

Power analysis of Genome Wide Association Studies based on simulations of phenotypes

V. Perduca, C. Sinoquet, R. Mourad, G. Nuel

IBC 2012, Kobe



Example

id	pheno	SNP1	SNP2
1	0	Aa	bb
2	0	aa	bB
3	1	AA	bB
4	0	aa	bb
5	1	Aa	BB
6	1	AA	BB
7	0	aa	bB

- ▶ pheno = status: 0 (control), 1 (case)

Example

id	pheno	SNP1	SNP2
1	0	Aa	bb
2	0	aa	bB
3	1	AA	bB
4	0	aa	bb
5	1	Aa	BB
6	1	AA	BB
7	0	aa	bB

SNP1	A	a
0	1	7
1	5	1

$$\Rightarrow p = 0.03$$

SNP2	B	b
0	2	6
1	5	1

$$\Rightarrow p = 0.1$$

- ▶ pheno = status: 0 (control), 1 (case)
- ▶ H_0 = no association, H_1 = association

Designing Genome-Wide Association Studies: Sample Size, Power, Imputation, and the Choice of Genotyping Chip

Chris C. A. Spencer¹, Zhan Su², Peter Donnelly¹, Jonathan Marchini^{1*}

Department of Statistics, University of Oxford, Oxford, United Kingdom

Power computed empirically

- ▶ phenotype: Y_i
- ▶ genotype: X_i

Power computed by simulating under

- ▶ H_1 : assumption of a disease model $\pi_i = \mathbb{P}(Y_i = 1|X_i)$
- ▶ H_0 : $\pi_i = \pi$ for all i

Simulations under H_0

Constraint: sample must have exactly n_1 cases and n_0 controls as in the original data

H_0

Phenotype shuffling

Example

id	pheno	Sim 1	SNP1	SNP2
1	0	1	Aa	bb
2	0	0	aa	bB
3	1	0	AA	bB
4	0	1	aa	bb
5	1	1	Aa	BB
6	1	0	AA	BB
7	0	0	aa	bB

Simulations under H_0

Constraint: sample must have exactly n_1 cases and n_0 controls as in the original data

H_0

Phenotype shuffling

Example

id	pheno	Sim 1	Sim 2	SNP1	SNP2
1	0	1	1	Aa	bb
2	0	0	1	aa	bB
3	1	0	0	AA	bB
4	0	1	0	aa	bb
5	1	1	0	Aa	BB
6	1	0	1	AA	BB
7	0	0	0	aa	bB

Simulations under H_0

Constraint: sample must have exactly n_1 cases and n_0 controls as in the original data

H_0

Phenotype shuffling

Example

id	pheno	Sim 1	Sim 2	Sim 3	SNP1	SNP2
1	0	1	1	0	Aa	bb
2	0	0	1	0	aa	bB
3	1	0	0	1	AA	bB
4	0	1	0	0	aa	bb
5	1	1	0	0	Aa	BB
6	1	0	1	1	AA	BB
7	0	0	0	1	aa	bB

Simulations under H_1

Constraint \mathcal{C} : sample must have exactly n_1 cases and n_0 controls

H_1 : $\pi_i = \mathbb{P}(\text{pheno } Y_i = 1 \mid \text{geno } X_i)$

One solution:

$$\mathbb{P}(X_i | Y_i) = \frac{\mathbb{P}(Y_i | X_i) \mathbb{P}(X_i)}{\mathbb{P}(Y_i)}$$

Problems:

- ▶ $\mathbb{P}(X)$: genotype model must take into account LD structure!
- ▶ need for extra data (e.g. reference panel of haplotypes from HapMap)
- ▶ $X \gg Y$

This strategy is implemented in **HAPGEN**

- ▶ Limited disease model: no epistasis, no gene-environment interactions...

Simulations under H_1 : alternative solution

Constraint \mathcal{C} : sample must have exactly n_1 cases and n_0 controls

H_1 : $\pi_i = \mathbb{P}(\text{pheno } Y_i = 1 \mid \text{geno } X_i)$

$Y_i \sim \mathcal{B}(\pi_i)$ **but** how to sample under the constraint?

