- Test SNPs are equally distant sampled along the chromosome
- Compute 6 similarity matrices containing same number of SNPs at growing distance to test SNPs.



[Lippert et al., 2011, Listgarten et al., 2012]

イロト イポト イヨト イヨト

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.
- Proximal contamination on IBD phenotype [WTCCC, 2007]
- Test SNPs are equally distant sampled along the chromosome
- Compute 6 similarity matrices containing same number of SNPs at growing distance to test SNPs.



[Lippert et al., 2011, Listgarten et al., 2012]

イロト イポト イヨト イヨト

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.



- 4 同 6 4 日 6 4 日 6

- SNPs that are in close proximity on the chromosome are highly correlated.
- Having a SNP in the similarity matrix that is close to the SNP tested is equivalent to conditioning on this SNP in the null-model.
- Correct by removing a sliding window around test-SNP from the similarity matrix.
- Correction is computed efficiently by subtracting a low-rank term.



Power study

- 3000 individuals
- Two populations (6:4)
- Balding-Nichols model
- 100000 non-causal SNPs
- 100 causal SNPs $(\sigma_g^2 = \sigma_e^2 = 0.1)$
- Test-set of 5100 SNPs



[Listgarten et al., 2012]

(日) (同) (三) (三)

Inflammatory bowel disease [WTCCC, Nature 2007]

Algorithm parameters				Algorithm performance					
Name	SNP selection method	#SNPs in matrix	Avoid prox conta m	λ_{GC}	False Positives	True Positives	Runtime (min) without speedup	Runtime (min) with speedup	Memory use (GB)
FaST-LMM-Select	Select	310	yes	1.08	0	100	1.3 x 10 ³	45	<1
FaST-LMM all	All	All	yes	1.09	2	108	4.0 x 10 ⁶	4567	86
FaST-LMM orig 310	Equi-spaced	310	yes	1.26	15	128	1.1 x 10 ³	6	<1
FaST-LMM orig 4K	Equi-spaced	4000	yes	1.17	8	114	2.1 x 10 ⁵	30	2
Traditional	All	All	no	0.97	2	64	42	NA	45

SNPs considered True Positive if:

- Reported in WTCCC paper [WTCCC, Nature 2007]
- Reported in meta analysis [Franke et al., Nat Gen 2010]
- ► In major histocompatibility complex (MHC) region

[Burton et al., 2007]

소리가 소문가 소문가 소문가 ...

C. Lippert & O. Stegle

Linear models III: Improved linear mixed models

- Relationship between linear mixed models and multi-variate modelling using linear regression
- Computing genetic similarities on associated SNPs improves correction FaST-LMM select
- Inclusion of SNPs in LD to the SNP tested causes deflated tests (proximal contamination)
- Efficient exclusion of SNPs in LD from the similarity matrix feasible using low-rank updates

$$\begin{split} & \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}\right).\\ & \propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta}+\tilde{\mathbf{X}}\boldsymbol{\theta};\sigma_{e}^{2}\boldsymbol{\delta}\mathbf{I}\right)\cdot\mathcal{N}\left(\boldsymbol{\theta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}\right)\mathsf{d}\boldsymbol{\theta}. \end{split}$$

- Relationship between linear mixed models and multi-variate modelling using linear regression
- Computing genetic similarities on associated SNPs improves correction FaST-LMM select
- Inclusion of SNPs in LD to the SNP tested causes deflated tests (proximal contamination)
- Efficient exclusion of SNPs in LD from the similarity matrix feasible using low-rank updates

$$\begin{split} & \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathsf{T}}+\sigma_{e}^{2}\mathbf{I}\right).\\ & \propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta}+\tilde{\mathbf{X}}\boldsymbol{\theta};\sigma_{e}^{2}\boldsymbol{\delta}\mathbf{I}\right)\cdot\mathcal{N}\left(\boldsymbol{\theta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}\right)\mathsf{d}\boldsymbol{\theta}. \end{split}$$

- Relationship between linear mixed models and multi-variate modelling using linear regression
- Computing genetic similarities on associated SNPs improves correction FaST-LMM select
- Inclusion of SNPs in LD to the SNP tested causes deflated tests (proximal contamination)
- Efficient exclusion of SNPs in LD from the similarity matrix feasible using low-rank updates

$$\begin{split} & \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathrm{T}}+\sigma_{e}^{2}\mathbf{I}\right).\\ & \propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta}+\tilde{\mathbf{X}}\boldsymbol{\theta};\sigma_{e}^{2}\delta\mathbf{I}\right)\cdot\mathcal{N}\left(\boldsymbol{\theta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}\right)\mathrm{d}\boldsymbol{\theta}. \end{split}$$

- Relationship between linear mixed models and multi-variate modelling using linear regression
- Computing genetic similarities on associated SNPs improves correction FaST-LMM select
- Inclusion of SNPs in LD to the SNP tested causes deflated tests (proximal contamination)
- Efficient exclusion of SNPs in LD from the similarity matrix feasible using low-rank updates

$$\begin{split} & \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta};\sigma_{g}^{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\mathrm{T}}+\sigma_{e}^{2}\mathbf{I}\right).\\ & \propto \int \mathcal{N}\left(\mathbf{y}|\mathbf{x}\boldsymbol{\beta}+\tilde{\mathbf{X}}\boldsymbol{\theta};\sigma_{e}^{2}\delta\mathbf{I}\right)\cdot\mathcal{N}\left(\boldsymbol{\theta}|\mathbf{0};\sigma_{g}^{2}\mathbf{I}\right)\mathsf{d}\boldsymbol{\theta}. \end{split}$$



References I

- P. Burton, D. Clayton, L. Cardon, N. Craddock, P. Deloukas, A. Duncanson, D. Kwiatkowski, M. McCarthy, W. Ouwehand, N. Samani, et al. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*, 447(7145):661–678, 2007.
- H. M. Kang, N. A. Zaitlen, C. M. Wade, A. Kirby, D. Heckerman, M. J. Daly, and E. Eskin. Efficient control of population structure in model organism association mapping. *Genetics*, 107, 2008.
- H. M. Kang, J. H. Sul, S. K. Service, N. A. Zaitlen, S. Y. Kong, N. B. Freimer, C. Sabatti, and E. Eskin. Variance component model to account for sample structure in genome-wide association studies. *Nat. Genet.*, 42:348–354, Apr 2010.
- C. Lippert, J. Listgarten, Y. Liu, C. Kadie, R. Davidson, and D. Heckerman. Fast linear mixed models for genome-wide association studies. *Nature Methods*, 8(10):838;835, 10 2011. doi: 10.1038/nmeth.1681.
- J. Listgarten, C. Lippert, C. Kadie, R. Davidson, E. Eskin, and D. Heckerman. Improved linear mixed models for genome-wide association studies. *Nature Methods*, 9(6):525–526, 2012.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