

## Annotated R code for a deer parentage analysis

This uses the example data and methods described in:

Dodds, K G, McEwan, J C, Brauning, R, Van Stijn, T C, Rowe, S J, McEwan, K M and Clarke, S M (2019) Exclusion and genomic relatedness methods for assignment of parentage using genotyping-by-sequencing data. *bioRxiv* 582585. doi:10.1101/582585

Command	Description
<code>genofile &lt;- "HapMap.hmc.txt.gz"</code>	Name of data file (example file can be downloaded from <a href="https://gsajournals.figshare.com/s/7ca45accf6ae82047c86">https://gsajournals.figshare.com/s/7ca45accf6ae82047c86</a> )
<code>pedfile &lt;- "DeerPedGBS.csv"</code>	"pedigree file" for parentage analysis (can be downloaded from <a href="https://gsajournals.figshare.com/s/7ca45accf6ae82047c86">https://gsajournals.figshare.com/s/7ca45accf6ae82047c86</a> )
<code>groupsfile &lt;- "Ped-Groups.csv"</code>	"groups file" for parentage analysis (can be downloaded from <a href="https://gsajournals.figshare.com/s/7ca45accf6ae82047c86">https://gsajournals.figshare.com/s/7ca45accf6ae82047c86</a> )
<code>sampdepth.thresh &lt;- 0.3</code>	Remove results with mean depth less than 0.3 (used in GBSsummary)
<code>cex.pointsize &lt;- 1.2</code>	Increase default size of text in graphical output by 20%
<code>functions.only &lt;- TRUE</code>	Do not run the KGD code when sourcing it in
<code>sink("GBSParDeerOut.txt")</code>	Redirect screen output to specified file
<code>source("GBS-Chip-Gmatrix.R")</code>	Load the main KGD functions (not run automatically because <i>functions.only</i> is TRUE). Obtained from <a href="https://github.com/AgResearch/KGD">https://github.com/AgResearch/KGD</a> .
<code>readGBS()</code>	Read the input file of "reference" and "alternate" allele counts for each sample and SNP. The default format is used ("uneak").
<code>outlevel &lt;- 1</code>	Reduce the amount of QC output
<code>GBSsummary()</code>	Run the main function for QC and setting up structures for further analyses. 12 samples are removed due to mean sample depth < 0.3. 2646 SNPs with mean depth <0.1 or with MAF=0 are removed.
<code>breed &lt;- read.table(text=seqID,sep="_", stringsAsFactors=FALSE)[,1]</code>	<i>seqID</i> contains the sample identifiers from the input file. In this example the identifiers are of the form <breed>_<number>, where <breed> is one of R (Red Deer) or W (Wapiti). This instruction extracts the text before the "_" into a character vector called <i>breed</i> .
<code>fcolo &lt;- c("darkblue","darkred") [match(breed,c("W","R"))]</code>	Set up a vector of colours ( <i>fcolo</i> ) to be used (darkblue for Wapiti, darkred for Red Deer)
<code>snpsubset &lt;- which(HWdis &gt; -0.05)</code>	Vector containing the positions of SNPs that pass the filter of Hardy Weinberg disequilibrium ( <i>HWdis</i> ) > -0.05. See finplot.png for a depiction of <i>HWdis</i> , MAF and mean SNP depth for the SNPs.
<code>GHW &lt;- calcG(npc=4, snpsubset=snpsubset,</code>	Main function for estimating relatedness. Here the

