

Course roadmap

AlphaSimR

Day 1: Simulation of breeding programmes **BASICS**

Day 2: ... Quantitative genetics

Day 3: ... Estimation with linear mixed models

Day 4: ... Spatial variation & GxE interactions

Day 5: ... Ancestral recombination graphs

Day 5 agenda – Ancestral Recombination Graphs

- 09:00-10:30 Spatial modelling, ARG
- 10:30-11:00 Refreshments break
- 11:00-12:30 ARG
- 12:30-13:30 Lunch break
- 13:30-15:00 ARG / ARG Practicals
- 15:00-15:30 Refreshments break
- 15:30-17:00 ARG Practicals / Open-end



THE UNIVERSITY
of EDINBURGH



Spatial modelling improves genetic evaluation in smallholder breeding programs

Gregor Gorjanc, Chris Gaynor, Jon Bancic, Daniel Tolhurst

UNE, Armidale

2024-02-09



Learning objectives

Separating genetic and environmental effects is a critical component of any quantitative genetics model

- Showcase challenge and solution for modelling data from smallholder settings
- Aside: APY approximation

Modelling environmental effects

- Naïve model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{a} + \mathbf{e}^*$$
$$\mathbf{a} \sim N(\mathbf{0}, \mathbf{A}\sigma_a^2)$$
$$\mathbf{e}^* \sim N(\mathbf{0}, \mathbf{E}\sigma_{e^*}^2)$$

- Model contemporary group effect, say, herd, herd-season, ...

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{a} + \mathbf{W}\mathbf{h} + \mathbf{e}$$
$$\mathbf{a} \sim N(\mathbf{0}, \mathbf{A}\sigma_a^2)$$
$$\mathbf{h} \sim N(\mathbf{0}, \mathbf{H}\sigma_h^2)$$
$$\mathbf{e} \sim N(\mathbf{0}, \mathbf{E}\sigma_e^2)$$

Extensive literature on “fixed” vs. “random” treatment due to data structure or views/opinions:

- unbalanced designs & bias
- ability to estimate effects
- ...

Challenge with small contemporary groups



JDS
Communications®
2021; 2:366–370

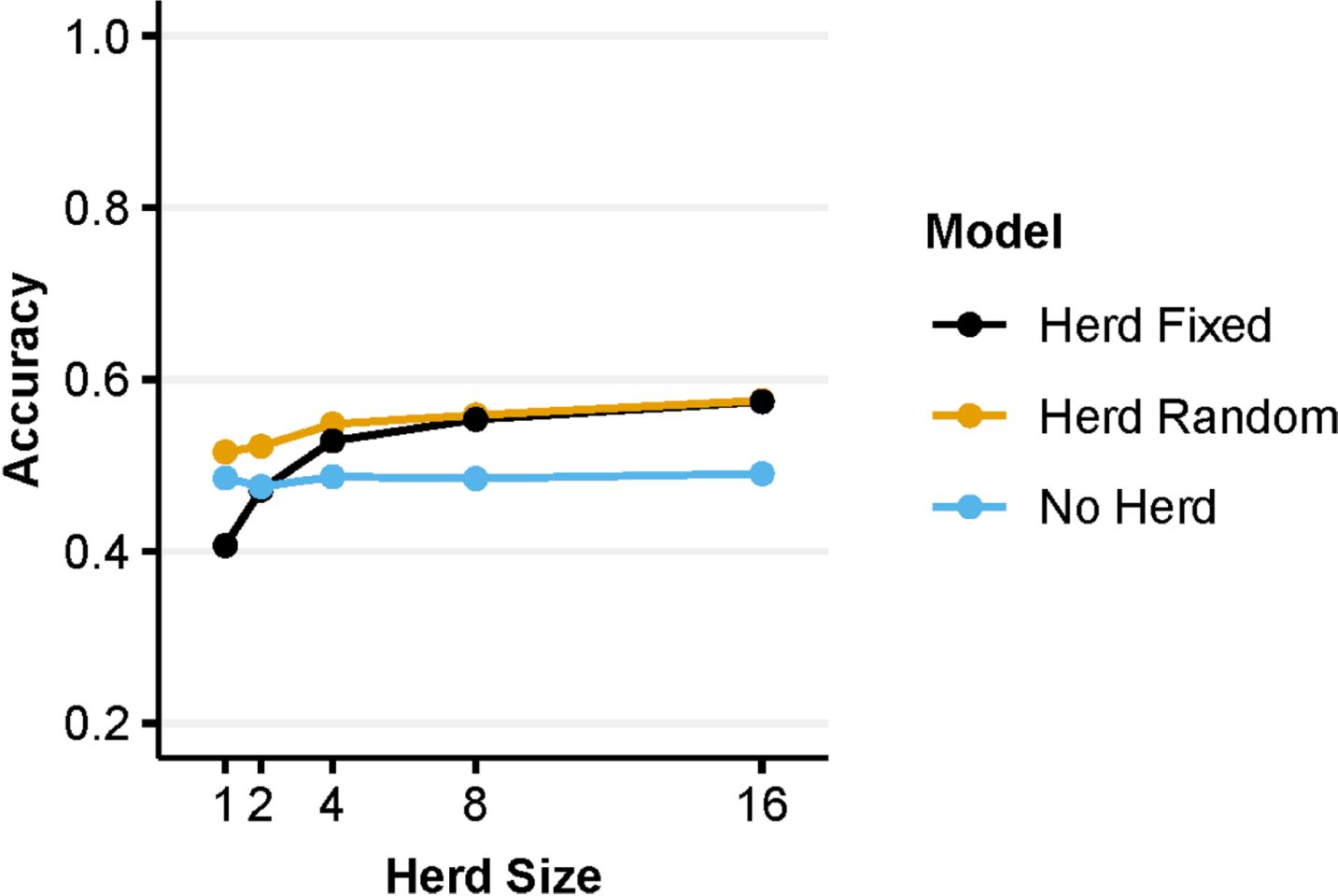
<https://doi.org/10.3168/jdsc.2021-0092>

Short Communication
Genetics

Genomic evaluations using data recorded on smallholder dairy farms in low- to middle-income countries

Owen Powell,^{1*}  Raphael Mrode,^{2,3}  R. Chris Gaynor,¹ Martin Johnsson,^{1,4}  Gregor Gorjanc,¹ 
and John M. Hickey¹

Challenge with small contemporary groups



Environmental/Spatial modelling



- A solution?
 - borrow information from neighbours (spatial model) and/or
 - measure key environmental indicators (location covariates)

Selle et al. *Genet Sel Evol* (2020) 52:69
<https://doi.org/10.1186/s12711-020-00588-w>

