

GCA: An R package for genetic connectedness analysis using pedigree and genomic data

Haipeng Yu^{1*} and Gota Morota^{1*}

¹Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, Blacksburg, VA, USA 24061

Running title: Genetic connectedness analysis software

ORCID: 0000-0002-8923-9733 (HY) and 0000-0002-3567-6911 (GM).

* Corresponding author:

Haipeng Yu
Department of Animal and Poultry Sciences
Virginia Polytechnic Institute and State University
175 West Campus Drive
Blacksburg, Virginia 24061 USA.
E-mail: haipengyu@vt.edu

Gota Morota
Department of Animal and Poultry Sciences
Virginia Polytechnic Institute and State University
175 West Campus Drive
Blacksburg, Virginia 24061 USA.
E-mail: morota@vt.edu

Abstract

Background: Genetic connectedness is a critical component of genetic evaluation as it assesses the comparability of predicted genetic values across management units. Genetic connectedness also plays an essential role in quantifying the linkage between reference and validation sets in whole-genome prediction. Despite its importance, there is no user-friendly software tool available to calculate connectedness statistics.

Results: We developed the GCA R package to perform genetic connectedness analysis for pedigree and genomic data. The software implements a large collection of various connectedness statistics as a function of prediction error variance or variance of unit effect estimates. The GCA R package is available at GitHub and the source code is provided as open source.

Conclusion: The GCA R package allows users to easily assess the connectedness of their data. It is also useful to determine the potential risk of comparing predicted genetic values of individuals across units or measure the connectedness level between training and testing sets in genomic prediction.

Keywords: genetic connectedness, prediction error of variance, variance of unit effect estimates

Background

Genetic connectedness quantifies the extent to which estimated breeding values can be fairly compared across management units or contemporary groups [1, 2]. Genetic evaluation is known to be more robust when the connectedness level is high enough due to sufficient sharing of genetic material across groups. In such scenarios, the best linear unbiased prediction minimizes the risk of uncertainty in ranking of individuals. On the other hand, limited or no sharing of genetic material leads to less reliable comparisons of genetic evaluation methods [3]. High-throughput genetic variants spanning the entire genome available for a wide range of agricultural species have now opened up an opportunity to assess connectedness using genomic data. A recent study showed that genomic relatedness strengthens the measures of connectedness across units compared with the use of pedigree relationships [4]. The concept of genetic connectedness was later extended to measure the connectedness level between reference and validation sets in whole-genome prediction. In general, it was observed that increased connectedness led to increased prediction accuracy of genetic values evaluated by a cross-validation [5]. Comparability of total genetic values across units by accounting for additive as well as non-additive genetic effects has also been investigated [6].

Despite the importance of connectedness, there is no user-friendly software tool available that offers computation of a comprehensive list of connectedness statistics. Therefore, we developed a genetic connectedness analysis R package, GCA, which measures the connectedness between individuals across units using pedigree and genomic data. The objective of this article is to describe a large collection of connectedness statistics implemented in the GCA package, overview the software architecture, and present several examples using simulated data.

Connectedness statistics

A list of connectedness statistics supported by the GCA R package is shown in Figure 1. These statistics can be classified into core functions derived from either prediction error variance (PEV) or variance of unit effect estimates (VE). PEV-derived metrics include prediction error variance of differences (PEVD), coefficient of determination (CD), and prediction error correlation (r). Further, each metric based on PEV can be summarized at the unit level as the average PEV of all pairwise differences between individuals across units, average PEV within and across units, or using a contrast vector. VE-derived metrics include variance of differences in unit effects (VED), coefficient of determination of VED (CDVED), and connectedness rating (CR). For each VE-derived metric, three correction factors accounting for the number of fixed effects can be applied. These include no correction (0), correcting for one fixed effect (1), and correcting for two or more fixed effects (2). Thus, a combination of core functions, metrics, summary functions, and correction factors uniquely characterizes connectedness statistics. Further, the overall connectedness statistic can be obtained by calculating the average of the pairwise connectedness statistics across units.

Core functions

Prediction error variance (PEV)

A PEV matrix is obtained from Henderson's mixed model equations (MME) by assuming a standard linear mixed model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$, where \mathbf{y} , \mathbf{b} , \mathbf{u} , and $\boldsymbol{\epsilon}$ refer to a vector of phenotypes, systematic effects, random additive genetic effects, and residuals, respectively [7]. The \mathbf{X} and \mathbf{Z} are incidence matrices associating systematic effects and genetic values to observations, respectively. The joint

distribution of random effects is as given below.

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} \sim N \left[\begin{pmatrix} \mathbf{X}\mathbf{b} \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbf{Z}\mathbf{K}\sigma_u^2\mathbf{Z}' + \mathbf{I}\sigma_\epsilon^2 & \mathbf{Z}\mathbf{K}\sigma_u^2 & \mathbf{I}\sigma_\epsilon^2 \\ \mathbf{K}\mathbf{Z}'\sigma_u^2 & \mathbf{K}\sigma_u^2 & 0 \\ \mathbf{I}\sigma_\epsilon^2 & 0 & \mathbf{I}\sigma_\epsilon^2 \end{pmatrix} \right],$$

where \mathbf{K} is a relationship matrix, σ_u^2 is the additive genetic variance, and σ_ϵ^2 is the residual variance.

The corresponding MME is as given below.

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \mathbf{K}^{-1}\lambda \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{Z}'\mathbf{y} \end{bmatrix},$$

where $\lambda = \frac{\sigma_\epsilon^2}{\sigma_u^2}$ is the ratio of variance components. The inverse of the MME coefficient matrix derived from this model is as given below.

$$\begin{aligned} \mathbf{C}^{-1} &= \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \mathbf{K}^{-1}\lambda \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \mathbf{C}^{11} & \mathbf{C}^{12} \\ \mathbf{C}^{21} & \mathbf{C}^{22} \end{bmatrix}. \end{aligned}$$

Then the PEV of \mathbf{u} is derived as shown in Henderson [7].

$$\begin{aligned} \text{PEV}(\mathbf{u}) &= \text{Var}(\hat{\mathbf{u}} - \mathbf{u}) \\ &= \text{Var}(\mathbf{u}|\hat{\mathbf{u}}) \\ &= (\mathbf{Z}'\mathbf{M}\mathbf{Z} + \mathbf{K}^{-1}\lambda)^{-1}\sigma_\epsilon^2 \\ &= \mathbf{C}^{22}\sigma_\epsilon^2, \end{aligned}$$

where $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is the absorption (projection) matrix for fixed effects and \mathbf{C}^{22} is the lower right quadrant of the inverse of coefficient matrix. Note that $\text{PEV}(\mathbf{u})$ can be viewed as the posterior variance of \mathbf{u} .

