

Quantitative Genomics and Genetics 2018 Project

Darya Akimova

May 8, 2018

All of the provided data files were imported successfully and the quality of the data was accessed as follows to ensure that

```
# is there any missing data?  
anyNA(list(geno, pheno, covars, snp_info, gene_info))  
  
[1] FALSE  
  
# are there approximately equal numbers of people in each covariate group?  
table(covars$Population)
```

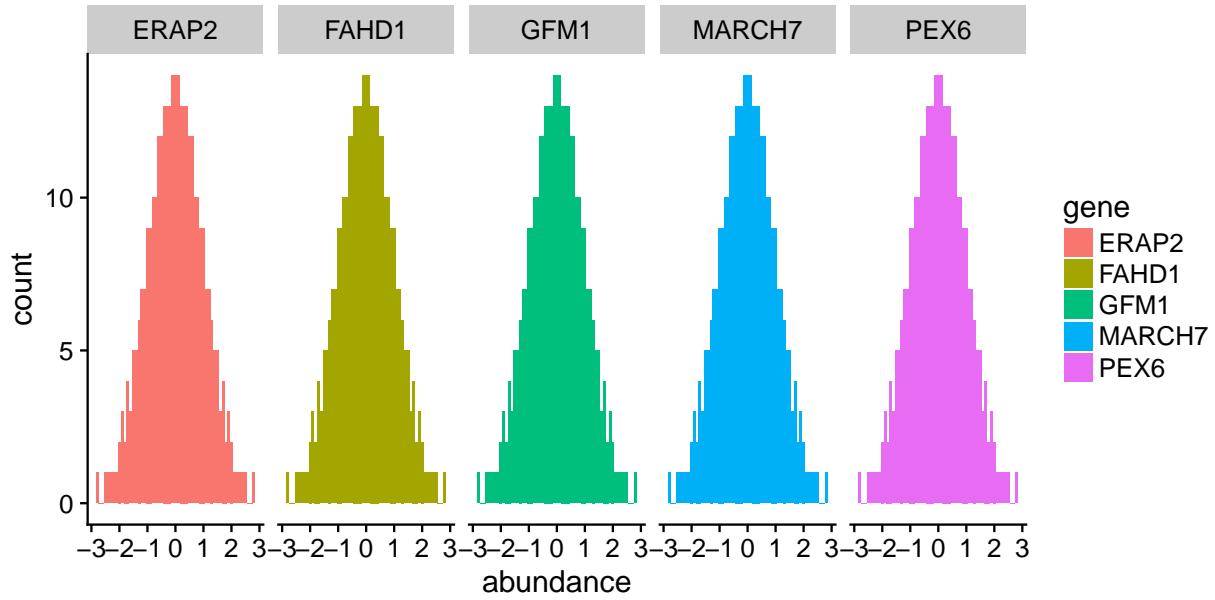
```
CEU FIN GBR TSI  
78 89 85 92  
table(covars$Sex)
```

```
FEMALE    MALE  
181       163  
  
# does the coding of the genotypes makes sense across all the data?  
table(as.matrix(geno))
```

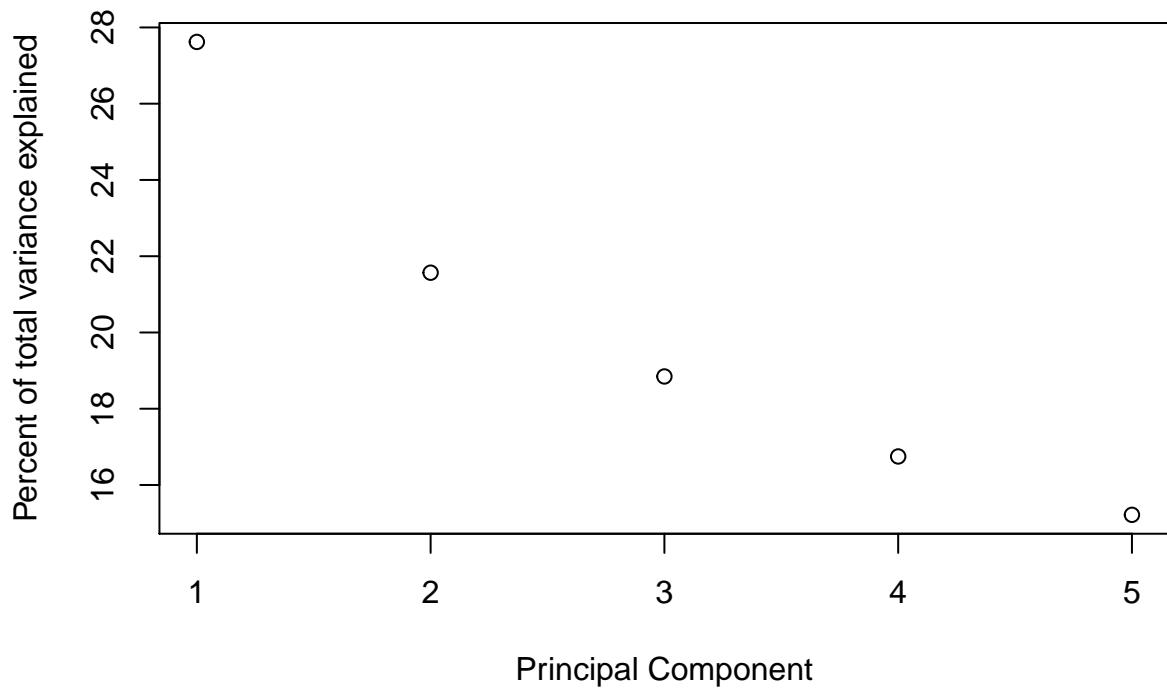
```
0      1      2  
8181444 5811217 3207339  
  
# calculate allele frequency  
geno_sums <- map_dbl(geno, function(x) sum(x) / (nrow(geno) * 2))  
# any genotypes need to be removed because of MAF < 5%?  
any(geno_sums < 0.05 | geno_sums > 0.95)
```

