

# First Estimates

Lucas Ramalho Anderson

21/10/2020

```
library ( dplyr )
library ( ggplot2 )

tempDir = "/scratch/genevol/users/lucas/"
saveDir = "/raid/genevol/users/lucas/heritability/plots/"
```

## Introduction

### Step 1

After the removal of monomorphisms, filtration of the desired samples and after obtaining the list of non-correlated snp's per chromosome (considering correlation value of  $\sqrt{0.1}$ ), it is now desired to calculate the GRM matrix.

```
# Read file with all chromosomes
# allChrFile = SeqArray::seqOpen ( paste0 ( tempDir , "allChr.gds" ) )
# List of all genes of interest (after pruning)
listGenes = readRDS ( paste0 ( tempDir , "fullPrunedList.rds" ) )

# GRM - calculated as defined in CGTA
# grm_obj = SNPRelate::snpgdsGRM( allChrFile , snp.id = listGenes , method = "GCTA")

# Estimating through "gaston" package
altReadSnps = gaston::read.vcf( paste0 ( tempDir , "allChr.vcf.gz" ) )

## ped stats and snps stats have been set.
## 'p' has been set.
## 'mu' and 'sigma' have been set.

# setting "p" parameter - correction with mean "p" and std sqrt(2p(1-2p))
gaston::standardize( altReadSnps ) <- "p"
grm_matrix = gaston::as.matrix ( altReadSnps )
# grm_scaled = scale( grm_matrix , center = T , scale = T )
# grm_scaled = readRDS (paste0(tempDir , "scaledMatrixBk.rds"))

# manual_GRM = ( 1 / nrow ( grm_scaled ) ) * grm_scaled %*% t ( grm_scaled )
# GRM matrix calculation (GCTA)
grm_alt_p = gaston::GRM ( altReadSnps , which.snps = listGenes )

## Warning in which.snps & is.autosome(x@snps$chr): longer object length is not a
## multiple of shorter object length
```

```

# transform matrix into dataframe (3 columns - col1 = samples each row, col2 = samples each column , c
dfGrm = reshape2::melt(grm_alt_p)

# indexing with numeric values each sample (columns and rows)
dfGrm$sampLines = rep ( seq ( 1 , nrow ( grm_alt_p ) ) , nrow ( grm_alt_p ) )
dfGrm$sampCols = sort ( rep ( seq ( 1 , nrow ( grm_alt_p ) ) , nrow ( grm_alt_p ) ) )

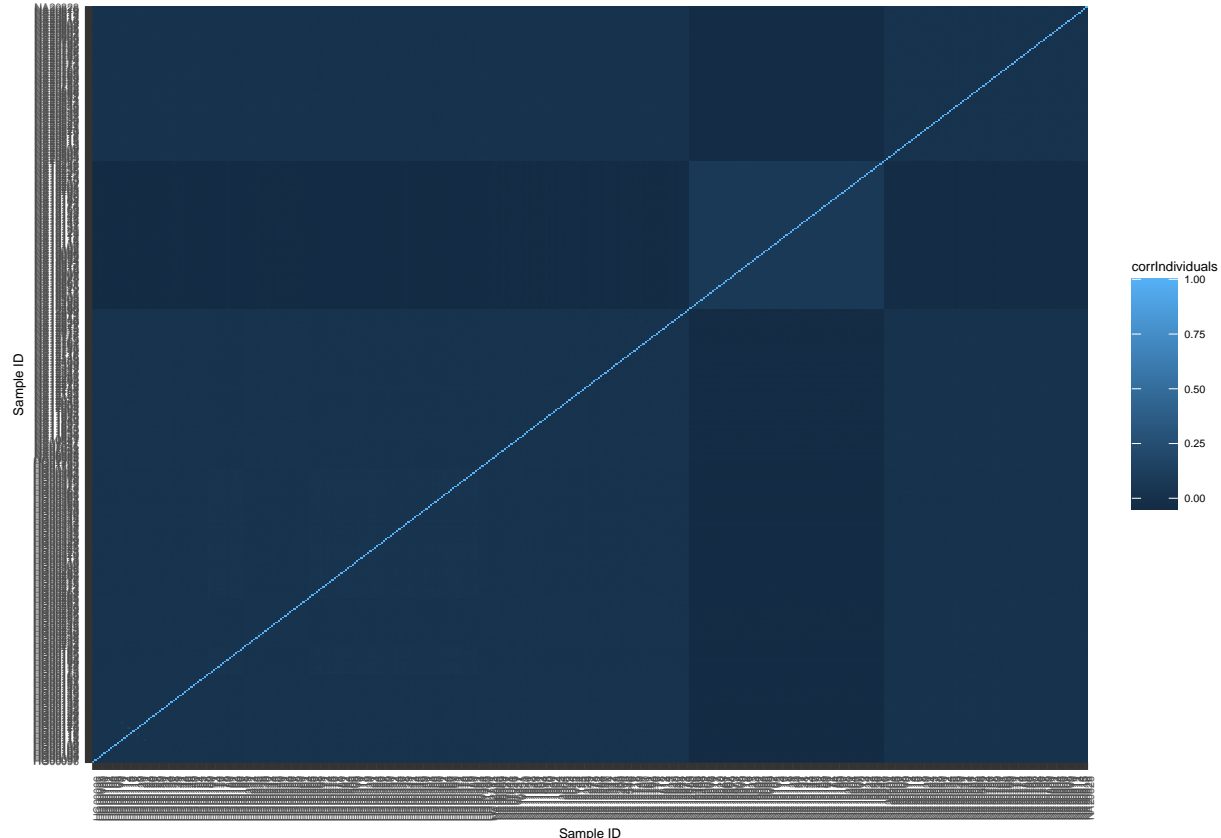
# To calculate the correlation between individuals, the calculation  $A_{ij}/\sqrt{A_{ii}}\sqrt{A_{jj}}$  will be d
# dataframe with only diag. values
dfGrmDiag = dfGrm[dfGrm$sampLines==dfGrm$sampCols,]
# sqrt of those values
dfGrmDiag = dfGrmDiag %>% mutate ( sqrtVal = sqrt ( value ) , sqrtVal2 = sqrt ( value ) )

# merging each  $A_{ii}$  for each row and col
dfGrmM = merge ( dfGrm , dfGrmDiag[ ,c ( "sqrtVal" , "sampLines" ) ] , on = c ( "sampLines" ) )
dfGrmM2 = merge ( dfGrmM , dfGrmDiag[ ,c ( "sqrtVal2" , "sampCols" ) ] , on = c ( "sampCols" ) )

# Calculating  $A_{ij}/(\sqrt{A_{ii}}\sqrt{A_{jj}})$ 
dfGrmFinal = dfGrmM2 %>% mutate ( corrIndividuals = value / ( sqrtVal * sqrtVal2 ) ) %>% arrange ( Var1

# plot heatmap - correlation between individuals
dfGrmFinal %>% ggplot( aes ( x = Var1 , y = Var2 , fill = corrIndividuals ) ) +
geom_tile() +
theme( axis.text.x = element_text(angle = 90, hjust = 1) , text = element_text (size = 5) ) +
labs ( x = "Sample ID" , y = "Sample ID" )

