

Modern genetic evaluation systems: A Python-based programming approach

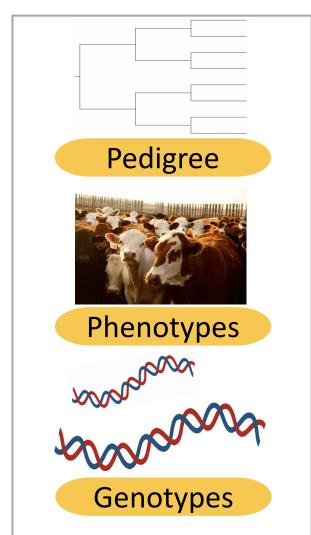
<u>Kristin Lee</u>, Gordon Vander Voort, Ricardo Ventura, Flavio Schenkel, Angela Cánovas

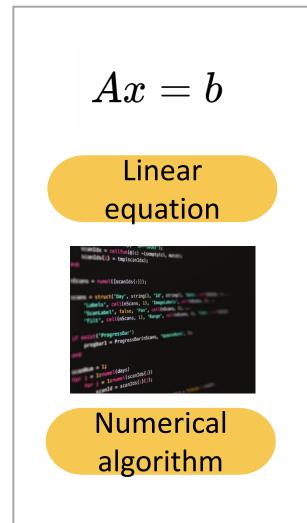
Centre for Genetic Improvement of Livestock, Department of Animal Biosciences, University of Guelph, Guelph, Ontario, Canada

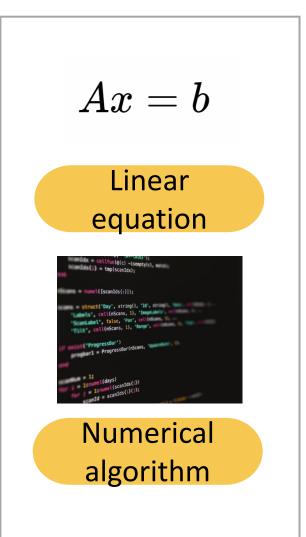


Introduction

Genetic evaluation systems











Introduction

Traditional genetic evaluation systems are developed in C or Fortran

Benefits

- Lower-level languages
 - Efficiency
 - Numerical & scientific computing
 - Performance optimization
 - Explicit memory management



Limitations

- Less flexibility
- Slower development
 - Higher maintenance
 - → High costs
 - → Affect innovation

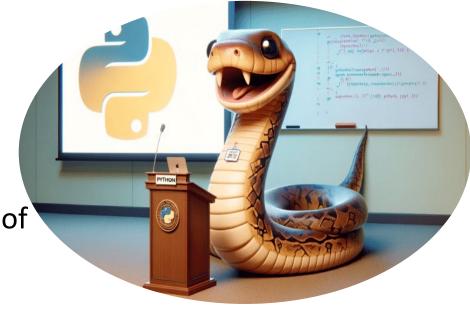
Objective

Python could be used to develop an adaptable software that can produce accurate breeding values within practical timeframes

Benefits

- Higher-level language
- Faster development
- Easier maintenance

→ Easy & fast integration of new features



Limitations

- Slower performance
- Higher memory requirements
- → Modern computers

Data

Simulated purebred angus beef cattle:

- AlphaSimR
- Sample size: N = 976,400
- Pedigree: 15 generations
- Phenotypes: Birth weight, weaning gain, post-weaning gain
- Masked phenotypes for final simulated generation



Data

Genotypes:

- 48,981 SNP markers
- Random sampling from final simulated generation
- Sample Size: 10,000
- 3 replicates of sample size



Linear model

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} + \begin{bmatrix} \mathbf{Z_1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z_2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{Z_3} \end{bmatrix} \begin{bmatrix} \mathbf{a_1} \\ \mathbf{a_2} \\ \mathbf{a_3} \end{bmatrix} + \begin{bmatrix} \mathbf{e_1} \\ \mathbf{e_2} \\ \mathbf{e_3} \end{bmatrix}$$

y = vector of phenotype records

 μ = phenotypic mean

Z = incidence matrix relating phenotypes to animal effects

a = vector of random additive genetic effects

e = vector of residual effects

Assumptions

$$\operatorname{Var}\begin{bmatrix}\mathbf{a}\\\mathbf{e}\end{bmatrix} = \begin{bmatrix}\mathbf{G} \otimes \mathbf{Rel} & 0\\ 0 & \mathbf{R} \otimes \mathbf{I}\end{bmatrix}$$

