Modifications of BIC

GWAS Simulation QT

Case Control studies

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー わえぐ

Out look 00 0000 0000

# Model Selection Procedures for Genome Wide Association Studies

#### Florian Frommlet

Department of Medical Statistics, Medical University Vienna

Liege, May 2013



Int	ro	lu	ct	io	n
	00		0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Contents

- $1. \ {\sf Introduction\ into\ GWAS}$
- 2. Model selection with Modifications of BIC
- 3. Simulation study for quantitative traits
- 4. Different approaches for case control studies
- 5. Outlook

Introduction	
00000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Out look 00 0000 00

## Genetic association studies

#### General purpose

#### Detect genomic regions which are associated with some trait Traits might be e.g. quantitative or dichotomous

### Genetic markers

- Traditional: Genes that encode certain phenotypes
- Today: DNA sequence information e.g. SNPs, Copy number variation, etc

Introduction	
00000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 0000

## Genetic association studies

#### General purpose

Detect genomic regions which are associated with some trait Traits might be e.g. quantitative or dichotomous

#### Genetic markers

- Traditional: Genes that encode certain phenotypes
- Today: DNA sequence information e.g. SNPs, Copy number variation, etc

Introduction	
00000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

## Genetic association studies

#### General purpose

Detect genomic regions which are associated with some trait Traits might be e.g. quantitative or dichotomous

#### Genetic markers

- Traditional: Genes that encode certain phenotypes
- Today: DNA sequence information e.g. SNPs, Copy number variation, etc.

Introduction	
00000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

## Genetic association studies

#### General purpose

Detect genomic regions which are associated with some trait Traits might be e.g. quantitative or dichotomous

#### Genetic markers

- Traditional: Genes that encode certain phenotypes
- Today: DNA sequence information e.g. SNPs, Copy number variation, etc.

Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies 0000000 0000000000 Out look 00 0000

# SNPs as genetic markers

## Single Nucleotide Polymorphism

#### SNP: Point mutation

Humans: Some 20 million SNPs known, figure increasing rapidly

## SNP Arrays

Affymetrix 6: ca. 1 million SNPs Latest Illumina: ca. 5 million SNPs



Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies 0000000 0000000000 Out look 00 0000

# SNPs as genetic markers

Single Nucleotide Polymorphism

SNP: Point mutation

Humans: Some 20 million SNPs known, figure increasing rapidly

### **SNP** Arrays

Affymetrix 6: ca. 1 million SNPs Latest Illumina: ca. 5 million SNPs



Introduction	
000000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

イロト 不得 とうき イヨト

э.

Out look 00 0000 00

# Technology

#### SNP arrays

- Approximately 10 years around
- Similar to RNA Micro-Arrays
- Both variants of SNP on array (say A and a)
  ⇒ Different intensities for genotypes AA, Aa, aa

#### First step

Image segmentation similar to RNA Micro-Arrays

Introduction	
000000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

・ロト ・ 御 ト ・ ヨ ト ・ ヨ ト ・ ヨ ・

Out look 00 0000 00

# Technology

#### SNP arrays

- Approximately 10 years around
- Similar to RNA Micro-Arrays
- Both variants of SNP on array (say A and a)
  - $\Rightarrow$  Different intensities for genotypes AA, Aa, aa

#### First step

Image segmentation similar to RNA Micro-Arrays

Introduction	
000000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 0000

# Technology

#### SNP arrays

- Approximately 10 years around
- Similar to RNA Micro-Arrays
- Both variants of SNP on array (say A and a)
  - $\Rightarrow$  Different intensities for genotypes AA, Aa, aa

#### First step

Image segmentation similar to RNA Micro-Arrays

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Technology

#### Good for calling

#### Bad for calling

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



#### Quality measures are needed

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Technology

#### Good for calling

#### Bad for calling

(日) (個) (E) (E) (E)



#### Quality measures are needed

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Downstream Analysis

 $Y \leftarrow X_1, \ldots, X_p$ 

Inte

• Y .... quantitative (e.g. height) or categorical (e.g. disease status)

ション ふゆ アメリア メリア しょうくの

- $X_j \in \{-1,0,1\}$  for different genotypes
- n observations

Typical:  $n>10^3$ ,  $p>10^5$ 

Question: Which  $X_j$  are associated with Y?

State of the art analysis: Single marker tests

- Test statistic for each SNP (ANOVA,  $\chi^2$ , etc.)
- Multiple testing correction (Bonferroni, FDR control, permutation tests, ...)

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Downstream Analysis

 $Y \leftarrow X_1, \ldots, X_p$ 

Inte

• Y .... quantitative (e.g. height) or categorical (e.g. disease status)

ション ふゆ アメリア メリア しょうくの

- $X_j \in \{-1,0,1\}$  for different genotypes
- n observations

Typical:  $n>10^3$ ,  $p>10^5$ 

### Question: Which $X_j$ are associated with Y?

State of the art analysis: Single marker tests

- Test statistic for each SNP (ANOVA,  $\chi^2$ , etc.)
- Multiple testing correction (Bonferroni, FDR control, permutation tests, ....)

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Downstream Analysis

 $Y \leftarrow X_1, \ldots, X_p$ 

• Y .... quantitative (e.g. height) or categorical (e.g. disease status)

ション ふゆ アメリア メリア しょうくの

- $X_j \in \{-1,0,1\}$  for different genotypes
- n observations

Typical:  $n > 10^3$ ,  $p > 10^5$ 

#### Question:

Intr

```
Which X_i are associated with Y?
```

State of the art analysis: Single marker tests

- Test statistic for each SNP (ANOVA,  $\chi^2$ , etc.)
- Multiple testing correction (Bonferroni, FDR control, permutation tests, ...)

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

#### Model Index vector $M = [j_1, \dots, j_{k_M}]$

#### Quantitative Trait: Linear regression

$$\mathcal{M}: Y = X_M \beta_M + \epsilon, \quad X_M = [X_{j_1}, \dots, X_{j_{k_M}}]$$

うして ふゆう ふほう ふほう うらつ

Columns of design matrix  $X_j$ :

- (-1,0,1) additive effects
- (1,0,1) dominance effects

Case control studies: Logistic regression

- 1. How to evaluate what is a good model?
- 2. How to find a good model?

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Model

Index vector  $M = [j_1, \ldots, j_{k_M}]$ 

Quantitative Trait: Linear regression

$$\mathcal{M}: Y = X_M \beta_M + \epsilon, \quad X_M = [X_{j_1}, \dots, X_{j_{k_M}}]$$

うして ふゆう ふほう ふほう うらつ

Columns of design matrix  $X_j$ :

- (-1,0,1) additive effects
- (1,0,1) dominance effects

Case control studies: Logistic regression

- 1. How to evaluate what is a good model?
- 2. How to find a good model?

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Model

Index vector  $M = [j_1, \ldots, j_{k_M}]$ 

Quantitative Trait: Linear regression

$$\mathcal{M}: Y = X_M \beta_M + \epsilon, \quad X_M = [X_{j_1}, \dots, X_{j_{k_M}}]$$

Columns of design matrix  $X_j$ :

- (-1,0,1) additive effects
- (1,0,1) dominance effects

#### Case control studies: Logistic regression

- 1. How to evaluate what is a good model?
- 2. How to find a good model?

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Model

Index vector  $M = [j_1, \ldots, j_{k_M}]$ 

Quantitative Trait: Linear regression

$$\mathcal{M}: Y = X_M \beta_M + \epsilon, \quad X_M = [X_{j_1}, \dots, X_{j_{k_M}}]$$

Columns of design matrix  $X_j$ :

- (-1,0,1) additive effects
- (1,0,1) dominance effects

Case control studies: Logistic regression

- 1. How to evaluate what is a good model?
- 2. How to find a good model?

Modifications of BIC

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# Model selection criteria

# Classical Maximum likelihood $L_M$ with penalties based on model size

 $-2 \log L_M + \text{Penalty} \cdot k_M$ 

Examples: AIC, BIC, RIC, Mallows C, etc.

AIC ... Penalty = 2, BIC ... Penalty =  $\log n$ 

More recent LASSO:  $L_1$  — Penalty Elastic Net:  $L_1$  and  $L_2$  — Penalty etc.

Modifications of BIC ••••• ••••• GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000

# Model selection criteria

# Classical Maximum likelihood $L_M$ with penalties based on model size

 $-2 \log L_M + \text{Penalty} \cdot k_M$ 

Examples: AIC, BIC, RIC, Mallows C, etc.

AIC ... Penalty = 2, BIC ... Penalty =  $\log n$ 

More recent LASSO: L<sub>1</sub>— Penalty Elastic Net: L<sub>1</sub> and L<sub>2</sub>— Penalty etc.

Modifications of BIC ••••• ••••• GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 000

# Model selection criteria

# Classical Maximum likelihood $L_M$ with penalties based on model size

 $-2 \log L_M + \text{Penalty} \cdot k_M$ 

Examples: AIC, BIC, RIC, Mallows C, etc.

AIC ... Penalty = 2, BIC ... Penalty =  $\log n$ 

More recent LASSO:  $L_1$  – Penalty Elastic Net:  $L_1$  and  $L_2$  – Penalty etc.

