

Pedigree file and the inverse of the relationship matrix in honey bees

Pim Brascamp, 19 April 2023

For comments and questions please contact pim.brascamp@wur.nl

This note describes two R-programs. The first produces a pedigree file for honey bees consisting of colonies, dams and sires. Colonies are groups of workers, dams are individual queens and sires are groups of drone-producing queens (DPQ). The programs allow for the situation that the sire is only one queen, which in particular is relevant for single-drone insemination. The program also allows natural mating. The second program produces the numerator relationship matrix **A** and its inverse.

For theory see Brascamp and Bijma (GSE, 2014, DOI:10.1186/s12711-014-0053-9), Brascamp *et al.* (Apidologie, 2016, DOI: 10.1007/s13592-016-0427-9) and a erratum to the Apidologie paper (Apidologie, 2018, DOI: 10.1007/s13592-018-0573-3). The erratum concerns the interpretation of the variance component for colony's worker effect. A paper going into this issue more extensively was published in GSE (Brascamp and Bijma GSE, 2019, DOI: 10.1186/s12711-019-0510-6.). The program computing the inverse of **A** does not include the inversion for very large datasets as developed by Bernstein et al. (2018, J Anim Breed Genet. 2018;135:323–332.)

The programs and files are the result of continuous development. The programs show traces of that and another consequence is sequences of items in the files are sometimes no longer logical. In 2015 a FORTRAN-version of AMD-AINV was written by Jérémie Vandenplas which was extremely much faster than the R-version. Since then the R-version was extended with additional mating options. On the basis of Jérémie's work recently Tristan Kistler wrote a FORTRAN-version for his simulation work using the current R-version.

The current relationship matrix in **A** assumes a Poisson distribution for the probability that two full-sibs descend from the same drone, and for the probability that they descend from the same drone-producing queen. As pointed out by Manuel Du this leads to probabilities larger than 1 if NS or ND equal 1, where NS equals the number of DPQ constituting a sire and ND equals the number of drones a queen on average fertilizes. We will look into this in the future.

In Brascamp and Bijma, 2019, we show that the estimates of genetic parameters depend upon the assumptions about the base population. If the assumptions don't represent the real situation, estimates of genetic parameters will be biased considerably. We distinguished the situation that drone-producing queens are unrelated in the base population or that their additive genetic relationship has equilibrium values because the breeding program already was running some time preceding the data to be analysed. The default values for NS and ND that are input in the computer program setting up the pedigree relate to this base population. Sometimes, however, with more complicated (and probably more realistic) breeding programs

it is not obvious what's taken to be the base population. As the levels of the estimated genetic parameters are influenced considerably by the assumptions about the base population it is important to explicitly pay attention to this issue when publishing results.

Both programs sometimes produce printed warnings.

Both programs produce log-files (.txt). Sometimes a program may run very short. This may relate to an early error, reported in the log-file. Apart from input-errors, reported problems generally concern sequence problems in the pedigree file or that the product of **A** and its inverse is not the identity matrix. If you can't solve these problems please let me know (pim.brascamp@wur.nl) and I can try to help you out.

There are two possible assumptions for the relatedness between the drone-producing-queens constituting base sires. Either they are assumed to be unrelated, or their relatedness is assumed to be in equilibrium. The assumption is embedded input in the second program.

The programs deal with the following situations or combinations of these:

1. A sire is identified by the dam of the full-sib group of drone-producing-queens that constitutes the sire. A particular dam may have different sires, as e.g. the case when she produced a sire for a mating station but also for instrumental insemination.
2. A sire is identified itself. There are a number of options for this case:
 - a. The sire is a single sire producing one drone (single-drone insemination) or various drones for insemination. The program accommodates the situation that a particular single sire produces various numbers of drones.
 - b. The sire is a group of open mating drones produced by unknown DPQ. In that case the additive genetic relationship between drone-producing queens is taken to be zero.

Alfons Willam is gratefully acknowledged because our collaboration led to the first version of the programs. Also, Ralph Büchler, Arista Bee Research (BartJan Fernhout, Bart Barten, Mari van Iersel and Tieme Wanders) and Florence Phocas and Tristan Kistler are gratefully acknowledged for let me work with their data such that I could extend and check the programs for various assumptions in real data. Tristan Kistler and Benjamin Basso used the programs extensively giving rise to some changes especially related to open mating.

1. pedigree.r (version 19).