GSE Genetics
Selection
Evolution

RESEARCH ARTICLE

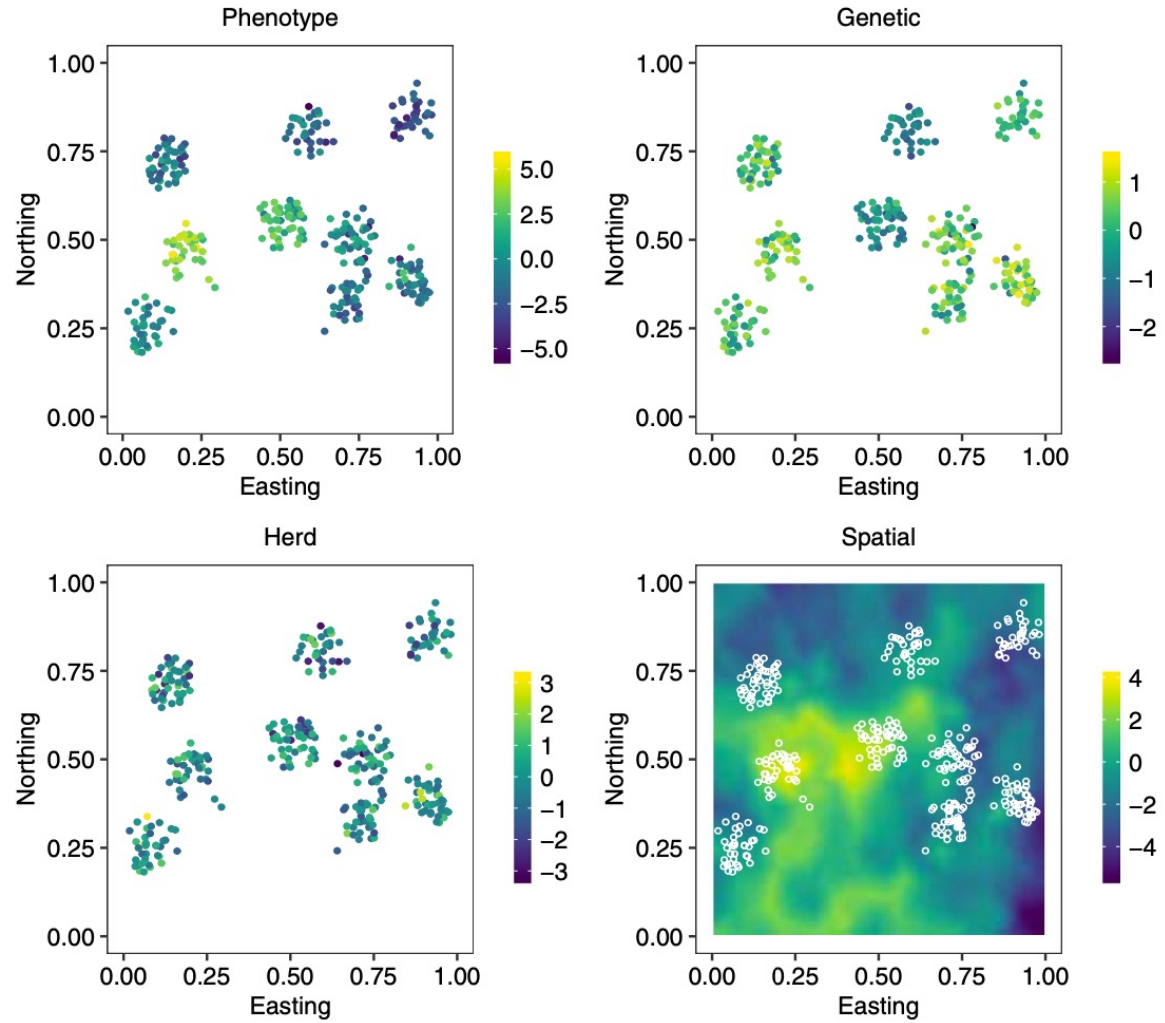
Open Access

Spatial modelling improves genetic evaluation in smallholder breeding programs



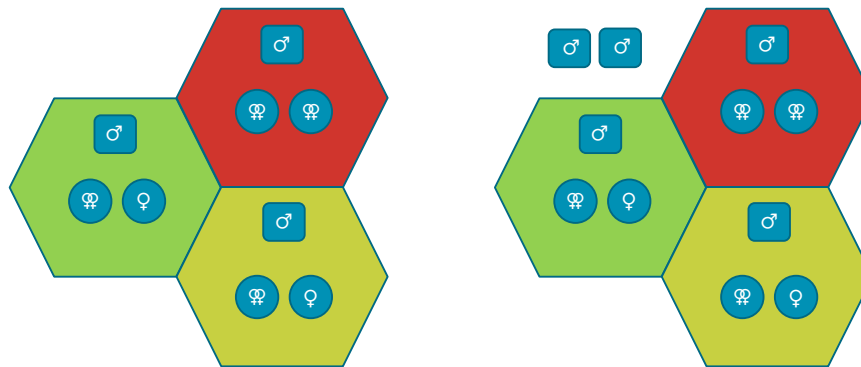
Maria L. Selle^{1*} , Ingelin Steinsland¹, Owen Powell², John M. Hickey² and Gregor Gorjanc²