Solutions:

1. **Rejection algorithm:** draw $Y \sim P(Y|X)$ until \mathcal{C} is true \Rightarrow waiting time in $O(1/P(\mathcal{C}))$

Simulations under H_1 : alternative solution

Constraint \mathcal{C} : sample must have exactly n_1 cases and n_0 controls

H_1 : $\pi_i = \mathbb{P}(\text{pheno } Y_i = 1 \mid \text{geno } X_i)$

$Y_i \sim \mathcal{B}(\pi_i)$ **but** how to sample under the constraint?

Solutions:

1. **Rejection algorithm:** draw $Y \sim P(Y|X)$ until \mathcal{C} is true \Rightarrow waiting time in $O(1/P(\mathcal{C}))$
2. **MCMC:** start from Y such as \mathcal{C} true, then perform moves that preserves $\mathcal{C} \Rightarrow$ many iterations to allow good mixing (slow)

Simulations under H_1 : alternative solution

Constraint \mathcal{C} : sample must have exactly n_1 cases and n_0 controls

H_1 : $\pi_i = \mathbb{P}(\text{pheno } Y_i = 1 \mid \text{geno } X_i)$

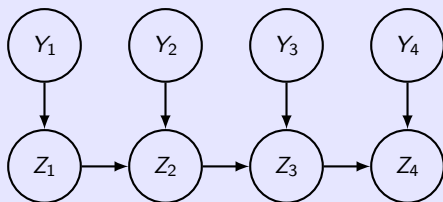
$Y_i \sim \mathcal{B}(\pi_i)$ but how to sample under the constraint?

Solutions:

1. **Rejection algorithm:** draw $Y \sim P(Y|X)$ until \mathcal{C} is true \Rightarrow waiting time in $O(1/P(\mathcal{C}))$
2. **MCMC:** start from Y such as \mathcal{C} true, then perform moves that preserves $\mathcal{C} \Rightarrow$ many iterations to allow good mixing (slow)
3. **Constrained backward sampling algorithm:** our contribution!

Our backward sampling: formalism

- ▶ $Z_i := \#$ cases among inds $1, \dots, i = Y_1 + \dots + Y_i = Z_{i-1} + Y_i$
- ▶ $\mathcal{C} = \{\sum_i^n Y_i = n_1\} = \{Z_n = n_1\}$, where $n = n_0 + n_1$



- ▶ $\mathbb{P}(Y_{1:n}, Z_{1:n}) = \mathbb{P}(Z_1|Y_1) \prod_{i=1}^n \mathbb{P}(Y_i) \prod_{j=2}^n \mathbb{P}(Z_j|Z_{j-1}, Y_j)$

⇒ A (very simple) BN!

⇒ Idea: **adapting BN message propagation algorithms** for sampling $\mathbb{P}(Y_1, \dots, Y_n | \mathcal{C})$.

Backward sampling

- ▶ Problem is solved by sampling the Heterogeneous Markov Chain:

$$\mathbb{P}(Y_1, \dots, Y_n | \mathcal{C}) = \mathbb{P}(Y_1 | \mathcal{C}) \cdot \mathbb{P}(Y_2 | Z_1, \mathcal{C}) \cdot \dots \cdot \mathbb{P}(Y_n | Z_{n-1}, \mathcal{C})$$

Definition (Backward quantities)

For $i = 1, \dots, n$:

$$B_i(m) = \mathbb{P}(Z_n = n_1 | Z_i = m) = \mathbb{P}(\mathcal{C} | Z_i = m).$$

Theorem

1.

$$B_{i-1}(m) = \pi_i B_i(m+1) + (1 - \pi_i) B_i(m)$$

2.

$$\mathbb{P}(Y_i = 1 | Z_{i-1} = m, \mathcal{C}) = \frac{\pi_i B_i(m+1)}{B_{i-1}(m)}$$