Variance of unit effect estimates (VE)

An alternative option for the choice of core function is to use VE, which is based on the variance-covariance matrix of estimated management unit or contemporary group effects. Kennedy and Trus [8] argued that mean PEV over unit (PEV_{Mean}) defined as the average of PEV between individuals within the same unit can be approximated by $\text{VE} = \text{Var}(\hat{b})$, that is

$$\begin{aligned} \text{VE0} &= \text{Var}(\hat{b}) \\ &= [\mathbf{X}'\mathbf{X} - \mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{K}^{-1}\lambda)^{-1}\mathbf{Z}'\mathbf{X}]^{-1}\sigma_{\epsilon}^2 \\ &\approx \text{PEV}_{\text{Mean}} \end{aligned} \tag{1}$$

Holmes et al. [9] pointed out that the agreement between PEV_{Mean} and VE0 depends on a number of fixed effects other than the management group fitted in the model. They proposed exact ways to derive PEV_{Mean} as a function of VE and suggested addition of a few correction factors. When unit effect is the only fixed effect included in the model, the exact PEV_{Mean} can be obtained as given below.

$$\text{VE1} = \text{PEV}_{\text{Mean}} = \text{Var}(\hat{b}) - \sigma_{\epsilon}^2(\mathbf{X}'\mathbf{X})^{-1}, \tag{2}$$

where $\mathbf{X}'\mathbf{X}^{-1}$ is a diagonal matrix with i th diagonal element equal to $\frac{1}{n_i}$, and n_i is the number of records in unit i . Thus, the term $\sigma_{\epsilon}^2(\mathbf{X}'\mathbf{X})^{-1}$ corrects the number of records within units. Accounting for additional fixed effects beyond unit effect when computing PEV_{Mean} is given by the following

equation.

$$\begin{aligned}
 \text{VE2} &= \text{PEV}_{\text{Mean}} & (3) \\
 &= \text{Var}(\hat{b}_1) - \sigma_\epsilon^2 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \\
 &+ (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \text{Var}(\hat{b}_2) \mathbf{X}_2' \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \\
 &+ (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \text{Cov}(\hat{b}_2, \hat{b}_1) \\
 &+ \text{Cov}(\hat{b}_1, \hat{b}_2) \mathbf{X}_2' \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1}, & (4)
 \end{aligned}$$

where \mathbf{X}_1 and \mathbf{X}_2 represent incidence matrices for units and other fixed effects, respectively, and \hat{b}_1 and \hat{b}_2 refer to the estimates of unit effects and other fixed effects, respectively [9]. This equation is suitable for cases in which there are two or more fixed effects fitted in the model.

Connectedness metrics

Below we describe connectedness metrics implemented in the GCA package. These metrics are the function of PEV or VE described earlier (Figure 1).

Prediction error variance of difference (PEVD)

A PEVD metric measures the prediction error variance difference of breeding values between individuals from different units [8]. The PEVD between two individuals i and j is expressed as shown below.

$$\begin{aligned}
 \text{PEVD}(\hat{u}_i - \hat{u}_j) &= [\text{PEV}(\hat{u}_i) + \text{PEV}(\hat{u}_j) - 2\text{PEC}(\hat{u}_i, \hat{u}_j)] \\
 &= (\mathbf{C}_{ii}^{22} - \mathbf{C}_{ij}^{22} - \mathbf{C}_{ji}^{22} + \mathbf{C}_{jj}^{22}) \sigma_\epsilon^2 \\
 &= (\mathbf{C}_{ii}^{22} + \mathbf{C}_{jj}^{22} - 2\mathbf{C}_{ij}^{22}) \sigma_\epsilon^2, & (5)
 \end{aligned}$$

where PEC_{ij} is the off-diagonal element of the PEV matrix corresponding to the prediction error covariance between errors of genetic values.

Individual average PEVD: When pairwise PEVD is first computed at the individual level using equation (5), these estimates need to be aggregated and summarized at the unit level. A calculation of summary PEVD can be traced back to Kennedy and Trus [8] as the average of PEVD between individuals across two units.

$$\text{PEVD}_{i'j'} = \frac{1}{n_{i'} \cdot n_{j'}} \sum \text{PEVD}_{i'j'},$$

where $n_{i'}$ and $n_{j'}$ are the total number of records in units i' and j' , respectively and $\sum \text{PEVD}_{i'j'}$ is the sum of all pairwise differences between the two units. We refer to this summary method as individual average. A flow diagram illustrating the computational procedure is shown in Figure 2A.

Group average PEVD: The second summary method applies equation (5) after calculating PEV_{Mean} of i' th and j' th units and mean prediction error covariance (PEC_{Mean}) between i' th and j' th units.

$$\text{PEVD}_{i'j'} = \overline{\text{PEV}}_{i'i'} + \overline{\text{PEV}}_{j'j'} - 2\overline{\text{PEC}}_{i'j'}, \quad (6)$$

where $\overline{\text{PEV}}_{i'i'}$, $\overline{\text{PEV}}_{j'j'}$, and $\overline{\text{PEC}}_{i'j'}$ denote PEV_{Mean} in i' th and j' th units, and PEC_{Mean} between i' th and j' th units. We refer to this summary method as group average as illustrated in Figure 2B.

Contrast PEVD: The third summary method is PEVD of contrast between a pair of units.

$$\text{PEVD}(\mathbf{x}) = \mathbf{x}'\mathbf{C}^{22}\mathbf{x}\sigma_{\epsilon}^2,$$

where \mathbf{x} is a contrast vector involving $1/n_{i'}$, $1/n_{j'}$ and 0 corresponding to individuals belonging to i' th, j' th, and the remaining units. The sum of elements in \mathbf{x} equals to zero. A flow diagram showing a computational procedure is shown in Figure 2C.

Coefficient of determination (CD)

A CD metric measures the precision of genetic values and can be interpreted as the square of the correlation between the predicted and the true difference in the genetic values or the ratio of posterior and prior variances of genetic values \mathbf{u} [10]. A notable difference between CD and PEVD is that CD penalizes connectedness measurements when across units include individuals that are genetically too similar [4, 5]. A pairwise CD between individuals i and j is given by the following equation.

$$\begin{aligned} \text{CD}_{ij} &= \frac{\text{Var}(\hat{\mathbf{u}})}{\text{Var}(\mathbf{u})} \\ &= \frac{\text{Var}(\mathbf{u}) - \text{Var}(\mathbf{u}|\hat{\mathbf{u}})}{\text{Var}(\mathbf{u})} \\ &= 1 - \frac{\text{Var}(\mathbf{u}|\hat{\mathbf{u}})}{\text{Var}(\mathbf{u})} \\ &= 1 - \lambda \frac{\mathbf{C}_{ii}^{22} + \mathbf{C}_{jj}^{22} - 2\mathbf{C}_{ij}^{22}}{\mathbf{K}_{ii} + \mathbf{K}_{jj} - 2\mathbf{K}_{ij}}, \end{aligned}$$

where \mathbf{K}_{ii} and \mathbf{K}_{jj} are i th and j th diagonal elements of \mathbf{K} , and \mathbf{K}_{ij} is the relationship between i th and j th individuals [11].