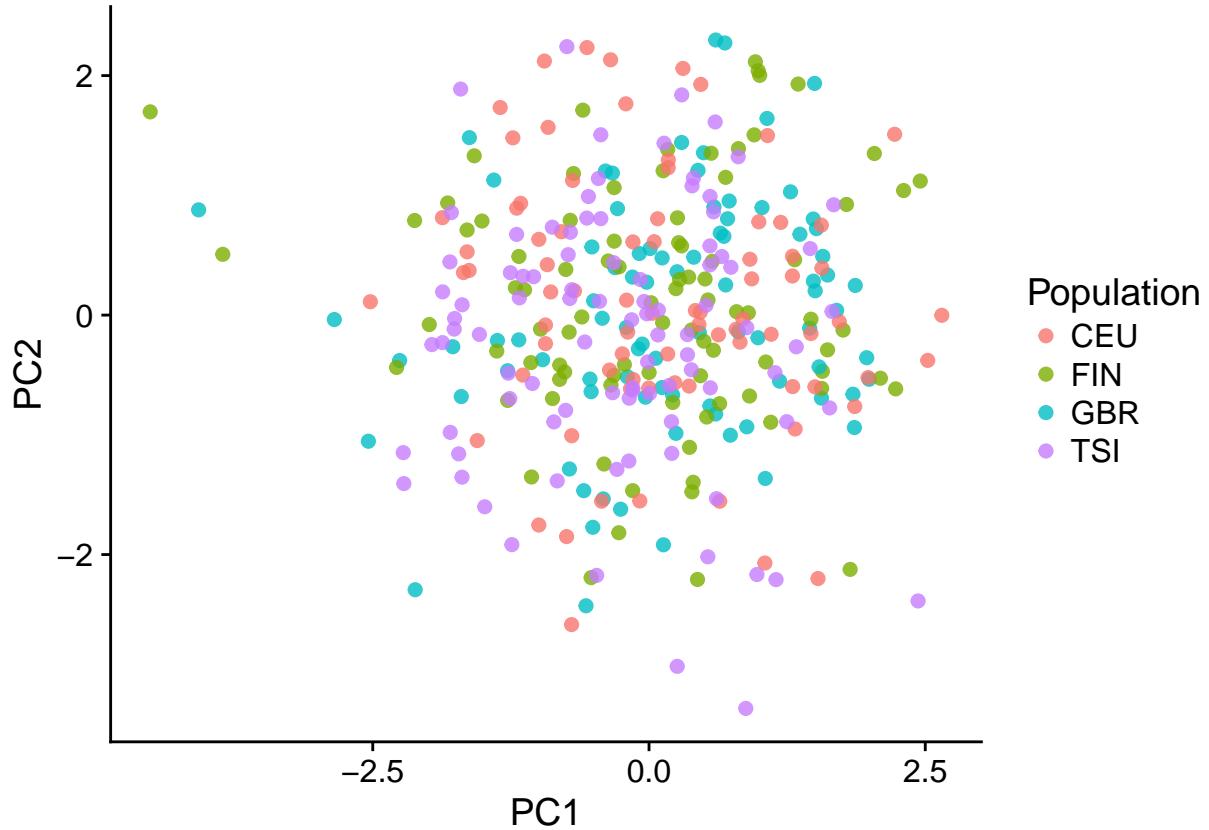
```
[1] FALSE  
  
# no genotypes need to be removed because of MAF < 5%  
# are there any genotypes without associated phenotypes, or vice versa?  
all.equal(rownames(geno), rownames(pheno))
```

