

RESEARCH

Grouping of genomic markers in populations with family structure

Dörte Wittenburg*, Michael Doschoris and Jan Klosa

*Correspondence:

wittenburg@fhn-dummerstorf.de

Leibniz Institute for Farm Animal
Biology, Institute of Genetics and
Biometry, 18196 Dummerstorf,
Germany

Full list of author information is
available at the end of the article

Abstract

Background: Linkage and linkage disequilibrium (LD) between genome regions cause dependencies among genomic markers. Due to family stratification in populations with non-random mating in livestock or crop, the standard measures of population LD such as r^2 may be biased. Grouping of markers according to their interdependence needs to account for the actual population structure in order to allow proper inference in genome-based evaluations.

Results: Given a matrix reflecting the strength of association between markers, groups are built successively using a greedy algorithm; largest groups are built first. As an option, a representative marker is selected for each group. We provide an implementation of the grouping approach as a new function to the R package `hscovar`. This package enables the calculation of the theoretical covariance between biallelic markers for half- or full-sib families and the calculation of representative markers. In case studies, we have shown that the number of groups comprising dependent markers was smaller and representative SNPs were spread more uniformly over the investigated chromosome region when the family stratification was respected compared to a population-LD approach. In a simulation study, we observed that sensitivity and specificity of a genome-based association study improved if selection of representative markers took family structure into account.

Conclusions: Chromosome segments which frequently recombine in the underlying population can be identified from the matrix of pairwise dependence between markers. Representative markers can be exploited, for instance, for dimension reduction prior to a genome-based association study or the grouping structure itself can be employed in a grouped penalization approach.

Keywords: Single nucleotide polymorphism; Covariance matrix; Clustering; TagSNP; Group lasso; SNP-BLUP

¹**Background**

²Genomic markers are an invaluable source for characterizing genetic variety and to
³elucidate the relationship between genetic and phenotypic variation in breeding pop-
⁴ulations. Dependencies among genomic markers are caused by linkage and linkage
⁵disequilibrium (LD) between genome regions. Though this condition complicates
⁶investigations on which genetic variants are truly associated with trait expression
⁷[1], dependencies can be advantageous for grouping of markers. For example, clus-
⁸tering based on a greedy algorithm [2], hierarchical clustering (e.g., [3]) or grouping
⁹via interval-graph modeling [4] exploit the presence of LD blocks which are regions
¹⁰of particularly high correlation. To allow for proper inferences of such approaches, a
¹¹suitable measure for the strength of dependence is needed. For instance, measuring
¹²LD in terms of r^2 [5] is a natural choice but it is meaningful only for popula-
¹³tions without stratification. In livestock and crop breeding, however, populations
¹⁴are often characterized by strong family stratification due to non-random mating
¹⁵of selected individuals. As examples, large paternal half-sib families are typical for
¹⁶cattle populations whereas chicken or fish populations consist of full-sib families.
¹⁷In plant breeding, maternal half-sib families are often produced in, for instance,
¹⁸wheat and clover. Then, linkage between markers within family leads to haplotype
¹⁹frequencies among progeny that are not conclusive for estimating r^2 . Hence, there
²⁰is need to promote measures of marker dependence which takes into account the
²¹particular family structure.

²²Especially in situations of ultra-dense panels of single nucleotide polymorphisms
²³(SNPs), it is often sufficient to investigate representative SNPs (“tagSNPs”) out
²⁴of each cluster. This subset can help identifying trait-associated genome regions
²⁵in genome-wide association studies and allows comparing genome characteristics
²⁶between ethnics/species/breeds (e.g., [2]). As the choice of tagSNPs is a consequence
²⁷of grouping, it is also influenced by the underlying population structure.

²⁸The objective of this paper is to exploit the family structure of a population for
²⁹specifying groups of associated markers. We generalize the grouping approach of
³⁰Carlson *et al.* [2] in order to allow binning of markers given a correlation matrix or
³¹any kind of similarity matrix with scaled entries in [0, 1]. We investigate three case
³²studies and a simulation study. For each case study, we visually inspect the correla-
³³tion matrix and link to the outcome of grouping. Usability for genome-based asso-

ciation studies is shown as one possible field of application. Results were compared¹
 to the commonly used population-LD approach which ignores family structure. We²
 provide a new function to the R package `hscovar` (available at CRAN) that enables³
 grouping of markers and selection of representative markers.⁴