Individual average CD: Individual average CD is derived from the average of CD between individuals across two units.

$$\begin{aligned}
 CD_{i'j'} &= 1 - \lambda \cdot \frac{\frac{1}{n_{i'} \cdot n_{j'}} \cdot \sum (\mathbf{C}^{22}_{i'i'} + \mathbf{C}^{22}_{j'j'} - 2\mathbf{C}^{22}_{i'j'})}{\frac{1}{n_{i'} \cdot n_{j'}} \cdot \sum (\mathbf{K}_{i'i'} + \mathbf{K}_{j'j'} - 2\mathbf{K}_{i'j'})} \\
 &= 1 - \frac{\frac{1}{n_{i'} \cdot n_{j'}} \cdot \sigma_e^2 \cdot \sum (\mathbf{C}^{22}_{i'i'} + \mathbf{C}^{22}_{j'j'} - 2\mathbf{C}^{22}_{i'j'})}{\frac{1}{n_{i'} \cdot n_{j'}} \cdot \sigma_u^2 \cdot \sum (\mathbf{K}_{i'i'} + \mathbf{K}_{j'j'} - 2\mathbf{K}_{i'j'})} \\
 &= 1 - \frac{\frac{1}{n_{i'} \cdot n_{j'}} \sum \text{PEVD}_{i'j'}}{\frac{1}{n_{i'} \cdot n_{j'}} \cdot \sigma_u^2 \cdot \sum (\mathbf{K}_{i'i'} + \mathbf{K}_{j'j'} - 2\mathbf{K}_{i'j'})} \\
 &= 1 - \frac{\sum \text{PEVD}_{i'j'}}{\sigma_u^2 \cdot \sum (\mathbf{K}_{i'i'} + \mathbf{K}_{j'j'} - 2\mathbf{K}_{i'j'})}.
 \end{aligned}$$

A flow diagram of individual average CD is shown in Figure 3A.

Group average CD: Similar to the group average PEVD statistic, PEV_{Mean} and PEC_{Mean} can be used to summarize CD at the unit level.

$$\begin{aligned}
 CD_{i'j'} &= 1 - \lambda \cdot \frac{\overline{\mathbf{C}^{22}_{i'i'}} + \overline{\mathbf{C}^{22}_{j'j'}} - 2\overline{\mathbf{C}^{22}_{i'j'}}}{(\overline{\mathbf{K}_{i'i'}} + \overline{\mathbf{K}_{j'j'}} - 2\overline{\mathbf{K}_{i'j'}})} \\
 &= 1 - \frac{\sigma_e^2 \cdot \overline{\mathbf{C}^{22}_{i'i'}} + \overline{\mathbf{C}^{22}_{j'j'}} - 2\overline{\mathbf{C}^{22}_{i'j'}}}{\sigma_u^2 \cdot (\overline{\mathbf{K}_{i'i'}} + \overline{\mathbf{K}_{j'j'}} - 2\overline{\mathbf{K}_{i'j'}})} \\
 &= 1 - \frac{\overline{\text{PEV}_{i'i'}} + \overline{\text{PEV}_{j'j'}} - 2\overline{\text{PEC}_{i'j'}}}{\sigma_u^2 \cdot (\overline{\mathbf{K}_{i'i'}} + \overline{\mathbf{K}_{j'j'}} - 2\overline{\mathbf{K}_{i'j'}})} \\
 &= 1 - \frac{\text{PEVD}_{i'j'}}{\sigma_u^2 \cdot (\overline{\mathbf{K}_{i'i'}} + \overline{\mathbf{K}_{j'j'}} - 2\overline{\mathbf{K}_{i'j'}})}. \tag{7}
 \end{aligned}$$

Here, $\overline{\mathbf{K}_{i'i'}}$, $\overline{\mathbf{K}_{j'j'}}$ and $\overline{\mathbf{K}_{i'j'}}$ refer to the means of relationship coefficients in units i' and j' , and the mean relationship coefficient between two units i' and j' , respectively. Graphical derivation of group average CD is illustrated in Figure 3B.

Contrast CD: A contrast of CD between any pair of units is given by [11]

$$\begin{aligned}
 \text{CD}(\mathbf{x}) &= 1 - \frac{\text{Var}(\mathbf{x}'\mathbf{u}|\hat{\mathbf{u}})}{\text{Var}(\mathbf{x}'\mathbf{u})} \\
 &= 1 - \lambda \cdot \frac{\mathbf{x}'\mathbf{C}^{22}\mathbf{x}}{\mathbf{x}'\mathbf{K}\mathbf{x}} \\
 &= 1 - \frac{\mathbf{x}'\mathbf{C}^{22}\mathbf{x} \cdot \sigma_e^2}{\mathbf{x}'\mathbf{K}\mathbf{x} \cdot \sigma_u^2} \\
 &= 1 - \frac{\text{PEVD}(\mathbf{x})}{\mathbf{x}'\mathbf{K}\mathbf{x} \cdot \sigma_u^2}.
 \end{aligned}$$

A flow diagram showing the computational procedure is shown in Figure 3C.

Prediction error correlation (r)

Prediction error correlation, known as pairwise r statistic, between individuals i and j is calculated from the elements of the PEV matrix [12].

$$r_{ij} = \frac{\text{PEC}(\hat{u}_i, \hat{u}_j)}{\sqrt{\text{PEV}(\hat{u}_i) \cdot \text{PEV}(\hat{u}_j)}}.$$

Individual average r: The summary method based on individual average calculates pairwise r for all pairs of individuals followed by averaging all r measures across units.

$$r_{i'j'} = \frac{1}{n_{i'} \cdot n_{j'}} \cdot \sum \frac{\text{PEC}(\hat{u}_{i'}, \hat{u}_{j'})}{\sqrt{\text{PEV}(\hat{u}_{i'}) \cdot \text{PEV}(\hat{u}_{j'})}}.$$

This summary method for r statistic was used in Yu et al. [4] and calculation steps are shown in Figure 4A.

Group average r: This is known as flock connectedness in the literature, which calculates the ratio of PEV_{Mean} and PEC_{Mean} [3]. This group average connectedness for r between two units i' and j' is

given by the following equation.

$$\begin{aligned}
 r_{i'j'} &= \frac{\overline{\text{PEC}}_{i'j'}}{\sqrt{\text{PEV}_{i'i'} \cdot \text{PEV}_{j'j'}}} \\
 &= \frac{1/n_{i'} \sum \text{PEC}_{i'j'} 1/n_{j'}}{\sqrt{(1/n_{i'})^2 \sum \text{PEV}_{i'i'} \cdot (1/n_{j'})^2 \sum \text{PEV}_{j'j'}}} \\
 &= \frac{\sum \text{PEC}_{i'j'}}{\sqrt{\sum \text{PEV}_{i'i'} \cdot \sum \text{PEV}_{j'j'}}}. \tag{8}
 \end{aligned}$$

A graphical derivation is presented in Figure 4B.