```



```

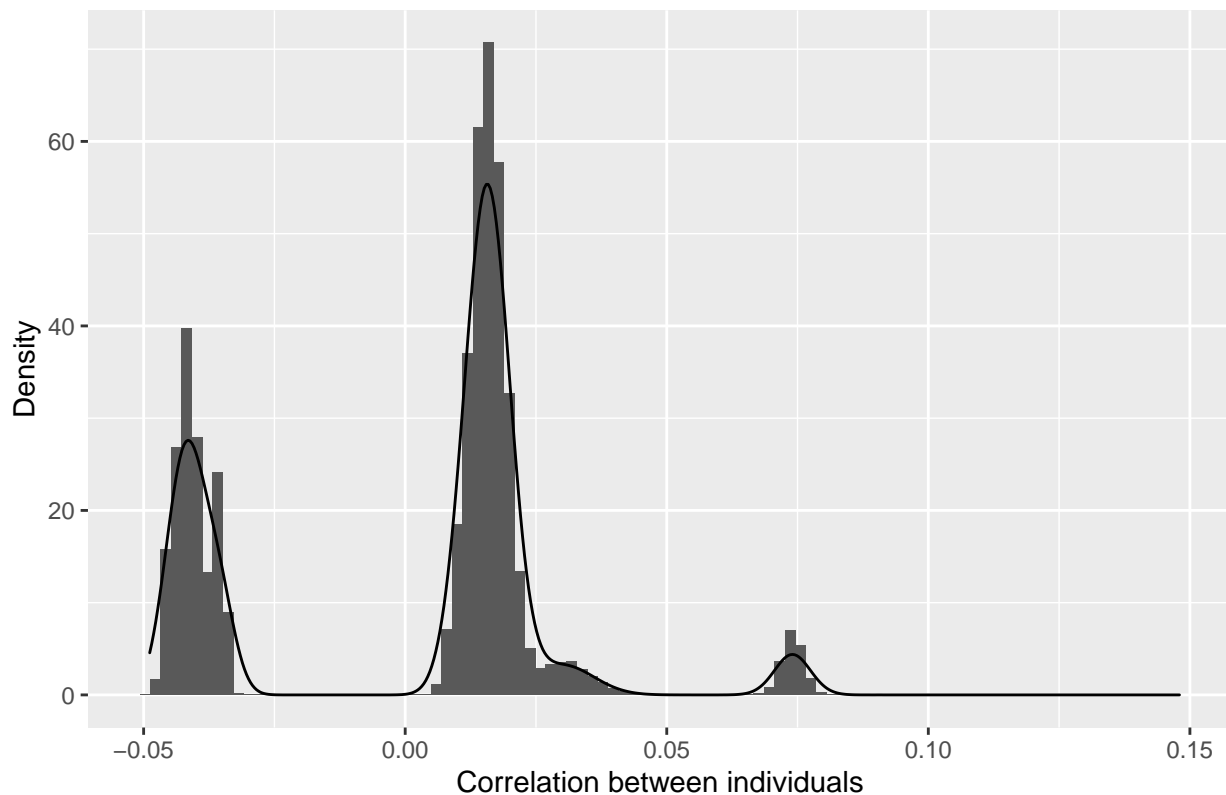
# It seems there are blocks with higher correlation between individuals between the samples

# Filter of all correlation values between individuals
dfUniqueCorr = dfGrmFinal %>% filter ( corrIndividuals < .9999 ) %>% distinct ( corrIndividuals , .keep

# Histogram and density of correlation values
dfUniqueCorr %>% ggplot ( aes ( x = corrIndividuals ) ) +
geom_histogram ( aes(y=..density..) , bins = 100 ) +
geom_density ( ) +
labs ( x = "Correlation between individuals" , y = "Density" , title = "Histogram of correlation between