Modifications of BIC

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

・ロト ・ 御 ト ・ ヨ ト ・ ヨ ト ・ ヨ ・

Out look 00 0000 0000

# Model selection for p > n

#### Classical theory for AIC and BIC

Derived for constant p, while  $n \to \infty$ 

Results for p > n no longer correct e.g. BIC no longer consistent

#### Problem

In case of sparsity and p > n BIC chooses too large models

Modifications of BIC

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 0000

# Model selection for p > n

#### Classical theory for AIC and BIC

Derived for constant p, while  $n \to \infty$ 

Results for p > n no longer correct e.g. BIC no longer consistent

#### Problem

In case of sparsity and p > n BIC chooses too large models

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Out look

# Schwarz BIC in case of sparsity

# Source of problem

BIC derived in Bayesian context

$$P(M|Y) = \frac{P(Y|M)\pi(M)}{P(Y)}$$

BIC ignores model prior  $\pi(M)$ , i.e. equivalent with uniform prior for all models  $\Rightarrow$  informative prior for model size

e.g. p models of size 1,  $\binom{p}{p/2}$  models of size p/2

If one expects only few causal SNPs

 $\Rightarrow$  BIC selects too large models

## Solution

Use model prior  $\pi(M)$  which takes into account p

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Schwarz BIC in case of sparsity

## Source of problem

BIC derived in Bayesian context

$$P(M|Y) = \frac{P(Y|M)\pi(M)}{P(Y)}$$

BIC ignores model prior  $\pi(M)$ , i.e. equivalent with uniform prior for all models  $\Rightarrow$  informative prior for model size

e.g. p models of size 1,  $\binom{p}{p/2}$  models of size p/2

If one expects only few causal SNPs  $\Rightarrow$  BIC selects too large models

Solution Use model prior  $\pi(M)$  which takes into account p

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Out look 00 0000 00

# Schwarz BIC in case of sparsity

## Source of problem

BIC derived in Bayesian context

$$P(M|Y) = \frac{P(Y|M)\pi(M)}{P(Y)}$$

BIC ignores model prior  $\pi(M)$ , i.e. equivalent with uniform prior for all models  $\Rightarrow$  informative prior for model size

e.g. p models of size 1,  $\binom{p}{p/2}$  models of size p/2

If one expects only few causal SNPs

 $\Rightarrow$  BIC selects too large models

#### Solution

Use model prior  $\pi(M)$  which takes into account p

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

**Case Control studies** 0000000 0000 0000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

Out look

# Schwarz BIC in case of sparsity

## Source of problem

BIC derived in Bayesian context

$$P(M|Y) = \frac{P(Y|M)\pi(M)}{P(Y)}$$

BIC ignores model prior  $\pi(M)$ , i.e. equivalent with uniform prior for all models  $\Rightarrow$  informative prior for model size

e.g. p models of size 1,  $\binom{p}{p/2}$  models of size p/2

If one expects only few causal SNPs

 $\Rightarrow$  BIC selects too large models

#### Solution

Use model prior  $\pi(M)$  which takes into account p

Modifications of BIC ○○○● ○○○○○ GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# First modification of BIC

# Original BIC [Schwarz (1978)]

 $BIC = -2\log L_M + k_M\log n$ 

# mBIC [Bogdan et al. (2004)]

Model prior 
$$\pi(M) = \omega^{k_M} \cdot (1-\omega)^{p-k_M}$$
 yields

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

- $\omega \dots$  Prior probability of causal SNPs  $\Rightarrow$  defines d Recommendation if no prior information:  $d = -2 \log 4$
- Orthogonal design ⇒ mBIC controls FWER (closely related to Bonferroni correction)

Intro	du	cti	on	
000	00	0		

Modifications of BIC ○○○● ○○○○○○ GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

First modification of BIC

# Original BIC [Schwarz (1978)]

 $BIC = -2\log L_M + k_M\log n$ 

mBIC [Bogdan et al. (2004)]

Model prior  $\pi(M) = \omega^{k_M} \cdot (1-\omega)^{p-k_M}$  yields

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

- $\omega \dots$  Prior probability of causal SNPs  $\Rightarrow$  defines d Recommendation if no prior information:  $d = -2 \log 4$
- Orthogonal design ⇒ mBIC controls FWER (closely related to Bonferroni correction)

Modifications of BIC

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

First modification of BIC

# Original BIC [Schwarz (1978)]

 $BIC = -2\log L_M + k_M\log n$ 

mBIC [Bogdan et al. (2004)]

Model prior  $\pi(M) = \omega^{k_M} \cdot (1-\omega)^{p-k_M}$  yields

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

- $\omega \dots$  Prior probability of causal SNPs  $\Rightarrow$  defines d Recommendation if no prior information:  $d = -2 \log 4$
- Orthogonal design ⇒ mBIC controls FWER (closely related to Bonferroni correction)

Introduction	Modifications of BIC	GWAS Simulation QT	Case Contro
000000	000●	00000000	0000000
	000000		0000

Control studies Outlook

First modification of BIC

# Original BIC [Schwarz (1978)]

 $BIC = -2\log L_M + k_M\log n$ 

mBIC [Bogdan et al. (2004)]

Model prior  $\pi(M) = \omega^{k_M} \cdot (1-\omega)^{p-k_M}$  yields

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

- $\omega \dots$  Prior probability of causal SNPs  $\Rightarrow$  defines d Recommendation if no prior information:  $d = -2 \log 4$
- Orthogonal design  $\Rightarrow$  mBIC controls FWER (closely related to Bonferroni correction)

Int	ro	d٢	1 C1	i	0	n	
	0	00	DC				

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look 00 0000 00

# FDR-controlling modifications of BIC

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

# mBIC2 [Frommlet et al. (2011)]

Model selection criterion which under orthogonality controls FDR

 $mBIC2 = -2 \log L_M + k_M [\log(np^2) + d] - 2 \log k_M!$ 

#### Properties

- Penalisation based on ideas of [Abramovich et al. (2006)]
- mBIC2 has certain optimality properties (as we will see)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
0000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

# FDR-controlling modifications of BIC

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

# mBIC2 [Frommlet et al. (2011)]

#### Model selection criterion which under orthogonality controls FDR

$$mBIC2 = -2 \log L_M + k_M [\log(np^2) + d] - 2 \log k_M!$$

・ロト ・ 日 ・ エ ヨ ・ ト ・ 日 ・ う へ つ ・

- Penalisation based on ideas of [Abramovich et al. (2006)]
- mBIC2 has certain optimality properties (as we will see)

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
0000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

# FDR-controlling modifications of BIC

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

# mBIC2 [Frommlet et al. (2011)]

Model selection criterion which under orthogonality controls FDR

$$mBIC2 = -2 \log L_M + k_M [\log(np^2) + d] - 2 \log k_M!$$

・ロト ・ 日 ・ エ ヨ ・ ト ・ 日 ・ う へ つ ・

- Penalisation based on ideas of [Abramovich et al. (2006)]
- mBIC2 has certain optimality properties (as we will see)
| duction | Modifications of BIC  | GWAS Simulation QT | Case Control studies       | Out look         |
|---------|-----------------------|--------------------|----------------------------|------------------|
| 0000    | 0000<br>000000<br>000 | 00000000           | 0000000<br>0000<br>0000000 | 00<br>0000<br>00 |

### FDR-controlling modifications of BIC

$$mBIC = -2\log L_M + k_M[\log(np^2) + d]$$

### mBIC2 [Frommlet et al. (2011)]

Model selection criterion which under orthogonality controls FDR

$$mBIC2 = -2 \log L_M + k_M [\log(np^2) + d] - 2 \log k_M!$$

ション ふゆ アメリア メリア しょうくの

#### Properties

- Penalisation based on ideas of [Abramovich et al. (2006)]
- mBIC2 has certain optimality properties (as we will see)

troduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 000000 000	00000000	0000000 0000 0000000	00 00000

# Ideas underlying mBIC2

#### Penalizing scheme by [Abramovich et al. (2006)]

$$\frac{RSS_M}{\sigma^2} + \sum_{i=1}^{k_M} q_N^2(\alpha i/2p) \tag{1}$$

(日) (四) (日) (日) (日)

with  $q_N(lpha)$  the (1-lpha) - quantile of standard normal

- Benjamini Hochberg (BH) corresponds to largest local minimum of (1)
- Corresponding step down procedure corresponds to smallest local minimum of (1)
- Approximation of penalty term using  $\alpha = n^{-1/2}$  yields mBIC2

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

## Ideas underlying mBIC2

Penalizing scheme by [Abramovich et al. (2006)]

$$\frac{RSS_M}{\sigma^2} + \sum_{i=1}^{k_M} q_N^2(\alpha i/2p) \tag{1}$$

with  $q_N(lpha)$  the (1-lpha) - quantile of standard normal

- Benjamini Hochberg (BH) corresponds to largest local minimum of (1)
- Corresponding step down procedure corresponds to smallest local minimum of (1)
- Approximation of penalty term using  $\alpha = n^{-1/2}$  yields mBIC2