Contrast r: A contrast of r is defined as below.

$$r(\mathbf{x}) = \mathbf{x}'\mathbf{r}\mathbf{x}.$$

A flow diagram illustrating a computational procedure is shown in Figure 4C.

Variance of differences in unit effects (VED)

A metric VED, which is a function of VE can be used to measure connectedness. All PEV-based metrics follow a two-step procedure in the sense that they first compute the PEV matrix at the individual level and then apply one of the summary methods to derive connectedness at the unit level. In contrast, VE-based metrics follow a single-step procedure such that we can obtain connectedness between units directly. Moreover, since the number of fixed effects is oftentimes smaller than the number of individuals in the model, the computational requirements for VED are expected to be lower [9]. Note that all VE-derived approaches can be classified based on the number of fixed effects to be corrected including no correction (0), correction for one fixed effect (1), and correction for two or more fixed effects (2) [9]. Below we discuss connectedness metrics that are derived from VED.

VED0: Using the summary method group average, the VED0 statistic [8] estimates PEVD alike connectedness with VE rather than PEV_{Mean} . We can obtain VED0 between two units i' and j' by

replacing PEV_{Mean} in equation (6) with VE0 defined in equation (1).

$$VED0_{i'j'} = VE0_{i'i'} + VE0_{j'j'} - 2VE0_{i'j'}, \quad (9)$$

VED1: A VED statistic that corrects for the presence of unit effect is obtained by replacing PEV_{Mean} in equation (6) with VE1 defined in equation (2). This corrects for the number of individuals in the units.

$$VED1_{i'j'} = VE1_{i'i'} + VE1_{j'j'} - 2VE1_{i'j'}, \quad (10)$$

VED2: Similarly, VED statistic based on VE2 is obtained by replacing PEV_{Mean} in equation (6) with VE2 defined in equation (3). This formula accounts for fixed effects other than unit effect.

$$VED2_{i'j'} = VE2_{i'i'} + VE2_{j'j'} - 2VE2_{i'j'}, \quad (11)$$

Coefficient of determination of VED

CDVED0: A CDVED0 statistic, which is a CD statistic based on VE0, is defined by replacing PEV_{Mean} in equation (7) with VE0. A pairwise CDVED0 between two units i' and j' is given by the following equation.

$$CDVED0_{i'j'} = 1 - \frac{VE0_{i'i'} + VE0_{j'j'} - 2VE0_{i'j'}}{\sigma_u^2 \cdot (\bar{\mathbf{K}}_{i'i'} + \bar{\mathbf{K}}_{j'j'} - 2\bar{\mathbf{K}}_{i'j'})}$$

CDVED1: CDVED1 is obtained by replacing PEV_{Mean} in equation (7) with VE1.

$$CDVED1_{i'j'} = 1 - \frac{VE1_{i'i'} + VE1_{j'j'} - 2VE1_{i'j'}}{\sigma_u^2 \cdot (\bar{\mathbf{K}}_{i'i'} + \bar{\mathbf{K}}_{j'j'} - 2\bar{\mathbf{K}}_{i'j'})}$$

CDVED2: Similarly, CDVED2 is obtained by replacing PEV_{Mean} in equation (7) with VE2.

$$CDVED2_{i'j'} = 1 - \frac{VE2_{i'i'} + VE2_{j'j'} - 2VE2_{i'j'}}{\sigma_u^2 \cdot (\bar{\mathbf{K}}_{i'i'} + \bar{\mathbf{K}}_{j'j'} - 2\bar{\mathbf{K}}_{i'j'})},$$

Connectedness rating (CR)

CR0: A CR statistic first proposed by Mathur et al. [13] is similar to equation (8). However, it uses variances and covariances of estimated unit effects. Specifically, we replace PEV_{Mean} with VE0, and CR0 between two units i' and j' is given by the following equation.

$$CR0_{i'j'} = \frac{VE0_{i'j'}}{\sqrt{VE0_{i'i'} \cdot VE0_{j'j'}}}.$$

CR1: A CR1 statistic is obtained by replacing PEV_{Mean} in equation (8) with VE1.

$$CR1_{i'j'} = \frac{VE1_{i'j'}}{\sqrt{VE1_{i'i'} \cdot VE1_{j'j'}}},$$

CR2: In the same manner, a CR2 statistic is obtained by replacing PEV_{Mean} in equation (8) with VE2.

$$CR2_{i'j'} = \frac{VE2_{i'j'}}{\sqrt{VE2_{i'i'} \cdot VE2_{j'j'}}},$$

Software Description

Overview of software architecture

The GCA R package is implemented entirely in R, which is an open source programming language and environment for performing statistical computing [14]. The package is hosted on a GitHub page accompanied by a detailed vignette document. Computational speed was improved by integrating C++ code into R code using the Rcpp package [15]. The initial versions of the algorithms and the R code were used in previous studies [4–6] and were enhanced further for efficiency, usability, and documentation in the current version to facilitate connectedness analysis. The GCA R package provides a comprehensive and effective tool for genetic connectedness analysis and whole genome prediction, which further contributes to the genetic evaluation and prediction.

Installing the GCA Package

The current version of the GCA R package is available at GitHub (<https://github.com/HaipengU/GCA>). The package can be installed using the devtools R package [16] and loaded into the R environment.

Box 1: Installing the GCA Package

```
install.packages("devtools")  
  
library(devtools)  
  
install_github('HaipengU/GCA')  
  
library(GCA)
```

Simulated data

We simulated a cattle data set using QMSim software [17] to illustrate the usage of GCA package. This data set is included in the package as an example data set. A total of 2,500 cattle spanning

five generations were simulated with pedigree and genomic information available for all individuals. We simulated 10,000 evenly distributed biallelic single nucleotide polymorphisms and 2,000 randomly distributed quantitative trait loci (QTL) across 29 pairs of autosomes with 100 cM per chromosome. A single phenotype with a heritability of 0.6 and a fixed covariate of sex were simulated. This was followed by simulating units using the k-medoid algorithm [18] coupled with the dissimilarity matrix derived from a numerator relationship matrix as shown in previous studies [4-6]. The data set was stored as an R object in the package.

Box 2: Loading the data

```
data(package = 'GCA')$results[, "Item"] # list all data files in the GCA package  
data(GCCattle) # load the data  
dim(cattle.pheno) # phenotype and fixed effects  
dim(cattle.W) # marker matrix
```

The genotype object is a $2,500 \times 10,000$ marker matrix. The phenotype object is a $2,500 \times 6$ matrix, including the columns of progeny, sire, dam, sex, unit, and phenotype.

Application of GCA Package

Below we show the usage of the main function `gca` followed by some specific examples using CD. Box 3 lists all input arguments for the `gca` function.