Spatial context



Simulation

- Smallholder breeding programme
- Connectedness scenarios



- Phenotype = **Location** + Herd + Genetics + Noise
0.40 0.25 0.10 0.25

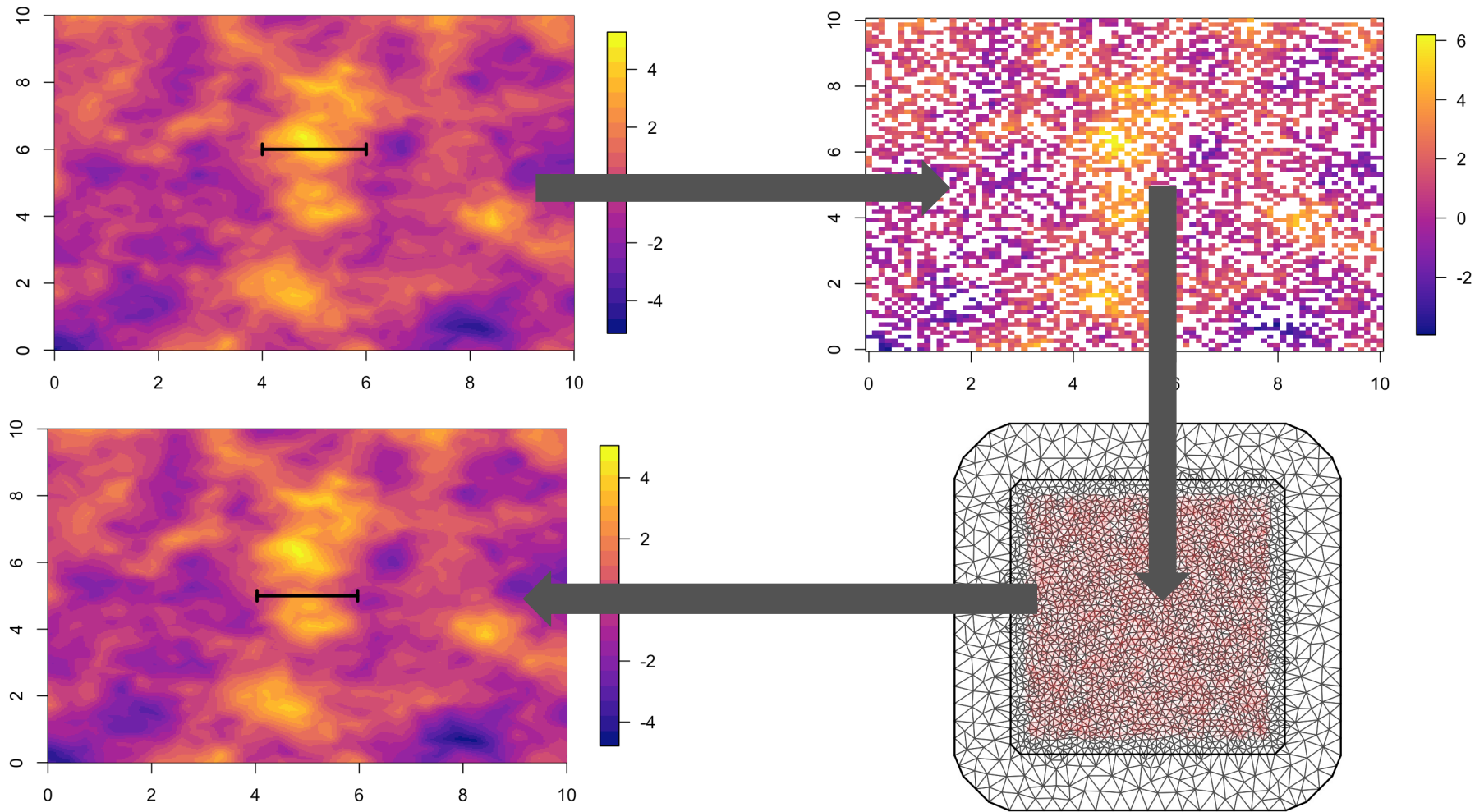
- **Location** = sum of 10 spatially varying covariates (rainfall, temp, ...)
- Observe 5 covariates with noise and 2 as binary indicators

Spatial modelling

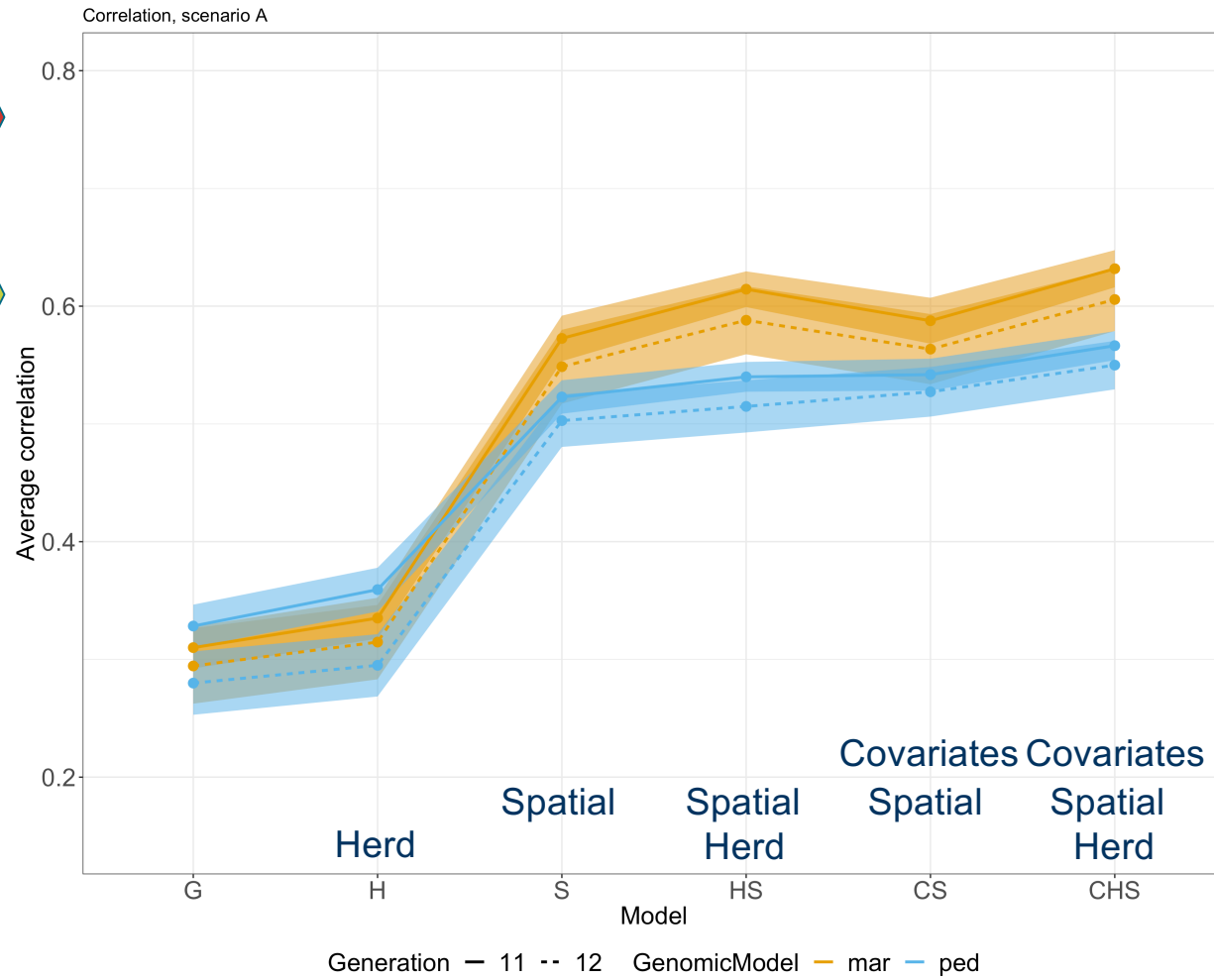
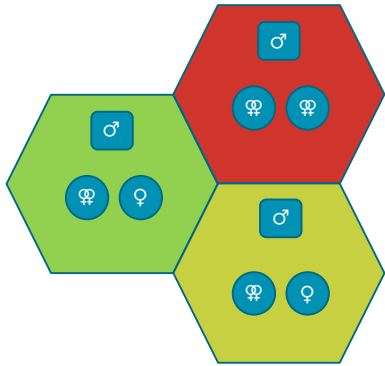
Phenotype = Location + Herd + Genetics + Noise