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

## Ideas underlying mBIC2

Penalizing scheme by [Abramovich et al. (2006)]

$$\frac{RSS_M}{\sigma^2} + \sum_{i=1}^{k_M} q_N^2(\alpha i/2p) \tag{1}$$

ション ふゆ アメリア メリア しょうくの

with  $q_N(lpha)$  the (1-lpha) - quantile of standard normal

- Benjamini Hochberg (BH) corresponds to largest local minimum of (1)
- Corresponding step down procedure corresponds to smallest local minimum of (1)
- Approximation of penalty term using  $\alpha = n^{-1/2}$  yields mBIC2

Modifications of BIC

GWAS Simulation QT

Case Control studies

Out look

# Optimality properties of mBIC2

#### Asymptotic Bayes optimality under sparsity (ABOS) Topic of Bogdan et al. (2011), Frommlet et al. (2013) for multiple testing

#### Essential idea for regression

• Two groups model for regressors:

$$P(\beta_i \neq 0) = \eta$$
, with  $\eta$  small

while *n* and *p* are large

Compare misclassification rate of procedure with optimal Bayes rule

Int	ro	lu	ct	io	n
	00		0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# Optimality properties of mBIC2

Asymptotic Bayes optimality under sparsity (ABOS) Topic of Bogdan et al. (2011), Frommlet et al. (2013) for multiple testing Essential idea for regression

• Two groups model for regressors:

$$P(eta_i
eq 0)=\eta, \hspace{0.2cm} ext{with} \hspace{0.2cm} \eta ext{ small}$$

while *n* and *p* are large

• Compare misclassification rate of procedure with optimal Bayes rule

Int	ro	lu	ct	io	n
	00		0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# Optimality properties of mBIC2

Asymptotic Bayes optimality under sparsity (ABOS) Topic of Bogdan et al. (2011), Frommlet et al. (2013) for multiple testing Essential idea for regression

• Two groups model for regressors:

$$P(eta_i
eq 0)=\eta, \hspace{0.2cm} ext{with} \hspace{0.2cm} \eta ext{ small}$$

while *n* and *p* are large

• Compare misclassification rate of procedure with optimal Bayes rule

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

Simulation for unknown  $\sigma$ ,  $\eta \propto p^{-1}$ 

Orthogonal design, p = n

**Misclassification rate** as a function of p,  $\eta(128) = 0.125$ 



Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

Simulation for unknown  $\sigma$ ,  $\eta \propto p^{-1/2}$ 

Orthogonal design, p = n

**Misclassification rate** as a function of *p*,  $\eta(128) = 0.125$ 



Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

Simulation for unknown  $\sigma$ ,  $\eta \propto p^{-1/4}$ 

Orthogonal design, p = n

**Misclassification rate** as a function of p,  $\eta(128) = 0.125$ 



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

Simulation for unknown  $\sigma$ ,  $\eta \propto p^{-1/8}$ 

Orthogonal design, p = n

**Misclassification rate** as a function of p,  $\eta(128) = 0.125$ 



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

Chen and Chen (2008): Prior  $\binom{p}{k_M}^{\kappa-1}$ 

$$EBIC = -2\log L_M + k_M\log n + 2\log \binom{p}{k_M}^{1-\kappa}$$

with  $0 \le \kappa \le 1$ .

•  $\kappa = 1 \Rightarrow \text{ original BIC}$ 

•  $\kappa = 0 \Rightarrow$  asymptotically equivalent with mBIC2

Chen and Chen (2008): Consistency results for EBIC

Under certain assumptions on design matrix for non-orthogonal case

うして ふゆう ふほう ふほう うらつ

Similar consistency results hold for mBIC2

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

Chen and Chen (2008): Prior  $\binom{p}{k_M}^{\kappa-1}$ 

$$EBIC = -2\log L_M + k_M\log n + 2\log \binom{p}{k_M}^{1-\kappa}$$

with  $0 \le \kappa \le 1$ .

- $\kappa = 1 \Rightarrow \text{ original BIC}$
- $\kappa = 0 \quad \Rightarrow \quad \text{asymptotically equivalent with mBIC2}$

#### Chen and Chen (2008): Consistency results for EBIC

• Under certain assumptions on design matrix for non-orthogonal case

うして ふゆう ふほう ふほう うらつ

Similar consistency results hold for mBIC2

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

Chen and Chen (2008): Prior  $\binom{p}{k_M}^{\kappa-1}$ 

$$EBIC = -2\log L_M + k_M\log n + 2\log \binom{p}{k_M}^{1-\kappa}$$

with  $0 \le \kappa \le 1$ .

- $\kappa = 1 \Rightarrow \text{ original BIC}$
- $\kappa = 0 \quad \Rightarrow \quad \text{asymptotically equivalent with mBIC2}$

#### Chen and Chen (2008): Consistency results for EBIC

• Under certain assumptions on design matrix for non-orthogonal case

うして ふゆう ふほう ふほう うらつ

• Similar consistency results hold for mBIC2

roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
0000	0000 0000000 •00	00000000	0000000 0000 0000000	00 0000 00

Chen and Chen (2008): Prior  $\binom{p}{k_M}^{\kappa-1}$ 

$$EBIC = -2\log L_M + k_M\log n + 2\log \binom{p}{k_M}^{1-\kappa}$$

with  $0 \le \kappa \le 1$ .

- $\kappa = 1 \Rightarrow \text{ original BIC}$
- $\kappa = 0 \quad \Rightarrow \quad \text{asymptotically equivalent with mBIC2}$

#### Chen and Chen (2008): Consistency results for EBIC

- Under certain assumptions on design matrix for non-orthogonal case
- Similar consistency results hold for mBIC2

Intro duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 000000 000	00000000	0000000 0000 0000000	00 0000 00

# Comparison of criteria for known $\sigma$

Orthogonal design, p = n = 64



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 00•	00000000	0000000 0000 0000000	00000

# Comparison of criteria for unknown $\sigma$ Orthogonal design, p = n = 64



・ロト 4回ト 4 E ト 4 E ト E - のQ(0)

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Simulation study for GWAS

#### Frommlet et al. (2011 b)

Real SNP Data: POPRES from dbGaP

- 309790 SNPs for 649 individuals (Caucasians)
- k = 40 **causal SNPs** chosen such that MAF between 0.3 and 0.5 Correlation between -0.12 and 0.1
- Simulation of quantitative trait under additive model M

$$Y_i = \sum_{j=1}^{40} eta_j X_{ij} + \epsilon_i, \qquad \epsilon_i \sim \mathcal{N}(0, 1)$$

 $eta_j$  equally spaced between 0.27 and 0.66

• 1000 simulation runs

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Simulation study for GWAS

### Frommlet et al. (2011 b)

Real SNP Data: POPRES from dbGaP

- 309790 SNPs for 649 individuals (Caucasians)
- k = 40 causal SNPs chosen such that MAF between 0.3 and 0.5 Correlation between -0.12 and 0.1
- Simulation of quantitative trait under additive model M

$$Y_i = \sum_{j=1}^{40} \beta_j X_{ij} + \epsilon_i, \qquad \epsilon_i \sim \mathcal{N}(0, 1)$$

 $eta_j$  equally spaced between 0.27 and 0.66

1000 simulation runs

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Simulation study for GWAS

#### Frommlet et al. (2011 b)

Real SNP Data: POPRES from dbGaP

- 309790 SNPs for 649 individuals (Caucasians)
- k = 40 causal SNPs chosen such that MAF between 0.3 and 0.5 Correlation between -0.12 and 0.1
- Simulation of quantitative trait under additive model M

$$Y_i = \sum_{j=1}^{40} eta_j X_{ij} + \epsilon_i, \qquad \epsilon_i \sim \mathcal{N}(0, 1)$$

 $\beta_j$  equally spaced between 0.27 and 0.66

1000 simulation runs

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Simulation study for GWAS

### Frommlet et al. (2011 b)

Real SNP Data: POPRES from dbGaP

- 309790 SNPs for 649 individuals (Caucasians)
- k = 40 causal SNPs chosen such that MAF between 0.3 and 0.5 Correlation between -0.12 and 0.1
- Simulation of quantitative trait under additive model M

$$Y_i = \sum_{j=1}^{40} eta_j X_{ij} + \epsilon_i, \qquad \epsilon_i \sim \mathcal{N}(0, 1)$$

 $\beta_j$  equally spaced between 0.27 and 0.66

1000 simulation runs

Intr	o du	cti	on
000		0	

GWAS Simulation QT

 Case Control studies
 Outlook

 0000000
 00

 0000
 0000

 0000
 0000

 000000
 0000

### Heritability

Total heritability:

$$H^2 = rac{{
m Var}\,\left( {X_M eta_M } 
ight)}{{1 + {
m Var}}\left( {X_M eta_M } 
ight)}$$

Individual heritability:

$$h_j^2 = \frac{\beta_j^2 \operatorname{Var}(X_j)}{1 + \operatorname{Var}(X_M \beta_M)} ,$$

#### Values in our simulation study

Total heritability:  $H^2pprox 0.81$ . Individual heritability:  $h_i^2$  between 0.006 and 0.037