Rel = relationship matrix

G = (co)variance matrices for genetic effects

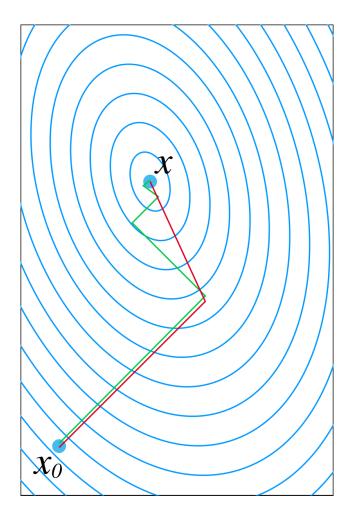
R = (co)variance matrices for residual effects

I = identity matrix

Numerical algorithm

Preconditioned conjugate gradient with iteration on data:

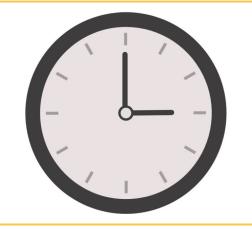
- Diagonal preconditioner
- Serial, multiprocessing, sparse
- Convergence criteria = $\frac{\|\mathbf{r_i}\|}{\|\mathbf{b}\|}$ = 1×10^{-5}



Software performance

Runtime analysis:

- The software was executed five times
- Average time elapsed of major components was reported



Server specifications:

- Operating system: Linux (64-bit)
- Memory (RAM): 96GB
- Processor: Intel(R) Xeon(R) Gold 6242 CPU @ 2.80GHz
- Cores: 8



Results verification

MiXBLUP:

- Commercial software
- Genetic evaluation system
- Pearson correlation
- Estimated breeding values



Software performance

| Component | Time (seconds) |
|------------------------------------|----------------|
| Read data | 1.79 |
| Preprocessing | 38.51 |
| Relationship matrix | |
| Mendelian sampling | 0.19 |
| Inbreeding | 15,020.00 |
| Inverse of the relationship matrix | 3.57 |
| Preconditioned conjugate gradient | |
| A) Serial processing | 2,688.15 |
| B) Multiprocessing | 933.24 |
| C) Sparse matrix operations | 770.12 |
| Write results | 4.09 |

Software performance

| Component | Time (seconds) | |
|------------------------------------|----------------|--|
| Read data | 1.79 | |
| Preprocessing | 38.51 | |
| Relationship matrix | | |
| Mendelian sampling | 0.19 | |
| Inbreeding | 15,020.00 | |
| Inverse of the relationship matrix | 3.57 | |
| Preconditioned conjugate gradient | | |
| A) Serial processing | 2,688.15 | |
| B) Multiprocessing | 933.24 | |
| C) Sparse matrix operations | 770.12 | |
| Write results | 4.09 | |

Software performance

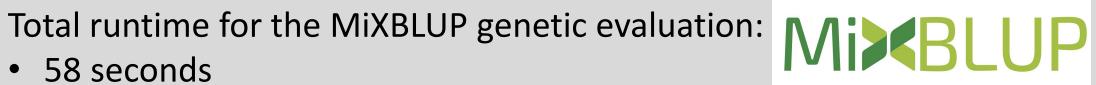
| Component | Time (seco | onds) | |
|------------------------------------|------------|----------------|-------|
| Read data | 1.79 | | |
| Preprocessing | 38.51 | | |
| Relationship matrix | | | |
| Mendelian sampling | 0.19 | | |
| Inbreeding | 15,020.00 | 4 hours 10 min | nutes |
| Inverse of the relationship matrix | 3.57 | | |
| Preconditioned conjugate gradient | | | |
| A) Serial processing | 2,688.15 | | |
| B) Multiprocessing | 933.24 | | |
| C) Sparse matrix operations | 770.12 | | |
| Write results | 4.09 | | |

Software performance

| Component | Time (seco | nds) |
|------------------------------------|------------|-----------------------|
| Read data | 1.79 | |
| Preprocessing | 38.51 | |
| Relationship matrix | | |
| Mendelian sampling | 0.19 | |
| Inbreeding | 15,020.00 | |
| Inverse of the relationship matrix | 3.57 | |
| Preconditioned conjugate gradient | | |
| A) Serial processing | 2,688.15 | 44 minutes 48 seconds |
| B) Multiprocessing | 933.24 | 15 minutes 33 seconds |
| C) Sparse matrix operations | 770.12 | 12 minutes 50 seconds |
| Write results | 4.09 | 12 |