- Location
 - Spatial $\mathbf{l} \sim N(\mathbf{0}, \text{Matern}(\mathbf{E}, \text{Kappa}, \text{Var}_L))$
 - Spatial + Covariates
- Herd $h \sim N(\mathbf{0}, \text{IVar}_H)$
- Genetics
 - Pedigree-based $\mathbf{a} \sim N(\mathbf{0}, \mathbf{A}_{\text{ped}} \text{Var}_A)$
 - Genome-based $\mathbf{a} \sim N(\mathbf{0}, \mathbf{A}_{\text{mar}} \text{Var}_A)$

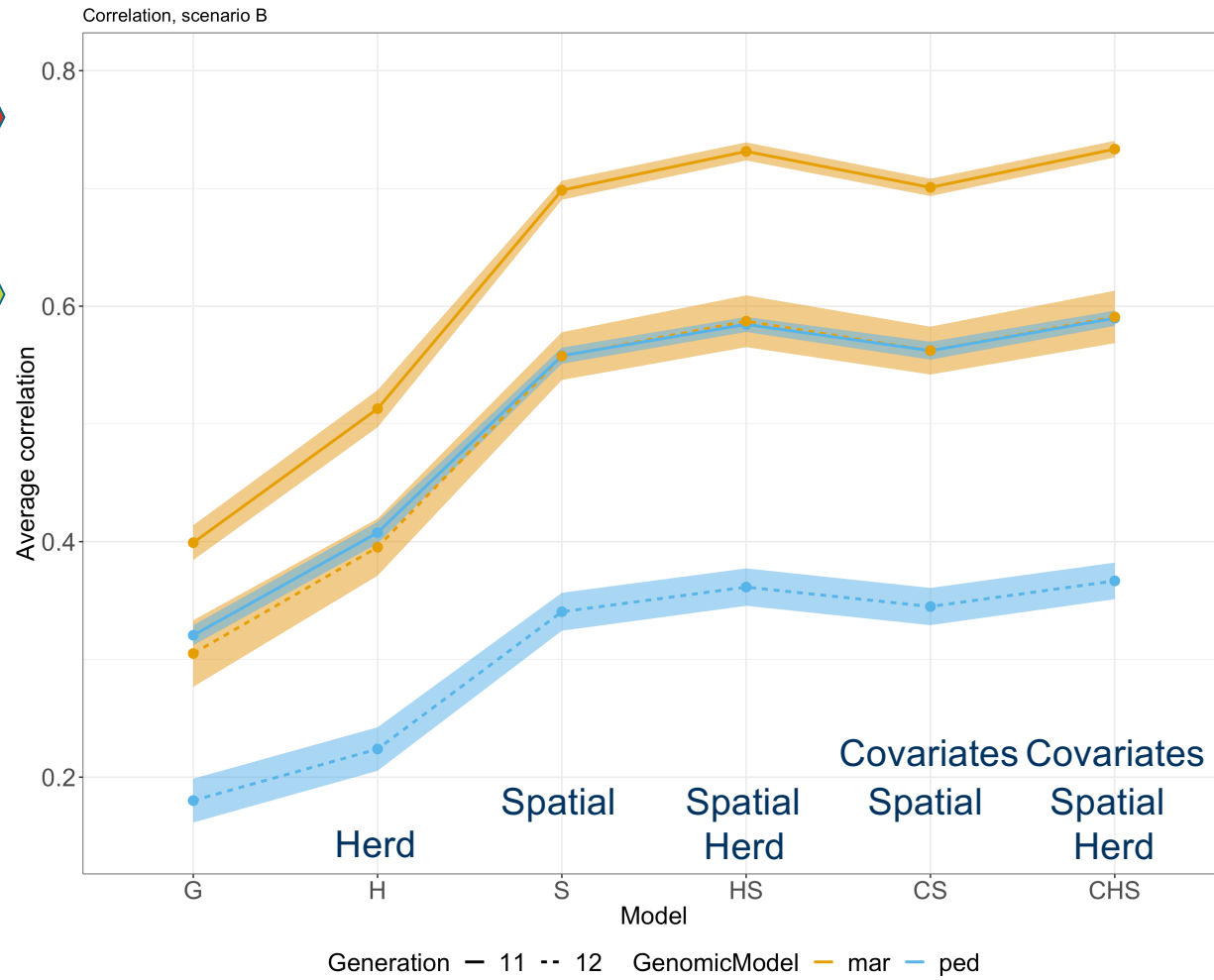
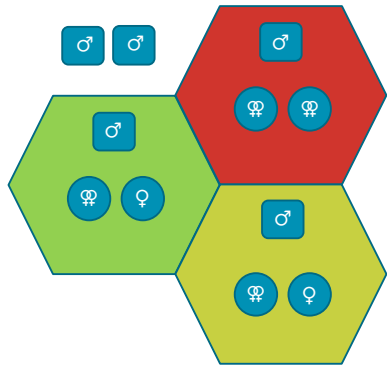
SPDE approach (similar idea to APY)



Accuracy of evaluation & prediction (simulation)



Accuracy of evaluation & prediction (simulation)