Int	ro	du	cti	on	
oc	0		0		

GWAS Simulation QT

 Case Control studies
 Outlook

 0000000
 00

 0000
 000

 0000
 000

 000000
 000

 000000
 000

### Heritability

Total heritability:

$$H^2 = rac{{{
m Var}}\left( {{X_M}{eta _M}} 
ight)}{{1 + {
m Var}}\left( {{X_M}{eta _M}} 
ight)}$$

Individual heritability:

$$h_j^2 = rac{eta_j^2 {\sf Var} \left( X_j 
ight)}{1 + {\sf Var} \left( X_M eta_M 
ight)} \; ,$$

Values in our simulation study

Total heritability:  $H^2pprox 0.81$ . Individual heritability:  $h_i^2$  between 0.006 and 0.037

Int	ro	du	cti	on	
oc	0		0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look 00 0000 00

### Heritability

Total heritability:

$$H^2 = rac{{{
m Var}}\left( {{X_M}{eta _M}} 
ight)}{{1 + {
m Var}}\left( {{X_M}{eta _M}} 
ight)}$$

Individual heritability:

$$h_j^2 = rac{eta_j^2 {\sf Var} \left( X_j 
ight)}{1 + {\sf Var} \left( X_M eta_M 
ight)} \; ,$$

#### Values in our simulation study

Total heritability:  $H^2 \approx 0.81$ . Individual heritability:  $h_i^2$  between 0.006 and 0.037

◆□> <圖> < => < => < => < => < <</p>

Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Definition of false positives and true positives

#### Problem Causal SNPs are known

Frequently strongly correlated SNPs are selected: Are these to be classified as true or false positives?

#### Common solution

Define threshold value C (E.g. C = 0.7 or C = 0.9). If correlation between detected SNP and causal SNP larger than C  $\Rightarrow$  Classification as true positive

Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Definition of false positives and true positives

#### Problem Causal SNPs are known

Frequently strongly correlated SNPs are selected: Are these to be classified as true or false positives?

#### Common solution

Define threshold value C (E.g. C = 0.7 or C = 0.9).

If correlation between detected SNP and causal SNP larger than C

 $\Rightarrow$  Classification as true positive

ntroduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	0000000	0000000 0000 0000000	00 0000 00

Comparison between model selection and single marker tests

#### Four methods:

	Per marker	Model selection
Control of FWER	Bonferroni	mBIC
Control of FDR	Benjamini Hochberg	mBIC2

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ ―臣 …の�?

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look

### False Discovery Rate

#### FDR and Power

Cutoffs define TP and FP for correlated marker





Power

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 0000000	00 0000 00

Explanation for low power

Correct model: 
$$Y_i = \beta_0 + \sum_{l \in M^*} \beta_l X_{il} + \epsilon_i$$

Test statistic:  $F_j = \frac{(n-2)MSS(X_j)}{RSS(X_j)}$ 

Non-centrality parameters:

$$\nu_{M,j} = \frac{\left(\sum_{l=1}^{k} \beta_l \operatorname{Cov} (X_j, X_l)\right)^2}{\sigma^2 \operatorname{Var} (X_j)}$$
  
$$\nu_{R,j} = \sum_{l \in M^* \setminus \{j\}} \sum_{r \in M^* \setminus \{j\}} \frac{\beta_l \beta_r}{\sigma^2} \left( \operatorname{Cov} (X_l, X_r) - \frac{\operatorname{Cov} (X_l, X_j) \operatorname{Cov} (X_r, X_j)}{\operatorname{Var} (X_j)} \right)$$

・ロト ・個ト ・ヨト ・ヨト 三日 -

 $\nu_{R,i}$  contains contribution of all other causal SNPs

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 0000000	00 0000 00

Explanation for low power

Correct model: 
$$Y_i = \beta_0 + \sum_{I \in M^*} \beta_I X_{iI} + \epsilon_i$$

Test statistic:  $F_j = \frac{(n-2)MSS(X_j)}{RSS(X_j)}$ 

#### Non-centrality parameters:

$$\nu_{M,j} = \frac{\left(\sum_{l=1}^{k} \beta_l \operatorname{Cov} (X_j, X_l)\right)^2}{\sigma^2 \operatorname{Var} (X_j)}$$
  
$$\nu_{R,j} = \sum_{l \in M^* \setminus \{j\}} \sum_{r \in M^* \setminus \{j\}} \frac{\beta_l \beta_r}{\sigma^2} \left( \operatorname{Cov} (X_l, X_r) - \frac{\operatorname{Cov} (X_l, X_j) \operatorname{Cov} (X_r, X_j)}{\operatorname{Var} (X_j)} \right)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

 $\nu_{R,i}$  contains contribution of all other causal SNPs

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 0000000	00 0000 00

Explanation for low power

Correct model: 
$$Y_i = \beta_0 + \sum_{I \in M^*} \beta_I X_{iI} + \epsilon_i$$

Test statistic:  $F_j = \frac{(n-2)MSS(X_j)}{RSS(X_j)}$ 

#### Non-centrality parameters:

$$\nu_{M,j} = \frac{\left(\sum_{l=1}^{k} \beta_l \operatorname{Cov} (X_j, X_l)\right)^2}{\sigma^2 \operatorname{Var} (X_j)}$$
  
$$\nu_{R,j} = \sum_{l \in M^* \setminus \{j\}} \sum_{r \in M^* \setminus \{j\}} \frac{\beta_l \beta_r}{\sigma^2} \left( \operatorname{Cov} (X_l, X_r) - \frac{\operatorname{Cov} (X_l, X_j) \operatorname{Cov} (X_r, X_j)}{\operatorname{Var} (X_j)} \right)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

 $\nu_{R,j}$  contains contribution of all other causal SNPs

duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000	0000 0000000 000	000000000	0000000 0000 0000000	00 0000 00

## Power for mBIC2 and BH

#### Bigger surprise:

Intro



▲ロト ▲圖 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● のへで

troduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 000000 000	000000000	0000000 0000 0000000	00 0000 00

### Explanation for problems of BH



Non-centrality parameter:

$$\sqrt{\nu_{M,j}} = \left| \frac{\beta_j}{\sigma} \sqrt{\operatorname{Var}(X_j)} + \frac{\sum_{l \neq j} \beta_l \operatorname{Cov}(X_j, X_l)}{\sigma \sqrt{\operatorname{Var}(x_j)}} \right|$$

Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆@▶ ◆臣▶ ◆臣▶ ─ 臣…

Out look 00 0000 0000

## Practical conclusions for GWAS analysis

In case of complex traits:

- Single marker tests (SMT) have low power
  - $\Rightarrow$  One aspect in the discussion about "missing heritability"
- SMT have difficulties to rank the importance of causal SNPs
   ⇒ Problem with replicability in GWAS
- For the same reason SMT systematically detect some false positives which are not correlated with any of the causal SNPs

Model selection approach helps to some extent

Modifications of BIC 0000 0000000 000 GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

## Practical conclusions for GWAS analysis

In case of complex traits:

- Single marker tests (SMT) have low power
  - $\Rightarrow$  One aspect in the discussion about "missing heritability"
- SMT have difficulties to rank the importance of causal SNPs
   ⇒ Problem with replicability in GWAS
- For the same reason SMT systematically detect some false positives which are not correlated with any of the causal SNPs

Model selection approach helps to some extent

Modifications of BIC

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Practical conclusions for GWAS analysis

In case of complex traits:

- Single marker tests (SMT) have low power
  - $\Rightarrow$  One aspect in the discussion about "missing heritability"
- SMT have difficulties to rank the importance of causal SNPs
  - $\Rightarrow$  Problem with replicability in GWAS
- For the same reason SMT systematically detect some false positives which are not correlated with any of the causal SNPs

Model selection approach helps to some extent
Modifications of BIC

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Practical conclusions for GWAS analysis

In case of complex traits:

- Single marker tests (SMT) have low power
  - $\Rightarrow$  One aspect in the discussion about "missing heritability"
- SMT have difficulties to rank the importance of causal SNPs
  - $\Rightarrow$  Problem with replicability in GWAS
- For the same reason SMT systematically detect some false positives which are not correlated with any of the causal SNPs

Model selection approach helps to some extent

Intr	o du	cti	on
000		0	

GWAS Simulation QT

Case Control studies • 000000 • 000

・ロト ・ 御 ト ・ ヨ ト ・ ヨ ト ・ ヨ ・

Out look 00 0000 00

# Logistic Regression for Case Control

Usual model  $Y_i$  is Bernoulli,  $P(Y_i = 1) = p_i$ , with

$$\log(p_i/(1-p_i))=eta_0+\sum_{i\in M}eta_jX_{ij}$$

#### We discuss three selected methods

- HLASSO: Hoggart, Balding (2008)
- GWASelect: He and Lin (2011)
- MOSGWA: Our own approach

Intr	o du	cti	on
000		0	

GWAS Simulation QT

Case Control studies • 000000 000 000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# Logistic Regression for Case Control