- `Kmatrix`: Genetic relationship matrix constructed from either pedigree or genomics.
- `Xmatrix`: Fixed effects incidence matrix excluding intercept. The first column of the `Xmatrix` should start with unit effects followed by other fixed effects if applicable.
- `sigma2a` and `sigma2e`: Estimates of additive genetic and residual variances, respectively.
- `MUScenario`: A vector of fixed factor units.
- `statistic`: Choice of connectedness statistic. Available options include

- 1 PEV-derived functions: PEVD_IdAve, PEVD_GrpAve, PEVD_contrast, CD_IdAve, CD_GrpAve, CD_contrast, r_IdAve, r_GrpAve, and r_contrast
 - 2 VE-derived functions: VED0, VED1, VED2, CDVED0, CDVED1, CDVED2, CR0, CR1, and CR2.
- NumofMU: Return either pairwise unit connectedness (Pairwise) or overall connectedness across all units (Overall).
 - Uidx: An integer indicating the last column number of units in the Xmatrix. This Uidx is required for VED2, CDVED2, and CR2 statistics. The default is NULL.
 - scale (logical): Should sigma2a be used to scale statistic (i.e., PEVD_IdAve, PEVD_GrpAve, PEVD_contrast, VED0, VED1, and VED2) so that connectedness is independent of measurement unit? The default is TRUE.
 - diag (logical): Should the diagonal elements of the PEV matrix (i.e., PEVD_GrpAve, CD_GrpAve, and r_GrpAve) or the K matrix (CDVED0, CDVED1, and CDVED2) be included? The default is TRUE.

Box 3: A list of input arguments for the `gca` function

```
gca(Kmatrix, Xmatrix, sigma2a, sigma2e, MUScenario,  
statistic, NumofMU, Uidx = NULL, scale = TRUE, diag = TRUE)
```

Example 1: Pairwise connectedness across units

The following example demonstrates the pairwise CD_IdAve across units with no additional fixed effect.

Box 4: Example of pairwise CD_IdAve across units

```
X_fixed <- model.matrix(~ -1 + factor(cattle.pheno$Unit)) # incidence matrix of units
G <- computeG(cattle.W) # genomic relationship matrix
sigma2a <- 0.6 # additive genetic variance
sigma2e <- 0.4 # residual variance
CD_IdAve <- gca(Kmatrix = G, Xmatrix = X_fixed, sigma2a = sigma2a, sigma2e = sigma2e,
MUScenario = as.factor(cattle.pheno$Unit), statistic = 'CD_IdAve', NumofMU = 'Pairwise')
```

Here, the 'X_fixed' is the incidence matrix of units with the intercept excluded. The 'G' is the first type of genomic relationship matrix in VanRaden [19]. The statistic 'CD_IdAve' calculates CD measures using individual average as a summary method. The option 'Pairwise' in the 'NumofMU' argument returns a square matrix containing pairwise connectedness measures across units.

Example 2: Overall connectedness across units

We present the calculation of overall CD_GrpAve measures across units by changing the argument of 'statistic' to 'CDGrpAve' in this example. The CD statistic is summarized at the unit level using PEV_{Mean} and PEC_{Mean} . Changing the argument 'NumofMU' to 'Overall' returns the average of all pairwise connectedness measures between units. The definitions of other arguments are identical as shown in Box 4.

Box 5: Example of overall CD_GrpAve across units

```
CD_GrpAve <- gca(Kmatrix = G, Xmatrix = X_fixed, sigma2a = sigma2a, sigma2e =
sigma2e, MUScenario = as.factor(cattle.pheno$Unit), statistic = 'CD_GrpAve', NumofMU =
'Overall')
```

Example 3: Pairwise connectedness across units with fixed effects of units and sex

The following example shows the pairwise connectedness of CDVED2 while correcting for two fixed effects, namely units and sex.

Box 6: Example of pairwise CDVED2 across units

```
X_fixed <- model.matrix(~ -1 + factor(cattle.pheno$Unit)
+ factor(cattle.pheno$Sex)) # incidence matrix of units and sex

G <- computeG(cattle.W) # genomic relationship matrix

sigma2a <- 0.6 # additive genetic variance

sigma2e <- 0.4 # residual variance

CDVED2 <- gca(Kmatrix = G, Xmatrix = X_fixed, sigma2a = sigma2a, sigma2e = sigma2e,
MUScenario = as.factor(cattle.pheno$Unit), statistic = 'CDVED2', NumofMU = 'Pairwise',
Uidx = 8)
```

This code returns CD measures based on VE2.

Conclusions

The GCA R package provides users with a comprehensive tool for analysis of genetic connectedness using pedigree and genomic data. The users can easily assess the connectedness of their data and be mindful of the uncertainty associated with comparing genetic values of individuals involving different management units or contemporary groups. Moreover, the GCA package can be used to measure the level of connectedness between training and testing sets in the whole-genome prediction paradigm. This parameter can be used as a criterion for optimizing the training data set. In summary, we contend that the availability of the GCA package to calculate connectedness allows breeders and geneticists to make better decisions on comparing individuals in genetic evaluations and inferring linkage between any pair of individual groups in genomic prediction.

Availability and implementation

The GCA R source code is provided as free and open source. The webpage <https://github.com/HaipengU/GCA> was created as a nexus of all genetic connectedness related functions and examples available in the GCA R package. The vignette is available at <https://haipengu.github.io/Rmd/Vignette.html>.

Declarations

Funding

This work was supported in part by Virginia Polytechnic Institute and State University startup funds to GM.

Authors' contributions

HY and GM developed the software tool and wrote the manuscript. GM supervised and directed the study. All authors read and approved the manuscript.

Ethics approval and consent to participate

Not applicable.

