

# Detecting past contraction in population size using haplotype homozygosity

C. Merle<sup>1,2&3</sup>, J-M. Marin<sup>1&3</sup>, F. Rousset<sup>3&4</sup> and R. Leblois<sup>2&3</sup>

1 - Institut de Montpelliérain Alexander Grothendieck (IMAG, UM); 2 - Centre de Biologie pour la Gestion des Populations (CBGP, INRA); 3 - Institut de biologie computationnelle (IBC);  
4 - Institut des Sciences de l'Evolution (ISEM, CNRS)

## Next generation genetic data

Classical inference methods of the demographic history (likelihood based methods (IS, MCMC), approximate bayesian methods and Site Frequency Spectrum), suitable for polymorphism data sets consisting in some loci and assuming the genealogies of different loci are independent, do not exploit the genetic recombination.

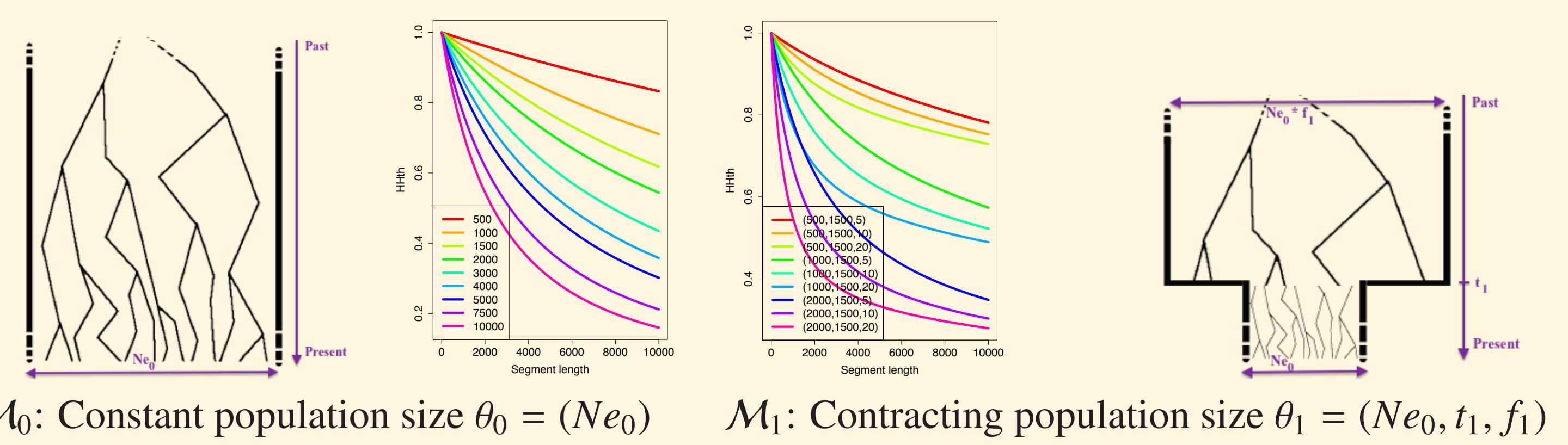
Take advantage of genome wide data with known genome

Pairwise alignment  
GT CATGAGCCGA GT  
GA CATGAGCCGA TT

Consider the dependency of genealogies of adjacent positions in the genome.

New inference approaches of demographic history based on the conserved sequence lengths in the pairwise alignment within a diploid genome : [3], [5], [4], [2], [1].

The patterns of linkage disequilibrium (LD) between polymorphic markers are shaped by the ancestral population history as we can see on haplotype homozygosity curves which measure the LD.