Usual model  $Y_i$  is Bernoulli,  $P(Y_i = 1) = p_i$ , with

$$\log(p_i/(1-p_i))=eta_0+\sum_{i\in M}eta_jX_{ij}$$

#### We discuss three selected methods

- HLASSO: Hoggart, Balding (2008)
- GWASelect: He and Lin (2011)
- MOSGWA: Our own approach

ntro du ction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	<b>000000</b> 0000 0000000	00 0000 00

## Hoggart, ..., Balding: Bioinformatics (2008)

More Bayesian approach using shrinkage priors on coefficients of logistic regression models

- Gaussian prior (implemented, but no selection)
- Double exponential prior (DE)  $\Rightarrow$  equivalent to Lasso
- Normal exponential gamma (NEG)  $\Rightarrow$  Hyper Lasso

## Densities of DE and NEG

$$DE(\beta|\xi) = \int_{\sigma^2=0}^{\infty} N(\beta|0,\sigma^2) Ga(\sigma^2|1,\xi^2/2) d\sigma^2 = \frac{\xi}{2} \exp\left(-\xi|\beta|\right)$$

$$NEG(\beta|\lambda,\gamma) = \int_{\Psi=0}^{\infty} \int_{\sigma^2=0}^{\infty} N(\beta|0,\sigma^2) Ga(\sigma^2|1,\Psi) Ga(\Psi|\lambda,\gamma^2) d\sigma^2 d\Psi$$

ntro du ction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	<b>000000</b> 0000 0000000	00 0000 00

## Hoggart, ..., Balding: Bioinformatics (2008)

More Bayesian approach using shrinkage priors on coefficients of logistic regression models

- Gaussian prior (implemented, but no selection)
- Double exponential prior (DE)  $\Rightarrow$  equivalent to Lasso
- Normal exponential gamma (NEG)  $\Rightarrow$  Hyper Lasso

## Densities of DE and NEG

$$DE(\beta|\xi) = \int_{\sigma^2=0}^{\infty} N(\beta|0,\sigma^2) Ga(\sigma^2|1,\xi^2/2) d\sigma^2 = \frac{\xi}{2} \exp\left(-\xi|\beta|\right)$$

$$NEG(\beta|\lambda,\gamma) = \int_{\Psi=0}^{\infty} \int_{\sigma^2=0}^{\infty} N(\beta|0,\sigma^2) Ga(\sigma^2|1,\Psi) Ga(\Psi|\lambda,\gamma^2) d\sigma^2 d\Psi$$

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	00000000000000000000000000000000000000	00 0000 00

# **NEG** priors

Logarithm of prior densities fixed to have the same density at the origin (Taken from Hoggart et al., 2008)

$$NEG(\beta|\lambda,\gamma) \propto \exp\left(rac{eta^2}{4\gamma^2}
ight) D_{-2\lambda-1}\left(rac{|eta|}{\gamma}
ight)$$

where *D* is parabolic cylinder function



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	000000 0000 0000000	00 0000 00

# **NEG** priors

Logarithm of prior densities fixed to have the same density at the origin (Taken from Hoggart et al., 2008)

$$\mathsf{NEG}(eta|\lambda,\gamma)\propto\exp\left(rac{eta^2}{4\gamma^2}
ight) \mathsf{D}_{-2\lambda-1}\left(rac{|eta|}{\gamma}
ight)$$

where D is parabolic cylinder function



◆□▶ ◆◎▶ ◆□▶ ◆□▶ ─ □

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies

・ロト ・ 御 ト ・ ヨ ト ・ ヨ ト

э.

Out look 00 0000 00

## **HLASSO**

## Optimisation algorithm

Not fully Bayesian, but searching only for posterior mode

$$\log p(\beta|X, Y) = \ell(\beta) - q(\beta) + const ,$$

with

$$\ell(\beta) := \log L(\beta|X, Y), \quad q(\beta) := -\log NEG(\beta|\lambda, \gamma)$$

HLASSO has rather efficient implementation to find optimum of  $\log p$  CLG algorithm (cyclic coordinate descent) with clever bounds to speed up

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー わえぐ

Outlook

## **HLASSO**

## Optimisation algorithm

Not fully Bayesian, but searching only for posterior mode

$$\log p(\beta|X, Y) = \ell(\beta) - q(\beta) + const$$
,

with

$$\ell(\beta) := \log L(\beta|X, Y), \quad q(\beta) := -\log NEG(\beta|\lambda, \gamma)$$

HLASSO has rather efficient implementation to find optimum of  $\log p$  CLG algorithm (cyclic coordinate descent) with clever bounds to speed up

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## CLG - Basic idea

Run iteratively and repeatedly through all coefficients with component-wise Newton

$$\beta_j^{new} = \beta_j - \frac{\frac{\partial}{\partial \beta_j} \ell(\beta) - q'(\beta_j)}{\frac{\partial^2}{\partial \beta_j^2} \ell(\beta) - q''(\beta_j)}$$

If 
$$eta_j^{\mathit{new}}\cdoteta_j<0$$
 then set  $eta_j^{\mathit{new}}=0$ 

## Specifically if $\beta_j = 0$

Consider both limits  $\beta_j = 0^+$  and  $\beta_j = 0^-$ No change of sign is equivalent to

$$\left|\frac{\partial}{\partial\beta_j}\ell(\beta)\right|_{\beta_j=0}>q'(\beta_j=0^+)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

ntroduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	<b>0000000</b> 0000 0000000	00 0000 00

### CLG - Basic idea

Run iteratively and repeatedly through all coefficients with component-wise Newton

$$\beta_j^{new} = \beta_j - \frac{\frac{\partial}{\partial \beta_j} \ell(\beta) - q'(\beta_j)}{\frac{\partial^2}{\partial \beta_j^2} \ell(\beta) - q''(\beta_j)}$$

If 
$$eta_j^{\mathit{new}}\cdoteta_j<0$$
 then set  $eta_j^{\mathit{new}}=0$ 

## Specifically if $\beta_j = 0$

Consider both limits  $\beta_j = 0^+$  and  $\beta_j = 0^-$ No change of sign is equivalent to

$$\left|\frac{\partial}{\partial\beta_j}\ell(\beta)\right|_{\beta_j=0}>q'(\beta_j=0^+)$$

・ロト ・ 日 ・ エ ヨ ・ ト ・ 日 ・ う へ つ ・

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Parameter tuning to control FWER

Asymptotic normality under null gives  $\hat{\beta}_j \sim \mathcal{N}\left(0, \frac{n_0+n_1}{n_0n_1}\right)$ From last relationship of previous slide one then gets relationship to determine type I error

$$q'(\beta_j = 0^+) = \sqrt{\frac{n_0 n_1}{n_0 + n_1}} \Phi^{-1}(1 - \alpha/2)$$

• Double Exponential:  $q'(eta_j=0^+)=\xi$ 

• NEG: (Careful: Typo in paper)

$$q'(\beta_j = 0^+) = const \cdot rac{2\lambda + 1}{\gamma}$$

・ロト ・ 日 ・ モート ・ 田 ・ うへで

Now actually two parameters to be fitted, Hoggart et al. use  $\lambda=$  0.05 (after trying other parameters) and the  $\gamma$  follows

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Parameter tuning to control FWER

Asymptotic normality under null gives  $\hat{\beta}_j \sim \mathcal{N}\left(0, \frac{n_0+n_1}{n_0n_1}\right)$ From last relationship of previous slide one then gets relationship to determine type I error

$$q'(\beta_j = 0^+) = \sqrt{\frac{n_0 n_1}{n_0 + n_1}} \Phi^{-1}(1 - \alpha/2)$$

• Double Exponential:  $q'(eta_j=0^+)=\xi$ 

• NEG: (Careful: Typo in paper)

$$q'(\beta_j = 0^+) = const \cdot \frac{2\lambda + 1}{\gamma}$$

・ロト ・ 日 ・ モート ・ 田 ・ うへで

Now actually two parameters to be fitted, Hoggart et al. use  $\lambda=$  0.05 (after trying other parameters) and the  $\gamma$  follows

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

## Parameter tuning to control FWER

Asymptotic normality under null gives  $\hat{\beta}_j \sim \mathcal{N}\left(0, \frac{n_0+n_1}{n_0n_1}\right)$ From last relationship of previous slide one then gets relationship to determine type I error

$$q'(\beta_j = 0^+) = \sqrt{\frac{n_0 n_1}{n_0 + n_1}} \Phi^{-1}(1 - \alpha/2)$$

- Double Exponential:  $q'(eta_j=0^+)=\xi$
- NEG: (Careful: Typo in paper)

$$q'(eta_j=0^+)=const\cdotrac{2\lambda+1}{\gamma}$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ の へ ()

Now actually two parameters to be fitted, Hoggart et al. use  $\lambda=0.05$  (after trying other parameters) and the  $\gamma$  follows

Int	ro d	u ct	ion	
	00	00		

GWAS Simulation QT

Case Control studies 000000● 000000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000 00

# **HLASSO**

## Parameter tuning to control FWER

<code>HLASSO</code> offers choice of parameter  $\alpha$  which corresponds to uncorrected significance level