Comparing the three algorithms

Validation on a toy dataset

- ▶ The three algorithms are **consistent**: by simulating phenotypes under H_1 with each method we obtain the same value of power
- ▶ **Backward outperforms the others**:

n	f_0	$\mathbb{P}(\mathcal{C})$	Rej	MCMC	Backward
20	0.2	$4.5 \cdot 10^{-3}$	0.4 s	7.1 m	0.05 s
20	0.1	$1.7 \cdot 10^{-5}$	1.5 m	7.1 m	0.05 s
20	0.07	$6.7 \cdot 10^{-7}$	38.5 m	7.3 m	0.05 s
20	0.05	$2.9 \cdot 10^{-8}$	11.2 h	7.2 m	0.1 s
40	0.2	$8.2 \cdot 10^{-5}$	17.4 s	7.2 m	0.1 s
100	0.2	$8.7 \cdot 10^{-10}$	NA	8.0 m	0.2 s
100	0.1	$5.8 \cdot 10^{-22}$	NA	7.9 m	0.2 s
100	0.01	$1.1 \cdot 10^{-69}$	NA	8.0 m	0.2 s

- ▶ Backward and HAPGEN consistent

Application

Dataset

Genotypes from 629 individuals from the 1000 Genomes Project. 314 cases. First 100,000 SNPs from Chr X. MAF > 5%. Total: **8,048 SNPs**.

Disease: additive model (β) with epistasis (η)

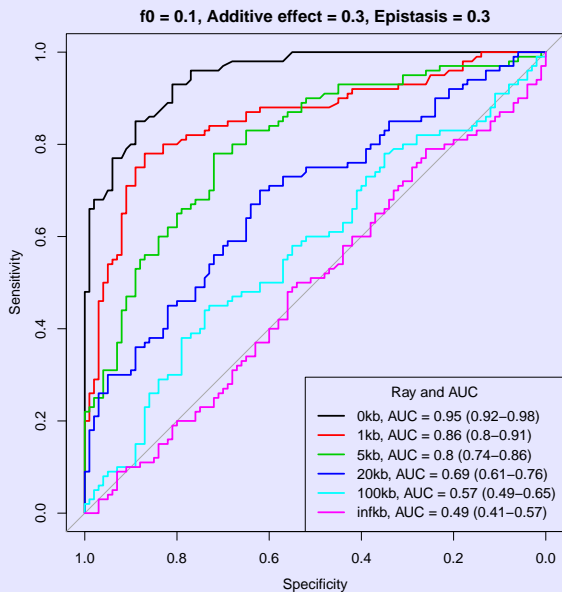
Two disease SNPs S_1 and S_2 (pos. 627,641 and 1,986,325) with no LD.

$$\pi_i = f_0 \times \text{RR} = f_0 \times \begin{cases} 1.0 + \beta \cdot X_i^{S_1} & \text{if } X_i^{S_2} = 0 \\ 1.0 + \beta \cdot X_i^{S_2} & \text{if } X_i^{S_1} = 0 \\ 1.0 + \eta + \beta \cdot (X_i^{S_1} + X_i^{S_2}) & \text{if } X_i^{S_1} \cdot X_i^{S_2} > 0 \end{cases}$$

The statistics

- ▶ For each SNP: trend p-values under H_0, H_1
- ▶ Intervals I_1, I_2 centered in S_1, S_2 with radius ρ , $\mathcal{R}_\rho = I_1 \cup I_2$
- ▶ $S := \max(-\log_{10}(\text{p-values SNPs in } \mathcal{R}_\rho))$

Results: varying the candidate region R_ρ



Results: varying the design

Role of the population size

n	AUC [95% CI]
629	0.49 [0.41, 0.57]
1258	0.78 [0.71, 0.84]
1887	0.92 [0.88, 0.96]
2516	0.93 [0.90, 0.97]

Table: $\rho = +\infty$, epistasis $\eta = 0.3$, additive effect $\beta = 0.3$, $f_0 = 0.1$.