6 **Methods**

The dependence between pairs of SNPs, each with two alleles A and B, can be⁷
 expressed in terms of a covariance or correlation matrix. It has already been shown⁸
 in the literature how to calculate the theoretical covariance between markers in⁹
 a population consisting of half-sib families [1]. It requires a genetic map, haplo-¹⁰
 types of the common parent and LD information (or haplotype frequencies) of the¹¹
 population the individual parent comes from. This approach can be extended to be¹²
 applicable to full-sib families by adding the paternal and maternal contribution into¹³
 a single covariance matrix; the derivation is summarized in Additional file 1. Hence,¹⁴
 a covariance matrix can be derived for any family structure, and this constitutes¹⁵
 the input of the following grouping approach.¹⁶

18 **Grouping of markers**

We generalized the strategy of Carlson *et al.* [2] for binning markers and selecting¹⁹
 representatives to be applicable to any symmetric matrix which reflects a measure²⁰
 of dependence between markers and has entries scaled in [0, 1]. In particular, we²¹
 considered the correlation matrix R . The idea is that SNPs which are associated to²²
 each other are assigned to groups. Groups are built one after the other - the largest²³
 group at first. For $b = 1, 2, \dots$, the b -th group is identified by searching for the SNP²⁴
 that has most occurrences of absolute correlation to other SNPs larger than a given²⁵
 threshold t . More precisely, let S_b denote the set of SNP indices which have not²⁶
 been binned yet. Then, for each SNP $k \in S_b$, the set of highly associated SNPs is²⁷
 determined as (Step A)²⁸

$$29 \quad C_k = \{l \mid l \in S_b : |R_{k,l}| > t\}, \quad 29$$

31 and a set with highest cardinality (operator #) is chosen,³¹
 32

$$33 \quad c \in \arg \max_k \#C_k. \quad 33$$

¹Thus, C_c constitutes the b -th group, and a tagSNP is selected from this group. SNP¹
² c has strong correlation with any other SNP in C_c but it can happen that also other²
³SNPs of C_c fulfill this criterion. Hence, a set of candidates is given by (Step B) ³

$$T = \{k \mid k \in C_c \wedge \forall l \in C_c : |R_{k,l}| > t\}.$$

⁴
⁵
⁶
⁷If more than one candidate remains, then the $|T|/2$ -th SNP becomes the repre-⁷
⁸sentative of group b . A next round of iteration is started using $S_{b+1} = S_b \setminus C_c$ until^{8,9} S_{b+1}
⁹ $= \emptyset$. The number of groups only depends on the threshold t . Similar to [2],^{9,10} $t =$
¹⁰ 0.8 is a suitable value. In an extreme case (with t approaching 1), each SNP^{10,11} builds a
¹¹single group, yielding a complexity of this algorithm of $O(p^2)$, with p the^{11,12} total
¹²number of SNPs. This approach is implemented as function tagSNP in the¹² ¹³R
¹³package hscovar. A graphical representation of this algorithm is shown in Figure¹³

¹⁴¹. ¹⁴

¹⁵ Evaluation ¹⁵

¹⁶ It is an obvious choice to compare the family approach with a population-LD ap-¹⁶
¹⁷proach. Such an approach requires the population frequency of the different hap-¹⁷
¹⁸lotypes ($f_{A-A}, f_{A-B}, f_{B-A}, f_{B-B}$) from which LD between markers is computed ¹⁸
¹⁹in terms of r^2 according to [5]. For any marker pair k, l with allele frequencies ¹⁹
²⁰ $f_k = f_{A-A} + f_{A-B}$ and $f_l = f_{A-A} + f_{B-A}$, we have ²⁰

$$r_{k,l}^2 = \frac{(f_{A-A}f_{B-B} - f_{A-B}f_{B-A})^2}{f_k(1-f_k)f_l(1-f_l)}.$$

²¹
²²
²³
²⁴ The LD matrix can be computed from progeny genotypes using the function ld
²⁵from the R package snpStats version 1.38.0 [6]. This function undertakes phasing
²⁶of genotypes using a maximum-likelihood approach [7]. Based on the LD matrix
²⁷containing r^2 , representative markers can be derived as described above. ²⁷
²⁸Family and population-LD approach were compared using the Calinski-Harabasz
²⁹(CH) index [8] which measures the cluster quality with respect to inter- and intra-
³⁰cluster distances. The method with higher CH index performed better. For this, the
³¹function calinhara from the R package fpc version 2.2-9 was applied to the groups
³²obtained; the quality referred to distances based on the genotype matrix centered
³³within family and scaled (as described below). Furthermore, we present the pure ³³

number of groups and highlight those groups with group size of at least three. This¹
 also helped visualizing the location of corresponding representative markers.²

Selecting a representative set of markers is a natural tool for dimension reduc-³
 tion prior to genomic evaluations in order to reduce the impact of multicollinearity⁴
 among predictor variables. Representative SNPs capture cumulative effects of the⁵
 corresponding LD blocks on trait expression. We investigated a SNP-BLUP ap-⁶
 proach, which is widely used in genomic evaluations (e.g., [9]), and thereby demon-⁷
 strate one possible application of the suggested approach. Representative markers⁸
 selected from the family or from the population-LD approach were employed as⁹
 predictor variables in a regression model¹⁰

