# Parameter estimation in the ancestral regression with missing parental or grandparental genotypes

## Rodolfo J. C. Cantet,  Natalia S. Forneris

**Departamento de Producción Animal,
Facultad de Agronomía
Universidad de Buenos Aires;
INPA-CONICET**

# Introduction

- Animal model of Henderson-Quaas-Bulmer is a Data Generating Process with a Markovian and p.d. covariance matrix.

- Cantet et al (2017) proposed an individual based model ("ancestral regression", AR) that is autoregressive (causal), Markovian and easy to fit, but no parameter estimation.

- Model parameters (the individual´s $\beta_S$ and $\beta_D$) are *identifiable only* with information from a dense set of SNPs to estimate the sufficient statistics.

- Bayesian estimators require the distribution of $\beta_S$ & $\beta_D$.

- Goal: To present a Gibbs sampler for $\beta_S$ and $\beta_D$.

# Ancestral regression

$$a_{\mathrm{X}} = 0.5\ a_{\mathrm{S}} + 0.5\,a_{\mathrm{D}} + \beta_{\mathrm{S}}\left(a_{\mathrm{SS}} - a_{\mathrm{DS}}\right) + \beta_{\mathrm{D}}\left(a_{\mathrm{SD}} - a_{\mathrm{DD}}\right) + \phi_{\mathrm{X}}$$

$$\beta_{\mathrm{S}} = \frac{\Sigma_{\mathrm{X,PGP}} - \Sigma_{\mathrm{PGP,R}}\,\Sigma_{\mathrm{R}}^{-1}\,\Sigma_{\mathrm{R,PGP}}}{\Sigma_{\mathrm{PGP}} - \Sigma_{\mathrm{PGP,R}}\,\Sigma_{\mathrm{R}}^{-1}\,\Sigma_{\mathrm{R,PGP}}} \qquad \beta_{\mathrm{D}} = \frac{\Sigma_{\mathrm{X,MGP}} - \Sigma_{\mathrm{MGP,R}}\,\Sigma_{\mathrm{R}}^{-1}\,\Sigma_{\mathrm{R,MGP}}}{\Sigma_{\mathrm{MGP}} - \Sigma_{\mathrm{MGP,R}}\,\Sigma_{\mathrm{R}}^{-1}\,\Sigma_{\mathrm{R,MGP}}}$$

where the sufficient statistics are

$$\Sigma_{\mathrm{X,PGP}} = 0.5\left(\Sigma_{\mathrm{X,SS}} - \Sigma_{\mathrm{X,DS}}\right) \qquad \Sigma_{\mathrm{X,MGP}} = 0.5\left(\Sigma_{\mathrm{X,SD}} - \Sigma_{\mathrm{X,DD}}\right)$$

with corresponding variances

$$\Sigma_{\mathrm{PGP}} = 0.25\left(\Sigma_{\mathrm{SS}} + \Sigma_{\mathrm{DS}} - \Sigma_{\mathrm{SS,DS}}\right) \qquad \Sigma_{\mathrm{MGP}} = 0.25\left(\Sigma_{\mathrm{SD}} + \Sigma_{\mathrm{DD}} - \Sigma_{\mathrm{SD,DD}}\right)$$

# Path coefficient view of the Ancestral Regression and ssBLUP

# Covariance matrix of BV under AR

$$B_{X(i)} = \begin{bmatrix} 0...\beta_S & -\beta_S & ...0...\beta_D & -\beta_D & 0.5 & 0.5 \end{bmatrix} => B$$

## Autoregressive causal model

$$a = B\,a + \phi \quad \Rightarrow \quad \Sigma = \left(I - B\right)^{-1} D \left(I - B'\right)^{-1}$$

- The distribution of **a** is MVN (proved elsewhere), such that $\Sigma^{-1}$ from an autoregressive structure, follows an inverted Wishart, and the betas (in **B**) are standard normal (Roverato, 2000).

# Covariances between BV under AR

A = Ancestor,  S = Sire ,  SS = Sire of Sire , DS = Dam of Sire ,
D = Dam,    SD = Sire of Dam,  DD = Dam of Dam.