```

Histogram of correlation between distinct individuals



```

# The correlation blocks are bolder in this plot

# Readind file with HLA expressions and ancestry information
hlaExp = readr::read_tsv("/raid/genevol/heritability/hla_expression.tsv")

## Parsed with column specification:
## cols(
##   subject_id = col_character(),
##   continental_pop = col_character(),
##   population = col_character(),
##   sex = col_character(),
##   gene_name = col_character(),
##   NumReads = col_double(),
##   TPM = col_double()
## )

```

```

# Ancestry of all samples
ancestry = unique ( hlaExp[ , c ( "subject_id" , "continental_pop" )] )

# Merging ancestry info with correlation dataframe
check = merge ( dfUniqueCorr , ancestry , by.x = c ( "Var1" ) , by.y = c ( "subject_id" ) )
check2 = merge ( check , ancestry , by.x = c ( "Var2" ) , by.y = c ( "subject_id" ) )

tableAncestry = unique ( check[,c("continental_pop" , "Var1")] ) %>% select ( continental_pop ) %>% tab
knitr::kable( tableAncestry )

```

Ancestry	Freq	relFreq
AFR	87	19.59%
EUR	357	80.41%

```

# Approximately 20% of the 444 individuals are African, while the other 80% are European

# Checking the amount of comparisons between individuals with same ancestry and different ones
checkFin = check2 %>% mutate ( ancestries = ifelse ( continental_pop.x == continental_pop.y , continent

tableComparisons = table ( checkFin$ancestries ) %>% as.data.frame() %>% mutate ( freqRel = Freq/ sum (

knitr::kable ( tableComparisons )