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

with $\mathbf{y} = (y_1, \dots, y_n)^T$ the phenotype vector, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_\tau)^T$ the vector of genetic¹⁴
 effects captured by τ tagSNPs, the corresponding design matrix \mathbf{X} with dimensions¹⁵
 $16n \times \tau$ including the genotype codes in terms of major allele counts. The columns of¹⁶
 \mathbf{X} and the vector \mathbf{y} were centered within family and scaled to obtain an empirical¹⁷
 variance of one. The residual errors were assumed to be independently and normally¹⁸
 distributed. For convenience, no other effects were assumed. We used the R package¹⁹
 asreml version 3.0 [10] to estimate the vector of regression coefficients as²⁰

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$$

where the shrinkage parameter λ was estimated via AI-REML. Significance of the²⁴
 k -th SNP effect was tested by a t -like test statistic as in [1],²⁵

$$T_k = \frac{\hat{\beta}_k}{SD(\hat{\beta}_k)}.$$

Significance was reported if $T_k \geq q_{1-\alpha/2}$ or $T_k < q_{\alpha/2}$ using the $1 - \alpha/2$ and²⁹
 $\alpha/2$ quantile of the standard normal distribution. The SNP-BLUP approach was³⁰
 evaluated in terms of sensitivity (i.e., true-positive rate) and specificity (i.e., $1 -$ ³¹
 false-positive rate) over a range of type-I error α . We additionally verified the im-³²
 pact of threshold $t \in \{0.5, 0.6, 0.7, 0.8\}$ on grouping and its consequences on the³³
 performance of the SNP-BLUP approach.

¹**Data**

²The data sets used for studying dependencies between SNP markers differed in SNP²
³density and family structure. They covered a range of mean inter-marker distances³
⁴from 0.003 cM to 0.23 cM. The case study of mouse data was based on low-density⁴
⁵genotypes of full-sib families; progeny and parents were genotyped. The case study of⁵
⁶cattle data comprised medium-density genotypes of half-sib families with genotyped⁶
⁷progeny only. Furthermore, medium-density SNP data were available for full-sib⁷
⁸families in maize. High-density genotype data of half-sib families were generated⁸
⁹in simulations. For evaluation, the SNP-BLUP approach was applied to simulated⁹
¹⁰data only. Unless otherwise stated, computations were done using R version 4.0.3¹⁰
¹¹[11] and $t = 0.8$; all scripts are included as Additional files. 11

12

¹³**Mouse data**

¹⁴Genotype data of a heterogeneous stock of mice were available from [https://wp.¹⁴
¹⁵cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/](https://wp.cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/). We investigated chro-¹⁵
¹⁶mosome 17 because it harbors the highly recombining MHC region which affects sev-¹⁶
¹⁷eral immunological traits [12]. The chromosome data consisted of genotypes at 394¹⁷
¹⁸SNPs of 2002 individuals. After filtering for individual call rate $\geq 90\%$, 1998 geno-¹⁸
¹⁹typed individuals remained comprising 1759 progeny, 120 fathers and 119 mothers.¹⁹
²⁰In total, 138 full-sib families (family size ranged from 1 to 47) could be identified.²⁰
²¹The SNP call rate was $\geq 90\%$. All genotype data were phased with Beagle version²¹
²²5.1 [13] and parental haplotypes were selected to set up the correlation matrix. As-²²
²³suming a 1:1 relationship between physical (Build37 genome assembly) and genetic²³
²⁴distance of adjacent markers, the genetic length was 91 cM. SNP alleles were coded²⁴
²⁵in terms of the major allele in the given sample. The population-LD matrix was²⁵
²⁶calculated from progeny genotypes. 26

27

²⁷**Cattle data**

²⁸Genotype data of Holstein cattle were available from RADAR [https://dx.doi.²⁸
²⁹org/10.22000/280](https://dx.doi.org/10.22000/280). The data comprised 50K SNP-chip data of five half-sib fam-²⁹
³⁰ilies with $n = 265$ progeny in total; the family size ranged from 32 to 106. A 30
³¹chromosome window containing 300 SNPs was selected from BTA1. Based on the 31
³²physical ordering of markers according to the genome assembly ARS-UCD1.2, this 32
³³region corresponded to 20.59 – 39.44 Mbp. The haplotypes of sires were imputed from 33