 $\Rightarrow$  Choosing for example lpha= 0.05/p works

## Difference between LASSO and HLASSO

- Lighter tails of DE distribution ⇒ more shrinkage ⇒ Correlated SNPs tend to enter model to explain full effect of causal SNP
- NEG prior has heavier tails  $\Rightarrow$  less shrinkage

Int	ro	du	ct	io	n
	00		0		

GWAS Simulation QT

Case Control studies 000000● 000000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Outlook

**HLASSO** 

## Parameter tuning to control FWER

<code>HLASSO</code> offers choice of parameter  $\alpha$  which corresponds to uncorrected significance level

 $\Rightarrow$  Choosing for example lpha= 0.05/p works

## Difference between LASSO and HLASSO

- Lighter tails of DE distribution ⇒ more shrinkage
   ⇒ Correlated SNPs tend to enter model to explain full effect of causal SNP
- NEG prior has heavier tails  $\Rightarrow$  less shrinkage

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies	Out look
0000000	00
•000	0000
0000000	00

・ロト ・ 四ト ・ 日ト ・ 日 ・

# GWASelect

#### He and Lin: Bioinformatics (2011)

## Rough outline of algorithm

- 1. Sure Independence Screening (SIS)
- 2. Lasso for logistic regression
- 3. 'Pruning' of correlated SNPs

These steps are iterated based on conditional score tests  $\Rightarrow$  ISIS

# Stability selection

Meinshausen, Bühlmann (2010) Derforms the sub-le arreadure on 50 a

Perform the whole procedure on 50 subsamples

(randomly select 50% cases and 50 % controls)

 $\Rightarrow$  compute selection probabilities

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies	Out look
0000000	00
0000	0000
0000000	00

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ●

# GWASelect

#### He and Lin: Bioinformatics (2011)

## Rough outline of algorithm

- 1. Sure Independence Screening (SIS)
- 2. Lasso for logistic regression
- 3. 'Pruning' of correlated SNPs

#### These steps are iterated based on conditional score tests $\quad\Rightarrow\quad$ ISIS

# Stability selection

Meinshausen, Bühlmann (2010)

Perform the whole procedure on 50 subsamples

(randomly select 50% cases and 50 % controls)

 $\Rightarrow$  compute selection probabilities

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies	Out look
0000000	00
0000	0000
0000000	00

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

# GWASelect

#### He and Lin: Bioinformatics (2011)

## Rough outline of algorithm

- 1. Sure Independence Screening (SIS)
- 2. Lasso for logistic regression
- 3. 'Pruning' of correlated SNPs

These steps are iterated based on conditional score tests  $\Rightarrow$  ISIS

## Stability selection

Meinshausen, Bühlmann (2010)

Perform the whole procedure on 50 subsamples

(randomly select 50% cases and 50 % controls)

 $\Rightarrow$  compute selection probabilities

Modifications of BIC 0000 0000000 GWAS Simulation QT

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# GWASelect

#### First iteration

Preselection based on marginal tests (Cochran Armitage trend test), SIS theory: Consider  $0.9n/(4 \log n)$  SNPs with largest test statistic

#### Lasso

Model selection using Lasso on selected set of SNPs (cyclic coordinate decent for optimization) For 'dynamic' GWASelect tuning parameter  $\lambda$  determined by 5-fold cross validation

#### Pruning of Correlated SNPs

Remove SNPs from model which have pairwise correlation |R| < 0.8

Int	rod	luo	ati	on	
	000	00	С		

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
000000	00

▲ロト ▲圖ト ▲ヨト ▲ヨト ヨー のへで

# GWASelect

#### First iteration

Preselection based on marginal tests (Cochran Armitage trend test), SIS theory: Consider  $0.9n/(4 \log n)$  SNPs with largest test statistic

#### Lasso

Model selection using Lasso on selected set of SNPs (cyclic coordinate decent for optimization) For 'dynamic' GWASelect tuning parameter  $\lambda$  determined by 5-fold cross validation

#### Pruning of Correlated SNPs

Remove SNPs from model which have pairwise correlation |R| < 0.8

Int	rod	luo	ati	on	
	000	00	С		

GWAS Simulation QT

Case Control studies	O ut look
000000	00
0000	0000
0000000	00

▲ロト ▲圖ト ▲ヨト ▲ヨト ヨー のへで

# GWASelect

#### First iteration

Preselection based on marginal tests (Cochran Armitage trend test), SIS theory: Consider  $0.9n/(4 \log n)$  SNPs with largest test statistic

#### Lasso

Model selection using Lasso on selected set of SNPs (cyclic coordinate decent for optimization) For 'dynamic' GWASelect tuning parameter  $\lambda$  determined by 5-fold cross validation

## Pruning of Correlated SNPs

Remove SNPs from model which have pairwise correlation |R| < 0.8

Modifications of BIC

GWAS Simulation QT

 Case Control studies
 Outlook

 000000
 00

 000000
 00000

 000000
 00000

▲ロト ▲圖ト ▲ヨト ▲ヨト ヨー のへで

# GWASelect

### Second and third iteration

After pruning t SNPs left in model, say  $X_1, \ldots, X_t$ Interested to consider all influence of other SNPs conditional on SNPs already in the model, specifically for all  $X_r, r > t$ 

$$\log(p_i/(1-p_i)) = \beta_0 + \sum_{j=1}^t \beta_j X_{ij} + \gamma X_r$$

#### we would like to know if $\gamma=\mathbf{0}$

LRT too timeconsuming, but **Scoretest** very fast alternative  $\Rightarrow$  SIS step based on Score test statistics Keep  $0.05n/(4 \log n)$  best SNPs Then again Lasso and pruning of correlated SNPs

Modifications of BIC

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
0000000	00

# GWASelect

#### Second and third iteration

After pruning t SNPs left in model, say  $X_1, \ldots, X_t$ Interested to consider all influence of other SNPs conditional on SNPs already in the model, specifically for all  $X_r, r > t$ 

$$\log(p_i/(1-p_i)) = \beta_0 + \sum_{j=1}^t \beta_j X_{ij} + \gamma X_r$$

we would like to know if  $\gamma=0$  LRT too timeconsuming, but  ${\bf Scoretest}$  very fast alternative

⇒ SIS step based on Score test statistics Keep 0.05 $n/(4 \log n)$  best SNPs Then again Lasso and pruning of correlated SNPs

Modifications of BIC

GWAS Simulation QT

 Case Control studies
 Outlook

 000000
 00

 000000
 00000

 000000
 00000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

# GWASelect

### Second and third iteration

After pruning t SNPs left in model, say  $X_1, \ldots, X_t$ Interested to consider all influence of other SNPs conditional on SNPs already in the model, specifically for all  $X_r, r > t$ 

$$\log(p_i/(1-p_i)) = \beta_0 + \sum_{j=1}^t \beta_j X_{ij} + \gamma X_r$$

we would like to know if  $\gamma = 0$ LRT too timeconsuming, but **Scoretest** very fast alternative  $\Rightarrow$  SIS step based on Score test statistics Keep  $0.05n/(4 \log n)$  best SNPs Then again Lasso and pruning of correlated SNPs

Intro	du	cti	on
000	oc	0	

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
0000000	00

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

## **GWAS**elect

#### Software

- Stand alone program, difference between GWASelect and d-GWASelect
- In d-GWASelect only parameter to choose is threshold from stability selection

Recommended values: between 0.1 and 0.2

ntroduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

# MOSGWA for Case Control

## Based on criterion mBIC2

$$mBIC2 = -2 \log L_M + k_M [\log(np^2/4)] - 2 \log k_m!$$

#### Specific issue with logistic regression

 $p > n \Rightarrow$  problem of **complete separation** occurs even for large  $n \Rightarrow$  we use Firth correction

$$L^*(\beta \mid Y, X) := L(\beta) \cdot |I(\beta)|^{1/2}$$

 $|I(\beta)|^{1/2} \dots$  Jeffreys prior,  $I(\beta) = -E\left(\frac{\partial^2 \beta}{\partial \beta_r \partial \beta_s} \log L(\beta)\right)^2$ Neither HLASSO nor GWASelect have this problem, as they implicitl penalize too large  $\beta_i$ 

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

# MOSGWA for Case Control

Based on criterion mBIC2

$$mBIC2 = -2 \log L_M + k_M [\log(np^2/4)] - 2 \log k_m!$$

#### Specific issue with logistic regression

 $p > n \Rightarrow$  problem of **complete separation** occurs even for large  $n \Rightarrow$  we use Firth correction

$$L^*(\beta \mid Y, X) := L(\beta) \cdot |I(\beta)|^{1/2}$$

 $|I(\beta)|^{1/2}$  ... Jeffreys prior ,  $I(\beta) = -E\left(\frac{\partial^2\beta}{\partial\beta_r\partial\beta_s}\log L(\beta)\right)^2$ 

Neither HLASSO nor GWAS elect have this problem, as they implicitly penalize too large  $\beta_j$ 