<code>sfx="RWHW")</code>	<p>filtered SNPs are used. A PCA and heatmap are requested along with the standard output. <i>npc</i>=4 principal components are output and plotted. The output object contains a genomic relatedness matrix using the KGD method (<i>G5</i>) as well as an object containing the PCA results (<i>PC</i>). Output graphics include "RWHW" as part of the name.</p> <p>Setting <i>npc</i>=-4 will omit the heatmap (which can use a lot of CPU time with larger datasets).</p> <p>Setting <i>npc</i>=0 will omit the PCA, although it is recommended to include the PCA for diagnostic purposes if possible.</p> <p>The plot GRWHWdiagdepth.png shows estimated self-relatedness as a function of mean sample depth. Within a set of samples sequenced using the same protocol, we would not expect to see a relationship. Sometimes a negative relationship is observed which may be due to non-optimised lab protocols.</p>
<code>G5 &lt;- GHW\$G5</code>	Extract the KGD GRM from <i>GHW</i> .
<code>GCheck &lt;- "G5"</code>	Specify the GRM to use in the parentage analysis.
<code>set.seed(230985)</code>	Specify a seed so that the same bootstrap results are generated if the code is rerun.
<code>source("GBSPedAssign.R")</code>	Run the parentage analysis using the script obtained from <a href="https://github.com/AgResearch/KGD">https://github.com/AgResearch/KGD</a> . The main results are in a structure called BothMatches and written to the file BothMatches.csv. Various other diagnostics files are produced.
<code>dir.create("W")</code>	Create a subfolder for a Wapiti-only analysis
<code>setwd("W")</code>	Set the work directory to the Wapiti folder
<code>indW &lt;- which(breed=="W")</code>	Create a vector of positions of Wapiti animals in the data.
<code>pW &lt;- calcp(indsubset=indW)</code>	Calculate allele frequencies (based on allele counts) using only the Wapiti data.
<code>snpsubset &lt;- which(HWdis &gt; -0.05 &amp; pW &gt; 0 &amp; pW &lt; 1)</code>	Find the SNPs that pass the Hardy-Weinberg threshold (based on the full dataset) and that are not monomorphic in the Wapiti data.
<code>GHWW &lt;- calcG(snpsubset, indsubset=indW, sfx="W", puse=pW, calclevel=1)</code>	Calculate the GRM for the Wapiti only data, using the Wapiti allele frequencies. No PCA or heatmap is requested, and some other diagnostics are suppressed (by <i>calclevel</i> =1).
<code>G5W &lt;- GHWW\$G5</code>	Extract the KGD GRM
<code>seqIDW &lt;- seqID[indW]; if(length(GHWW\$samp.removed) &gt; 0 ) seqIDW &lt;- seqIDW[-GHWW\$samp.removed]</code>	Obtain the seqID vector for the Wapiti. Ensure that this list corresponds to the individuals in <i>G5W</i> (sometimes individuals are removed to ensure that all pairs have exceed the minimum specified co-call rate in <i>calcG</i> (in this case the default threshold of 0 was used). The reduced set of SNPs may have resulted in a few pairs with no SNPs in common, although this does not usually happen (especially if low depth samples are removed initially).
<code>GCheck &lt;- "G5W"</code>	Specify the GRM to use in the parentage analysis.
<code>puse &lt;- pW</code>	Specify the allele frequencies to use in the parentage analysis
<code>indsubset &lt;- indW</code>	Specify the individuals to use in the parentage analysis (Wapiti animals)
<code>rm(minr4inb)</code>	Remove the <i>minr4inb</i> (minimum parent relatedness to