¹progeny genotypes using the R package *hsphase* version 2.0.2 [14]. Maternal LD¹
²and paternal recombination rates between SNP pairs were estimated according to²
³Hampel *et al.* [15]. However, we used a 1:1 relationship between physical (Mbp)³
⁴and genetic positions (cM) for convenience; the genetic length of this window was⁴
⁵19 cM. Sire haplotypes and maternal LD were also part of the RADAR data set.⁵
⁶SNP alleles were coded in terms of the major maternal allele among progeny.⁶

⁸Maize data⁸

⁹Raw marker data were available from NCBI GEO database under Accession Num-⁹
¹⁰ber GSE50558, accompanied with physical coordinates corresponding to the genome¹⁰
¹¹assembly B73. The data set contained two maize panels, Flint and Dent, for which¹¹
¹²about 50K SNPs have been assessed in order to estimate recombination activity in¹²
¹³different maize populations [16]. We arbitrarily chose the Flint panel and chromo-¹³
¹⁴somes 2 for further analysis. In this panel, 13 full-sib families have been obtained by¹⁴
¹⁵crossing an inbred “central” line and several inbred “founder” lines. Double haploid,¹⁵
¹⁶(DH) lines have been derived from the F1 plants. This procedure allowed for study-¹⁶
¹⁷ing maternal meioses only. A cM:Mbp ratio of 0.80 was reported for Flint [16]; the¹⁷
¹⁸genetic length of chromosome 2 was approximately 188 cM. SNP genotypes of DH¹⁸
¹⁹progeny being heterozygous were set to missing value. After filtering the data for¹⁹
²⁰SNP and individual call rate $\geq 90\%$, $n = 1\,248$ out of 1 262 DH progeny and 1 447²⁰
²¹out of 2 030 SNPs remained. Rarely missing marker information of DH progeny were²¹
²²imputed by sampling the homozygous genotypes according to their frequencies. Af-²²
²³terwards loci with minor allele frequency less than 5% were discarded, yielding 956²³
²⁴SNPs. As haplotypes of DH progeny were given with certainty, the population-LD²⁴
²⁵matrix was set up directly using the squared Spearman correlation between SNPs²⁵
²⁶based on haploid data. The family approach solely considered the female part of²⁶
²⁷the covariance between SNPs. The haplotypes of F1 individuals were inferred from²⁷
²⁸the marker data of inbred lines. SNP alleles were coded according to central-line²⁸
²⁹origin.²⁹

³⁰Simulated data³⁰

³¹The simulation study resembled the population structure of a dairy cattle popula-³¹
³²tion. The setup of simulation design is fully described in [1]. Briefly, we considered³²
³³ $N = 1, 5, 10$ sires of half siblings. The overall number of progeny was $n = 1\,000$ ³³

¹equally partitioned into half-sib families. Quantitative traits were simulated which¹
²were influenced by 2 and 5 QTLs with equal effect sizes. QTLs contributed 30 %²
³to the trait variation (i.e., heritability 0.3). In total, 300 SNPs were simulated on³
⁴a chunk of DNA with 1 cM length. The data were generated using the R package⁴
⁵AlphaSimR version 0.13.0 [17]. SNP alleles were recoded in terms of the major allele⁵
⁶in the founder population. The simulation was repeated 100 times. For assessing⁶
⁷the SNP-BLUP approach, a window of 0.05 cM to both sides of a simulated QTL⁷
⁸was accepted as true-positive result. 8

¹⁰**Results and discussion** 10

¹¹**Case studies** 11

¹²We have shown applicability of the suggested software tool to empirical data. Espe-
¹³cially for the mouse data consisting of 138 genotyped full-sib families, the correlation
¹⁴matrix gave a clear representation of genomic regions with high or low interdepen-
¹⁵dence. In contrast to a population-LD approach, which did not account for family
¹⁶stratification, it was also possible to identify regions with positive or negative re-
¹⁷lationship. The population-LD approach exposed a wide region (13.82 – 21.13 Mbp)
¹⁸that had a strong association with the entire chromosome 17 shown as a red band
¹⁹in Figure 2b. With the family approach, this region revealed only high positive in-
²⁰terdependence with almost no impact on the remaining chromosome, see Figure 2a.
²¹Also Carlson *et al.* [2] reported that population stratification may generate artifac-
²²tual LD and hence makes an LD-selection algorithm sensitive. Moreover, a highly
²³fragmented region appeared in the range of 27.59 – 45.88 Mbp that overlaps MHC re-
²⁴gions 1 and 2 (https://www.ncbi.nlm.nih.gov/assembly/GCF_000001635.18/).
²⁵This was caused by high variation of parental haplotypes. Though the total num-
²⁶ber of groups was smaller with the family approach than with the population-LD
²⁷approach (83 vs. 98), the number of large groups, i.e., with group size ≥ 3 , was
²⁸almost equal (49 vs. 51), see Table 1. The CH index was about two times higher
²⁹with the family approach. Figure 2 shows that tagSNPs of groups containing at
³⁰least three markers were distributed more uniformly over the chromosome with the
³¹family approach than with the population-LD approach. The median of distances
³²between all tagSNPs was 0.70 cM and 0.52 cM in the family and population-LD
³³approach, respectively. 33