Real application

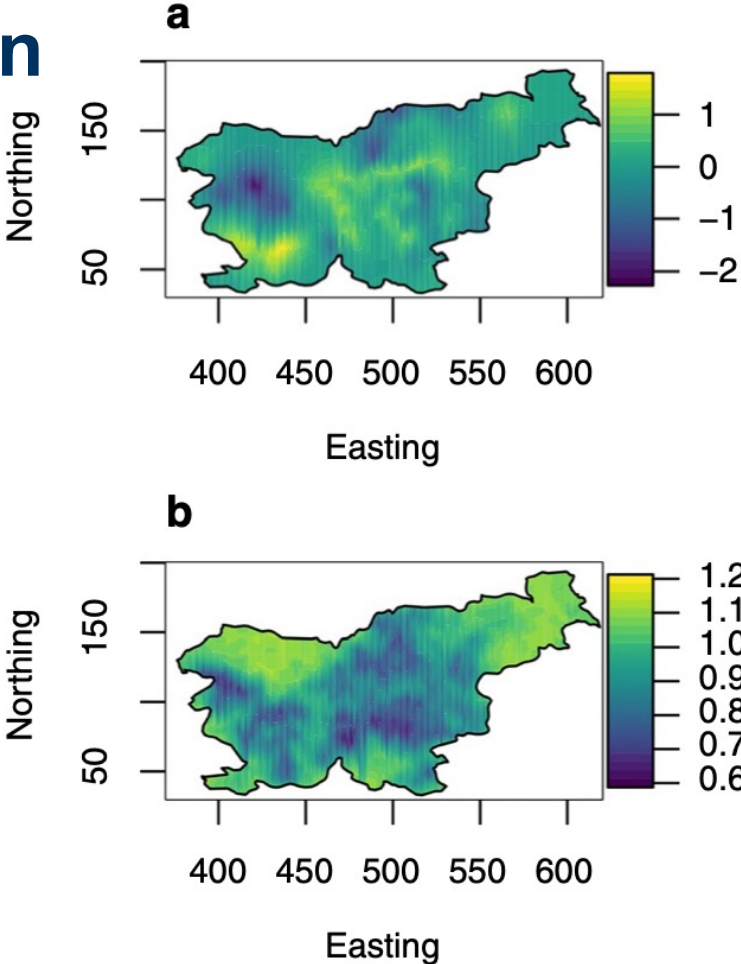
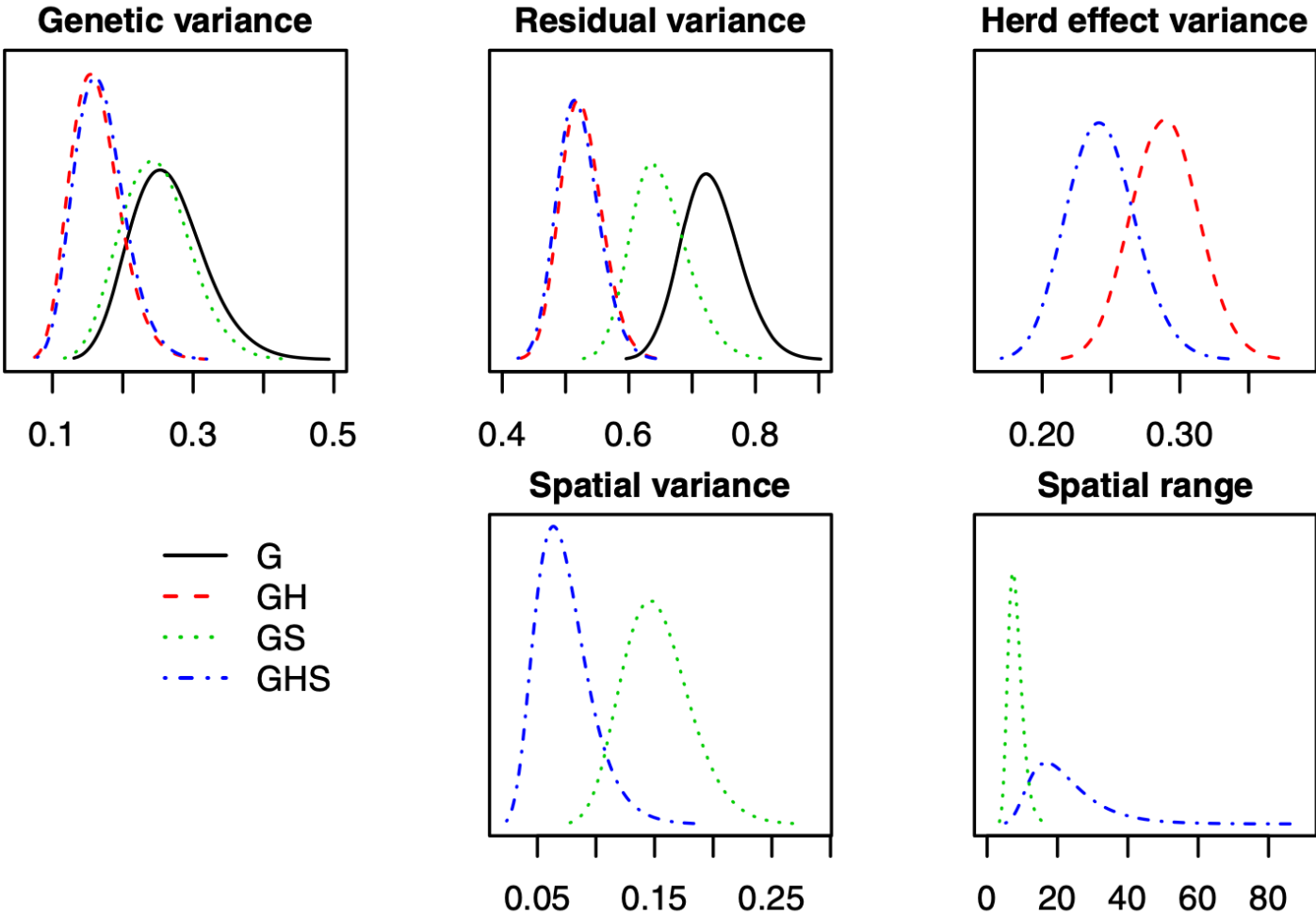


Fig. 5 Posterior mean (a) and standard deviation (b) of the estimated spatial effect (in units of posterior spatial standard deviation) from model GHS fitted to the real data—the axis units are in km

Estimated variance components



Does it matter? Yes!