```

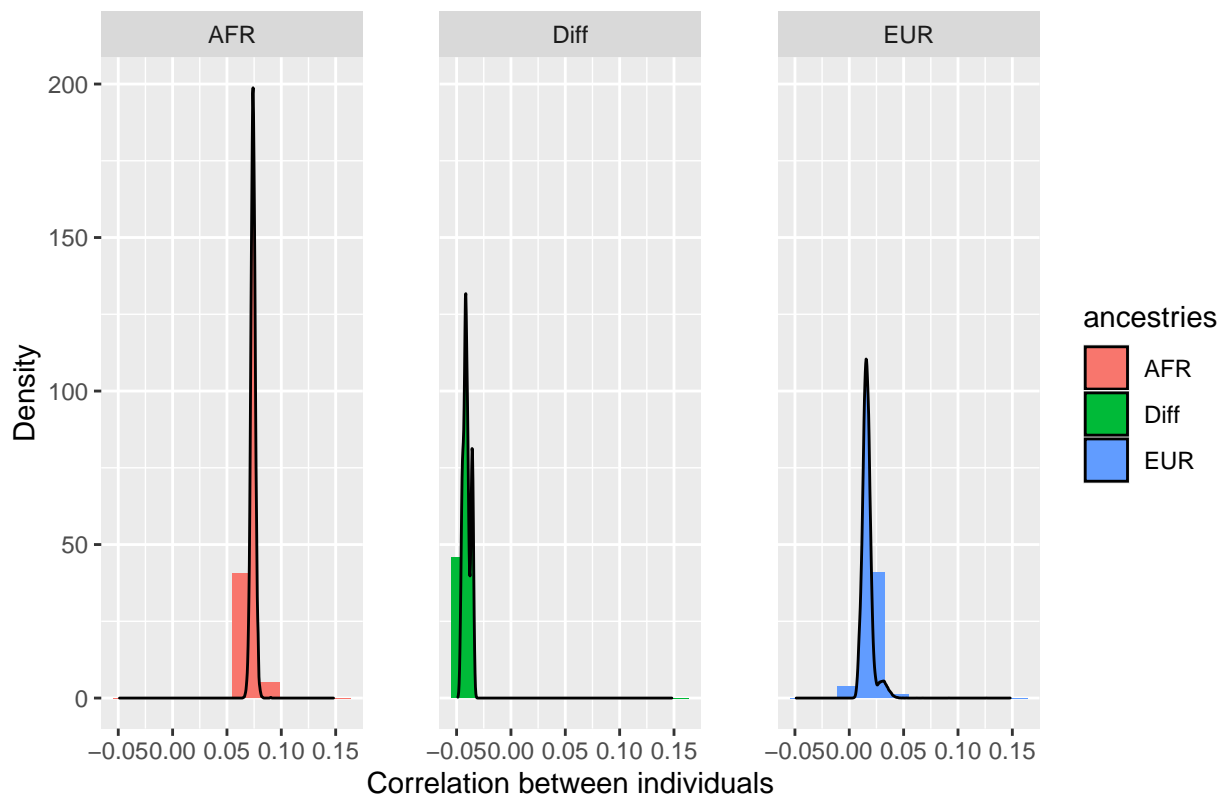
Ancestry	NumComparisons	freqRel
AFR	3741	0.0378682
Diff	31146	0.3152748
EUR	63903	0.6468570

```

checkFin %>% ggplot ( aes ( x = corrIndividuals , fill = ancestries ) ) +
geom_histogram ( aes(y=..density..) , bins = 10 ) +
geom_density ( ) +
facet_wrap ( ~ancestries ) +
theme(panel.spacing = unit (2, "lines")) +
labs ( x = "Correlation between individuals" , y = "Density" , title = "Histogram of correlation between

```

Histogram of correlation between distinct individuals



```
# display individuals with correlation greater than 10% in the sample
listGreatCorr = checkFin[ ( checkFin$corrIndividuals > .1 ) & ( checkFin$corrIndividuals < .999 ) , ] %>%
```

```
listGreatCorr
```

```
##      Var2      Var1 sampCols sampLines      value      sqrtVal      sqrtVal2
## 1 HG00120 HG00116      21        17 0.09009086 0.7819213 0.7784470
## 2 HG00240 HG00238      77        75 0.07855438 0.7790581 0.7979304
##      corrIndividuals continental_pop.x continental_pop.y ancestries
## 1          0.1480092                EUR                EUR                EUR
## 2          0.1263675                EUR                EUR                EUR
```

```
grm = grm_alt_p
```

```
# rownames ( grm ) = altReadSnps
```

```
# colnames ( grm ) = altReadSnps$sample.id
```

```
eigenValuesGrm = eigen ( grm )
```

```
dfEigen = eigenValuesGrm$values %>%
```

```
as.data.frame ( ) %>%
```

```
mutate ( order = 1:n() ) %>%
```