Input is a file input-pedigree.txt, the format of which is given in Appendix 1. Usually the input will not be in that form and has to be produced tailor-made. The program produces several output files of which log-pedigree.txt reports whether the pedigree seems to be okay, especially with respect to the sequence of the individuals in the pedigree. There are two requirements. Firstly, parents should precede offspring, but specific for honey bees,

aunts should precede cousins. Usually, the program takes care of this. If changes have to be made in the input file, it may be recommendable to let the changes be done by the program. Look for *file specific issues causing sequence problems* and add a script with the changes to be made. It often boils down to changing years of birth (column 5).

The test location for a colony can also be added (search for *this concerns the test location*). Sometimes this simply is column 3 of input-pedigree.txt, but often a combination of year of observation (column 1), bee breeder (column 2) and column 3.

Among further output there are four files that are input for the second program: pedigree.txt, numbers.txt, ident.txt and blocks.txt. In situation 2a, single sires, also singlesire.txt is input for the second program. The formats of these files are given in Appendix 1 as well. The first file is the pedigree file, yet to be completed by inbreeding coefficients, mendelian sampling terms and, for sires, the additive genetic relationships between drone-producing queens. The second summarizes some numbers, the third the identifications of the individuals in the pedigree and the fourth the start and finish sequence numbers of individuals with the same dam (and therefore sire). The last file contains the identifications of the single sires, the numbers of drones used, and the numbers of drone-producing queens and drones of its sire.

Apart from the external input (input-pedigree.txt) there is internal input, defined in the program itself. This should be checked and –if necessary changed- before running the program.

- Default values for NS (the number of drone-producing queens) and ND (the number of drones fertilizing a queen) should be given. Usually this is something like 8-20 and e.g. 12. These values are taken for matings in the base population and also in later generations if there is a record without mate specified. Special attention should be given to NS in case of open mating. See for this Appendix 2.
- There is a parameter codetest that equals 1 if a thorough test is carried out on the proper sequence of the pedigree. A second part of that test takes a lot of time but usually doesn't detect errors. For initial runs it may be useful to put codetest unequal 1 to limit the sequence test to the less time-consuming part, but ultimately codetest should be put to 1 to be pretty sure that inversion of the relationship matrix doesn't create problems.

Other variations in input are accommodated in pedigree-input.txt and are described below.

1a, 2a and 4a are input

- a. The input for the pedigree file usually consists of the identifications of the queen (1a) in the colony, its dam (2a) and the dam (4a) of the drone-producing queens (1b). Drones of those drone-producing queens are mated to the 1a. In addition the years of birth are to be added and, in case of colonies with observations, the year of observation.

- b. In some applications a 4a may have more than one 1b, for example a mating on an island and instrumental insemination. In that case, an additional column with an additional single digit needs to be entered per 4a in column 18 to distinguish between the 1b. It may happen that beforehand it is not realized that more than one 1b exists per 4a. NS and ND are input for each record in pedigree-input.txt and the program tests whether for a 4a always the same NS and ND are present in the input file. If that is not the case there is either a typo or different 1bs exist per 4a, and column 18 can be used if different 1bs exist per 4a. NS and ND for example likely are different for island mating and instrumental insemination.

Insemination from one drone-producing queen (1a, 2a, 1b and 4a are input)

In some applications, only one drone-producing queen is used. In that case the identification of 1b has to be input as well. Otherwise it is zero. Note that in this case the identification for 4a should also be given, to distinguish from the case where queens are open mated.

Open mating (1a, 2a and 1b are input)

It may be that the 4a is not known and should be zero because queens are open mated. For each record with an open-mated queen, a 1b has to be added. There are different options for this, discussed in Appendix 2.

Deeper pedigree

It may be that the input file not only contains records that have observations, but also deeper pedigree information. For deeper pedigree information without observation the year of observation needs to be zero. Note that also for those records the numbers of DPQ and drones of the queen's mate should be given.

Sequence problems

In some cases the program doesn't arrive at a sequence suitable to derive **A**. In that case enter in the program an additional line (or lines) that makes specific changes to the modified year of birth of a specific queen, usually the 2a (See in the program the paragraph file specific issues causing sequence problems). The output file pedigree_ident.txt may facilitate to detect the reason of a sequence error. It may, however be more efficient to make these changes in the program that produces the input file for pedigree.r.

