

A Mendelian Sampling model for genetic prediction

**Clement Carre, Fabrice Gamboa, David Cros, Gregor
Gorjanc, Eduardo Manfredi**

Additive models for prediction

Individual:

$$y = \mu + u + e, \quad \text{var}(u) = A\sigma_u \text{ and } \text{var}(e) = I\sigma_e$$

Marker:

$$y = \mu + Wm + e, \quad \text{var}(m) = I\sigma_m \text{ and } \text{var}(e) = I\sigma_e$$

$$y = \mu + u^* + e, \quad \text{var}(u^*) = WW'f(\sigma_m) \text{ and } \text{var}(e) = I\sigma_e$$

Individual + marker:

$$y = \mu + u + Wm + e$$

Individual + marker models: when?

When the individual variance σ_u is larger than the variance due to markers $f(\sigma_m)$

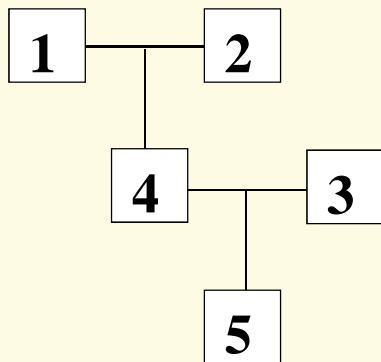
- Theory: elusive relationship between σ_u and σ_m (Gianola et al., 2009)
- Reports on missing heritability (Yang et al., 2011)
- Improvement of prediction accuracy in some applications (de los Campos et al., 2009; Duchemin et al., 2012)
- Genome complexities: LD is subjected to noise and to data conditions (marker coverage, reference and target populations), no-nucleotide polymorphism

So, revisit the model: $y = \mu + u + Wm + e$

A Mendelian Sampling Model (1/4)

Individual effects of descendants are a function of effects of base individuals and Mendelian Sampling (MS) (Quaas, 1976):

$$\mathbf{u}_d = \mathbf{M}^{-1} \begin{bmatrix} \mathbf{u}_b \\ \mathbf{s} \end{bmatrix}, \text{ with } \mathbf{M} = \mathbf{I} - 1/2 \mathbf{P}$$



$$\mathbf{u}_1 = \mathbf{u}_1$$

$$\mathbf{u}_5 = \left(\frac{\mathbf{u}_1}{4} + \frac{\mathbf{u}_2}{4} + \frac{\mathbf{s}_4}{2} \right) + \frac{\mathbf{u}_3}{2} + \mathbf{s}_5$$

$$\mathbf{M}^{-1} =$$

	1	2	3	4	5
1	1				
2		1			
3			1		
4	1/2	1/2		1	
5	1/4	1/4	1/2	1/2	1

A Mendelian Sampling Model (2/4)

Replace random Mendelian Sampling s by realized Mendelian Sampling r :

$$\mathbf{r} = [\mathbf{M}_{db} \quad \mathbf{M}_{dd}] \begin{bmatrix} \mathbf{W}_b \\ \mathbf{W}_d \end{bmatrix} \mathbf{m}$$

	1	2	3	4	5
1	1				
2		1			
3			1		
4	-1/2	-1/2		1	
5			-1/2	-1/2	1

	Snp1	Snp2	Snp3
1	0	1	1
2	2	1	0
3	1	1	1
4	1	2	0
5	2	1	1

$$\mathbf{r}_4 = \left[1 - \frac{0+2}{2}\right] \mathbf{m}_1 + \left[2 - \frac{1+1}{2}\right] \mathbf{m}_2 + \left[0 - \frac{1+0}{2}\right] \mathbf{m}_3$$

Mendelian sampling model (3/4)

Using parts (1) and (2) write individual effects of descendants as a function of effects of base individuals and MS:

$$\mathbf{u}_d = \mathbf{D}\mathbf{u}_b + (\mathbf{W}_d - \mathbf{D}\mathbf{W}_b)\mathbf{m}$$

$$\text{with } \mathbf{D} = -\mathbf{M}_{dd}^{-1} \mathbf{M}_{db}$$

Mendelian Sampling model (4/4)

Then, several models for phenotypes of descendants can be proposed:

Disjoint model:

$$y_d = \mu + Du_b + (W_{1d} - DW_{1b})m_1 + W_{2d}m_2 + e$$

For m_1 and m_2 , subsets of m (e.g. markers of QTL and rest of markers)

Embedded model:

$$\begin{aligned} y_d &= \mu + Du_b + (W_d - DW_b)m + W_d m + e \\ &= \mu + Du_b + (2W_d - DW_b)m + e \end{aligned}$$

Testing the embedded model

1. Simulated data (QMSim ; Sargolzaei and Frenkel, 2009)

Demography : G0 (n=25), G1 to G3 (training; n=120), G4 (target ; n=40)

Random mating; family size=1

Genome: 2 chromosomes (2M) with 10 biallelic QTL each

Trait: phenotypic variance=1 and genetic variance (QTL + infinitesimal effects)=0.4

Conditions studied: 200 vs 2000 SNP/chromosome ; variance explained by qtl: 10, 50 and 90% (200 replicates/condition)

Predictive ability: correlation between known simulated values (phenotypes and genetic values) and their corresponding predictions

Simulated scenario	Training data (%)^a	Test data(%)^a
QTL variance 10%		
200 SNP markers	103	112
2000 SNP markers	102	108
QTL variance 50%		
200 SNP markers	100	100
2000 SNP markers	100	98
QTL variance 90%		
200 SNP markers	99	95
2000 SNP markers	97	97

^a(%) is 100 times the ratio between the average accuracy under the Mendelian segregation model and the average accuracy under the marker model

Discussion

First test of the embedded model suggests advantages (NS) of the MS model when markers capture a small proportion of genetic variance. On the other hand, the embedded model did not improve the accuracy of the reference model when markers do capture a high proportion of genetic variance (NS).

Some simplifying assumptions that we used at this stage may limit the potential contribution of the embedded model: the assumption of uncorrelated effects of base individuals, and also the use of known variances.

The MS model gave an extra accuracy of +16% over the marker model in the analysis of a QTLMAS simulated data set (Lund et al., 2009), with a complex genetic model and large family size (n=10)

Conclusions

The decomposition of individual effects into base individual effects and mendelian sampling can be revisited to integrate realized MS instead of random MS.

This opens new possibilities for modelling: by splitting markers into groups or by embedding marker and individual effects. Also, the model allows the easy integration of genetic groups.

Ongoing work is on variance components estimation of the embedded and disjoint models, including the incomplete marker data situation, and the analyses of real data.