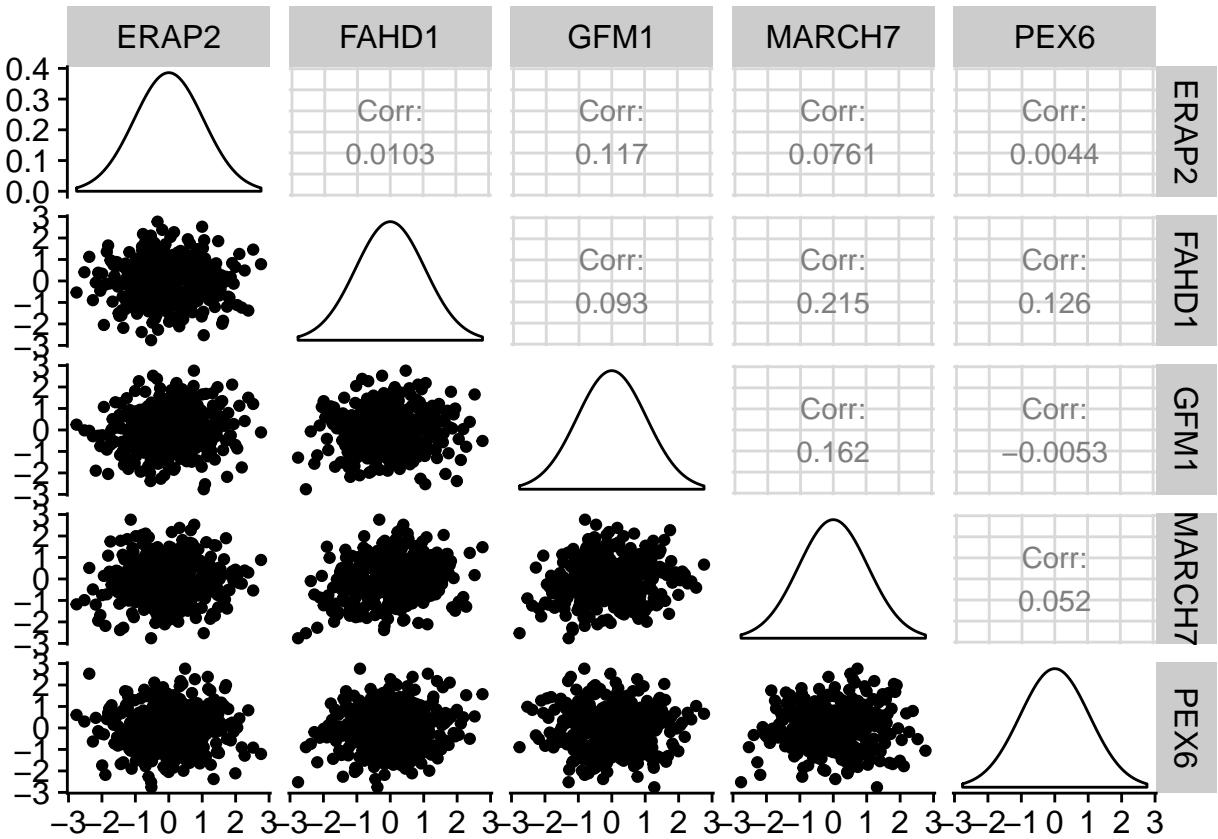
```
[1] TRUE
```

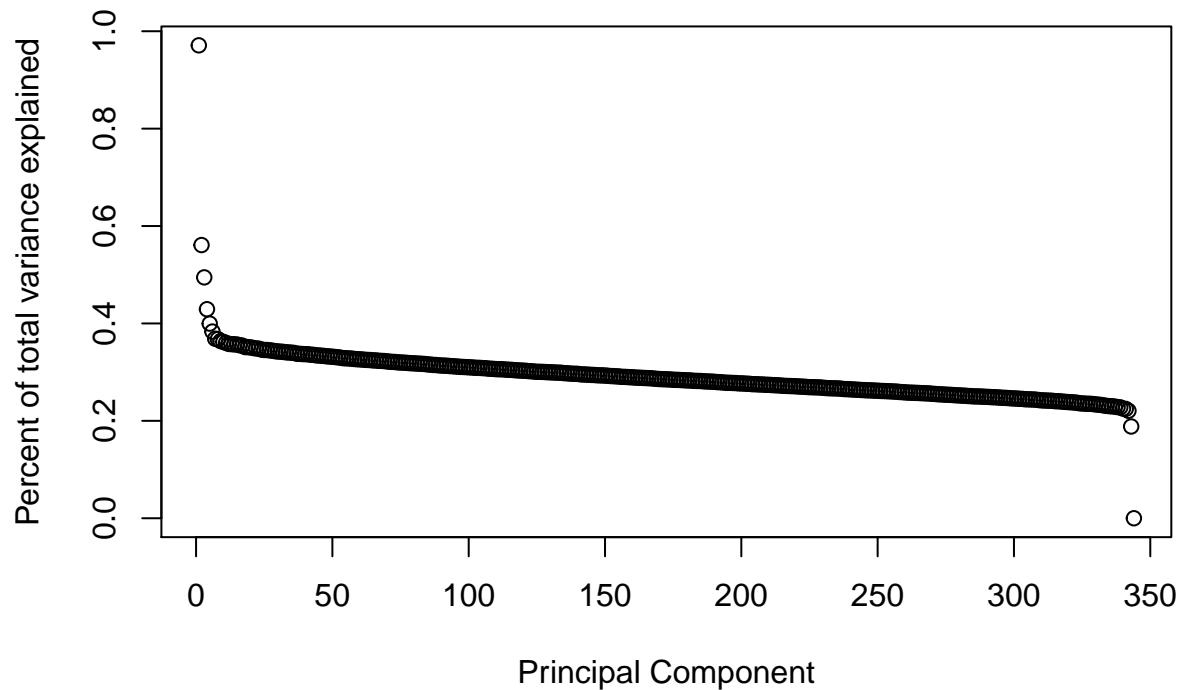


A PCA analysis was performed on the phenotype data