Pedigree.r is programmed such that all members in the pedigree have two or zero parents. If a parent is missing the program generates it. For dams it is a unique dam for each member and for sires it is a base sire with the specified assumptions, i.e. whether the drone-producing queens are unrelated or that their additive genetic relationship is in equilibrium, and NS and ND as specified in the beginning of the program. There may be

reasons to deviate from that (e.g. in case of open mating) and then the sire should be added as with open mating.

2. AMD-AINV.r (version 19)

Input are the four files given above. Furthermore, in the program itself, four parameters should be defined.

- *equi*, being zero or one. In case of zero the drone-producing queens in base sires are taken to be unrelated. If one, the additive genetic relationship between drone-producing-queens is assumed to be in equilibrium. This is likely to be the case if a selection program is already running a couple of generations before the start of the data that are analysed.
- *pinv*, being zero or one. In case of zero the inverse (AINV) of the numerator relationship matrix is computed utilizing mendelian sampling (co)variances. It is recommended to do this once because it is a fairly reliable test whether everything (i.e. the pedigree file but also the computer program related to special perhaps unforeseen cases) is okay.
- In case of one AINV is computed with the R-function 'solve'. It is likely that the first option (*pinv*=0) allows larger datasets than the latter, but the latter option (*pinv*=1) is (considerably) faster. So, particularly in the case of a simulation where several runs with different AINV need to be analysed, the latter option is to be preferred. *priF*, being zero or one. This concerns writing of inbreeding coefficients of possible future matings and if *priF*=0 these will not be computed nor written. If *priF*=1, an additional parameter should be entered, *yip*, that tells for how many years until the last year the coefficients should be written. The purpose of it is to give an indication of whether a particular planned mating doesn't result in too high inbreeding and that's only useful for queens that were born rather recently. Three coefficients are written: Mating of queens reared from queen A with drones of DPQ reared from queen B and *vv*, queens reared from queen A mated with drones of queen B, and queens reared from queen B mated to drones from queen A.

The program produces a log file: log-AMD-AINV.txt. Of special relevance is the result of the multiplication of *A* and its inverse that should be identity. Output file is AINV.giv which can be used as input for ASReml. Further output files are pedigree_complete.txt, ident_complete.txt and F-future.txt. In case of single sires also singlesires_ass.txt. See singlesires.txt in Appendix 1.

AMD-AINV.r assume that all individuals in the pedigree have two or zero known parents. In case an individual has one known parent, pedigree.r creates the second, that itself has no known parents. Animals with unknown parents form the base population.

The difference between pedigree.txt and pedigree_complete.txt is that for each individual the inbreeding coefficient and the mendelian sampling variance is added, and for sires also the additive genetic relationship between two drone-producing queens. Furthermore, the diagonal element of A-inverse is added, and the identification in column 1 (equals the sequence number in column 16) is replaced by the real identification in case of queens and fabricated identification in the case of colonies (groups of workers), sires (groups of drone-producing queens) or base dams.

Another difference is that in column 12 the diagonal element of AINV is given. This information is useful for the calculation of accuracies of EBVs.

- 3. Y.txt.** This is input file for ASReml and contains the identification of the colony (if worker effect of workers is to be included in the statistical model) and also the dam of workers (1a, if also the queen effect of the dam of workers is to be included in the statistical model) and furthermore fixed effect(s) and observations. The identifications are sequence numbers 2 (column 17) as in pedigree_complete.txt. For flexibility the creation of Y.txt has been brought outside of AMD-AINV.r. Therefore, a specific program needs to be written to produce Y.txt. Apart from the input file (the file that was used to produce input-pedigree.txt) also pedigree_complete.txt is needed, to link identifications to sequence numbers 2 and also to the sequence number 2 of the dam in the colony. A small program Y.r is available to serve as a starting point.
- 4. EBVs.** Often EBVs are a desired result of the analyses. In ASReml these are in the file with extension .sln. The identifications are the sequence numbers 2. Usually, the colonies are identified by the real identification of the 1a (the dam of workers) and therefore again pedigree_complete.txt is needed as input for a program to list EBVs.