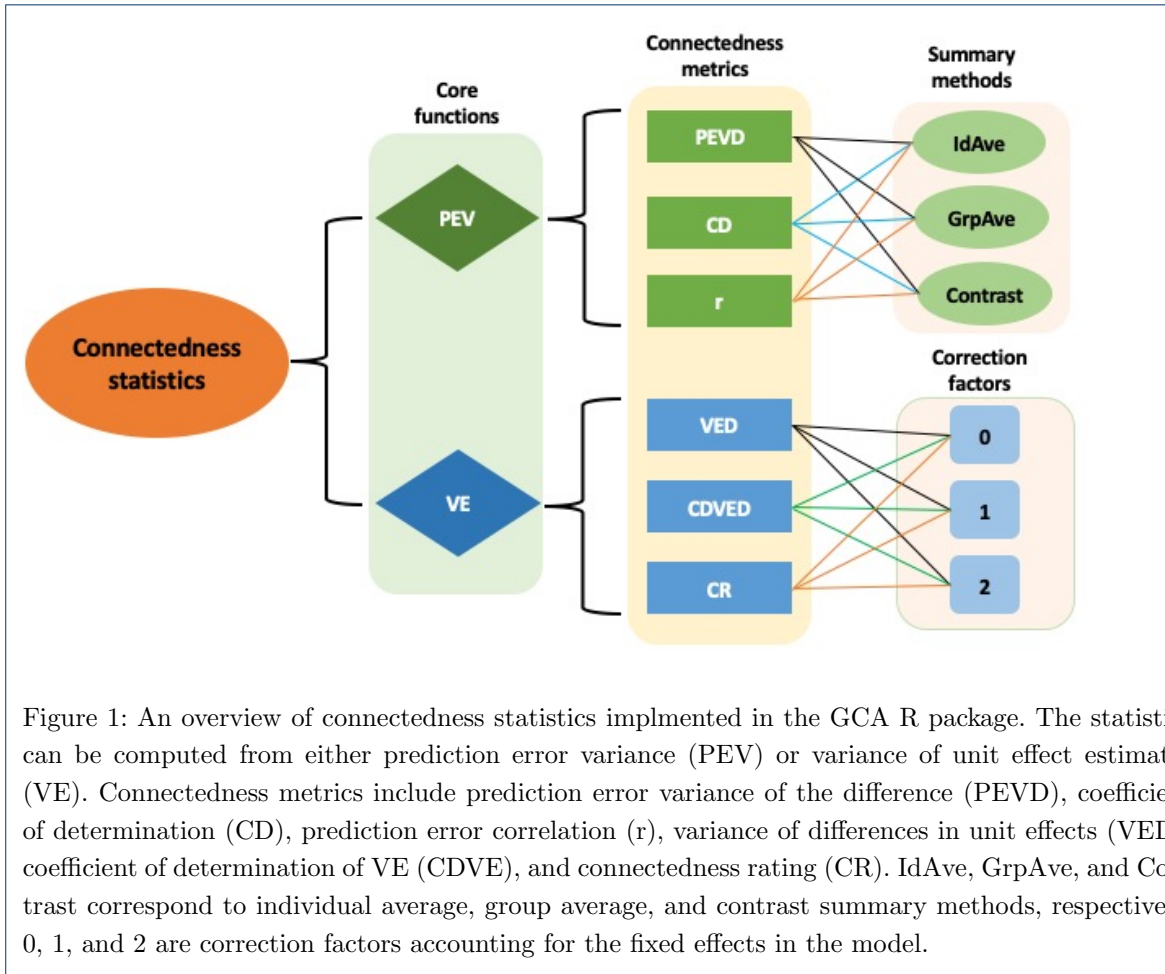
Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Figures



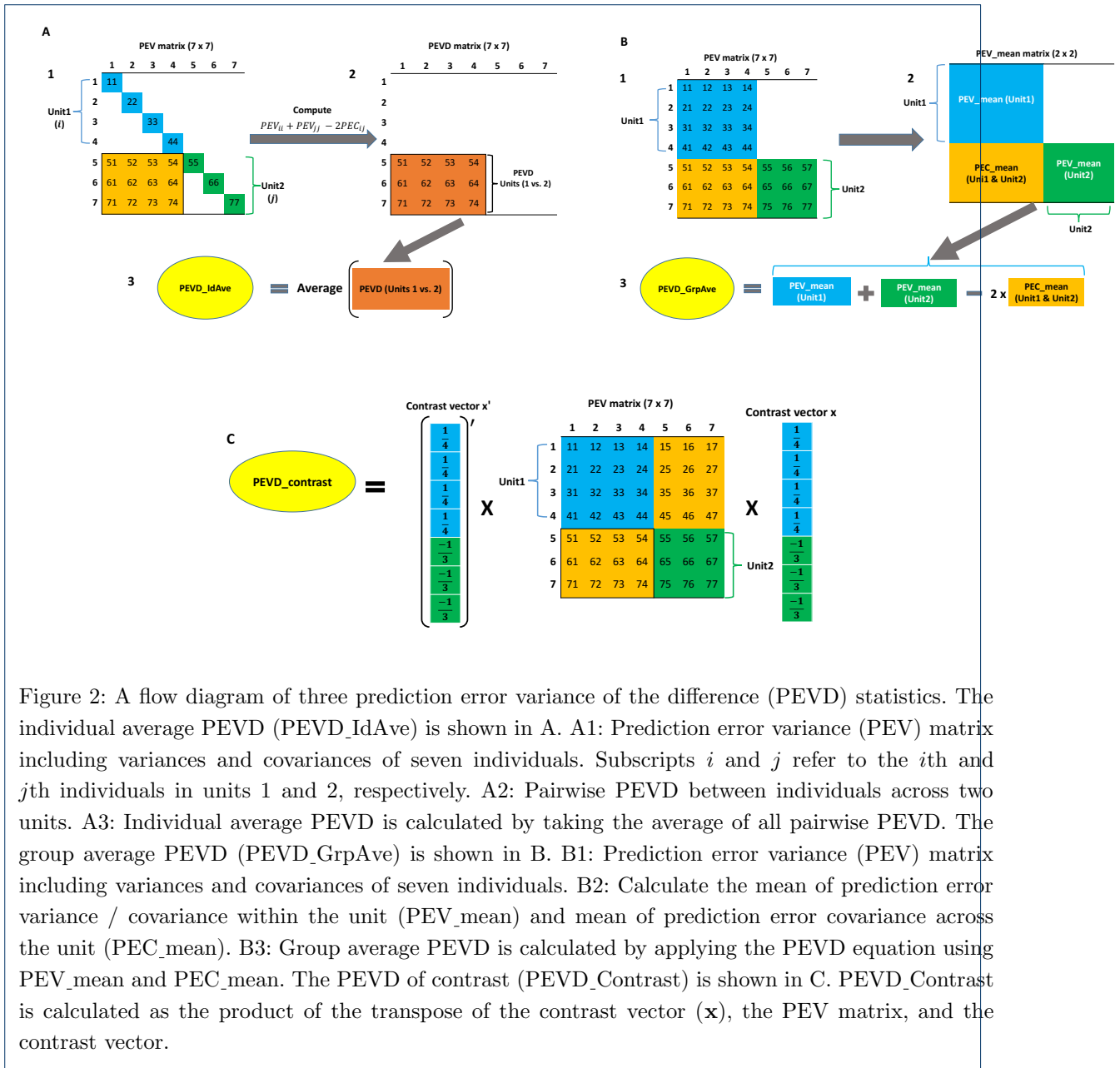


Figure 2: A flow diagram of three prediction error variance of the difference (PEVD) statistics. The individual average PEVD (PEVD_IdAve) is shown in A. A1: Prediction error variance (PEV) matrix including variances and covariances of seven individuals. Subscripts i and j refer to the i th and j th individuals in units 1 and 2, respectively. A2: Pairwise PEVD between individuals across two units. A3: Individual average PEVD is calculated by taking the average of all pairwise PEVD. The group average PEVD (PEVD_GrpAve) is shown in B. B1: Prediction error variance (PEV) matrix including variances and covariances of seven individuals. B2: Calculate the mean of prediction error variance / covariance within the unit (PEV_mean) and mean of prediction error covariance across the unit (PEC_mean). B3: Group average PEVD is calculated by applying the PEVD equation using PEV_mean and PEC_mean. The PEVD of contrast (PEVD_Contrast) is shown in C. PEVD_Contrast is calculated as the product of the transpose of the contrast vector (\mathbf{x}'), the PEV matrix, and the contrast vector.