## Covariance between an ancestor and X

$$\Sigma_{A,X} = 0.5\left(\Sigma_{A,S} + \Sigma_{A,D}\right) + \beta_S\left(\Sigma_{A,SS} - \Sigma_{A,DS}\right) + \beta_D\left(\Sigma_{A,SD} - \Sigma_{A,DD}\right)$$

### Covariances between two animals,
### neither of whom is an ancestor of the other

$$\begin{bmatrix} \Sigma_A & \Sigma_A B_Y{}' \\ B_X \Sigma_A & B_X \Sigma_A B_Y{}' \end{bmatrix}$$

# Inbreeding under AR

S = Sire ,    SS = Sire of Sire , DS = Dam of Sire ,
D = Dam,    SD = Sire of Dam,  DD = Dam of Dam.

$$F_{X(AR)} = 0.5\left[\Sigma_{S,D} + \beta_S\left(\Sigma_{SS,D} - \Sigma_{DS,D}\right) + \beta_D\left(\Sigma_{SD,S} - \Sigma_{DD,S}\right)\right.$$

$$\left. + \beta_S\,\beta_D\left(\Sigma_{SS,SD} - \Sigma_{SS,DD} - \Sigma_{DS,SD} + \Sigma_{DS,DD}\right)\right]$$

# Distribution of sufficient statistics and parameters

❑ Jimenez' Thesis: Simulation to obtain the empirical distribution. Pig and beef cattle genotypes to validate the distributions obtained.

❑ Cov(Grandparent, Individual) => Beta.

❑ **Sufficient statistics**: $\Sigma_{X, PGP}$

= Cov(Grandsire, X) − Cov(Grand-dam,X) => Normal

❑ $\beta_S$, $\beta_D$ => Normal.     $F_{AR}$ => Exponential.

# Estimating equations

$$L\,B_X{}' = \beta = \begin{bmatrix} \beta_S \\ \beta_D \end{bmatrix}$$

$$L' = \begin{bmatrix} 0.5 & 0 \\ -0.5 & 0 \\ 0 & 0.5 \\ 0 & -0.5 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

If all grandparents and parents are genotyped, betas are the solution of the system with 2 equations:

$$L\,\Sigma_A\,L'\,\beta = L\,\Sigma_{AX}$$

$$\mathrm{Var}\left(\phi_X\right) = \left[1 + F_{X(AR)} - L\,B_X\,\Sigma_A\,B_X{}'\,L'\right]\sigma_A^2$$

# Gibbs sampling algorithm for $\beta_S$ & $\beta_D$

1) Estimation of IBD relationships (data) with an algorithm that accounts for a) Pedigree, b) Inbreeding, c) LD: Han & Abney (2011, Genetic Epidemiology 35: 557-567).

2) Build up the reduced system:

$$L \, \Sigma_A \, L' \, \beta \; = \; L \, \Sigma_{AX}$$

3) Sampling $\beta_S$ and $\beta_D$ from

$$\beta_t \sim T \, N\left(\beta_{t-1}, (0.25)^2\right), \quad \beta_t \in [-0.25, 0.25]$$

Algorithm of Foulley (2000, GSE 32:631–635).

# IBD genomic relationships under AR for Sultán when dams are not genotyped



**Sultana Sumaj 566 Comodoro**

$$\beta_S = 0.075 \pm 0.036$$

$$\beta_D = 0.086 \pm 0.036$$

$$F_{AR} = 0.037$$

| $rel_{AR}$ | SS | DS | SD | DD |
|---|---|---|---|---|
| Sultán | <u>0.347</u> | 0.188 | 0.23 | <u>0.27</u> |
| $F_{ped} = 0.031$ | Sire | 0.57 | Dam | 0.51 |

# Missing data: work in progress

- Patterns of missing data are variable. Most problematic: **missing grand-dams**, because cov(Grand-dam, X) is part of the sufficient statistics (SuSt). Fortunately, the SuSt are Normal.

- Estimate the missing covariance using available genotypes from relatives **no more than 2 meioses apart** (*U* = Uncles, grand-uncles). *U* are many in pigs, but not enough in cattle.

- Collaterals (FS): $\Sigma_{FS} = 0.5 + 2\left(\beta_{SX}\,\beta_{SY} + \beta_{DX}\,\beta_{DY}\right)$

- HMM to estimate missing covariances through **IBD segment sharing** with different SNP panels, or even a small number of microsatellites used for paternity.