Appendix 1. Formats of various files

input-pedigree.txt

- 1 year of observation (zero if it concerns deeper pedigree info without observations)
- 2 bee breeder
- 3 testing station (often the combination of 1, 2 and 3 is HYS in the statistical model)
- 4 year of birth 1a
- 5 identification 1a
- 6 year of birth 2a
- 7 identification 2a
- 8 year of birth 4a
- 9 identification 4a
- 10 not used
- 11 not used
- 12 not used
- 13 not used
- 14 NS (number of drone-producing queens in the 1b mated to the 1a)
- 15 ND (average number of drones produced by the 1b mated to the 1a)
- 16 year of birth of 1b (if one single queen or open mating, otherwise zero)
- 17 identification of 1b (if known or fabricated in case of open mating, otherwise zero)
- 18 extra digit, with default is zero. This is only needed if a particular 4a occurs with different sets of drone-producing queens, that may differ for NS and ND.

numbers.txt

- 1 N of records in the input file
- 2 N of drone-producing queens in base sires (NS)
- 3 N of drones base sires on average contribute to a mating (ND)
- 4
- 5 N of colonies
- 6 N of records total (size of the complete pedigree file, and e.g. dimension of *A*.)
- 7 N of blocks of full sibs
- 8 Equals 0 if there are no single sire with NS=1, and 1 otherwise.

pedigree.txt, pedigree_complete.txt and pedigree_ident.txt

Table with in the first column pedigree.txt, in the second pedigree_complete.txt and in a third pedigree_ident.txt.

x item is present

xx sequence number 2

xxx identifications, for queens in colonies the original idents in input-pedigree.txt, for other queens, colonies, and sires, fabricated idents.

1	identification	x	xxx		xxx	colony, dam or sire
2	Sex (1=dam, 2=sire, 3=colony)	x	x		x	x
3	mate (for dams only)	x	x		xxx	mate
4	type of mating (for sires only)	x	x		xxx	mate's dam
5	year of birth modified	x	x		x	x
6	sequence number of dam	xx	xx		xxx	dam
7	sequence number of sire	xx	xx		xxx	sire
8	a_{ss} (for sires only)		x		xxx	sire's dam
9	F		x			
10	var MS		x			
11	testing station	x	x		x	
12	diagonal element of A		x			
13						
14						
15						
16	sequence number 1	x	x		x	
17	sequence number 2	xx	xx		xx	
18	year of birth	x	x		x	
19	NS (number of DPQ)	x	x		x	
20	ND (number of drones)	x	x		x	

Note that NS and ND for dams relate to those of their mate. For sires they relate to the numbers of the sires themselves. For single sires the situation is divergent. First of all ND for a single sire may vary among single sires and is therefore not informative. Furthermore, if the single sire appears as dam as well, the numbers relate to her mate.

Sequence number 2 is the sequence number of the individuals in the pedigree after sorting on year of birth modified in such a way that offspring of the same dam (and therefore sire) is in a block.

Column 12 gives the diagonal elements of A. These can be used to compute accuracies of EBVs for the case that the EBV only contains one random effect, for example worker effect. In that

case the accuracy equals $r_{A_i\hat{A}_i} = \sqrt{1 - \frac{s_i^2}{\sigma_{A_i}^2}}$, where s_i^2 is the square of seEffect from ASReml (more generally, it equals the square of the standard error of prediction, SEP) and $\sigma_{A_i}^2$ equals the diagonal element in column 12 multiplied by the estimate of σ_A^2 from ASReml. If the EBV contains for example worker and queen effect, then the inverse of the block in AINV referring to the worker and the queen is needed. ASReml can produce SEP in that case as well. The SEP of the worker group can be obtained by inclusion of the command line (personal communication Arthur Gilmour)

PREDICT worker a: b queen a: b ! ONLYUSE worker queen ! PARALLEL worker queen

directly after the line defining the model. In this equation a and b refer to sequence numbers in the pedigree file. (Note that we experienced a limit to the size of b – a that can be dealt with in a single run.)

ident.txt

This is an in-between file containing sorted on sequence number 1, after initially sorting the data on year of birth and identification of 1a.

- 1 true identification (for queens in the colonies) and fabricated for colonies, other queens (if present) and sires.

blocks.txt

The number of records is the 7th item in numbers.txt. For each record:

- 1 first sequence number 2 in the block
- 2 last sequence number 2 in the block

singlesires.txt

This file only is relevant in case of single sires with NS=1

- 1 sequence number 2 for the single sire
- 2 identification of the single sire
- 3 number of drones inseminated with

AINV.giv

- 1 sequence number 2 of the first member of the pedigree (i)
- 2 sequence number 2 of the second member of the pedigree (j)
- 3 element of AINV

The file is organized such that $i \geq j$. It only contains non-zero elements.