Final word



Weighted affectation for constrained sampling under H_1

- ▶ We modeled the problem as a (very simple) BN and worked out a message propagation-like algorithm
- ▶ We generalized the shuffle method by affecting the pheno of each individual i w.r.t. π_i under the constraint that the number of cases must be n_1

Backward vs concurrents

- ▶ Gold standard is HAPGEN but backward has several advantages:
 - ▶ no additional assumptions more than epidemiological ones
 - ▶ complete freedom in the choice of π_i (interactions, environment, prevalence, penetrance, etc)
 - ▶ fast (2 sec on a laptop for 2000 cases and 2000 controls)
- ▶ Rejection algorithm: cannot be used in practice
- ▶ MCMC: delicate to calibrate

References

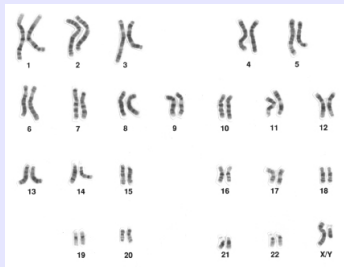
-  V. Perduca, C. Sinoquet, R. Mourad and G. Nuel; *Alternative methods for H1 simulations in Genome Wide Association Studies*. Hum Hered, 2012;73:95-104. Free Access
-  R package `wafect` 1.2 available on CRAN
`> vignette('wafect-tutorial')`

3 Assessing the power of GWAs

Given a GWA study method, it is crucial to assess its statistical power to detect susceptibility variants. Power can be estimated empirically by simulating disease (case and control) phenotypes. We illustrate how to assess the statistical power of GWA studies using `wafect` for phenotype simulations. In particular we will proceed as follows:

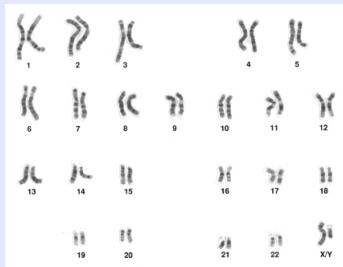
Appendix

Genetic background: Single Nucleotide Polymorphisms



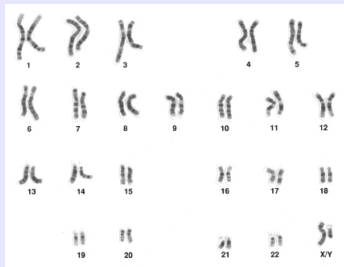
Ind	DNA
1	AGTTCCATCATGGTAAGC AGTTCCATTATGGTAAGC
2	AGTTCCATTATGGTAAGC AGTTCCATCATGGTAAGC
3	AGTTCCATTATGGTAAGC AGTTCCATTATGGTAAGC
4	AGTTCCATCATGGTAAGC AGTTCCATCATGGTAAGC

Genetic background: Single Nucleotide Polymorphisms



Ind	DNA
1	AGTTCCATCATGGTAAGC AGTTCCATTATGGTAAGC
2	AGTTCCATTATGGTAAGC AGTTCCATCATGGTAAGC
3	AGTTCCATTATGGTAAGC AGTTCCATTATGGTAAGC
4	AGTTCCATCATGGTAAGC AGTTCCATCATGGTAAGC

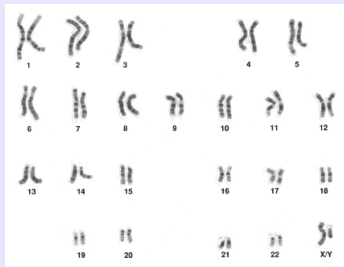
Genetic background: Single Nucleotide Polymorphisms



Ind	DNA	gen
1	AGTTCCATCATGGTAAGC AGTTCCATTATGGTAAGC	CT
2	AGTTCCATTATGGTAAGC AGTTCCATCATGGTAAGC	TC
3	AGTTCCATTATGGTAAGC AGTTCCATTATGGTAAGC	TT
4	AGTTCCATCATGGTAAGC AGTTCCATCATGGTAAGC	CC