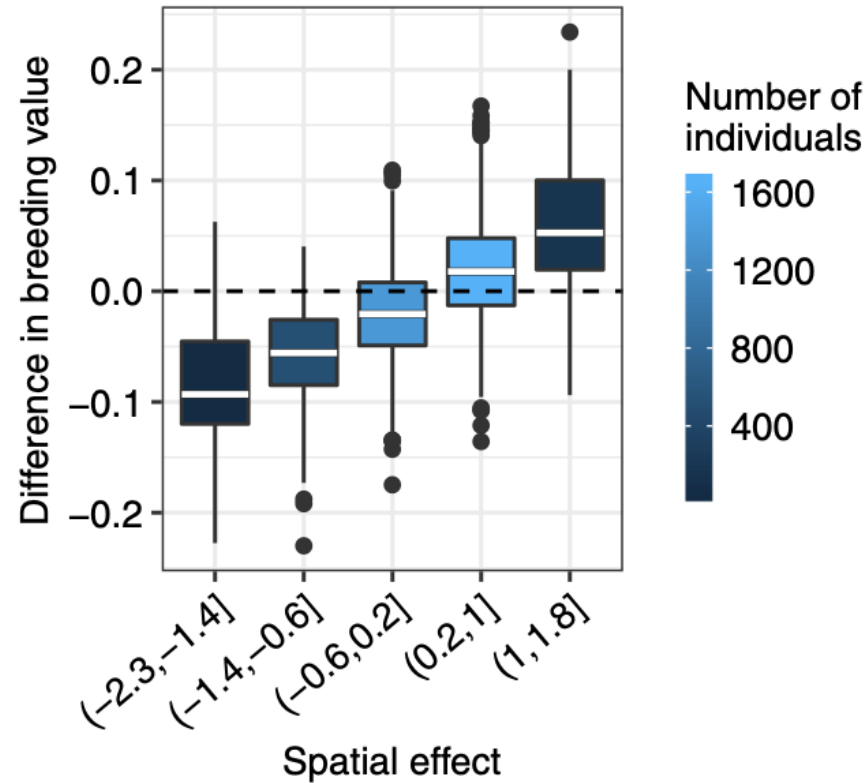


Fig. 6 The difference in estimated breeding values (in units of posterior genetic standard deviation) between models GH and GHS by the estimated spatial effect (in units of posterior spatial standard deviation) from model GHS fitted to the real data

Questions?!

APY approximation

- Number of genotyped individuals is growing rapidly!!!!
- Fitting genome-based models is becoming a challenge
 - Marker model scales with number of markers
 - Individual model scales with number of individuals
- Many solutions proposed to manage the scale
- APY (Algorithm for Proven and Young) is one of them
- Showcasing it due to connection to Day 3 and spatial models

Idea

- Pedigree-based model

$$\text{Var}(\mathbf{a}|\mathbf{T}) = \mathbf{TDT}^T \sigma_a^2 = \mathbf{A}\sigma_a^2$$

- if we use MME, we need inverse of \mathbf{A} , sparse with pedigree data
- sparse because of pedigree structure (conditioning on parents, Mendelian sampling terms are independent of parents)

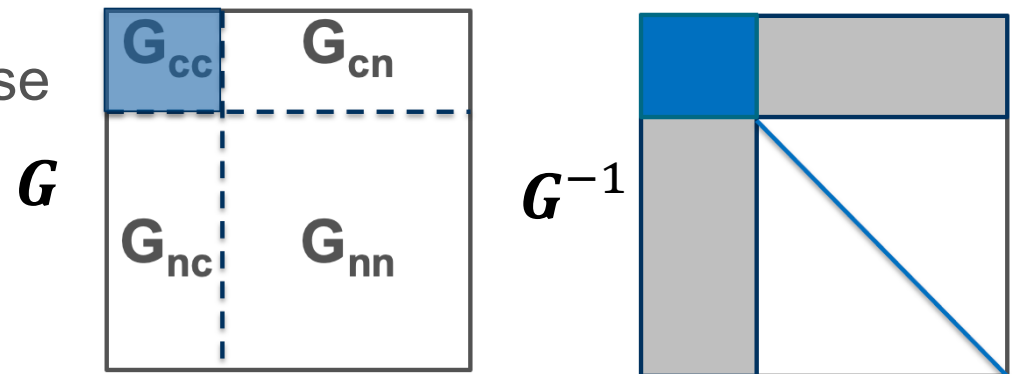
- Genome-based model

$$\text{Var}(\mathbf{a}|\mathbf{W}) = \mathbf{WW}^T \sigma_\alpha^2 = \mathbf{G}\sigma_\alpha^2$$

- if we use MME, we need inverse of \mathbf{G} , dense with genomic data
- dense because genomic data informs what we share or don't share (no independence)

Idea

- With limited number of markers and large number of ind., we soon incur linear dependencies \rightarrow we can explain genome of some ind. with genome of other ind. & what is not explained is independent
- APY
 - Pick some “core” individuals (=locations/knots in spatial community) & “non-core” individuals
 - Split \mathbf{G} between these two
 - Inverse \mathbf{G} can be made sparse



(Misztal, 2016)

How?

$$\text{Var}(\mathbf{a}|\mathbf{W}) = \mathbf{W}\mathbf{W}^T \sigma_\alpha^2 = \mathbf{G}\sigma_\alpha^2$$

$$\mathbf{a}_c = \mathbf{W}_c \boldsymbol{\alpha} + \mathbf{r}_c$$

$$\boldsymbol{\alpha} \sim N(\mathbf{0}, \mathbf{I}\sigma_\alpha^2)$$

$$\mathbf{r}_c \sim N(\mathbf{0}, \mathbf{I}m\sigma_\alpha^2)$$

$$\begin{aligned} \mathbf{a}_n &= \mathbf{W}_n \boldsymbol{\alpha} + \mathbf{r}_n \\ &= f(\mathbf{a}_c) + \mathbf{r}_n \end{aligned}$$

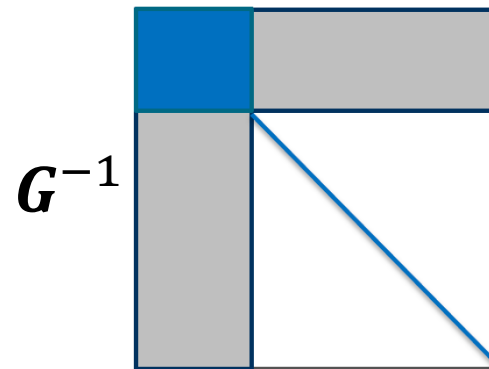
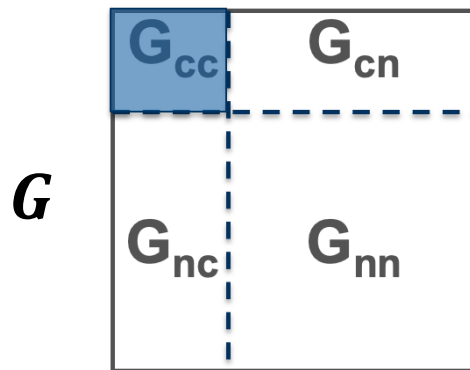
$$E(\mathbf{a}_n|\mathbf{a}_c) = E(\mathbf{a}_n) + \text{Cov}(\mathbf{a}_n, \mathbf{a}_c)\text{Var}(\mathbf{a}_c)^{-1}(\mathbf{a}_c - E(\mathbf{a}_c))$$

$$\text{Var}(\mathbf{a}_n|\mathbf{a}_c) = \text{Var}(\mathbf{a}_n) - \text{Cov}(\mathbf{a}_n, \mathbf{a}_c)\text{Var}(\mathbf{a}_c)^{-1}\text{Cov}(\mathbf{a}_c, \mathbf{a}_n)$$

How?