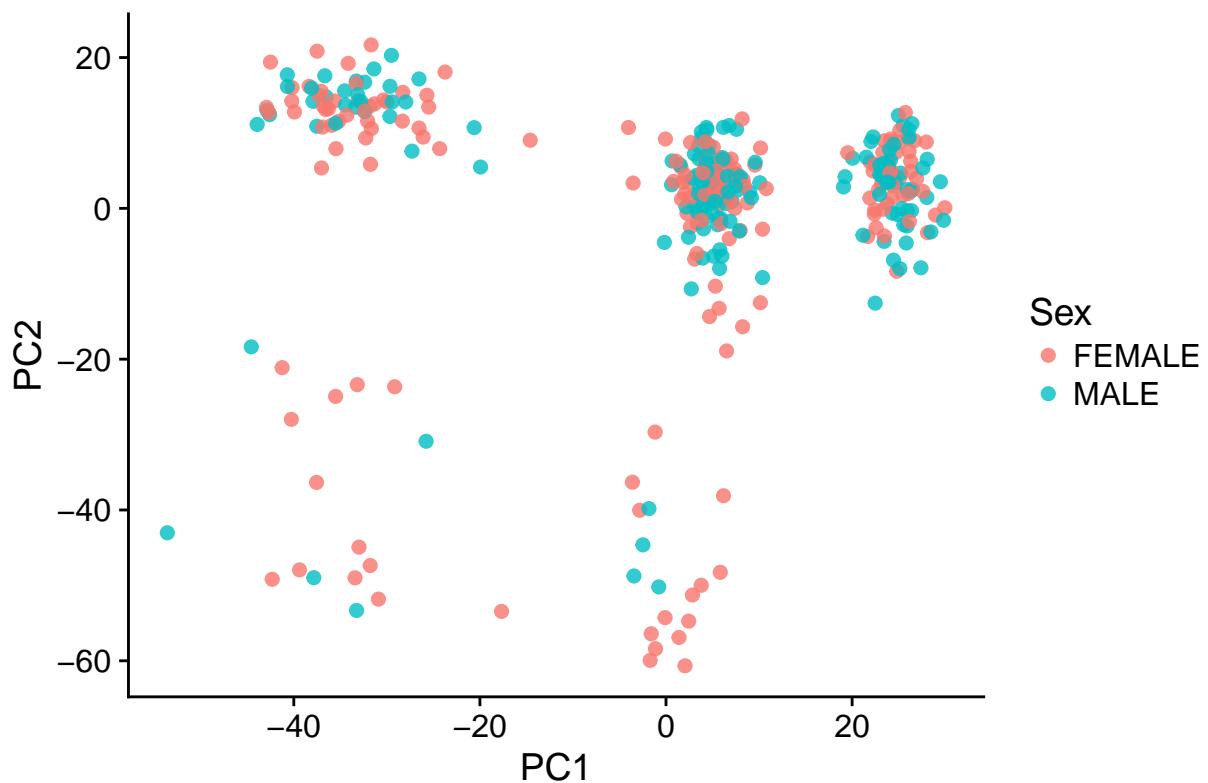




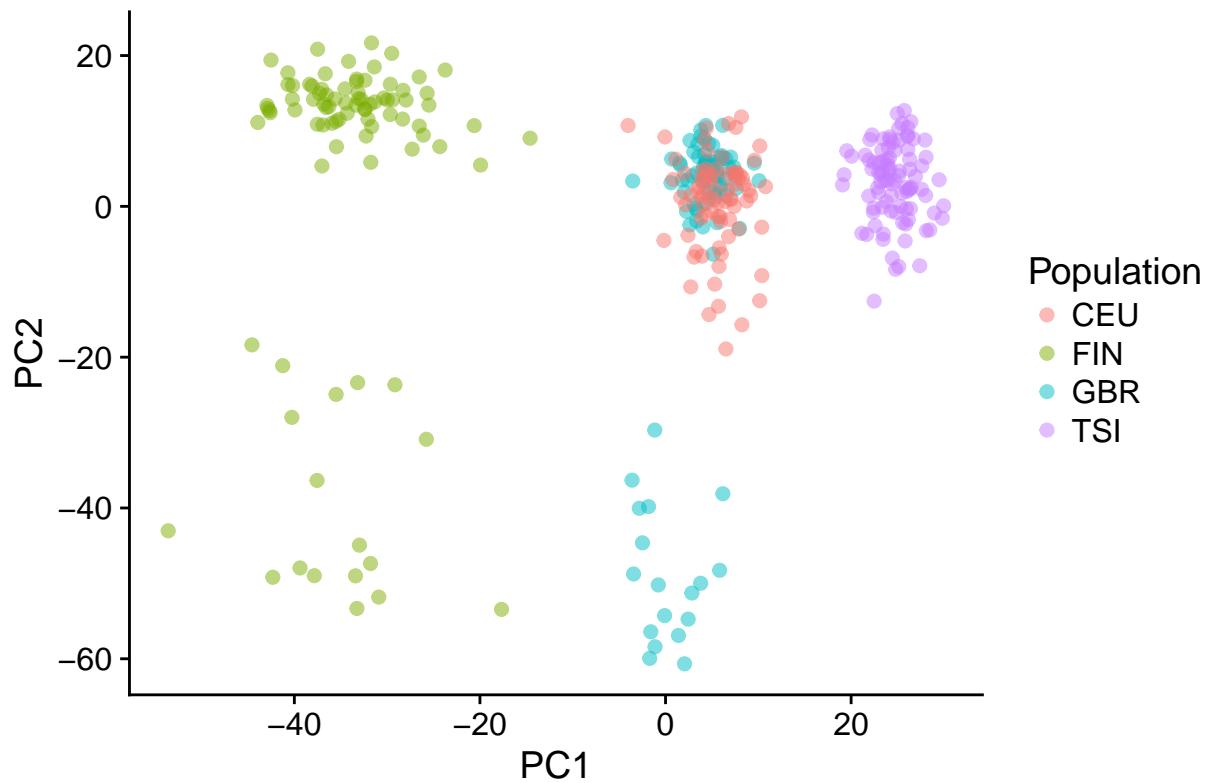


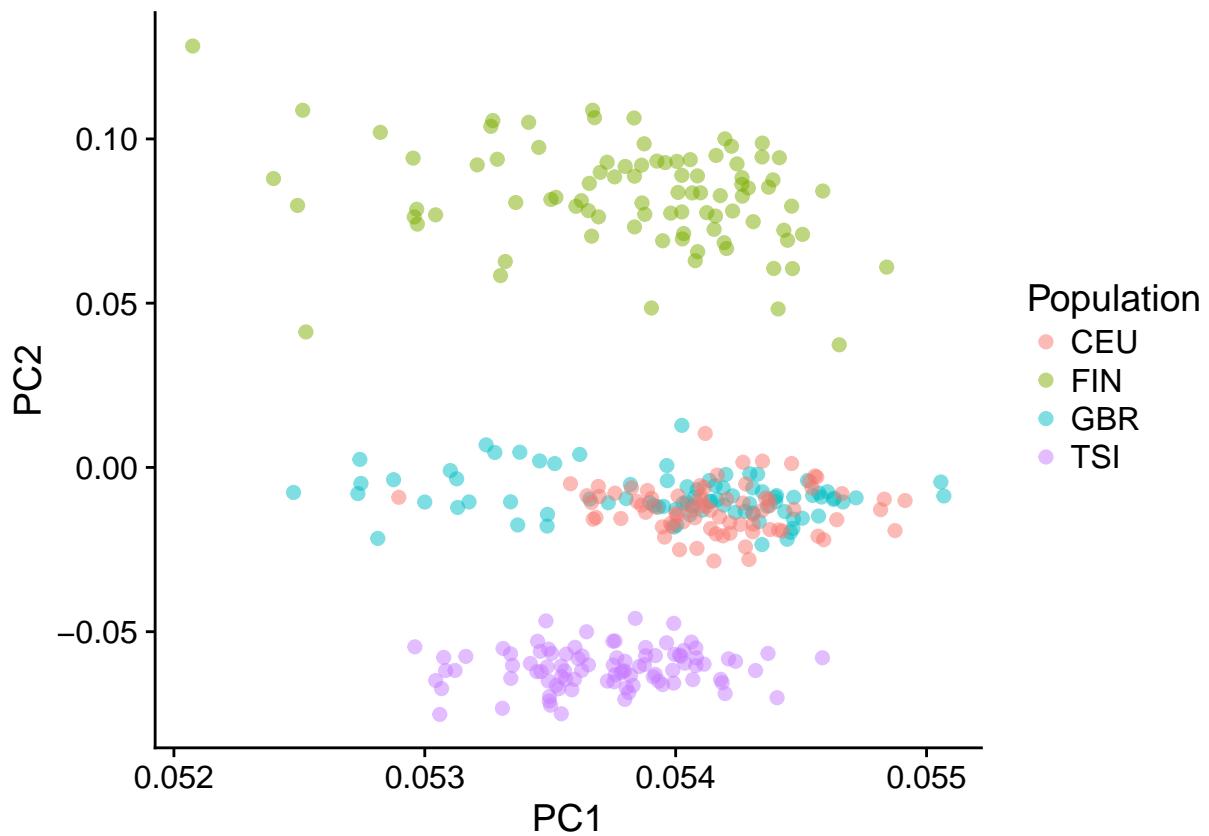
Genotype plots:

