

# Selection on expected maximum haploid breeding values can increase genetic gain in recurrent genomic selection

Dominik Müller\*· Pascal Schopp\* and Albrecht E. Melchinger\*

\*Institute of Plant Breeding, Seed Sciences and Population Genetics, University of Hohenheim, 70599 Stuttgart, Germany

January 28<sup>th</sup>, 2018

Correspondence: melchinger@uni-hohenheim.de

## Supplemental File 2

### Approximation of EMBV using a normal distribution

Using independent assortment of chromosomes and conclusions following from the central limit theorem, breeding values of a quantitative trait asymptotically follow a normal distribution. Let  $\sigma_i^2$  be the segregation variance in the progeny of a candidate  $i$ , defined as the variance of the breeding values of a population of DH lines developed from it. Let  $GEBV_i$  be the GEBV of the candidate. Then, for the breeding values of these DH lines, we have  $Y_i \sim \mathcal{N}(GEBV_i, \sigma_i^2)$ . If we denote with  $Y_{i(N_G)}$  the largest order statistic (maximum) of a sample of  $N_G$  DH lines, the definition of the EMBV implies that  $EMBV = E(Y_{(N_G)})$ . By the properties of the expectation and of order statistics, it is easy to see that

$$EMBV_i = GEBV_i + E(X_{(N_G)}) \sigma_i, \quad (1)$$

where  $E(X_{(N_G)})$  is the expectation of the largest order statistic of  $N_G$  random variables from  $\mathcal{N}(0, 1)$  (Figure S2-1). This demonstrates that the EMBV can be interpreted as the sum of the GEBV and the square root of the segregation variance, weighted by  $E(X_{(N_G)})$ .

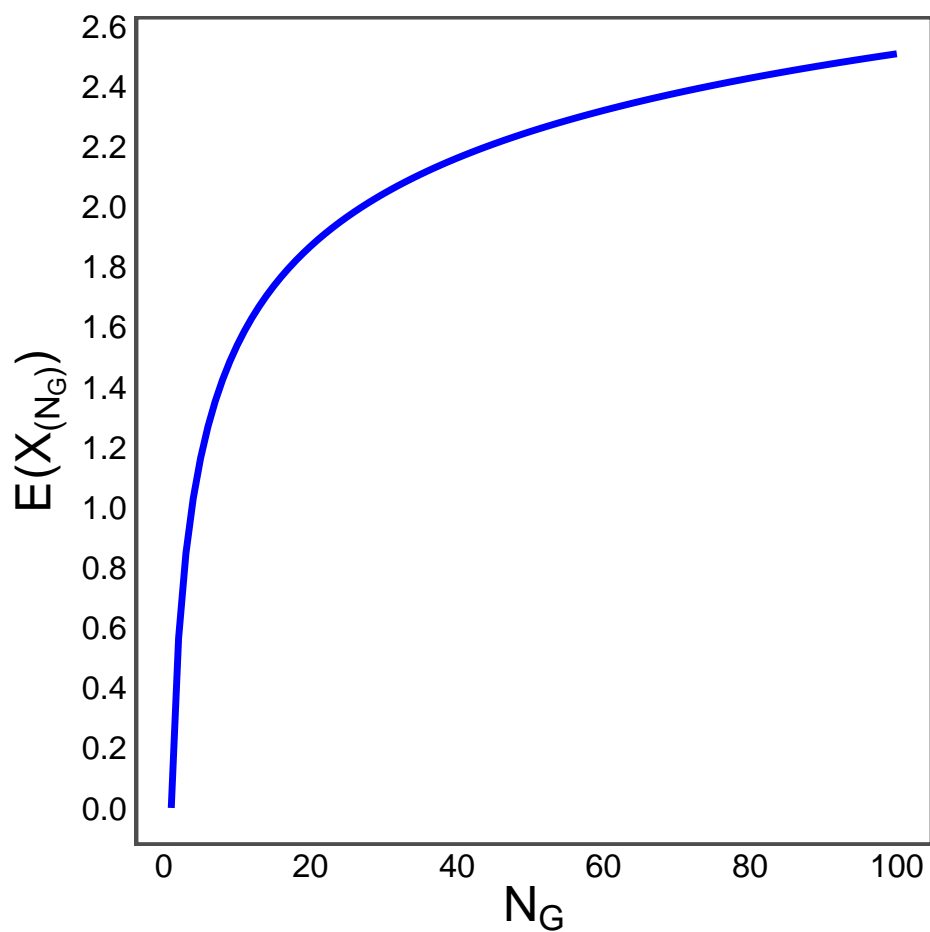
14 Estimates of the segregation variance can be obtained by repeatedly simulating DH progenies  
 15 and calculating the sample variance of their breeding values. Alternatively, it could be calculated  
 16 analytically using the well-known equation from Lynch and Walsh (1998)

$$\sigma_i^2 = 2 \sum_l^M \alpha_l^2 p_l(1 - p_l) + 2 \sum_l^M \sum_{l' \neq l}^M \alpha_l^2 \alpha_{l'}^2 D_{ll'}, \quad (2)$$