ション ふゆ く 山 マ チャット しょうくしゃ

nt roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

# MOSGWA for Case Control

Based on criterion mBIC2

$$mBIC2 = -2 \log L_M + k_M [\log(np^2/4)] - 2 \log k_m!$$

#### Specific issue with logistic regression

 $p > n \Rightarrow$  problem of **complete separation** occurs even for large  $n \Rightarrow$  we use Firth correction

$$L^*(\beta|Y,X) := L(\beta) \cdot |I(\beta)|^{1/2}$$

ション ふゆ く 山 マ チャット しょうくしゃ

 $|I(\beta)|^{1/2}$  ... Jeffreys prior ,  $I(\beta) = -E\left(\frac{\partial^2 \beta}{\partial \beta_r \partial \beta_s} \log L(\beta)\right)^2$ Neither HLASSO nor GWASelect have this problem, as they implicitly penalize too large  $\beta_i$ 

Int	ro	d٢	1 C1	i	0	n	
	0	00	DC				

GWAS Simulation QT

・ロッ ・雪 ・ ・ ヨ ・ ・

э

# MOSGWA, Model Search

#### Major issue for practical application $\Rightarrow 2^p$ potential models

# Strategies

Computation of ML much more time consuming than for linear regression  $\Rightarrow$  Even more important than for linear regression to keep models small

- Preselection of markers using marginal tests (Compare SIS)
- Heuristic greedy search procedures (as described in the next slide)
- Genetic Algorithms (as described later)

tro duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

Major issue for practical application  $\Rightarrow 2^p$  potential models

## Strategies

Computation of ML much more time consuming than for linear regression

 $\Rightarrow$   $% \left( E_{1},E_{2},E_{1},E_{2},E_{1},E_{2},E$ 

▲ロト ▲圖ト ▲ヨト ▲ヨト ヨー のへで

- Preselection of markers using marginal tests (Compare SIS)
- Heuristic greedy search procedures (as described in the next slide)
- Genetic Algorithms (as described later)

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

Major issue for practical application  $\Rightarrow 2^p$  potential models

## Strategies

Computation of ML much more time consuming than for linear regression

 $\Rightarrow$   $% \left( {{\rm Even more \ important \ than \ for \ linear \ regression \ to \ keep \ models \ small \ } \right)$ 

ション ふゆ アメリア メリア しょうくの

#### • Preselection of markers using marginal tests (Compare SIS)

- Heuristic greedy search procedures (as described in the next slide)
- Genetic Algorithms (as described later)

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

Major issue for practical application  $\Rightarrow 2^p$  potential models

## Strategies

Computation of ML much more time consuming than for linear regression

 $\Rightarrow$   $\;$  Even more important than for linear regression to keep models small

うして ふゆう ふほう ふほう うらう

- Preselection of markers using marginal tests (Compare SIS)
- Heuristic greedy search procedures (as described in the next slide)
- Genetic Algorithms (as described later)

it roduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 0000 000000	00 0000 00

Major issue for practical application  $\Rightarrow 2^p$  potential models

## Strategies

Computation of ML much more time consuming than for linear regression

 $\Rightarrow$   $\;$  Even more important than for linear regression to keep models small

ション ふゆ く 山 マ チャット しょうくしゃ

- Preselection of markers using marginal tests (Compare SIS)
- Heuristic greedy search procedures (as described in the next slide)
- Genetic Algorithms (as described later)

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Model Search strategy

## Step 1

- Preselection and sorting based on marginal tests (CAT)
- Fast stepwise search (to be specified)

# Step 2

• Preselection and sorting based on Score tests conditional on model obtained in Step 1

ション ふゆ アメリア メリア しょうくの

- Idea: For logistic regression Score test much faster than LRT
- Fast stepwise search (to be specified)

Described strategy good, but tends to get stuck in too small models

 $\Rightarrow$  Start with search using milder criterion

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	000000 000 000000	00 0000 00

# Model Search strategy

## Step 1

- Preselection and sorting based on marginal tests (CAT)
- Fast stepwise search (to be specified)

# Step 2

 Preselection and sorting based on Score tests conditional on model obtained in Step 1
 Idea: For logistic regression Score test much faster than LRT

ション ふゆ く 山 マ チャット しょうくしゃ

• Fast stepwise search (to be specified)

Described strategy good, but tends to get stuck in too small models

 $\Rightarrow$  Start with search using milder criterion
Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	000000 000 000000	00 0000 00

# Model Search strategy

## Step 1

- Preselection and sorting based on marginal tests (CAT)
- Fast stepwise search (to be specified)

# Step 2

 Preselection and sorting based on Score tests conditional on model obtained in Step 1
Idea: For logistic regression Score test much faster than LRT

・ロト ・ 日 ・ エ ヨ ・ ト ・ 日 ・ うらつ

• Fast stepwise search (to be specified)

Described strategy good, but tends to get stuck in too small models

 $\Rightarrow$  Start with search using milder criterion

troduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	000000 000 000000	00 0000 00

# Model Search strategy

#### Fast stepwise search

Starting point: Regressors sorted by test statistic lterate the following three steps till no further improvement

- Directed forward search: Find the first regressor in sorted list which decreases mBIC2 and add to model
- Exchange step: See if substituting any regressor in model with candidate SNP foe exchange decreases mBIC2 *Candidates*: neighboring SNPs or in first step SNPs preselected with CAT

ション ふゆ く 山 マ チャット しょうくしゃ

• Backward step: Routine backward elimination

Fast enough to deal with full GWAS data sets

ntro duction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
00000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

# Model Search strategy

#### Fast stepwise search

Starting point: Regressors sorted by test statistic lterate the following three steps till no further improvement

- Directed forward search: Find the first regressor in sorted list which decreases mBIC2 and add to model
- Exchange step: See if substituting any regressor in model with candidate SNP foe exchange decreases mBIC2 *Candidates*: neighboring SNPs or in first step SNPs preselected with CAT

ション ふゆ く 山 マ チャット しょうくしゃ

• Backward step: Routine backward elimination

Fast enough to deal with full GWAS data sets

Modifications of BIC

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
0000000	00

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ●

# Simulation study

## Simulation study

#### Comparison of MOSGWA, GWASelect and HLasso Again SNP data from POPRES (dbGaP), now more than 4000 individuals

#### First simulation under global null

For four different sets of SNPs Chr1, Chr1 + Chr2, Chr1 - Chr4, Chr1 - Chr6

#### Second simulation

24 causal SNPs (uncorrelated, MAF > 0.3) Simulate 200 instances under logistic regression model

Effect  $\beta_j$  sizes between 0.2 and 0.26 Half of causal SNPs removed before search

Modifications of BIC

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
0000000	00

ション ふゆ く 山 マ チャット しょうくしゃ

## Simulation study

## Simulation study

#### Comparison of MOSGWA, GWASelect and HLasso Again SNP data from POPRES (dbGaP), now more than 4000 individuals

#### First simulation under global null

#### For four different sets of SNPs Chr1, Chr1 + Chr2, Chr1 - Chr4, Chr1 - Chr6

#### Second simulation

24 causal SNPs (uncorrelated, MAF > 0.3) Simulate 200 instances under logistic regression model

Effect  $\beta_j$  sizes between 0.2 and 0.26 Half of causal SNPs removed before search

Modifications of BIC

GWAS Simulation QT

Case Control studies	O ut look
0000000	00
0000	0000
0000000	00

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ●

## Simulation study

## Simulation study

Comparison of MOSGWA, GWASelect and HLasso Again SNP data from POPRES (dbGaP), now more than 4000 individuals

#### First simulation under global null

For four different sets of SNPs Chr1, Chr1 + Chr2, Chr1 - Chr4, Chr1 - Chr6

#### Second simulation

24 causal SNPs (uncorrelated, MAF > 0.3) Simulate 200 instances under logistic regression model

Effect  $\beta_j$  sizes between 0.2 and 0.26 Half of causal SNPs removed before search

Intr	o du	cti	on
000		0	

GWAS Simulation QT

## Simulation under global null

Average number of False Positives HLASSO with parameters 0.1/p, 0.2/p, 0.3/pGWASelect with parameters 0.1, 0.2, 0.3



MOSGWA and GWASelect



▲ロト ▲圖ト ▲画ト ▲画ト 三直 - の久で

troduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	O ut look
00000	0000 0000000 000	00000000	000000 000000 000000	00 0000 00

Simulation under Model ( $k^* = 24$ )

#### False positives and Power

Cutoffs define TP and FP for correlated marker

FΡ



Power



▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ \_ 圖 \_ のへで

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# MOSGWA

#### Model Selection for Genome Wide Association

Package

- Written in C++ by Bodenstorfer, Dolejsi, Ruhaltinger
- Aim: Professional software for genetic researchers
- Currently
  - SNP array data (PLINK format and HDF5)
  - Linear and Logistic Regression
  - Allows for inclusion of covariates
- First version on-line since this week! https://sourceforge.net/projects/mosgwa/

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# MOSGWA

#### Model Selection for Genome Wide Association

Package

- Written in C++ by Bodenstorfer, Dolejsi, Ruhaltinger
- Aim: Professional software for genetic researchers
- Currently
  - SNP array data (PLINK format and HDF5)
  - Linear and Logistic Regression
  - Allows for inclusion of covariates
- First version on-line since this week! https://sourceforge.net/projects/mosgwa/

Int	ro	du	cti	0	n
	0	oc	00		

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look

# MOSGWA

#### Model Selection for Genome Wide Association

Package

- Written in C++ by Bodenstorfer, Dolejsi, Ruhaltinger
- Aim: Professional software for genetic researchers
- Currently
  - SNP array data (PLINK format and HDF5)
  - Linear and Logistic Regression
  - Allows for inclusion of covariates
- First version on-line since this week! https://sourceforge.net/projects/mosgwa/

Int	ro d	цς	tio	n
	00		)	

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 0● 0000

# MOSGWA

- Better search strategies Specifically genetic algorithm for search
- Mixed models (relatively easy for QT)
- Sequencing data, including methods for rare SNPs
- Admixture mapping
- Logic regression for interactions, etc.