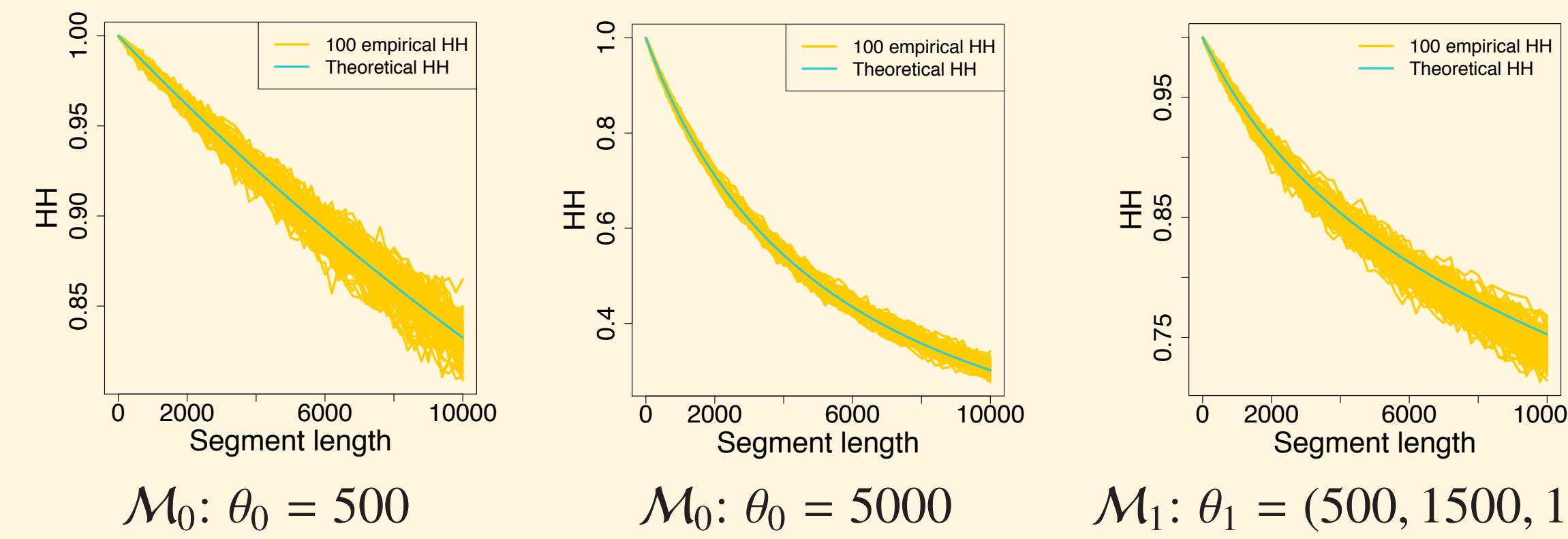


## Demographic history inference from genome wide sequence data

**Haplotype Homozygosity:**  $HH(i)$  : probability for  $i$  adjacent markers drawn at random in the whole genome sequence to be homozygote.

**Theoretical  $HH$**  of [3], denoted  $HH_{th}$  : coalescent based computation, assuming the mutation and recombination rate known and constant along the genome.

**Empirical  $HH$** , denoted  $\widehat{HH}$  computed from the observed data, is the segment proportion of at least  $i$  homozygotes adjacent markers.



We estimate  $\theta_0$  and  $\theta_1$  by :

$$\widehat{\theta}_0 \in \arg \min_{\theta_0} \sum_{i \in I} \left( \frac{HH_{th}(\theta_0, i) - \widehat{HH}(i)}{\widehat{HH}(i)} \right)^2$$

$$\widehat{\theta}_1 \in \arg \min_{\theta_1} \sum_{i \in I} \left( \frac{HH_{th}(\theta_1, i) - \widehat{HH}(i)}{\widehat{HH}(i)} \right)^2$$

Evaluation of  $i \mapsto HH_{th}(\theta_1, i)$  is time consuming  
→ optimization with the **blackbox** R package.