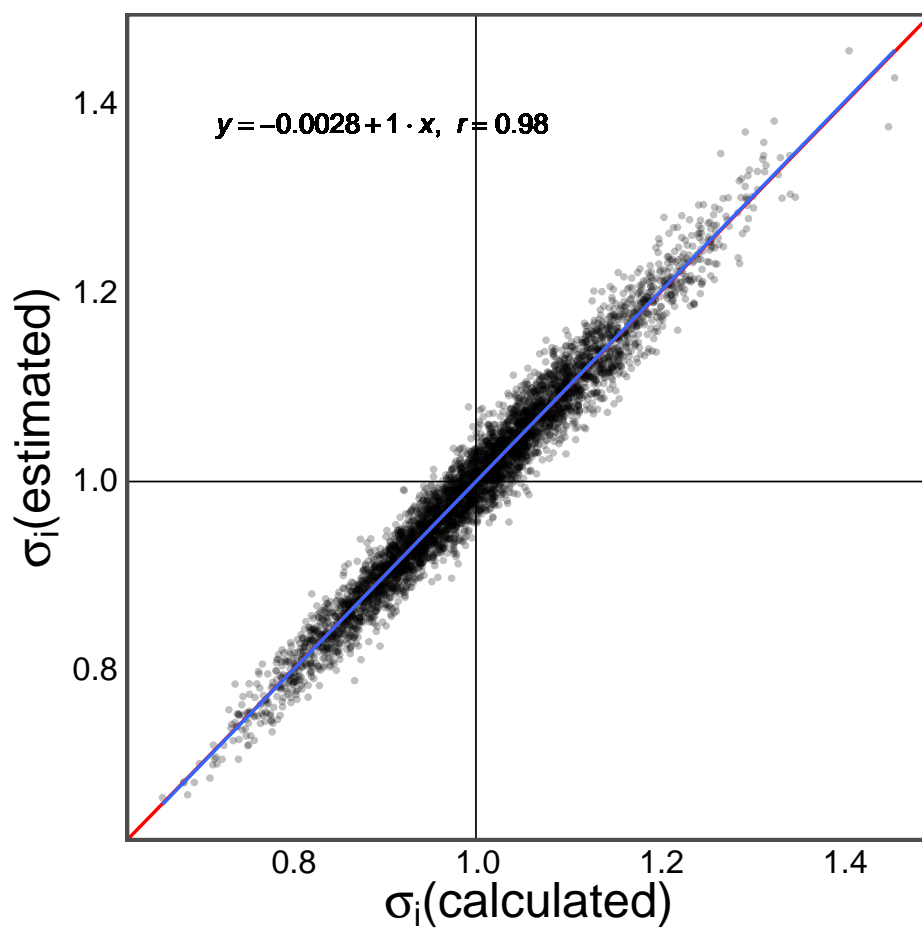
17 where  $D_{ll'}$  is the linkage disequilibrium (LD) between loci  $l$  and  $l'$  in the target population. As  
 18 the goal is to calculate  $\sigma_i^2$  for DH progeny, we computed  $D_{ll'}$  as  $D_{ll'} = (1 - 2r_{ll'})D'_{ll'}$ . Here,  $D'_{ll'}$  is  
 19 the LD in the parental gametes of the candidate itself, which is 0 if the candidate is homozygous at  
 20 either locus  $l$  or  $l'$ , 0.25 if linkage phases are identical and  $-0.25$  if linkage phases are reversed. The  
 21 recombination fraction  $r_{ll'}$  was computed from the genetic distance  $d_{ll'}$  between both loci using  
 22 Haldane's mapping function, which yields  $D_{ll'} = e^{-2d_{ll'}} D'_{ll'}$ . A similar derivation has been given by  
 23 Lehermeier et al. (2017).

24 In Figure S2-2, we compared estimated against calculated values of  $\sigma_i$  for a total of 1,000  
 25 individuals from the founder population ( $N_{chr} = 20$ ). The correlation was high ( $r \approx 0.98$ ) and the  
 26 regression of estimated on calculated values was unbiased ( $b = 1$ ). This was also true for other  
 27 chromosome numbers ( $N_{chr} = 5, 40$ ; data not shown). We conclude that an analytical solution  
 28 following equation 2 could lead, in conjunction with equation 1, to a rapid alternative to the  
 29 computationally heavy simulation-based computation of EMBV.