¹ For the cattle data consisting of five half-sib families, 239 SNPs were taken into¹
² account in the family approach. At 61 out of 300 SNPs, all sires were homozygous²
³ for the major allele, leading to zero entries on the diagonal of the covariance matrix.³
⁴ Thus, these loci were discarded when setting up the correlation matrix R . A clear⁴
⁵ distinction of regions with particularly high interdependence was not possible for⁵
⁶ any of the approaches (Fig. 3). In total, 11 groups with size ≥ 3 were found with⁶
⁷ both the family and population-LD approach (Tab. 1) but the representative SNPs⁷
⁸ of these groups were distributed more evenly over the chromosome window based⁸
⁹ on the family approach. The median of distances between all tagSNPs was 0.08 cM⁹
¹⁰ and 0.07 cM in the family and population-LD approach, respectively. When only¹⁰
¹¹ those SNPs were used in the population-LD approach that were considered in the¹¹
¹² family approach, the outcome of grouping differed: 239 SNPs were binned into 194¹²
¹³ groups (CH index 6.2); 8 out of them had group size of at least three SNPs vs. 300¹³
¹⁴ SNPs were binned into 210 groups (CH index 2.4); 11 groups with group size of at¹⁴
¹⁵ least three SNPs. 15

¹⁶ The correlation matrix corresponding to 13 full-sib families in maize is shown 16
¹⁷ in Figure 4. In three out of 956 SNPs, maternal SNP alleles of F1 plants were 17
¹⁸ missing; these loci were discarded in the family approach. In total, 953 SNPs have 18
¹⁹ been binned into 426 groups based on the correlation matrix but the CH index 19
²⁰ was about 50 % lower than with the population-LD approach where 956 SNPs were 20
²¹ grouped into 576 bins (Tab. 1). With both approaches, tagSNPs corresponding to 21
²² the bins of at least three SNPs were similarly distributed over the chromosome 22
²³ (family approach: 93 bins, population-LD approach: 70 bins) except a gap between 23
²⁴ SNP index 650 and 750 which was better covered with tagSNPs from the family 24
²⁵ approach. The median distance between all representative SNPs was 0.29 cM with 25
²⁶ the family and 0.16 cM with the population-LD approach. With both methods, a 26
²⁷ block of strong (positive) association among SNPs appeared in the region of 85.04 27
²⁸ to 95.31 Mbp which is in the vicinity of the functional centromere [18]. Two F1 28
²⁹ plants were the driving factor: they were completely heterozygous in this window. 29
³⁰ 30

³¹ Simulation study 31

³² The average number of groups over all repetitions, and groups with at least three 32
³³ SNPs are listed in Table 1; results are given for $t = 0.8$ and varying number of 33

¹half-sib families and QTLs. Grouping of markers appeared rather robust based on¹
²the family approach, about 60 groups have been built, but the number of groups²
³strongly varied with the population-LD approach suggesting a dependence on family³
⁴size. Large families needed more groups. Note that the number of rows in \mathbf{X} was⁴
⁵fixed ($n = 1\ 000$). The number of groups containing at least three markers was⁵
⁶rather constant between methods and for different family sizes and QTLs. The⁶
⁷CH index was at least two times larger with the family approach than with the⁷
⁸population-LD approach. The population-LD approach becomes competitive with⁸
⁹decreasing threshold t and performed better than the family approach with respect⁹
¹⁰to the CH index when $t \leq 0.6$ (see Additional file 7).¹⁰

¹¹ A SNP-BLUP approach was applied to simulated data in order to evaluate sen-¹¹
¹²sitivity and specificity if only tagSNPs were used as predictor variables in linear¹²
¹³regression. As an example, results based on one half-sib family and two simulated¹³
¹⁴QTLs are shown as ROC curve in Figure 5; the number of groups was almost equal¹⁴
¹⁵for this scenario. As expected, considering all SNPs simultaneously yielded highest¹⁵
¹⁶accuracy in terms of true-positive and false-positive rate. However, if the aim was to¹⁶
¹⁷use filtered data in SNP-BLUP, then tagSNPs obtained from the family approach¹⁷
¹⁸was the second best choice. Or in other words, using only one fifth of available¹⁸
¹⁹genotypic information led to almost the same accuracy of genome-based association¹⁹
²⁰studies as using all genotypic information. The choice of t had no influence on which²⁰
²¹method performed best, see Figure 5a for $t = 0.8$ and Figure 5b for $t = 0.5$. ROC²¹
²²curves looked very similar for all investigated scenarios of simulation though the²²
²³number of groups obtained from the population-LD approach increased with de-²³
²⁴creasing number of families. Hence, a direct relationship between number of groups²⁴
²⁵and sensitivity/specificity seems not to exist.²⁵