$$\begin{aligned} E(\mathbf{a}_n | \mathbf{a}_c) &= \text{Cov}(\mathbf{a}_n, \mathbf{a}_c) \text{Var}(\mathbf{a}_c)^{-1} \mathbf{a}_c \\ &= \mathbf{G}_{n,c} \mathbf{G}_{c,c}^{-1} \mathbf{a}_c \end{aligned}$$

$$\begin{aligned} \text{Var}(\mathbf{a}_n | \mathbf{a}_c) &= \text{Var}(\mathbf{a}_n) - \text{Cov}(\mathbf{a}_n, \mathbf{a}_c) \text{Var}(\mathbf{a}_c)^{-1} \text{Cov}(\mathbf{a}_c, \mathbf{a}_n) \\ &= \mathbf{G}_{n,n} - \mathbf{G}_{n,c} \mathbf{G}_{c,c}^{-1} \mathbf{G}_{c,n} \end{aligned}$$



(Miształ, 2016)

Core optimisation

Pocrnic et al. *Genetics Selection Evolution* (2022) 54:76
<https://doi.org/10.1186/s12711-022-00767-x>



GSE Genetics
Selection
Evolution

RESEARCH ARTICLE

Open Access

Optimisation of the core subset for the APY approximation of genomic relationships



Ivan Pocrnic^{1*} , Finn Lindgren², Daniel Tolhurst¹, William O. Herring³ and Gregor Gorjanc¹

Core optimisation - algorithms

Algorithm 1 Core subset optimisation using conditional covariance matrix \mathbf{C}

Require: n_c , \mathbf{k} , and $\mathbf{C}_0 = \mathbf{W}\mathbf{W}^\top$ ▷ Core subset size, Vector for core animals, and Covariance matrix

- 1: **for** i in 1 to n_c **do** ▷ Loop over the core subset
- 2: $k_i \leftarrow \operatorname{argmax}(\operatorname{diag}(\mathbf{C}_{i-1}))$ ▷ Find the i -th core animal
- 3: $\mathbf{e} \leftarrow 0$ ▷ Update the “selector” vector
- 4: $\mathbf{e}_{k_i} \leftarrow 1$
- 5: $\mathbf{C}_i \leftarrow \mathbf{C}_{i-1} - \mathbf{C}_{i-1}\mathbf{e}_{k_i}(\mathbf{e}_{k_i}^\top \mathbf{C}_{i-1} \mathbf{e}_{k_i})^{-1} \mathbf{e}_{k_i}^\top \mathbf{C}_{i-1}$ ▷ Update the covariance matrix
- 6: **end for**

Core optimisation - algorithms

Algorithm 1 Core subset optimisation using conditional covariance matrix \mathbf{C}

Require: n_c , \mathbf{k} , and $\mathbf{C}_0 = \mathbf{W}\mathbf{W}^\top$ \triangleright Core subset size, Vector for core animals, and Covariance matrix

- 1: **for** i in 1 to n_c **do** \triangleright Loop over the core subset
- 2: $k_i \leftarrow \text{argmax}(\text{diag}(\mathbf{C}_{i-1}))$ \triangleright Find the i -th core animal
- 3: $\mathbf{e} \leftarrow 0$ \triangleright Update the “selector” vector
- 4: $\mathbf{e}_{k_i} \leftarrow 1$
- 5: $\mathbf{C}_i \leftarrow \mathbf{C}_{i-1} - \mathbf{C}_{i-1}\mathbf{e}_{k_i}(\mathbf{e}_{k_i}^\top \mathbf{C}_{i-1}\mathbf{e}_{k_i})^{-1}\mathbf{e}_{k_i}^\top \mathbf{C}_{i-1}$ \triangleright Update the covariance matrix
- 6: **end for**

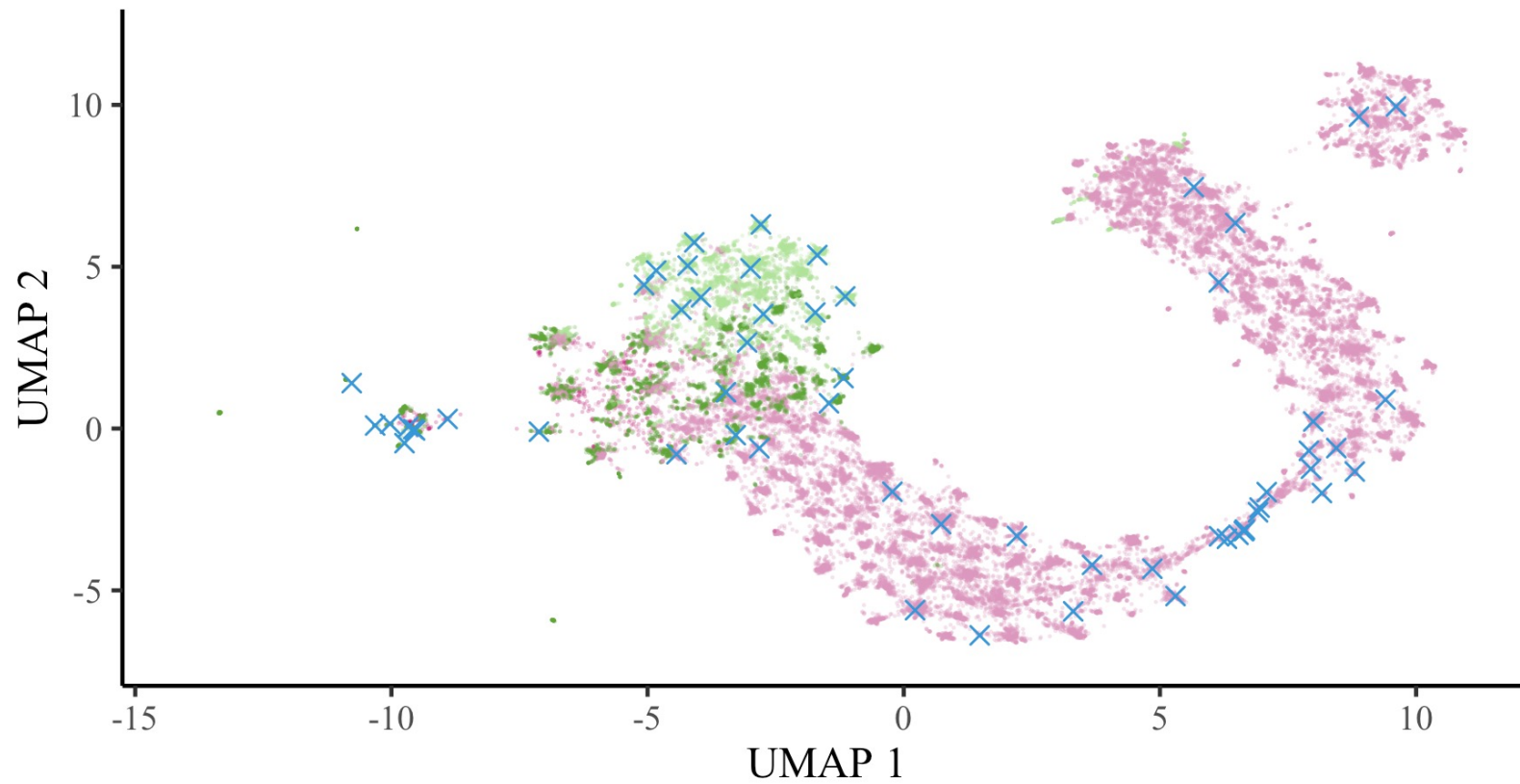
Algorithm 2 Core subset optimisation using conditional SNP genotype matrix \mathbf{W}