- ▶ Depending on its two **alleles**, for any given SNP there are three possible **genotypes**

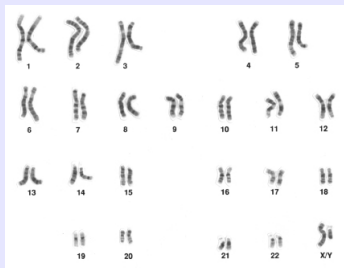
Genetic background: Single Nucleotide Polymorphisms



Ind	DNA	gen
1	AGTTCCATCATGGTAAGC AGTTCCATTATGGTAAGC	CT = 1
2	AGTTCCATTATGGTAAGC AGTTCCATCATGGTAAGC	TC = 1
3	AGTTCCATTATGGTAAGC AGTTCCATTATGGTAAGC	TT = 0
4	AGTTCCATCATGGTAAGC AGTTCCATCATGGTAAGC	CC = 2

- ▶ Depending on its two **alleles**, for any given SNP there are three possible **genotypes**
- ▶ The genotype of a SNP is coded with the number $i \in \{0, 1, 2\}$ of copies of the less frequent allele in the population, e.g. C

Genetic background: Single Nucleotide Polymorphisms



Ind	DNA	gen
1	AGTTCATCATGGTAAGC AGTTCATTATGGTAAGC	CT = 1
2	AGTTCATTATGGTAAGC AGTTCATCATGGTAAGC	TC = 1
3	AGTTCATTATGGTAAGC AGTTCATTATGGTAAGC	TT = 0
4	AGTTCATCATGGTAAGC AGTTCATCATGGTAAGC	CC = 2

- ▶ Depending on its two **alleles**, for any given SNP there are three possible **genotypes**
- ▶ The genotype of a SNP is coded with the number $i \in \{0, 1, 2\}$ of copies of the less frequent allele in the population, e.g. C
- ▶ SNPs are used as **markers** to identify the genomic regions associated with a phenotype (e.g. a disease)

GWAS

SNPs are used as **markers** to identify the genomic regions associated with a phenotype

Which **SNPs** across the genome are associated with a given disease?

1. Recruitment of n individuals: n_1 **cases** and n_0 **controls** ($n_1, n_0 \sim 10^3$)
2. High throughput genotyping of each individual with respect to all the SNPs ($\sim 10^5$)
3. For each SNP, test the association with the disease (e.g. χ^2 test):
 $H_0 = \text{no association}, H_1 = \text{association}$
4. **Choice of a statistics S** to analyze the signal
5. Correction for multiple testing

H_1 and H_0

For each individual i :

- ▶ $Y_i \in \{0, 1\}$: phenotype
- ▶ $X_i \in \{0, 1, 2\}^p$, $p = \#$ SNPs: genotype

H_1 : assumption of a disease model $\pi_i = \mathbb{P}(Y_i = 1|X_i)$

Example $p = 1$:

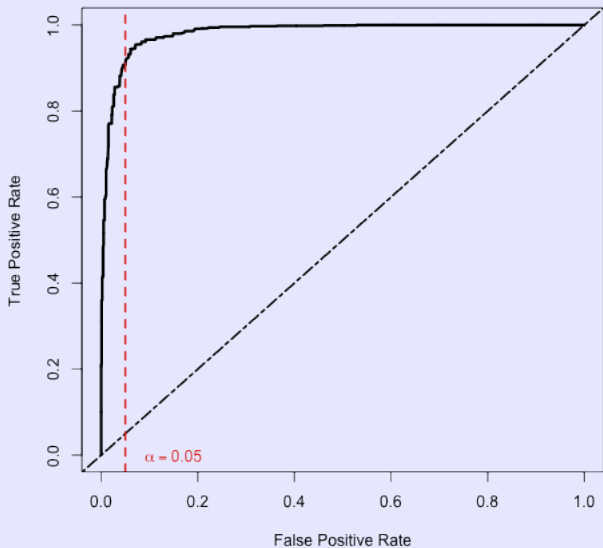
- ▶ $\pi_i = f_0$ if $X_i = 0$
- ▶ $\pi_i = f_1 = f_0 \cdot RR_1$ if $X_i = 1$
- ▶ $\pi_i = f_2 = f_0 \cdot RR_2$ if $X_i = 2$