²⁷Options for statistical-genetics approaches²⁷

²⁸ Instead of selecting representative markers for genome-based evaluations, employ-²⁸
²⁹ing the grouping structure itself can be a beneficial option. For instance, the group²⁹
³⁰assignment derived from the family approach can directly be considered in a group³⁰
³¹lasso approach [3, 19]. Then the effects of markers in a group of highly dependent³¹
³²markers will jointly be shrunk towards zero or enlarged with respect to the rele-³²
³³vance of this group for trait expression. Additional sparsity within group can be³³

¹achieved with a sparse-group lasso approach [20]. Grouped approaches shall be in-¹
²vestigated in more detail in future because they hold potential to cope with high²
³multicollinearity. Possible benefits will likely depend on characteristics of the sam-³
⁴ple, such as the number of families, SNP density, and population-genetic parameters,⁴
⁵e.g., heritability and heterozygosity. 5

⁶ In future research, the functionality of our package should be extended by grouping⁶
⁷methods based on LD blocks which can optionally put restrictions to the physical⁷
⁸distance between SNPs (similar to [21]). Other options for selecting tagSNPs (e.g.,⁸
⁹depending on allele frequency; [22]) will be verified. 9

10

10

11 **Conclusions** 11

¹²The extent of dependence among genomic markers is affected by the underlying¹²
¹³population structure. Representative markers can be selected more efficiently if the¹³
¹⁴corresponding matrix of pairwise dependencies takes this structure into account.¹⁴
¹⁵The correlation matrix for half- or full-sib families highlights regions of high depen-¹⁵
¹⁶dence between markers more precisely than the population-LD matrix. Additionally,¹⁶
¹⁷it reveals regions of positive or negative association among markers. We contributed¹⁷
¹⁸a new function tagSNP to the R package hscovar which is suited to samples from¹⁸
¹⁹livestock and crop populations with typical family stratification. The covariance ma-¹⁹
²⁰trix can be set up in a piecewise manner, either separately for each chromosome or²⁰
²¹based on other meaningful information. The resulting grouping structure can be ex-²¹
²²ploited in genome-based evaluations to handle the problem of high multicollinearity²²
²³between markers. 23