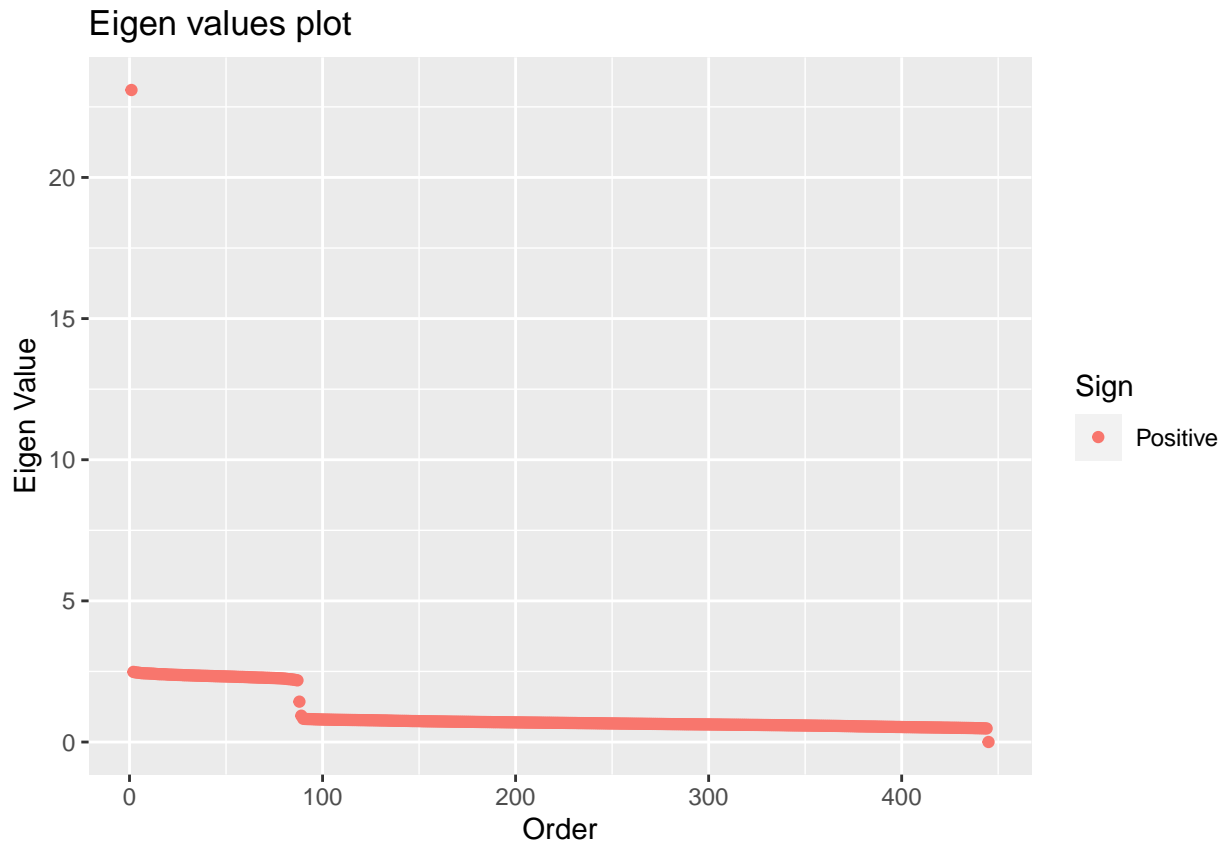
```
rename ( "Value" = '.' ) %>%
```

```
mutate ( neg = ifelse ( Value < 0 , "Negative" , "Positive" ) )
```

```
dfEigen %>% ggplot ( aes ( x = order , y = Value , colour = neg ) ) +
```

```
geom_point ( ) +
```

```
labs ( x = "Order" , y = "Eigen Value" , title = "Eigen values plot" , colour = "Sign" )
```