RR_1, RR_2 : relative risks; f_1, f_2 : penetrances

H_0 : $\pi_i = \pi$ for all i

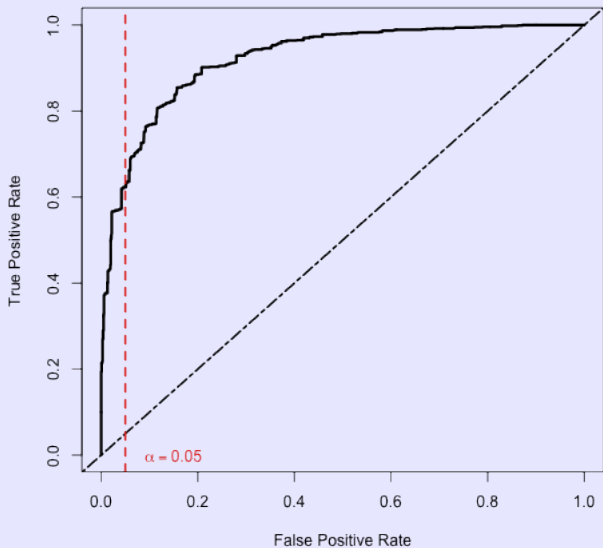
The observed genotype has no effect on the phenotype

Accuracy of a GWAS: ROC curves



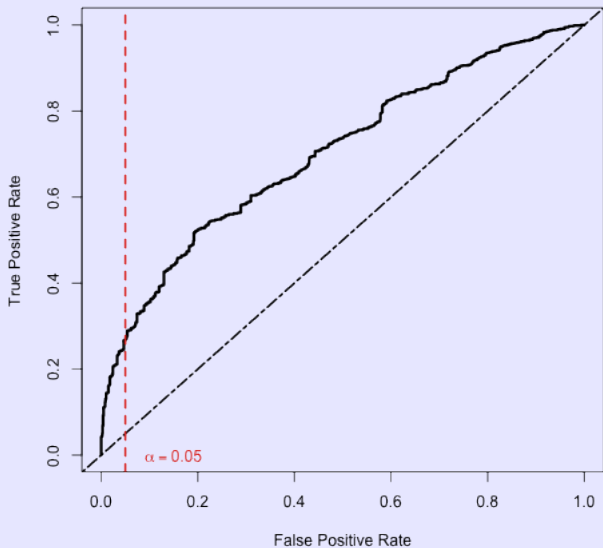
⇒ Good power

Accuracy of a GWAS: ROC curves



⇒ Lower power

Accuracy of a GWAS: ROC curves



⇒ Poor power

Estimating the AUC

Definition

AUC = area under the ROC curve

Qualitative interpretation of the AUC

AUC	0.5 – 0.6	0.6 – 0.7	0.7 – 0.8	0.8 – 0.9	0.9 – 1.0
quality	fail	poor	fair	good	excellent

In order to find empirically the ROC curve and its AUC we need to **sample the statistic distributions under H_0 and H_1** :

Proposition

If S_1, \dots, S_r is a sample under H_0 and T_1, \dots, T_r a sample under H_1 then

$$\widehat{\text{AUC}} = \frac{1}{r^2} \sum_{i,j} \mathbf{1}_{\{T_j \geq S_i\}} \quad \hat{\sigma}_{\max} = \sqrt{\frac{\widehat{\text{AUC}} \cdot (1 - \widehat{\text{AUC}})}{r}}$$

Sampling H_1 : Reject algorithm

Constraint: $\mathcal{C} = \{\sum_{i=1}^n Y_i = n_1\}$ must be fulfilled

Reject algorithm

1. draw $(Y_i)_{i=1\dots n}$
2. if \mathcal{C} holds then retain (Y_i) , else discard it and go back to 1

Problem: in practice, \mathcal{C} is a very rare event!