24

24

24 **Abbreviations** 24

²⁵**AI-REML:** Average information restricted maximum likelihood 25

²⁶**BLUP:** Best linear unbiased prediction 26

²⁶**BTA:** Bos taurus autosome 26

²⁷**cM:** CentiMorgan 27

²⁸**DNA:** Deoxyribonucleic acid 28

²⁸**FPR:** False positive rate 28

²⁹**GWAS:** Genomewide association study 29

³⁰**LD:** Linkage disequilibrium 30

³⁰**Mbp:** Mega base pairs 30

³¹**MHC:** Major histocompatibility complex 31

³¹**QTL:** Quantitative trait locus 31

³²**ROC:** Receiver operating characteristic 32

³²**SNP:** Single nucleotide polymorphism 32

³³**TPR:** True positive rate 33

1	Declarations	1
2	Ethics approval and consent to participate	2
	Not applicable.	
3		3
4	Consent for publication	4
	Not applicable.	
5		5
	Availability of data and materials	
6	All R scripts used are provided in Additional files 2–6. The R package <code>hscovar</code> version 0.4.0 is available at CRAN.	6
	The cattle data are accessible through RADAR https://dx.doi.org/10.22000/280 . The mouse data are	
7	obtainable from https://wp.cs.ucl.ac.uk/outbredmice/heterogeneous-stock-mice/ . Maize data are available	7
8	at NCBI GEO under Accession Number GSE50558.	8
9	Competing interests	9
	The authors declare that they have no competing interests.	
10		10
	Funding	
11	The project was funded by the German Research Foundation (DFG, WI 4450/2-1). The funder had no role in the	11
12	design of the study and collection, analysis, and interpretation of data and in writing the manuscript.	12
13	Authors' contributions	13
	DW developed the theory, implemented the statistical methods, performed the analysis, and wrote the manuscript.	
14	MD and JK contributed to software development and improved the manuscript. All authors have read and approved	14
15	the final manuscript.	15
16	Acknowledgments	16
	We thank the Reviewers for their helpful comments.	
17		17
	References	
18	1. Wittenburg D, Bonk S, Doschoris M, Reyer H. Design of experiments for fine-mapping quantitative trait loci in	18
19	livestock populations. <i>BMC Genetics</i> . 2020;21:66.	19
20	2. Carlson CS, Eberle MA, Rieder MJ, Yi Q, Kruglyak L, Nickerson DA. Selecting a Maximally Informative Set of	20
	Single-Nucleotide Polymorphisms for Association Analyses Using Linkage Disequilibrium. <i>The American Journal</i>	
21	of Human Genetics. 2004;74(1):106–120.	21
22	3. Dehman A, Ambroise C, Neuvial P. Performance of a Blockwise Approach in Variable Selection Using Linkage	22
	Disequilibrium Information. <i>BMC Bioinformatics</i> . 2015;16:148.	
23	4. Kim SA, Cho CS, Kim SR, Bull SB, Yoo YJ. A New Haplotype Block Detection Method for Dense Genome	23
24	Sequencing Data Based on Interval Graph Modeling of Clusters of Highly Correlated SNPs. <i>Bioinformatics</i> .	24
25	2018;34(3):388–397.	25
26	5. Hill W, Robertson A. Linkage Disequilibrium in Finite Populations. <i>TAG Theoretical and Applied Genetics</i> .	26
	1968;38(6):226–231.	
27	6. Clayton D. <code>snpStats: SnpMatrix and XSnpmatrix Classes and Methods</code> ; 2017. R package version 1.34.0.	27
	Available from: http://bioconductor.org/packages/release/bioc/html/snpStats.html .	
28	7. Clayton D, Leung HT. An R package for analysis of whole-genome association studies. <i>Human heredity</i> .	28
29	2007;64(1):45–51.	29
30	8. Caliński T, Harabasz J. A dendrite method for cluster analysis. <i>Communications in Statistics-theory and</i>	30
	<i>Methods</i> . 1974;3(1):1–27.	
31	9. Koivula M, Strandén I, Su G, Mäntysaari EA. Different Methods to Calculate Genomic	31
	Predictions—Comparisons of BLUP at the Single Nucleotide Polymorphism Level (SNP-BLUP), BLUP at the	
32	Individual Level (G-BLUP), and the One-Step Approach (H-BLUP). <i>Journal of Dairy Science</i> .	32
	2012;95(7):4065–4073.	
33	10. Butler D, Cullis BR, Gilmour A, Gogel B. <i>ASReml-R Reference Manual</i> . The State of Queensland, Department	33
	of Primary Industries and Fisheries, Brisbane. 2009; Available from: https://asreml.kb.vsni.co.uk/wp-content/uploads/sites/3/2018/02/ASReml-R-3-Reference-Manual.pdf .	

- 1 11. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2020. Available 1
 2 from: <https://www.R-project.org/>. 2
 3 12. Valdar W, Solberg LC, Gauguier D, Burnett S, Klenerman P, Cookson WO, et al. Genome-Wide Genetic 3
 4 Association of Complex Traits in Heterogeneous Stock Mice. *Nature Genetics*. 2006 Aug;38(8):879–887. 4
 5 13. Browning SR, Browning BL. Rapid and Accurate Haplotype Phasing and Missing-Data Inference for 5
 6 Whole-Genome Association Studies by Use of Localized Haplotype Clustering. *American Journal of Human 6
 7 Genetics*. 2007 Nov;81(5):1084–1097. 7
 8 14. Ferdosi M, Kinghorn B, van der Werf J, Lee S, Gondro C. hspbase: An R Package for Pedigree Reconstruction, 8
 9 Detection of Recombination Events, Phasing and Imputation of Half-Sib Family Groups. *BMC Bioinformatics*. 9
 10 2014;15(1):172. 10
 11 15. Hampel A, Teuscher F, Gomez-Raya L, Doschoris M, Wittenburg D. Estimation of Recombination Rate and 11
 12 Maternal Linkage Disequilibrium in Half-Sibs. *Frontiers in Genetics*. 2018;9:186. 12
 13 16. Bauer E, Falque M, Walter H, Bauland C, Camisan C, Campo L, et al. Intraspecific variation of recombination 13
 14 rate in maize. *Genome Biol*. 2013;14:R103. 14
 15 17. Gaynor RC, Gorjanc G, Hickey JM. AlphaSimR: An R-package for Breeding Program Simulations. *G3: Genes 15
 16 Genomes Genetics*. 2020;. 16
 17 18. Schneider KL, Xie Z, Wolfruber TK, Presting GG. Inbreeding drives maize centromere evolution. *Proceedings 17
 18 of the National Academy of Sciences*. 2016;113(8):E987–E996. 18
 19 19. Yuan M, Lin Y. Model Selection and Estimation in Regression with Grouped Variables. *Journal of the Royal 19
 20 Statistical Society B*. 2006;68(1):49–67. 20
 21 20. Simon N, Friedman J, Hastie T, Tibshirani R. A Sparse-Group Lasso. *Journal of Computational and Graphical 21
 22 Statistics*. 2013;22(2):231–245. 22
 23 21. Gabriel SB, Schaffner SF, Nguyen H, Moore JM, Roy J, Blumenstiel B, et al. The Structure of Haplotype 23
 24 Blocks in the Human Genome. *Science*. 2002;296(5576):2225–2229. 24
 25 22. Wang S, He S, Yuan F, Zhu X. Tagging SNP-Set Selection with Maximum Information Based on Linkage 25
 26 Disequilibrium Structure in Genome-Wide Association Studies. *Bioinformatics*. 2017;33(14):2078–2081. 26
 27
 28
 29