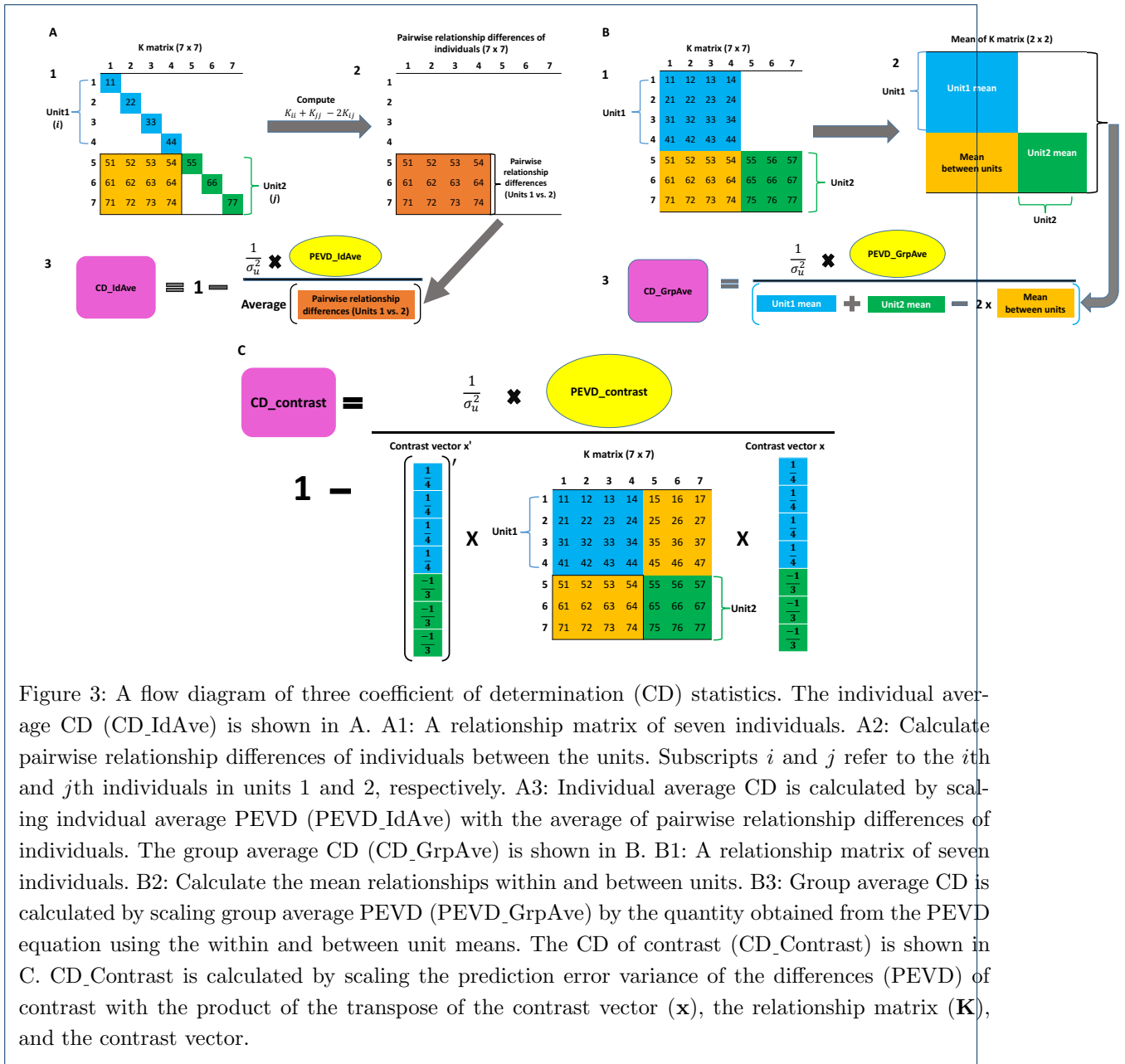


Figure 3: A flow diagram of three coefficient of determination (CD) statistics. The individual average CD (CD_IdAve) is shown in A. A1: A relationship matrix of seven individuals. A2: Calculate pairwise relationship differences of individuals between the units. Subscripts i and j refer to the i th and j th individuals in units 1 and 2, respectively. A3: Individual average CD is calculated by scaling individual average PEVD (PEVD_IdAve) with the average of pairwise relationship differences of individuals. The group average CD (CD_GrpAve) is shown in B. B1: A relationship matrix of seven individuals. B2: Calculate the mean relationships within and between units. B3: Group average CD is calculated by scaling group average PEVD (PEVD_GrpAve) by the quantity obtained from the PEVD equation using the within and between unit means. The CD of contrast (CD_Contrast) is shown in C. CD_Contrast is calculated by scaling the prediction error variance of the differences (PEVD) of contrast with the product of the transpose of the contrast vector (\mathbf{x}), the relationship matrix (\mathbf{K}), and the contrast vector.