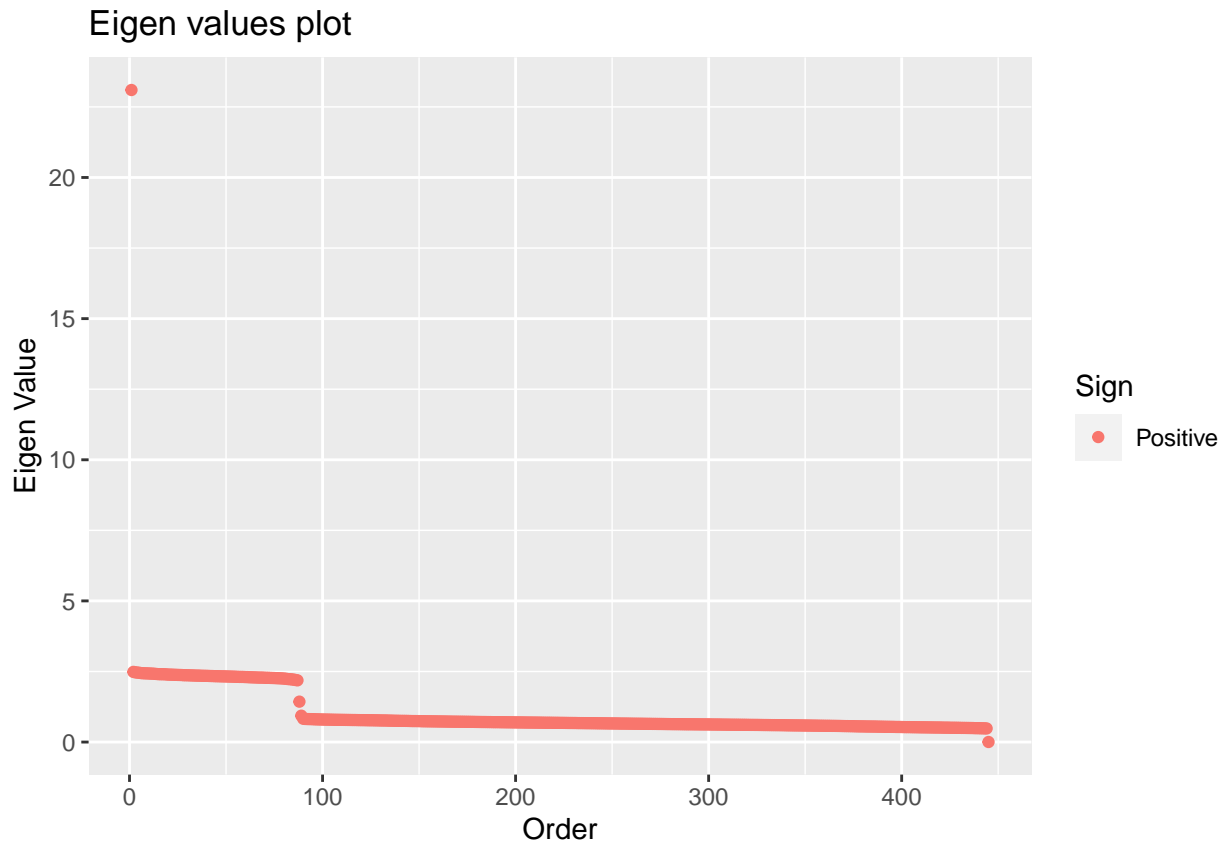
```
matrixCorrection = eigenValuesGrm$eigenvectors %*% diag( eigenValuesGrm$values + abs ( min ( eigenValuesGrm$
# rownames ( matrixCorrection ) = grm_obj$sample.id
# colnames ( matrixCorrection ) = grm_obj$sample.id

rownames ( matrixCorrection ) = rownames ( grm )
colnames ( matrixCorrection ) = colnames ( grm )

eigenCorr = eigen ( matrixCorrection )

dfEigenCorr = eigenCorr$values %>%
  as.data.frame ( ) %>%
  mutate ( order = 1:n() ) %>%
  rename ( "Value" = '.' ) %>%
  mutate ( neg = ifelse ( Value < 0 , "Negative" , "Positive" ) )

dfEigenCorr %>% ggplot ( aes ( x = order , y = Value , colour = neg ) ) +
  geom_point ( ) +
  labs ( x = "Order" , y = "Eigen Value" , title = "Eigen values plot" , colour = "Sign" )
```



```
expressionInterest = hlaExp %>% filter ( subject_id %in% colnames ( grm ) )

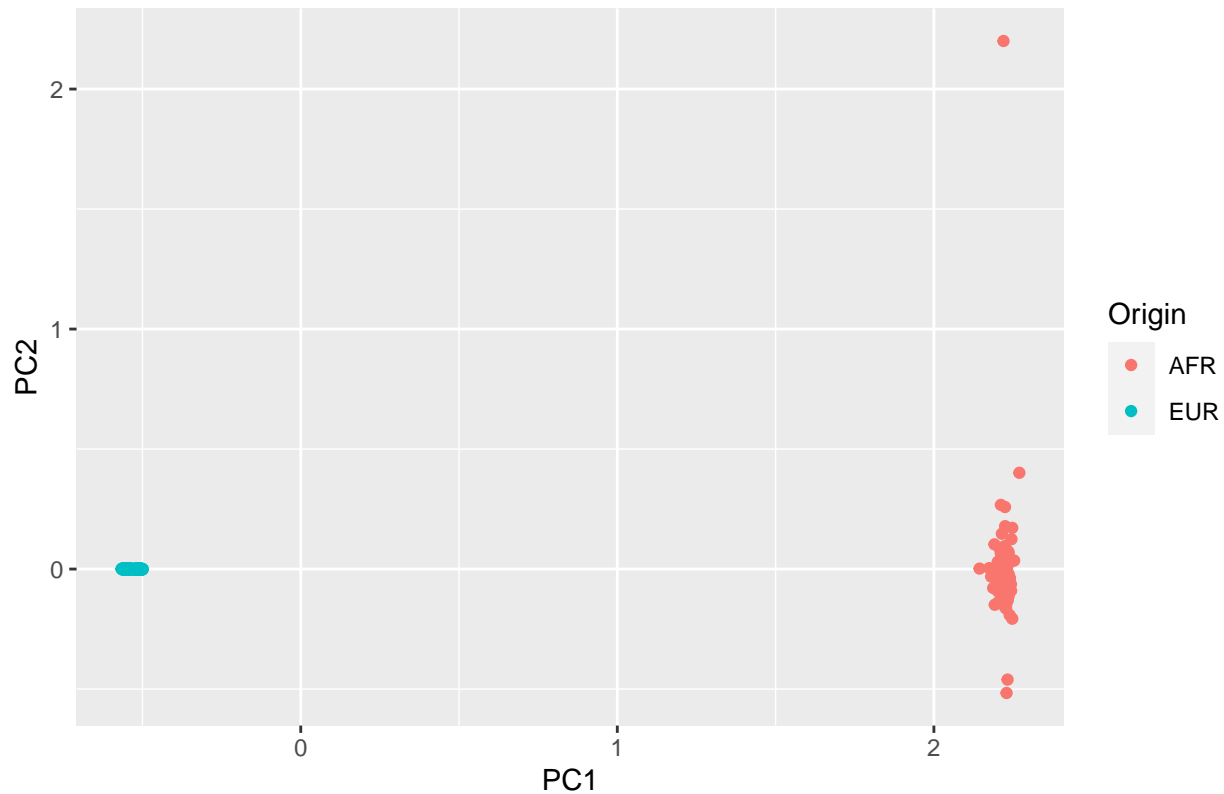
mainInfo = expressionInterest %>% distinct( subject_id , continental_pop ,population )
numEigen = 2
print ( paste0 ( "Total variation explained by the first ", numEigen , " eigen values: " , 100*round ( , 2 ) ) )

## [1] "Total variation explained by the first 2 eigen values: 5.61%"

vectors_ = eigenCorr$vectors[,1:numEigen]
calcScores = matrixCorrection %*% vectors_ %>%
  as.data.frame() %>%
  rename ( "PC1" = "V1" , "PC2" = "V2" ) %>%
  mutate ( subject_id = rownames ( matrixCorrection ) )
pcaPlot = merge ( mainInfo , calcScores )

pcaPlot %>% ggplot ( aes ( x = PC1 , y = PC2 , colour = continental_pop ) ) +
  geom_point ( ) +
  labs ( title = "PCA plot (first 2 dimensions)" , colour = "Origin" )
```