18 Figures 18

19
 20 **Figure 1** (a) Matrix highlights the SNP pairs with correlation larger than a certain threshold for a 20
 21 given SNP panel. The SNP involved most is marked by an arrow (Step A). (b) Correlations are 21
 22 considered within the corresponding SNP subset and a tagSNP is selected (Step B). (c) All SNPs 22
 23 associated with the tagSNP are removed from the remaining ungrouped SNP panel for the next 23
 24 round of iteration (Step A). Iterations continue until no ungrouped SNP is left. 24

25
 26 **Figure 2** Correlation (a) and LD matrix (b) for mouse data on chromosome 17. The red dots 26
 27 highlight representative SNPs of groups with at least three SNPs. In total, 394 SNPs were 27
 28 considered for the correlation and population-LD matrix. 28

29
 30 **Figure 3** Correlation (a) and LD matrix (b) for cattle data in a target region of chromosome 30
 31 1:20.6–39.4 Mbp. The red dots highlight representative SNPs of groups with at least three SNPs. 31
 32 In total, 300 SNPs were considered. The correlation matrix was set up at 239 SNPs and the 32
 33 population-LD matrix at 300 SNPs; missing values are filled in white color. 33

Figure 4 Correlation (a) and LD matrix (b) for maize data on chromosome 2. The red dots highlight representative SNPs of groups with at least three SNPs. In total, 956 SNPs were considered. The correlation matrix was set up at 953 SNPs and the population-LD matrix at 956 SNPs; missing values are filled in white color.

Figure 5 Sensitivity and specificity of testing SNP effects depending on threshold $t = 0.8$ (a) and $t = 0.5$ (b). ROC curves are based on 100 repeated simulations of genotypes and phenotypes in $N = 1$ half-sib family with 1 000 progeny (two QTL signals, heritability 0.3).

Table 1 Number of groups, number of groups with at least three SNPs and Calinski-Harabasz index (CH). Number of SNPs corresponds to the method applied (family or population LD). Average values of 100 repetitions are presented for simulated data and $t = 0.8$. Computing time for grouping markers was up to 0.2s except for the maize data which required 2s.

	Families	QTLs	Family approach				Population-LD approach			
			SNPs	Groups	≥ 3	CH	SNPs	Groups	≥ 3	CH
Simulation	10	2	283	59	10	80.7	300	21	9	40.0
	10	5	282	61	10	75.0	300	17	9	38.9
	5	2	282	61	10	71.9	300	29	8	32.6
	5	5	281	64	10	85.6	300	24	9	35.0
	1	2	281	59	9	61.4	300	56	6	20.5
	1	5	282	59	9	83.0	300	49	7	25.0
Mouse	138		394	83	49	38.8	394	98	51	20.6
Cattle	5		237	172	11	2.9	300	210	11	2.4
Maize	13		953	426	93	10.4	956	576	70	21.9

Tables

Additional Files

- Additional file 1 – PDF file summarizing the derivation of family-based correlation matrices
Given the population structure, half-sib families or full-sib families, the covariance matrix is analytically retrieved. Its computation using the R package `hscovar` is shown.
- Additional file 2 — TXT file containing R code for mouse data
Raw mouse data are processed and the matrix of correlation between markers is derived.
- Additional file 3 — TXT file containing R code for cattle data
Cattle data are processed and the matrix of correlation between markers is derived.
- Additional file 4 — TXT file containing R code for maize data
Raw maize data are processed and the matrix of correlation between markers is derived.
- Additional file 5 — TXT file containing R code for simulation study
With this script, genotype and phenotype data of half-sib families are simulated and a genome-based association analysis is carried out.
- Additional file 6 — TXT file containing R code for creating plots
Plots of correlation matrices and population-LD matrices are produced based on results with additional files 2–5.

¹ Additional file 7 — PDF file with additional tables	1
² Number of groups, number of groups with at least three SNPs and Calinski-Harabasz index for different simulation scenarios and varying threshold.	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33