#### Single-step

#### Ole F. Christensen

#### Aarhus University, Center for Quantitative Genetics and Genomics

GenSAP meeting 2014 in Korsør, day3

#### Single-step: theory

A joint model for genomic value and marker genotypes (coded 0, 1, 2):

$$g = \sum_{j} (M_j - 2\rho_j)\beta_j$$

$$\mathsf{E}[M_j] = 2\rho_j 1, \ \, \mathsf{Var}[M_j] = v_j A$$

(based on idea by Gengler et al. 2007 to infer missing genotypes).

- Assume *M* is multivariate normal
- Individuals with missing and observed genotypes

$$M = \left[ \begin{array}{c} M^{miss} \\ m^{obs} \end{array} \right], \quad A = \left[ \begin{array}{c} A_{11} & A_{12} \\ A_{21} & A_{22} \end{array} \right],$$

Single-step - extension of G to non-genotyped individuals

▶ Marginalisation (integration) of *M<sup>miss</sup>* gives

$$\begin{split} \mathsf{E}[g \mid m^{obs}] &= 0, \\ \mathsf{Var}[g \mid m^{obs}] &= \sigma_g^2 \begin{bmatrix} A_{12}A_{22}^{-1}G & A_{12}A_{22}^{-1}(G - A_{22})A_{22}^{-1}A_{21} + A_{11} \\ G & GA_{22}^{-1}A_{21} \end{bmatrix} \\ &= \sigma_g^2 H. \\ \mathsf{where} \ G &= \sum_j (m_j^{obs} - 2\rho_j 1) (m_j^{obs} - 2\rho_j 1)^T / \sum_j v_j \text{ and} \\ \sigma_g^2 &= \sum_j v_j \sigma_\beta^2. \end{split}$$

# Sparse inverse

$$H^{-1} = \left[ \begin{array}{cc} 0 & 0 \\ G^{-1} - A_{22}^{-1} & 0 \end{array} \right] + A^{-1}.$$

#### Single-step: a practical problem

$$G = \sum_{j} (m_j^{obs} - 2
ho_j 1) (m_j^{obs} - 2
ho_j 1)^T / \sum_{j} v_j$$

where  $\rho_j$  is allele frequency in base population,  $s = \sum_j v_j$  is a scaling.

- These are unknown!
- Allele frequencies: in principle both phenotypes and marker genotypes provide information about these. Inferring them is computationally challenging, and in addition there will be uncertainty in such estimation.
- ▶ Practical solution: use observed allele freq p̂<sub>j</sub> = 1<sup>T</sup> m<sup>obs</sup><sub>j</sub>/n and adjust G towards A<sub>22</sub>,

$$G_{adjust} = G\beta + 11^T \alpha$$

#### Single-step: SNP-model

$$g = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} A_{12}A_{22}^{-1}\sum_j(m_j^{obs} - 2\rho_j)\beta_j \\ \sum_j(m_j^{obs} - 2\rho_j)\beta_j \end{bmatrix} + \begin{bmatrix} \epsilon \\ 0 \end{bmatrix}$$
  
where  $Var(\epsilon) = A_{11} - A_{12}A_{22}^{-1}A_{21}$ .  
Substituting  $\rho = \hat{\rho} + (\rho - \hat{\rho})$ , then

$$g = \begin{bmatrix} A_{12}A_{22}^{-1} \\ 1 \end{bmatrix} \mu_g + \begin{bmatrix} A_{12}A_{22}^{-1}\sum_j (m_j^{obs} - 2\hat{\rho}_j 1)\beta_j \\ \sum_j (m_j^{obs} - 2\hat{\rho}_j 1)\beta_j \end{bmatrix} + \begin{bmatrix} \epsilon \\ 0 \end{bmatrix}$$
where  $\mu_j = 2\sum_j (\hat{\rho}_j - \rho_j)\beta_j$ 

where  $\mu_g = 2 \sum_j (\hat{\rho}_j - \rho_j) \beta_j$ .

•  $\mu_g$  related to genetic drift

## Single-step: genetic drift

• 
$$\mu_g = 2 \sum_j (\hat{\rho}_j - \rho_j) \beta_j$$
 (genetic drift)

- Fixed effect μ<sub>g</sub>, and inferred from phenotypes (Fernando et al. 2014, Vitezica et al. 2011)
- Proper prior on  $\mu_g$ ,

$$\mathsf{Var}(\mu_g 1 + (m^{obs} - 2\hat{\rho} 1^{\mathsf{T}})\beta) = \sigma_{\mu_g}^2 11^{\mathsf{T}} + \sigma_g^2 G$$

and infer  $\sigma_{\mu_g}^2$  based on matching to  $A_{22}$ , i.e. infer from marker genotypes (Vitezica et al. 2011).

Important point: adjustments like G<sub>adjust</sub> = Gβ + 11<sup>T</sup>α are related to genetic drift.

#### Adjusting G to A - some theory

(Powell et al. 2010, Vitezica et al. 2011, Meuwissen et al. 2011) •  $G_{adjust} = G(1 - \alpha/2) + 11^T \alpha$ 

- G reflects relationships relative to genotyped individuals.
- Idea: translate relationships such they are relative to base population
- Note: assumption of random assignment of alleles in each population, i.e. theory does not really support a population across generations.

Adjusting A to G - some theory

(Christensen, 2012)

- ► Alternatively: adjust A to G.
- A<sup>γ</sup> = A(1 − γ/2) + γ11<sup>T</sup> is relationship matrix with base individuals being related and inbreed. Computing formulas are as usual.
- Estimate γ from G<sub>1/2</sub> = (m<sup>obs</sup> − 1)(m<sup>obs</sup> − 1)<sup>T</sup>/s where s is unknown scaling parameter.
- A Bayesian derivation (prior on  $\rho_j$  and  $v_j$ ), very sketchy:

 $\mathsf{E}[M_j] = 2\mathsf{E}[\rho_j]\mathbf{1} = \mathbf{1}, \ \mathsf{Var}[M_j] = \mathsf{E}[v_j]A + 4\,\mathsf{Var}[\rho_j]\mathbf{1}\mathbf{1}^T \propto A(\gamma)$ 

#### Conclusion

- Genomic and pedigree relationships are incompatible. So adjustments are needed.
- Main problem is the allele frequencies.
- Three solutions:
  - Estimate genetic drift in model for phenotypes.
  - Adjust G:  $G_{adjust} = G\beta + 11^T \alpha$
  - Adjust A:  $A^{\gamma} = A(1 \gamma/2) + \gamma 11^{T}$

## Topics that I did not mention:

- Additional polygenic effect.
- Multiple breeds, genetic groups
- ► LDLA approach (Meuwissen et al. 2011).
- Computing !
- Computing !
- Computing !