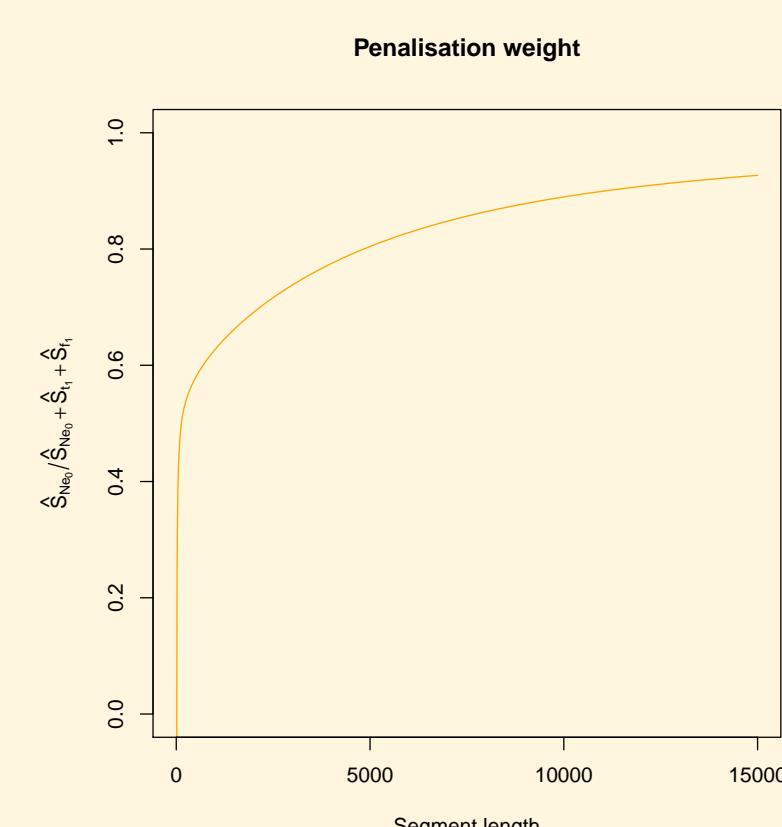
## Model choice criterion: embedded models

### Penalized mean square criterion

$$\arg \min_{j \in \{0,1\}} \sum_{i \in I} w_j(i) \left( \frac{HH_{th}(\theta_j, i) - \widehat{HH}(i)}{\widehat{HH}(i)} \right)^2$$

$$w_0(i) = \widehat{S}_{N_{e_0}}^{(1)}(i) / (\widehat{S}_{N_{e_0}}^{(1)}(i) + \widehat{S}_{t_1}^{(1)}(i) + \widehat{S}_{f_1}^{(1)}(i)),$$