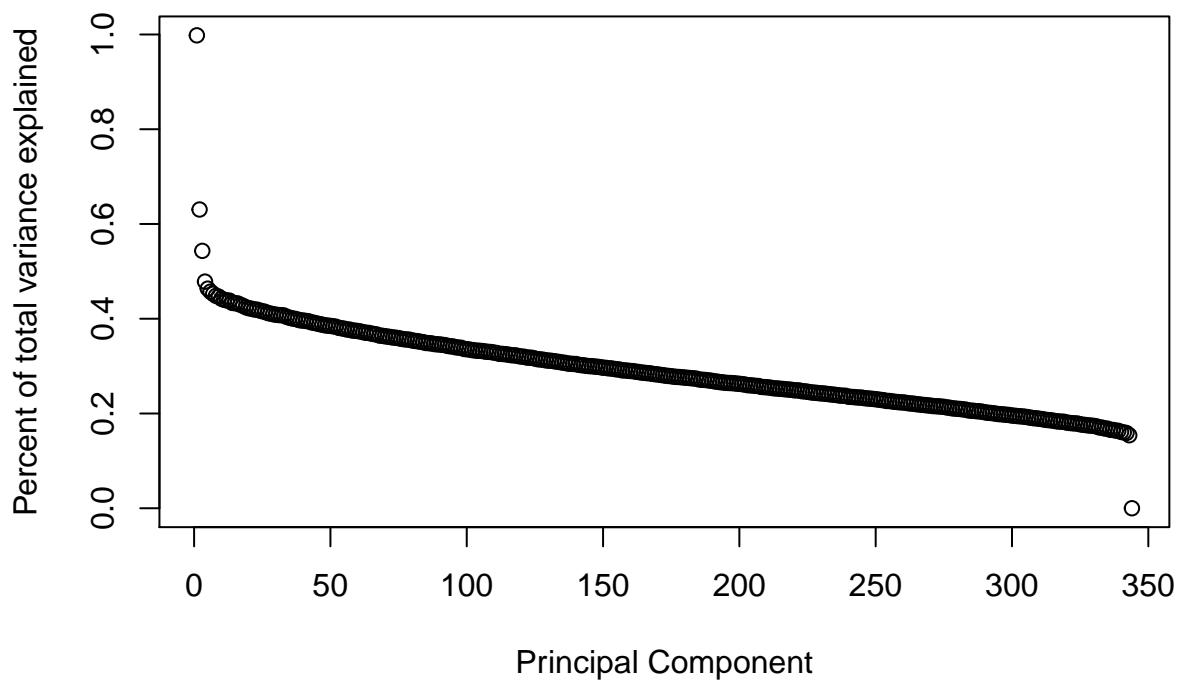
First two principal components, colored by Sex



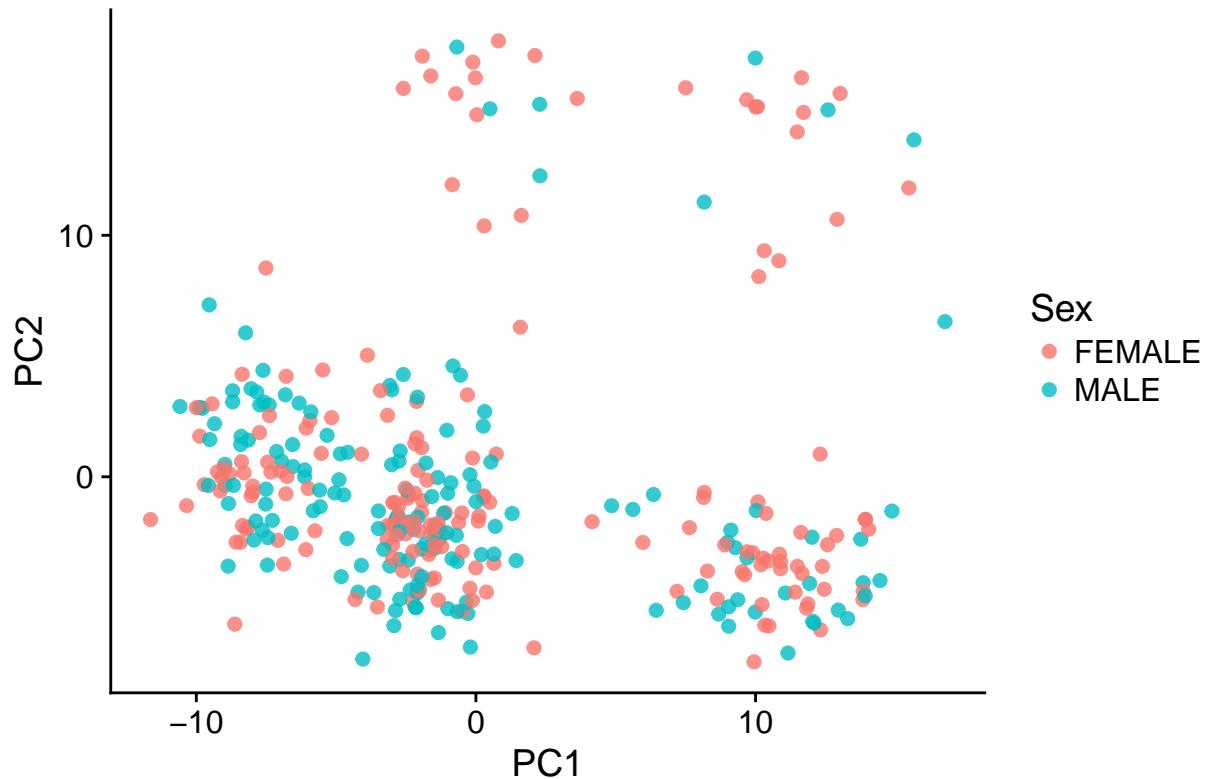
First two principal components, colored by Population