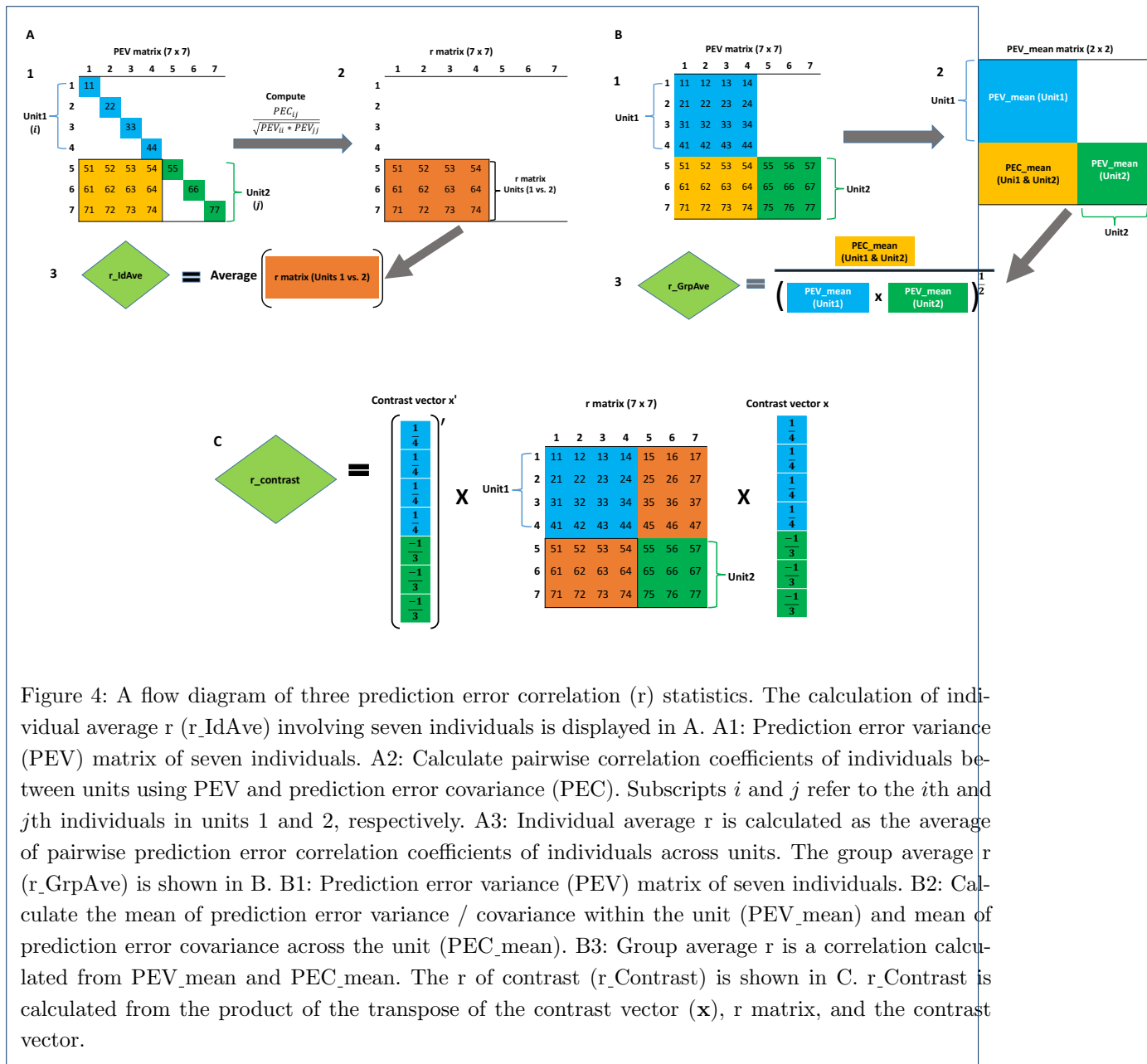


Figure 4: A flow diagram of three prediction error correlation (r) statistics. The calculation of individual average r (r_IdAve) involving seven individuals is displayed in A. A1: Prediction error variance (PEV) matrix of seven individuals. A2: Calculate pairwise correlation coefficients of individuals between units using PEV and prediction error covariance (PEC). Subscripts i and j refer to the i th and j th individuals in units 1 and 2, respectively. A3: Individual average r is calculated as the average of pairwise prediction error correlation coefficients of individuals across units. The group average r (r_GrpAve) is shown in B. B1: Prediction error variance (PEV) matrix of seven individuals. B2: Calculate the mean of prediction error variance / covariance within the unit (PEV_mean) and mean of prediction error covariance across the unit (PEC_mean). B3: Group average r is a correlation calculated from PEV_mean and PEC_mean. The r of contrast (r_Contrast) is shown in C. r_Contrast is calculated from the product of the transpose of the contrast vector (x'), r matrix, and the contrast vector.

Author details

References

1. Foulley, J., Schaeffer, L., Song, H., Wilton, J.: Progeny group size in an organized progeny test program of ai beef bulls using reference sires. *Canadian Journal of Animal Science* **63**(1), 17–26 (1983)
2. Foulley, J., Bouix, J., Goffinet, B., *et al.*: Connectedness in genetic evaluation. In: *Advances in Statistical Methods for Genetic Improvement of Livestock*, pp. 277–308. Springer, ??? (1990)
3. Kuehn, L., Notter, D., Nieuwhof, G., Lewis, R.: Changes in connectedness over time in alternative sheep sire referencing schemes. *Journal of Animal Science* **86**(3), 536–544 (2008)
4. Yu, H., Spangler, M.L., Lewis, R.M., Morota, G.: Genomic relatedness strengthens genetic connectedness across management units. *G3: Genes, Genomes, Genetics* **7**(10), 3543–3556 (2017)
5. Yu, H., Spangler, M.L., Lewis, R.M., Morota, G.: Do stronger measures of genomic connectedness enhance prediction accuracies across management units? *Journal of Animal Science* **96**(11), 4490–4500 (2018)
6. Momen, M., Morota, G.: Quantifying genomic connectedness and prediction accuracy from additive and non-additive gene actions. *Genetics Selection Evolution* **50**(1), 45 (2018)
7. Henderson, C.R.: *Applications of Linear Models in Animal Breeding*. University of Guelph, Third edition, Edited by Schaeffer LR. Guelph (1984)
8. Kennedy, B., Trus, D.: Considerations on genetic connectedness between management units under an animal model. *Journal of Animal Science* **71**(9), 2341–2352 (1993)
9. Holmes, J.B., Dodds, K.G., Lee, M.A.: Estimation of genetic connectedness diagnostics based on prediction errors without the prediction error variance–covariance matrix. *Genetics Selection Evolution* **49**(1), 29 (2017)
10. Laloë, D.: Precision and information in linear models of genetic evaluation. *Genetics Selection Evolution* **25**(6), 557 (1993)
11. Laloë, D., Phocas, F., Menissier, F.: Considerations on measures of precision and connectedness in mixed linear models of genetic evaluation. *Genetics Selection Evolution* **28**(4), 359 (1996)
12. Lewis, R., Crump, R., Simm, G., Thompson, R.: Assessing connectedness in across-flock genetic evaluations. *Proc. Brit. Soc. Anim. Sci* **121** (1999)
13. Mathur, P., Sullivan, B., Chesnais, J.: Measuring connectedness: concept and application to a large industry breeding program. In: *Proc. 7th World Congr. Genet. Appl. to Livest. Prod.*, vol. 19, p. 23 (2002)
14. R Core Team: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria (2019). R Foundation for Statistical Computing. <https://www.R-project.org/>
15. Eddelbuettel, D., François, R.: Rcpp: Seamless R and C++ integration. *Journal of Statistical Software* **40**(8), 1–18 (2011). doi:[10.18637/jss.v040.i08](https://doi.org/10.18637/jss.v040.i08)
16. Wickham, H., Hester, J., Chang, W., Hester, M.J.: Package 'devtools' (2019)
17. Sargolzaei, M., Schenkel, F.S.: Qmsim: a large-scale genome simulator for livestock. *Bioinformatics* **25**(5), 680–681 (2009)
18. Kaufman, P. L. & Rousseeuw: *Finding Groups in Data: an Introduction to Cluster Analysis*. John Wiley and Sons, New York (1990)
19. VanRaden, P.M.: Efficient methods to compute genomic predictions. *Journal of Dairy Science* **91**(11), 4414–4423 (2008)