Require: n_c , \mathbf{k} , and $\mathbf{W}_0 = \mathbf{W}$ \triangleright Core subset size, Vector for core animals, and SNP genotype matrix

- 1: **for** i in 1 to n_c **do** \triangleright Loop over the core subset
- 2: $k_i \leftarrow \text{argmax}(\text{diag}(\mathbf{W}_{i-1}\mathbf{W}_{i-1}^\top))$ \triangleright Find the i -th core animal
- 3: $\mathbf{w}_{k_i} \leftarrow \mathbf{W}_{i-1}[k_i,]$ \triangleright Conditional SNP genotypes of the animal i
- 4: $\mathbf{W}_i \leftarrow \mathbf{W}_{i-1} - \mathbf{W}_{i-1}\mathbf{w}_{k_i}^\top \mathbf{w}_{k_i} / (\mathbf{w}_{k_i}\mathbf{w}_{k_i}^\top)$ \triangleright Update the SNP genotype matrix
- 5: **end for**

Core optimisation – UMAP visualisation

Line ● BC1 ● BC2 ● F1 ● L1 × Conditional



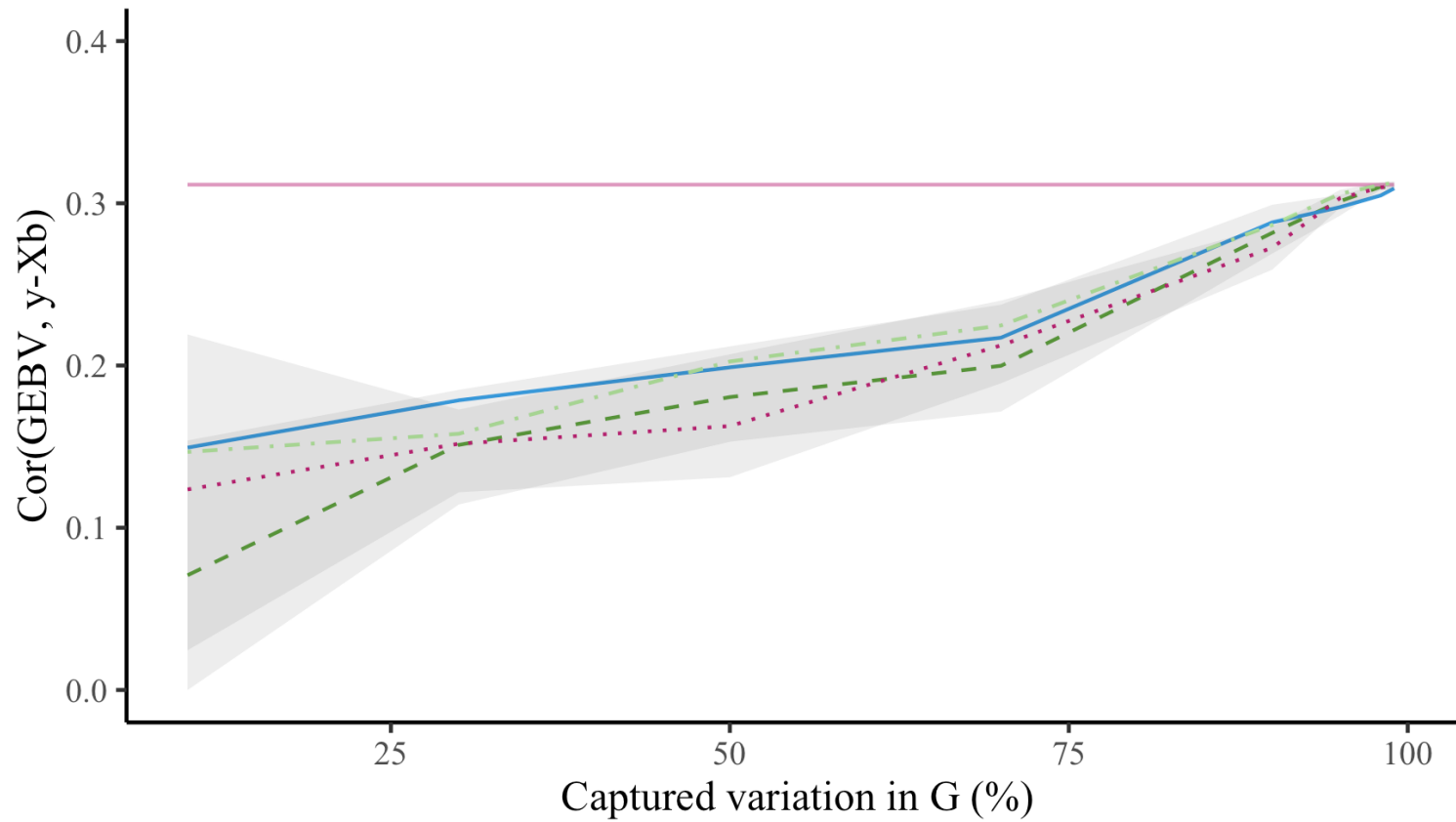
Core optimisation – UMAP visualisation

Line ● BC1 ● BC2 ● F1 ● L1 × Conditional



Impact on accuracy of prediction

Core selection Random Diagonal Weighted Conditional Full



Learning objectives

Separating genetic and environmental effects is a critical component of any quantitative genetics model

- Showcase challenge and solution for modelling data from smallholder settings
- Aside: APY approximation

Questions?!



THE UNIVERSITY
of EDINBURGH



Spatial modelling improves genetic evaluation in smallholder breeding programs

Gregor Gorjanc, Chris Gaynor, Jon Bancic, Daniel Tolhurst

UNE, Armidale

2024-02-09