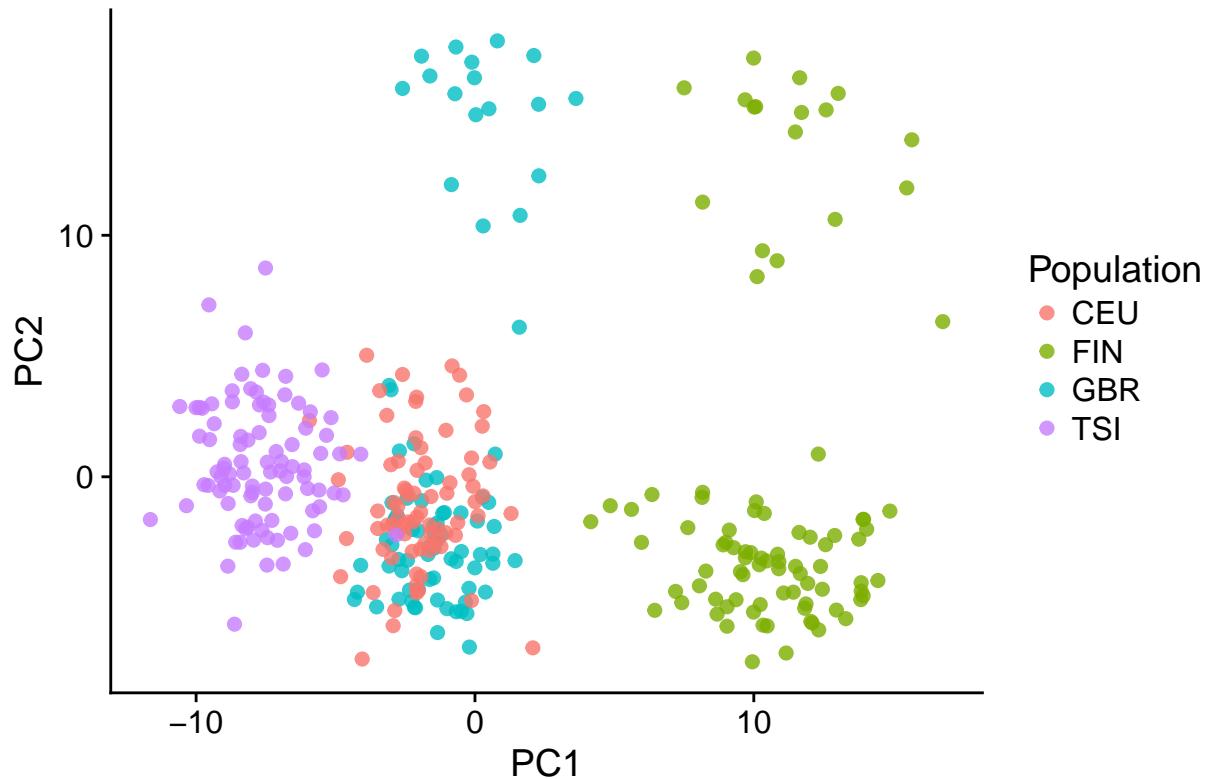




First two principal components, colored by Sex

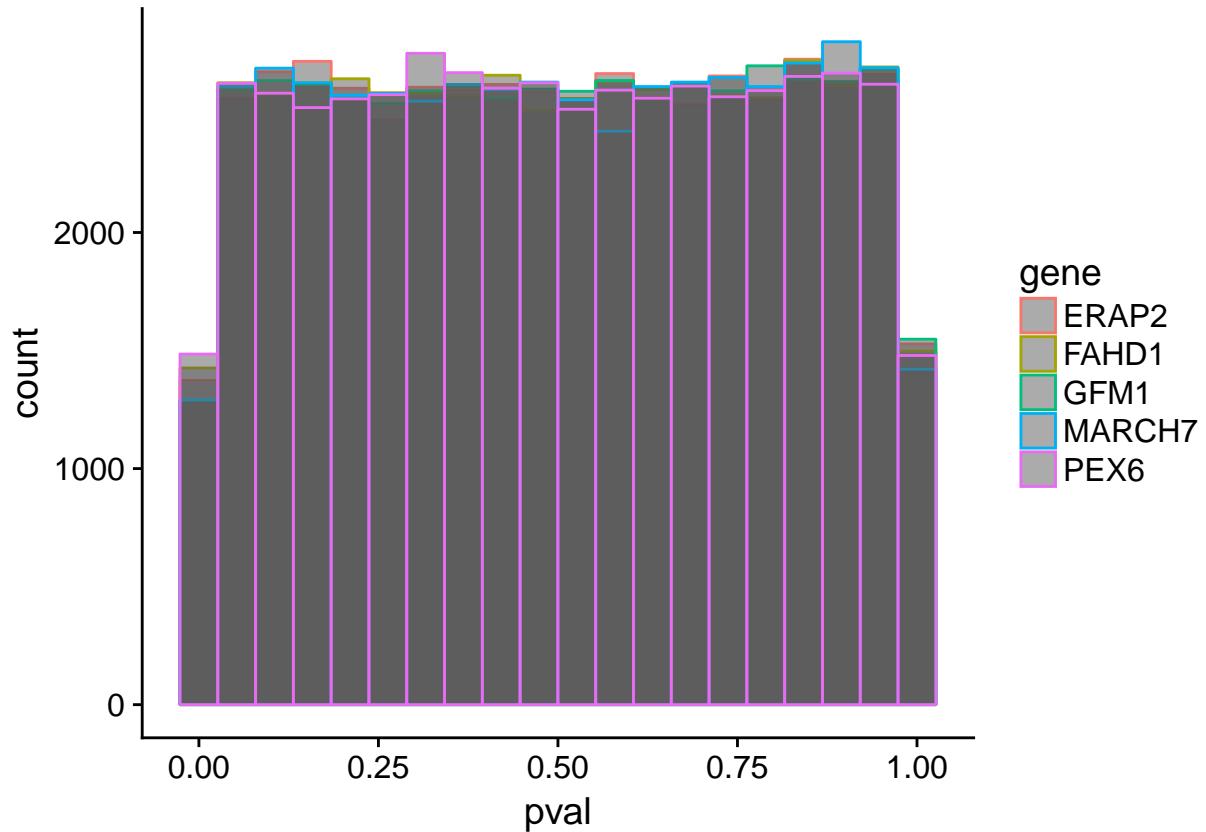


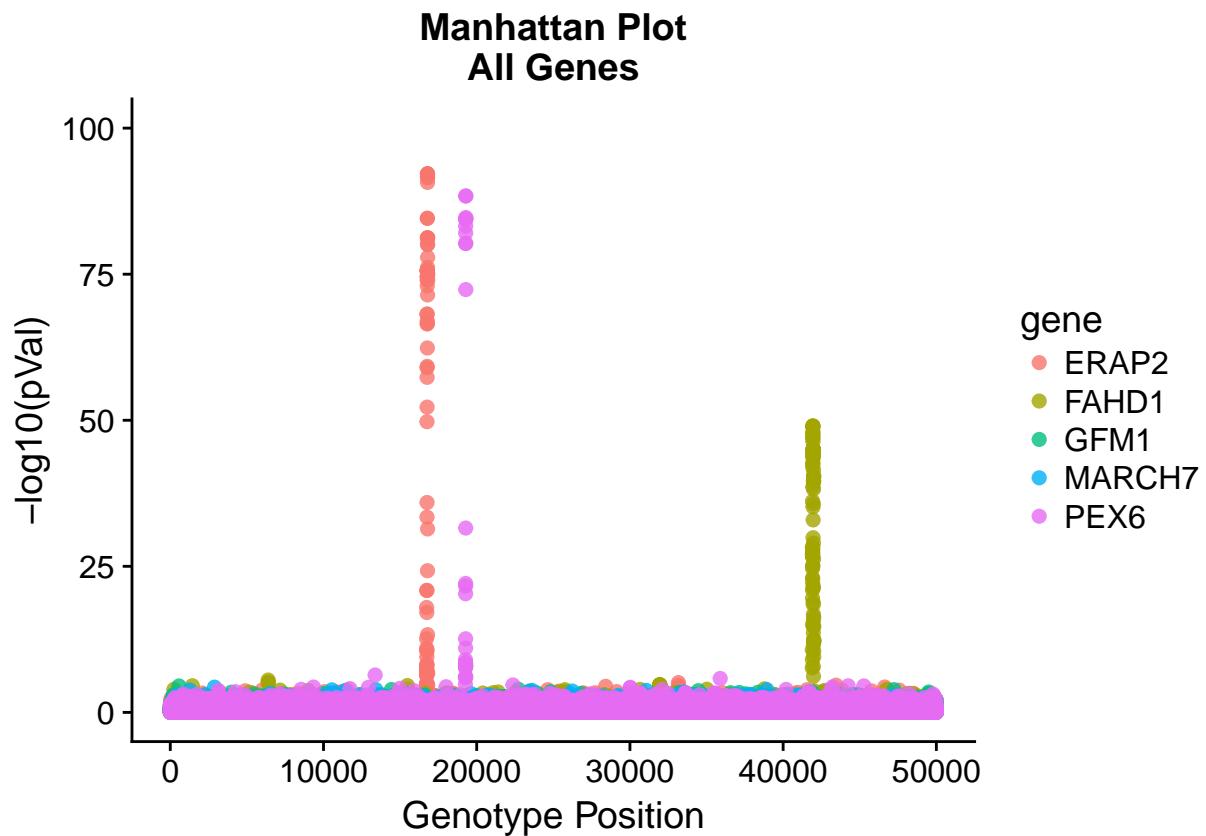
First two principal components, colored by Population