$$w_1(i) = 1.$$

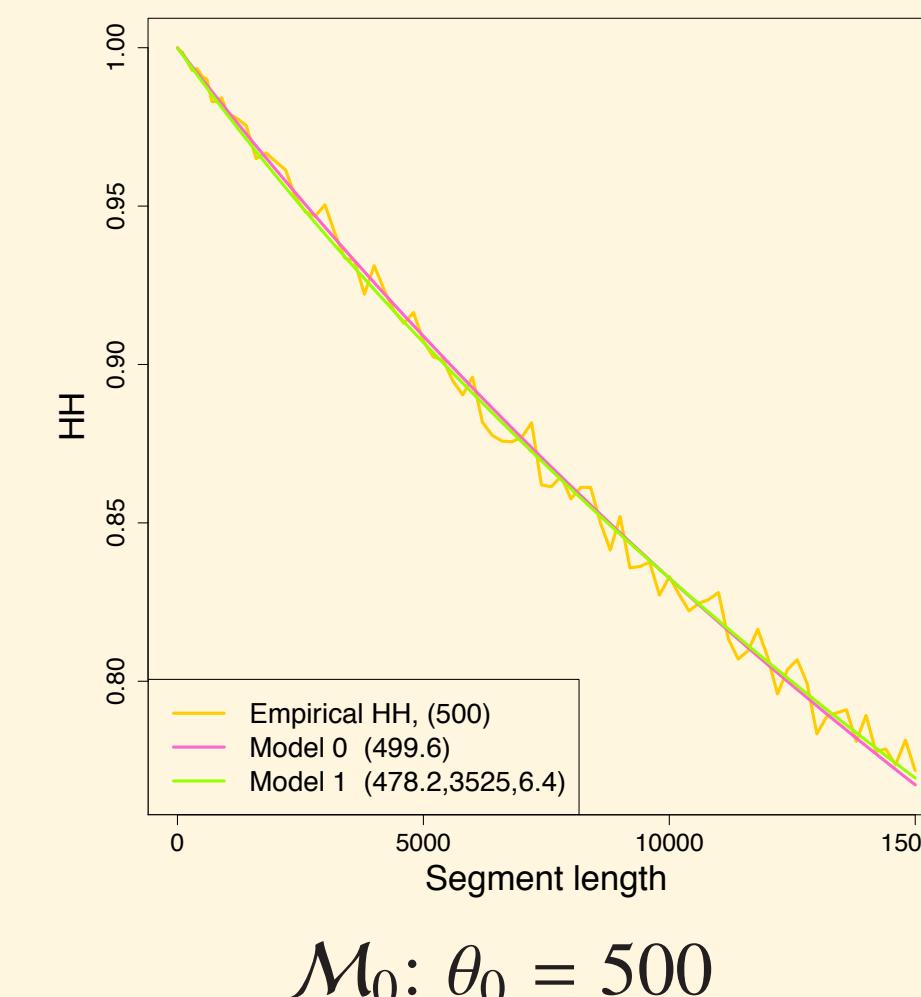


$\widehat{S}_{\phi}^{(1)}(i)$  : estimate of Sobol's sensitivity index of order one of a given parameter  $\phi$  computed for  $i$  adjacent markers under the more complex model (*sensitivity* R package).

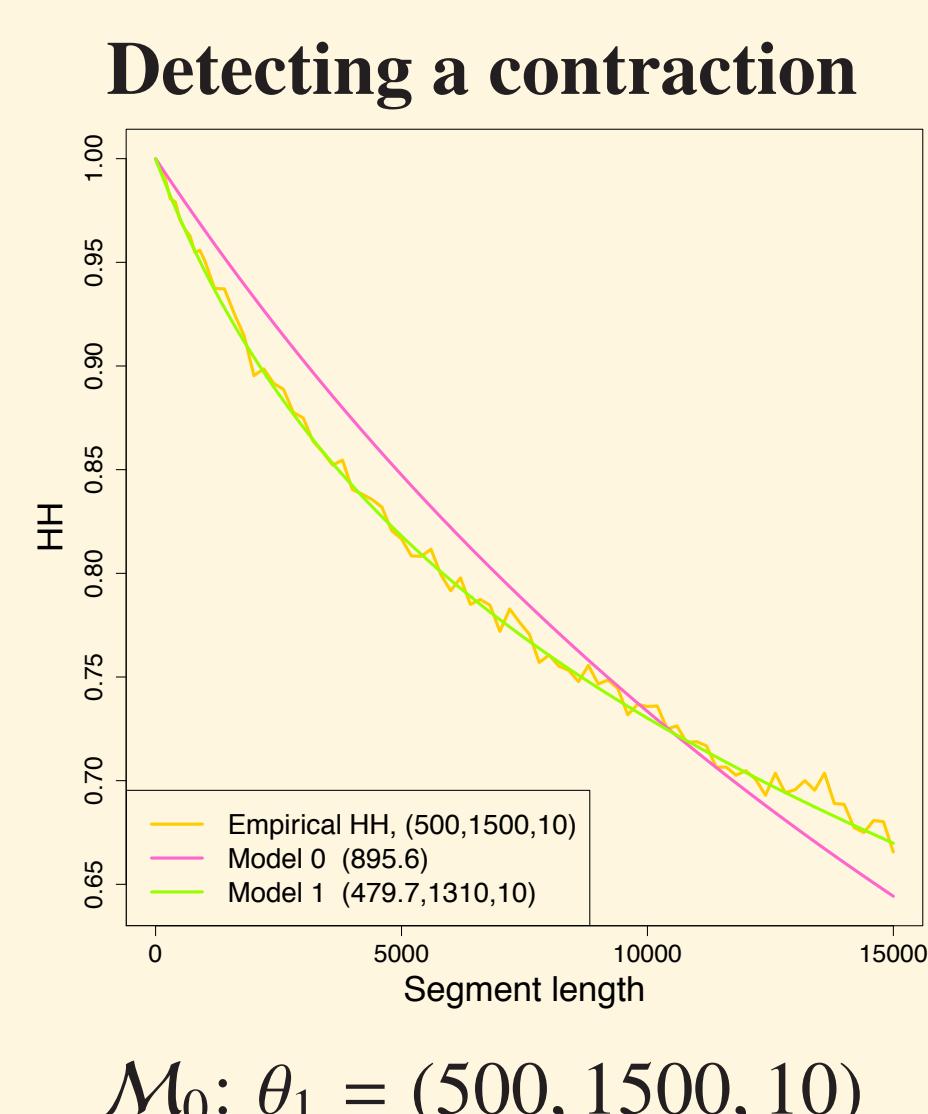
On a simulated data set with constant population size  $\theta_0 = 500$ , the non penalized mean square criterion would choose  $M_1$  whereas the penalized mean square criterion choose  $M_0$ .