PCA plot (first 2 dimensions)



```
# geom_text ( )
```

```
simpleModels = function ( exp_ , df ){

  dfFilter = df %>% filter ( gene_name == exp_ )

  fixed0 = lm ( TPM ~ 1 , data = dfFilter )
  sum0 = summary ( fixed0 )
  fixedEffectSigma = sum0$sigma^2

  mixedModel = coxme::lmeKin( dfFilter$TPM ~ 1 + (1|dfFilter$subject_id) , data=dfFilter, varlist=list(

  mixedEffectSigma = mixedModel$sigma^2
  sigmaA = as.numeric(mixedModel$vcov)

  # comparison = mixedEffectSigma/fixedEffectSigma

  # h = sigmaA / ( sigmaA + mixedEffectSigma)

  modelExpanded = coxme::lmeKin( dfFilter$TPM ~ 1 + dfFilter$PC1 + dfFilter$PC2 + (1|dfFilter$subject_id)

  mixedEffectSigmaExp <- modelExpanded$sigma^2
  # comparisonExp = modelExpanded/fixedEffectSigma
  sigmaAExp = as.numeric(modelExpanded$vcov)

  # hExp = sigmaAExp / (sigmaAExp + mixedEffectSigmaExp )
```



```

return ( c ( exp_ , fixedEffectSigma , mixedEffectSigma , sigmaA , mixedEffectSigmaExp , sigmaAExp ) )
}

listNames = unique ( expressionInterest$gene_name )
modelDf = merge ( expressionInterest , calcScores )

requiredInfo = NULL
for ( name_ in listNames ){

  requiredInfo = rbind ( requiredInfo , simpleModels ( exp_ = name_ ,df = modelDf ) )
}

finalDf = requiredInfo %>% as.data.frame ( ) %>% rename ("Gene" = "V1" ,
                                                         "fixedSigma" = "V2" ,
                                                         "residualMixedSigma" = "V3" ,
                                                         "randomEffectSigma" = "V4" ,
                                                         "residualMixedSigmaExp" = "V5" ,
                                                         "randomEffectSigmaExp" = "V6") %>%
mutate ( fixedSigma = as.numeric ( as.character ( fixedSigma ) ) ,
         residualMixedSigma = as.numeric ( as.character (residualMixedSigma)) ,
         randomEffectSigma = as.numeric ( as.character (randomEffectSigma)) ,
         residualMixedSigmaExp = as.numeric ( as.character (residualMixedSigmaExp)) ,
         randomEffectSigmaExp = as.numeric ( as.character (randomEffectSigmaExp))
         ) %>%
mutate ( comparisonNull = residualMixedSigma/fixedSigma ,
         comparisonNullExp = residualMixedSigmaExp/fixedSigma ,
         hSimple = randomEffectSigma / ( randomEffectSigma + residualMixedSigma ) ,
         hExpanded = randomEffectSigmaExp / ( randomEffectSigmaExp + residualMixedSigmaExp ))

finalDf %>% knitr::kable()