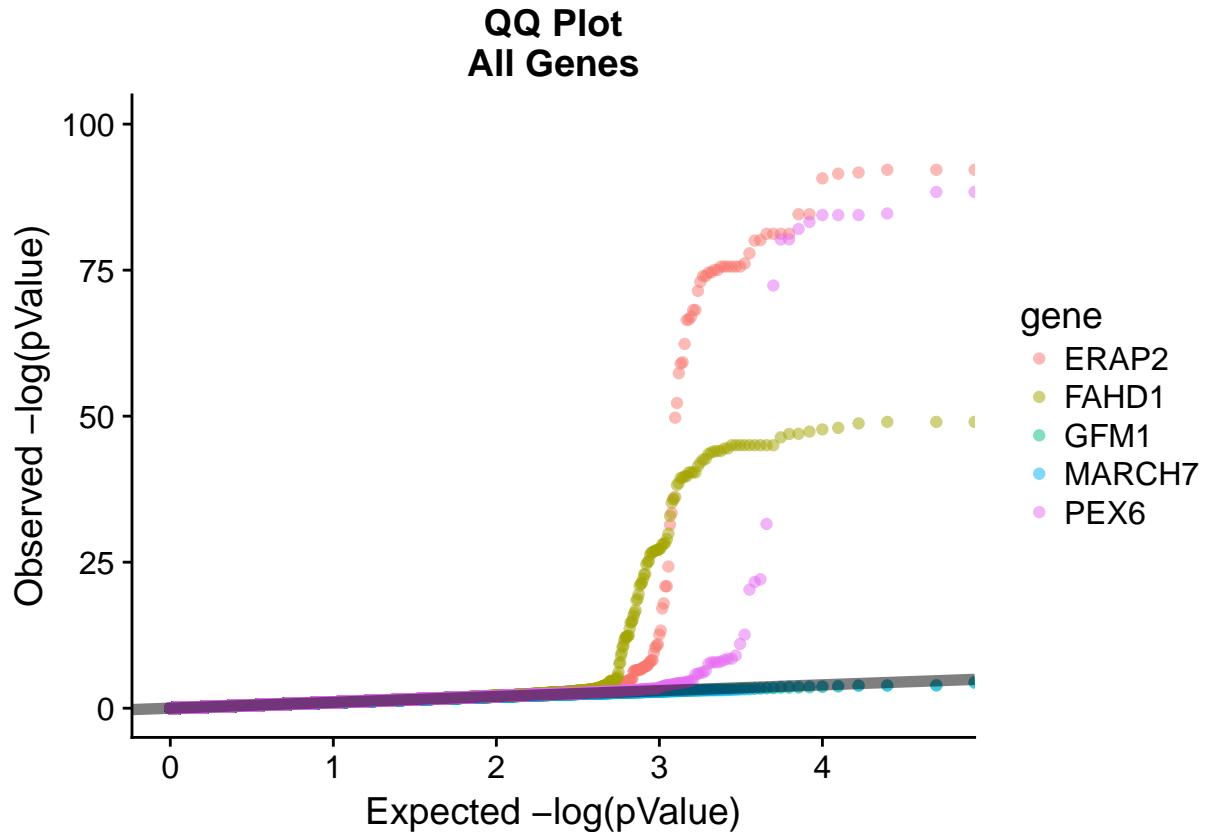
Based on the principal component analysis of the genomes, colored by the population of origin, it is possible to see that there is clearly population structure.

Test each covariate individually.





```
# A tibble: 3 x 2
# Groups:   gene [3]
  gene     n
  <chr> <int>
1 ERAP2    75
2 FAHD1    93
3 PEX6     31
```



Seems like there is a problem with just the base model. Based of the PCA analysis, it appears that the Population may be an important covariate to include in the models.