Theorem

Let $Z_j = \sum_{i=1}^j Y_i$. Then $\mathbb{P}(Z_i = m) = F_i(m)$, where

$$F_i(m) = F_{i-1}(m-1)\pi_i + F_{i-1}(m)(1-\pi_i),$$

with $F_0(m) = 0$ except for $F_0(0) = 1$.

In particular: $\mathbb{P}(\mathcal{C}) = \mathbb{P}(Z_n = n_1) = F_n(n_1)$.

Sampling H_1 : MCMC algorithm

MCMC

Start from a configuration $(Y_i)_{i=1\dots n}$ fulfilling the constraint and alternate two steps:

1. exchange Y_i and Y_j for two i, j s.t. $Y_i = 1$ and $Y_j = 0$
2. accept the move in 1 with rate

$$\alpha = \frac{(1 - \pi_i)\pi_j}{\pi_i(1 - \pi_j)}$$

The sequence of configurations that are generated is a Markov chain whose stationary distribution is the targeted distribution

Problem: delicate to choose the number of iterations needed for convergence (*burn-in*) and for ensuring independence of the samples

Forward and Backward quantities

Definition

$$F_i(m) = \mathbb{P}(Z_i = m)$$

$$B_i(m) = \mathbb{P}(Z_n = n_1 | Z_i = m) = \mathbb{P}(\mathcal{C} | Z_i = m)$$

Theorem

$$\mathbb{P}(\mathcal{C}) = F_n(n_1) = B_0(0)$$

$$\mathbb{P}(Y_i = 1 | \mathcal{C}) \propto \sum_m F_i(m) \pi_i B_i(m+1)$$

$$\mathbb{P}(Y_1 = 0 | \mathcal{C}) \propto \sum_m F_{i-1}(m) (1 - \pi_i) B_i(m)$$

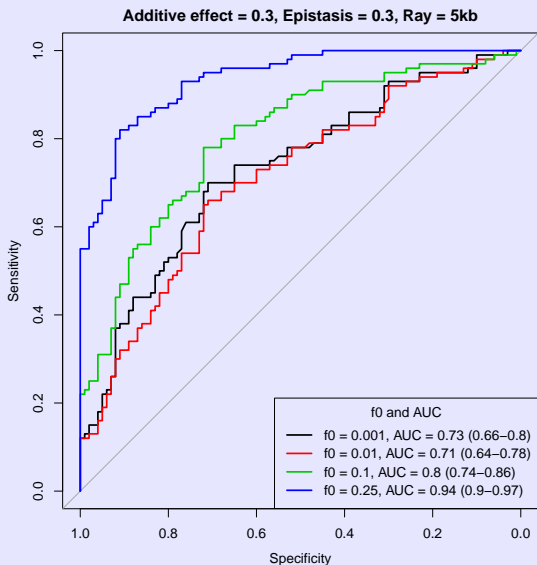
Comparing the three algorithms

AUC

n	f_0	Rej	MCMC	Backward
20	0.2	0.60 [0.53, 0.67]	0.58 [0.51, 0.65]	0.61 [0.54, 0.68]
20	0.1	0.59 [0.52, 0.66]	0.58 [0.51, 0.65]	0.58 [0.51, 0.65]
20	0.07	0.62 [0.55, 0.69]	0.54 [0.47, 0.61]	0.56 [0.49, 0.63]
20	0.05	0.44 [0.37, 0.51]	0.55 [0.48, 0.62]	0.53 [0.47, 0.60]
40	0.2	0.58 [0.50, 0.65]	0.54 [0.46, 0.61]	0.59 [0.52, 0.67]

Results: varying the design

Role of f_0



Results: varying the design

Role of the additive effect