```

Gene	fixedSigma	residualMixedSigma	randomEffectSigma	residualMixedSigmaExp	randomEffectSigmaExp	comparisonNull	comparisonNullExp	hSimple	hExpanded
HLA-A	180900.152	6073.64	51699.583	132663.39	43405.555	414 0.6969239	0.7333515	0.2908176	0.2465259
HLA-B	526383.047	1650.59	48732.587	511005.00	362.294720	0.8960216	0.9707854	0.0936475	0.0007085
HLA-C	127438.485	472.71	40355.059	88515.95	36479.149	516 0.6706982	0.6945783	0.3207166	0.2918446
HLA-DPA1	31587.292	9131.76	2066.406	30542.59	5.083768	0.9222622	0.9669267	0.0662349	0.0001664
HLA-DPB1	37625.102	5390.79	9591.940	33403.73	9.743842	0.6748363	0.8878043	0.2741908	0.0002916
HLA-DQA1	51654.464	1391.61	8644.165	48501.48	28.874006	0.8013173	0.9389602	0.1727597	0.0005950

Gene	fixed	Sigma <sup>2</sup> residual	Mixed	Sigma <sup>2</sup> Effect	Sigma <sup>2</sup> residual	Mixed	Sigma <sup>2</sup> Effect	Sigma <sup>2</sup> residual	Comparison	Null	SE	h	Expanded
HLA-DQB1	38524.9336	0.3333	2319.335	38005.50	27.971661	0.9353249	0.9865169	0.0604739	0.0007354				
HLA-DRA	390199.527	1874.60	15913.762	381119.37	32.196945	0.9530370	0.9767295	0.0410372	0.0000845				
HLA-DRB1	148507.265	5935.65	69781.856	83404.69	46950.398845	0.4439896	0.5616205	0.5141699	0.3601731				