F-future.txt

This file summarizes inbreeding coefficients (F) of possible matings using colonies of the last yip years. If there are two colonies x and y there are three possible matings.

- a. A virgin queen from x with drones of drone-producing queens from y and the reciprocal; both lead to the same F.
- b. A virgin queen from x with drones from y.
- c. Drones from x with a virgin queen from y.

These latter became relevant because in case of single-drone insemination the drone usually is produced by the queen in a tested colony and not by drone-producing queens reared from tested colonies.

- 1 sequence number 2 of the first colony (i)
- 2 sequence number 2 of the second colony (j)
- 3 real identification of the queen (dam) in colony i
- 4 real identification of the queen (dam) in colony j
- 5 F (%) for case a.
- 6 F (%) for case b.
- 7 F (%) for case c.

The file is organized such that $i \geq j$. The file only contains colonies born in the last three years.

Y.txt (for example)

- 1 sequence number 2 of the colony
- 2 sequence number 2 of the dam (if the statistical model needs it)
- 3 test location
- 4 Possible other fixed effects and observations

In earlier versions of pedigree-x.r the observations were part of input-pedigree.txt. Because the very variable numbers of observations that is no longer the case. There is a little example of a program Y.r that may serve as a starting point for Y.r.

Appendix 2

Open mating

Incorporation of open mating is still under development. This note describes the thoughts until now following from discussions with Tristan Kistler as he is looking into datasets with open mating.

In current programs pedigree.r and AMD-AINV.r the idea is to add artificial mates (groups of DPQ) in case of open mating. In pedigree.r these mates are converted into sires in the pedigree file. There are different options: For each mated queen add a mate, for queens in the reach of similar colonies producing drones add the same mate, or even add the same mate for all open matings. A fourth option is to declare the mate for open mating as missing, *assuming that all drones participating in mating are unrelated*.

To illustrate what happens look at a pedigree with two base dams a and b, two base open mating sires c and d and four descending queens:

e with parents a and c

f with parents a and c

g with parents b and c

h with parents a and d

In this pedigree e and f are fullsibs, e and h are maternal halfsibs and e and g are paternal halfsibs. It is assumed that the number of unrelated DPQ in case of open mating is s , such that the additive genetic variance of a sire equals $1/s$

We get the following additive genetic (co)variances.

	a	b	c	d	a	c	a	c	b	c	a	d
	a	b	c	d	e	f	g	h				
a	1	1	0	0	0.5	0.5	0	0.5				
b	1	1	0	0	0	0	0.5	0				
c	0	0	$1/s$	0	$1/2s$	$1/2s$	$1/2s$	0				
d	0	0	0	$1/s$	0	0	0	$1/2s$				
e	0.5	0	$1/2s$	0	1	$0.25 + 1/4s$	$1/4s$	0.25				
f	0.5	0	$1/2s$	0	$0.25 + 1/4s$	1	$1/4s$	0.25				
g	0	0.5	$1/2s$	0	$1/4s$	$1/4s$	1	0				
h	0.5	0	0	$1/2s$	0.25	0.25	0	1				

In this table the individuals a-h are given and in the top row the parents of e-h.

It follows that the additive genetic covariance

between FS equals $0.25 + 1/4s$

between mHS equals 0.25

between pHS equals $1/4s$

When the number of DPQ in case of open mating increases the additive genetic covariance between FS approaches 0.25 and between pHS approaches zero. These are also the values in case the mate would be declared missing, assuming that drones are unrelated.

This result implies that when the number of DPQ in case of open mating is very large, the options give very similar results.

Currently pedigree.r and AMD-AINV.r are programmed such that all members in a pedigree have zero or two parents. And if it happens to be one, an artificial base animal is added. In case an artificial sire is added it will obey the rules defined in pedigree.r and AMD-AINV.r with

respect to relatedness of DPQ in base sires and NS and ND. This will be changed in the future but for the time being when missing with unrelated drones is the aim, addition of one single open mating mate with a very large number of DPQ would give the desired results.

At first sight the 'best' solution for a real-life situation is to add 1b's that reflect the situation at hand. This may imply that for each test location (i.e. test location x year subclass) one 1b is added, with the number of DPQ that is felt reasonable for that test location. But if the NS given for that 1b is large, it probably doesn't matter much for practical application whether it is 250, 1000, or 1000,000, but this needs to be found out. In case of 1000,000 probably one and the same 1b for all open matings will give the same results as one 1b per test location, or one 1b per mating.