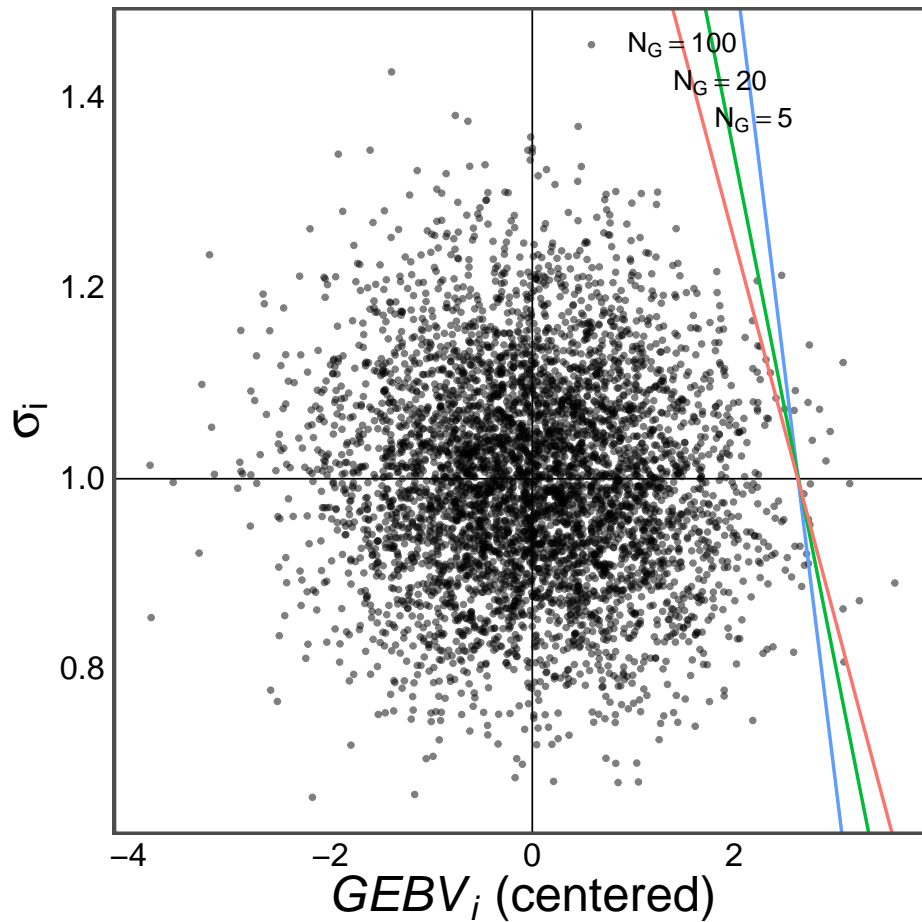
30 To illustrate how the selection of candidates depends on the value of  $N_G$ , we plotted the GEBVs  
 31 of the 1,000 individuals their value of  $\sigma_i$  (Figure S2-3). Lines were added that separate the top 20  
 32 candidates with highest EMBV from other individuals. These results show how with increasing  $N_G$ ,  
 33 candidates having a high GEBV are traded with candidates having a lower GEBV, but larger  $\sigma_i$ .

34 **Figures**

**Figure S2-1:** Expectation of the maximum order statistic  $E(X_{(N_G)})$  of a standard normal random variable  $X$ , depending on the sample size (number of gametes,  $N_G$ ).



**Figure S2-2:** Scatter plot of estimated versus calculated square root of segregation variance ( $\sigma_i$ ) for 1,000 individuals from the founder population with  $N_{chr} = 20$ .



**Figure S2-3:** Scatter plot of genomic estimated breeding values (GEBVs) and the square root of the estimated segregation variance ( $\sigma_i$ ) for 1,000 individuals from the founder population with  $N_{chr} = 20$ . Colored lines separate selected from unselected candidates for different values of  $N_G$ , given that in each case the top 20 candidates with the highest EMBV were selected.

35 **References**

- 36 Lehermeier, C., S. Teyssède, and C.-C. Schön (2017). “Genetic Gain Increases by Applying the  
37 Usefulness Criterion with Improved Variance Prediction in Selection of Crosses”. In: *Genetics*.  
38 Lynch, M. and B. Walsh (1998). *Genetics and Analysis of Quantitative Traits*. 1st edition. Sunderland:  
39 Sinauer Associates, page 980.