Int	ro d	цς	tio	n
	00		)	

GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 0● 0000

# MOSGWA

- Better search strategies Specifically genetic algorithm for search
- Mixed models (relatively easy for QT)
- Sequencing data, including methods for rare SNPs
- Admixture mapping
- Logic regression for interactions, etc.

Intro	duction	
000	000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Out look 0● 0000

## MOSGWA

- Better search strategies Specifically genetic algorithm for search
- Mixed models (relatively easy for QT)
- Sequencing data, including methods for rare SNPs
- Admixture mapping
- Logic regression for interactions, etc.

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Out look 0● 0000

## MOSGWA

- Better search strategies Specifically genetic algorithm for search
- Mixed models (relatively easy for QT)
- Sequencing data, including methods for rare SNPs
- Admixture mapping
- Logic regression for interactions, etc.

Intro	du	cti	on	
000	oc	0		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 0● 0000

## MOSGWA

- Better search strategies Specifically genetic algorithm for search
- Mixed models (relatively easy for QT)
- Sequencing data, including methods for rare SNPs
- Admixture mapping
- Logic regression for interactions, etc.

Modifications of BIC 0000 0000000 GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000

# Memetic algorithm

## Basic idea of genetic algorithm

- Work with population of models
- mBIC2 as measure of fitness of a model
- Use evolutionary dynamic to increase fitness of population Selection, recombination, mutation

## Memetic algorithm

Modifications of BIC 0000 0000000 GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look ○○ ●○○○

# Memetic algorithm

## Basic idea of genetic algorithm

- Work with population of models
- mBIC2 as measure of fitness of a model
- Use evolutionary dynamic to increase fitness of population Selection, recombination, mutation

### Memetic algorithm

Modifications of BIC 0000 0000000 GWAS Simulation QT

Case Control studies

◆□▶ ◆□▶ ★□▶ ★□▶ □ のQ@

Out look 00 0000

# Memetic algorithm

### Basic idea of genetic algorithm

- Work with population of models
- mBIC2 as measure of fitness of a model
- Use evolutionary dynamic to increase fitness of population Selection, recombination, mutation

#### Memetic algorithm

Modifications of BIC 0000 0000000 GWAS Simulation QT

Case Control studies

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Out look ○○ ●○○○

# Memetic algorithm

### Basic idea of genetic algorithm

- Work with population of models
- mBIC2 as measure of fitness of a model
- Use evolutionary dynamic to increase fitness of population Selection, recombination, mutation

### Memetic algorithm



Inti	ro d	цс	:ti	10	1
00	00	00	С		

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look 00 00●0

# Memetic algorithms

## MA particularly designed for our application

- Frommlet et al. (2012): Simulations for QTL mapping:
  - $\Rightarrow$  GA finds frequently better solution than stepwise search
- Implementation for GWAS finished, not yet integrated in MOSGWA

## Estimation of marker posteriors

MA gives many good solutions  $\Rightarrow$  model posteriors

$$\sum_{M \in \mathcal{M}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right) \;.$$

$$\mathsf{P}(j_r|Y) \approx \frac{\sum\limits_{M \in \mathcal{M}_r^{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right)}{\sum\limits_{M' \in \mathcal{M}^{Pool}} \exp\left(-\mathsf{mBIC}(M')/2\right)}$$

Int	ro	du	ct	io	n	
	0	00	0			

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look

# Memetic algorithms

## MA particularly designed for our application

- Frommlet et al. (2012): Simulations for QTL mapping:
  - $\Rightarrow \quad \mathsf{GA} \text{ finds frequently better solution than stepwise search}$
- Implementation for GWAS finished, not yet integrated in MOSGWA

## Estimation of marker posteriors

MA gives many good solutions  $\Rightarrow$  model posteriors

$$\sum_{M \in \mathcal{M}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right) \,.$$

$$\mathsf{P}(j_r|Y) \approx \frac{\sum\limits_{M \in \mathcal{M}_r^{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right)}{\sum\limits_{M' \in \mathcal{M}^{Pool}} \exp\left(-\mathsf{mBIC}(M')/2\right)}$$

Int	ro	du	ct	io	n	
	0	00	0			

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look ○○ ○○●○

# Memetic algorithms

## MA particularly designed for our application

- Frommlet et al. (2012): Simulations for QTL mapping:
  - $\Rightarrow \quad \mathsf{GA} \text{ finds frequently better solution than stepwise search}$
- Implementation for GWAS finished, not yet integrated in MOSGWA

# Estimation of marker posteriors

 $\mathsf{MA} \text{ gives many good solutions} \quad \Rightarrow \quad \mathsf{model posteriors}$ 

$$\sum_{M \in \mathcal{M}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right) \ .$$

$$\mathsf{P}(j_r|Y) \approx \frac{\sum\limits_{M \in \mathcal{M}_r^{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right)}{\sum\limits_{M' \in \mathcal{M}^{Pool}} \exp\left(-\mathsf{mBIC}(M')/2\right)}$$

Int	ro	du	ct	io	n	
	0	00	0			

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look

# Memetic algorithms

## MA particularly designed for our application

- Frommlet et al. (2012): Simulations for QTL mapping:
  - $\Rightarrow \quad \mathsf{GA} \text{ finds frequently better solution than stepwise search}$
- Implementation for GWAS finished, not yet integrated in MOSGWA

# Estimation of marker posteriors

 $\mathsf{MA} \text{ gives many good solutions} \quad \Rightarrow \quad \mathsf{model posteriors}$ 

$$\sum_{M \in \mathcal{M}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \mathsf{P}(Y|M) \cdot \pi(M) \approx \sum_{M \in \mathit{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right) \ .$$

$$\mathsf{P}(j_r|Y) \approx \frac{\sum\limits_{M \in \mathcal{M}_r^{Pool}} \exp\left(-\mathsf{mBIC}(M)/2\right)}{\sum\limits_{M' \in \mathcal{M}^{Pool}} \exp\left(-\mathsf{mBIC}(M')/2\right)}$$

Intro	duction	
000	000	

GWAS Simulation QT

Case Control studies 0000000 0000 0000000 Out look

# Marker posteriors from MA

#### Example from real data analysis

Morphological differences in Drosophila Simulans (Zeng et al. (2000))

Search over markers

With imputation (IM)



Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00 0000 00

### Literature

Schwarz (1978) Estimating the Dimension of a Model. Ann. Statist. 6(2), 461 - 464.

Benjamini, Hochberg (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. J. Roy. Statist. Soc. Ser. B. 57, 289 - 300.

**Bogdan, Ghosh, Doerge (2004)** Modifying the Schwarz Bayesian Information Criterion to locate multiple interacting quantitive trait loci. *Genetics*, **167**, 989 - 999.

Abramovich, Benjamini, Donoho, Johnstone (2006) Adapting to unknown sparsity by controlling the false discovery rate *Ann. Statist*.34, 584 - 653.

**Chen, Chen (2008)** Extended Bayesian Information criteria for model selection with large model spaces. *Biometrika* **95**, 759 - 771.

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Introduction	Modifications of BIC	GWAS Simulation QT	Case Control studies	Out look
000000	0000 0000000 000	00000000	0000000 0000 0000000	00000

### Literatur

**Baierl, Bogdan, Frommlet, Futschik, (2006)**. On locating multiple interacting quantitative trait loci in intercross designs. *Genetics* **173**, 1693 - 1703.

**Bogdan, Frommlet, Biecek, Cheng, Ghosh, Doerge, (2008)** Extending the Modified Bayesian Information Criterion (mBIC) to Dense Markers and Multiple Interval Mapping. *Biometrics* **64**, 1162 - 1169.

Bogdan, Chakrabati, Frommlet, Ghosh (2011) Asymptotic Bayes Optimality under sparsity of some multiple testing procedures. Ann. Statist. 39, 1551 - 1579

**Frommlet**, **Bogdan (2013)** Some optimality properties of FDR controlling rules under sparsity, *Electron. J. Statist.* **7**, 1328-1368.

Frommlet, Ruhaltinger, Twarog, Bogdan (2011) Modified versions of Bayesian Information Criterion for genome-wide association studies. CSDA, 56, 1038 - 1051

**Frommlet, Ljubic, Arnardottir, Bogdan (2012)** QTL Mapping Using a Memetic Algorithm with modifications of BIC as fitness function, *Stat. Appl. in Genet. and Molec. Biol.*, **11**(4), Article 2

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()