Software performance

• 58 seconds



Results verification

Pearson correlation coefficient (r) between estimated breeding values (EBV) of the Python and MiXBLUP genetic evaluation systems

| Trait | r _{EBV:EBV} |
|-------------------|----------------------|
| Birth weight | 1.0 |
| Weaning gain | 1.0 |
| Post-weaning gain | 1.0 |

Software performance

| Component | Time |
|-----------------------------------|-----------|
| Read genotypes | 56.06 |
| Preprocessing | 17.95 |
| Relationship matrix | |
| Genomic relationship matrix (G) | 29.49 |
| Blend & tune G | 5.06 |
| Invert G | 10.11 |
| Construct and invert A22 | 52,040.54 |
| Construct H-1 | 31.51 |
| Preconditioned conjugate gradient | |
| Sparse matrix operations | 2,367.50 |

Software performance

| Component | Time |
|-----------------------------------|-----------|
| Read genotypes | 56.06 |
| Preprocessing | 17.95 |
| Relationship matrix | |
| Genomic relationship matrix (G) | 29.49 |
| Blend & tune G | 5.06 |
| Invert G | 10.11 |
| Construct and invert A22 | 52,040.54 |
| Construct H-1 | 31.51 |
| Preconditioned conjugate gradient | |
| Sparse matrix operations | 2,367.50 |

Software performance

| Component | Time |
|-----------------------------------|-----------|
| Read genotypes | 56.06 |
| Preprocessing | 17.95 |
| Relationship matrix | |
| Genomic relationship matrix (G) | 29.49 |
| Blend & tune G | 5.06 |
| Invert G | 10.11 |
| Construct and invert A22 | 52,040.54 |
| Construct H-1 | 31.51 |
| Preconditioned conjugate gradient | |
| Sparse matrix operations | 2,367.50 |

Software performance

| Component | Time | |
|-----------------------------------|-----------|---------------------|
| Read genotypes | 56.06 | |
| Preprocessing | 17.95 | |
| Relationship matrix | | |
| Genomic relationship matrix (G) | 29.49 | |
| Blend & tune G | 5.06 | |
| Invert G | 10.11 | |
| Construct and invert A22 | 52,040.54 | 14 hours 27 minutes |
| Construct H-1 | 31.51 | |
| Preconditioned conjugate gradient | | |
| Sparse matrix operations | 2,367.50 | |

Software performance

| Component | Time |
|-----------------------------------|-----------|
| Read genotypes | 56.06 |
| Preprocessing | 17.95 |
| Relationship matrix | |
| Genomic relationship matrix (G) | 29.49 |
| Blend & tune G | 5.06 |
| Invert G | 10.11 |
| Construct and invert A22 | 52,040.54 |
| Construct H-1 | 31.51 |
| Preconditioned conjugate gradient | |
| Sparse matrix operations | 2,367.50 |

Software performance

| Component | Time | |
|-----------------------------------|-----------|----------------------|
| Read genotypes | 56.06 | |
| Preprocessing | 17.95 | |
| Relationship matrix | | |
| Genomic relationship matrix (G) | 29.49 | |
| Blend & tune G | 5.06 | |
| Invert G | 10.11 | |
| Construct and invert A22 | 52,040.54 | |
| Construct H-1 | 31.51 | |
| Preconditioned conjugate gradient | | |
| Sparse matrix operations | 2,367.50 | 39 minutes 27 second |

Software performance

Total runtime for the MiXBLUP genetic evaluation:

15 minutes 14 seconds



Results verification

Pearson correlation coefficient (r) between estimated breeding values for genotyped individuals of the Python & MiXBLUP genetic evaluation systems

| Trait | r _{EBV:EBV} |
|-------------------|----------------------|
| Birth weight | 0.99 |
| Weaning gain | 0.99 |
| Post-weaning gain | 0.99 |

Conclusions



Python for genetic evaluation system development:

- Simple development process
- Accurate
- Practical timeframes for small to medium sized breeding programs
- Components needing optimization identified

Acknowledgements





klee32@uoguelph.ca