	use for checking the inbreeding threshold). This ensures the default (no minimum) is used. <i>minr4inb</i> is set to the minimum parent pair relatedness during the parentage analysis. Removing the variable allows it to be reset.
pedfile <- "../DeerPedGBS.csv"	Specify the location of the pedigree file relative to the current working directory
groupsfile <- "../Ped-Groups.csv"	Also for the groups file
source("GBSPedAssign.R")	Run the pedigree analysis for the Wapiti data
MatchesW <- BothMatches	Place the parentage results into another data frame
write.csv(MatchesW,"BothMatchesW.csv",row.names=FALSE,quote=FALSE)	Rewrite the results to a file with a name different to the standard output name (e.g. this allows these results and the combined breed results to be opened together in Excel)
# Alt models	
uY <- which(MatchesW\$BothAssign=="Y")	Identify the offspring that were assigned both parents
bbopt <- optimize(ssbbmm,lower=0,upper=20,tol=0.001)	Fit the beta-binomial model to the assigned trios
depth2K <- depth2Kchoose (dmodel="bb", bbopt\$minimum)	Change the allele sampling model to the fitted beta-binomial model
mmstatsW.bb <- mismatch.2par (MatchesW\$IndivID, MatchesW\$BestFatherMatch, MatchesW\$BestMotherMatch)	Calculate mismatch rates using the fitted beta-binomial model
names(mmstatsW.bb) <- paste0(names(mmstatsW.bb),".bb")	Add ".bb" to the variable names of the beta-binomial mismatch rates
mpopt <- optimize(ssmpmm,lower=0.5,upper=0.9,tol=0.001)	Fit the modified p model to the assigned trios
depth2K <- depth2Kchoose (dmodel="modp", mpopt\$minimum)	Change the allele sampling model to the fitted modified p model
mmstatsW.mp <- mismatch.2par (MatchesW\$IndivID, MatchesW\$BestFatherMatch, MatchesW\$BestMotherMatch)	Calculate mismatch rates using the fitted modified p model
names(mmstatsW.mp) <- paste0(names(mmstatsW.mp),".mp")	Add ".mp" to the variable names of the modified p mismatch rates
MatchesW <- cbind(MatchesW, mmstatsW.bb, mmstatsW.mp)	Add the beta-binomial and modified p model mismatch rates to the Wapiti parentage results
write.csv(MatchesW,"BothMatchesW.csv",row.names=FALSE,quote=FALSE)	Write these results to a file (overwrites previous file)
depth2K <- depth2Kchoose (dmodel="modp") # back to default model	Reset the allele sampling model to the default model (can specify "modp" or "bb" with no parameter – in both cases reverts to the standard binomial model)
setwd("../")	Set the working directory to the original analysis folder.
dir.create("R")	Repeat the Wapiti analysis workflow for the Red Deer data
setwd("R")	
indR <- which(breed=="R")	
pR <- calcp(indsubset=indR)	
snpsubset <- which(HWdis > -0.05 & pR > 0 & pR < 1)	
GHWR <- calcG(snpsubset, indsubset=indR, sfx="R",puse=pR,calclevel=1)	
G5R <- GHWR\$G5	
seqIDR <- seqID[indR]; if(length(GHWR\$samp.removed) > 0 ) seqIDR <- seqIDR[-GHWR\$samp.removed]	

GCheck <- "G5R"	
puse <- pR	
indsubset <- indR	
rm(minr4inb)	
source("GBSPedAssign.R")	
MatchesR <- BothMatches	
write.csv(MatchesR,"BothMatchesR.csv",row.names =FALSE,quote=FALSE)	
# Alt models	
uY <- which(MatchesR\$BothAssign=="Y")	
bbopt <- optimize(ssbbmm,lower=0,upper=20, tol=0.001)	
depth2K <- depth2Kchoose (dmodel="bb", bbopt\$minimum)	
mmstatsR.bb <- mismatch.2par(MatchesR\$IndivID, MatchesR\$BestFatherMatch, MatchesR\$BestMotherMatch)	
names(mmstatsR.bb) <- paste0(names(mmstatsR.bb),".bb")	
mpopt <- optimize(ssmpmm,lower=0.5,upper=0.8, tol=0.001)	
depth2K <- depth2Kchoose (dmodel="modp", mpopt\$minimum)	
mmstatsR.mp <- mismatch.2par(MatchesR\$IndivID, MatchesR\$BestFatherMatch, MatchesR\$BestMotherMatch)	
names(mmstatsR.mp) <- paste0(names(mmstatsR.mp),".mp")	
MatchesR <- cbind(MatchesR, mmstatsR.bb, mmstatsR.mp)	
write.csv(MatchesR,"BothMatchesR.csv",row.names =FALSE,quote=FALSE)	
depth2K <- depth2Kchoose (dmodel="modp")	
setwd("../")	
sink()	Stop writing screen output to a file