On a simulated data set with contracting population size  $\theta_1 = (500, 1500, 10)$ , the penalized mean square criterion choose  $M_1$ .

### Avoid choosing the more complex model



$M_0: \theta_0 = 500$



$M_0: \theta_1 = (500, 1500, 10)$

## Numerical results

Data sets simulated under  $M_0, \theta_0 = N_{e_0} = 500$

$n_{pos}$	5000		10000		
	Locus	2000	200000	2000	200000
Data sets	100	1	100	1	
TPR	Adjustment	0.74	1	0.59	0
	Penalized Adj.	<b>0.94</b>	1	<b>0.98</b>	1
$\overline{N_{e_0}}$		494	494	502	501
		[478 – 510]	–	[496 – 509]	–

Data sets simulated under  $M_0, \theta_0 = N_{e_0} = 5000$

$n_{pos}$	5000		10000		
	Locus	2000	200000	2000	200000
Data sets	100	1	100	1	
TPR	Adjustment	0.41	0	0.54	0
	Penalized Adj.	<b>0.86</b>	1	<b>0.79</b>	0
$\overline{N_{e_0}}$		4998	5044	4915	4880
		[4960 – 5036]	–	[4881 – 4949]	–

Data sets simulated under  $M_1, \theta_1 = (N_{e_0}, t_1, f_1) = (500, 1500, 10)$

$n_{pos}$	10000		15000		20000		
	Locus	2000	200000	2000	200000	2000	200000
Data sets	100	1	100	1	100	1	
TPR	Penalized Adj.	<b>0.81</b>	1	<b>0.92</b>	1	<b>0.95</b>	1
$\overline{N_{e_0}}$		634	633	561	529	534	492
		[617 – 650]	–	[547 – 576]	–	[520 – 548]	–
$\overline{t_1}$		2573	2445	2031	1592	1881	1446
		[2402 – 2744]	–	[1918 – 2144]	–	[1727 – 2036]	–
$\overline{f_1}$		12.4	11.5	12.3	9.1	12.2	9.49
		[11.8 – 13.0]	–	[11.7 – 12.8]	–	[11.6 – 12.7]	–

Holstein data set

$n_{pos}$	10000		
	Locus	2000	200000
Data sets	100	1	
CDR	Penalized Adj.	<b>0.93</b>	1
$\overline{N_{e_0}}$		6690	6924
		[6533 – 6847]	–
$\overline{t_1}$		5965	6532
		[5559 – 6372]	–
$\overline{f_1}$		4.23	4.0
		[3.90 – 4.57]	–

## References

- Sharon R Browning and Brian L Browning. Accurate non-parametric estimation of recent effective population size from segments of identity by descent. *The American Journal of Human Genetics*, 97(3):404–418, 2015.
- Kelley Harris and Rasmus Nielsen. Inferring demographic history from a spectrum of shared haplotype lengths. *PLoS genetics*, 9(6):e1003521, 2013.
- I. MacLeod, T. Meuwissen, B. Hayes, and M. Goddard. A novel predictor of multilocus haplotype homozygosity: comparison with existing predictors. *Genetics research*, 91(6):413–426, 2009.
- I. MacLeod, D. Larkin, H. Lewin, B. Hayes, and M. Goddard. Inferring Demography from Runs of Homozygosity in Whole-Genome Sequence, with Correction for Sequence Errors. *Molecular Biology and Evolution*, 30(9):2209–2223, 2013.
- Pier Francesco Palamara, Todd Lencz, Ariel Darvasi, and Itsik Pe'er. Length distributions of identity reveal fine-scale demographic history. *The American Journal of Human Genetics*, 91(5):809–822, 2012.

## Acknowledgements