But first, let's test the relationship of each phenotype with each covariate:

```
## # A tibble: 1 x 2
##   gene   p.value
##   <chr>    <dbl>
## 1 ERAP2 0.00799

## # A tibble: 1 x 2
##   gene   p.value
##   <chr>    <dbl>
## 1 ERAP2 0.00914

## [1] FALSE

## # A tibble: 1 x 2
##   gene   p.value
##   <chr>    <dbl>
## 1 ERAP2 0.00787

## # A tibble: 3 x 6
##   gene   term           estimate std.error statistic p.value
##   <chr> <chr>          <dbl>     <dbl>      <dbl>    <dbl>
## 1 ERAP2 (Intercept) 0.399     0.140      2.86  0.00454
## 2 ERAP2 as.numeric(factor(Population~ -0.124     0.0474     -2.63  0.00901
## 3 ERAP2 factor(Sex)MALE -0.172     0.105      -1.64  0.103

## # A tibble: 1 x 2
```

```

##   gene p.value
##   <chr> <dbl>
## 1 ERAP2 0.00678

## # A tibble: 5 x 6
##   gene   term          estimate std.error statistic p.value
##   <chr> <chr>        <dbl>     <dbl>      <dbl>    <dbl>
## 1 ERAP2 (Intercept)  0.273     0.123     2.21    0.0278
## 2 ERAP2 factor(Population)FIN -0.207     0.152    -1.37    0.172
## 3 ERAP2 factor(Population)GBR -0.0606    0.153    -0.397   0.692
## 4 ERAP2 factor(Population)TSI -0.458     0.150    -3.06    0.00241
## 5 ERAP2 factor(Sex)MALE     -0.172     0.106    -1.63    0.105

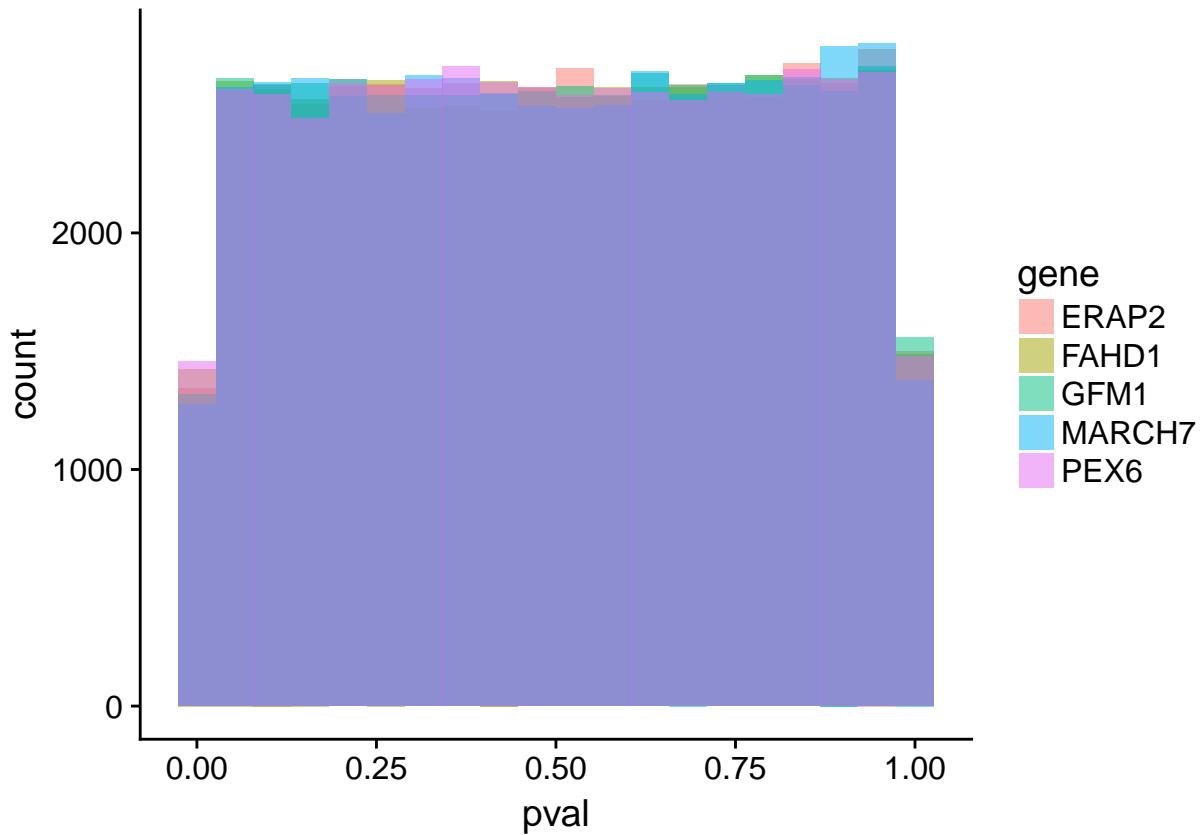
## [1] FALSE

## # A tibble: 1 x 2
##   gene p.value
##   <chr> <dbl>
## 1 ERAP2 0.00326

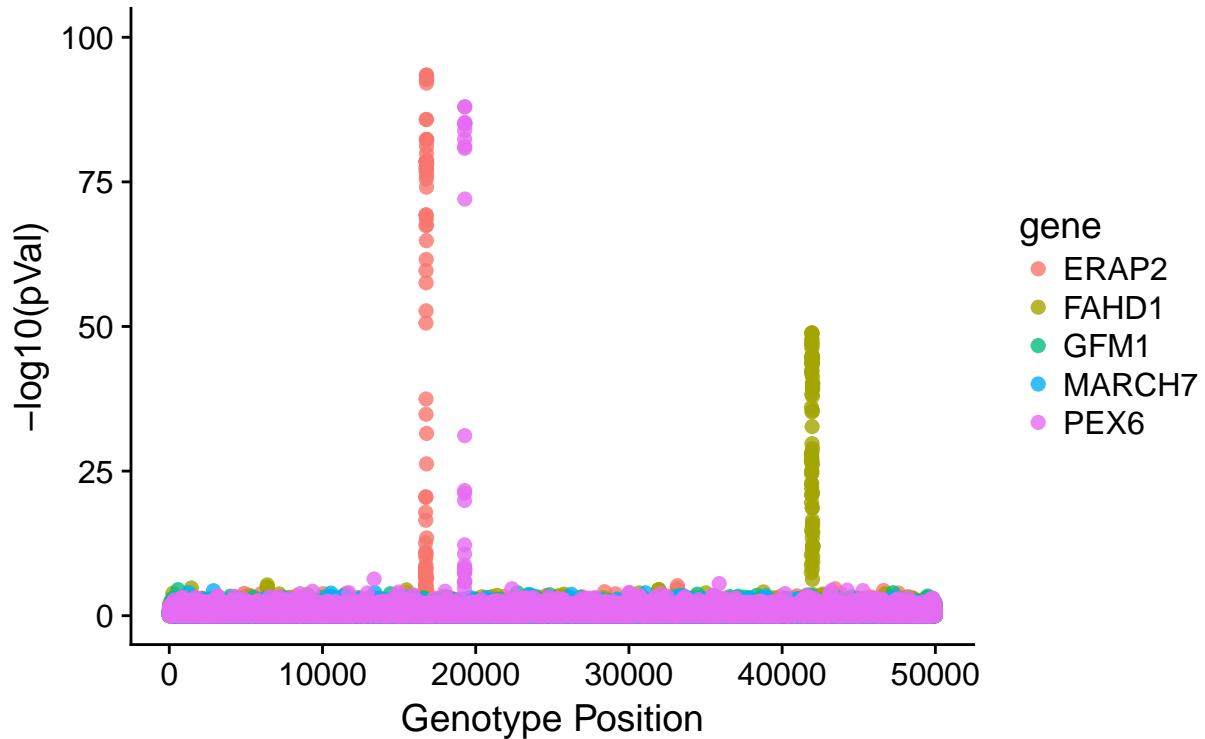
```

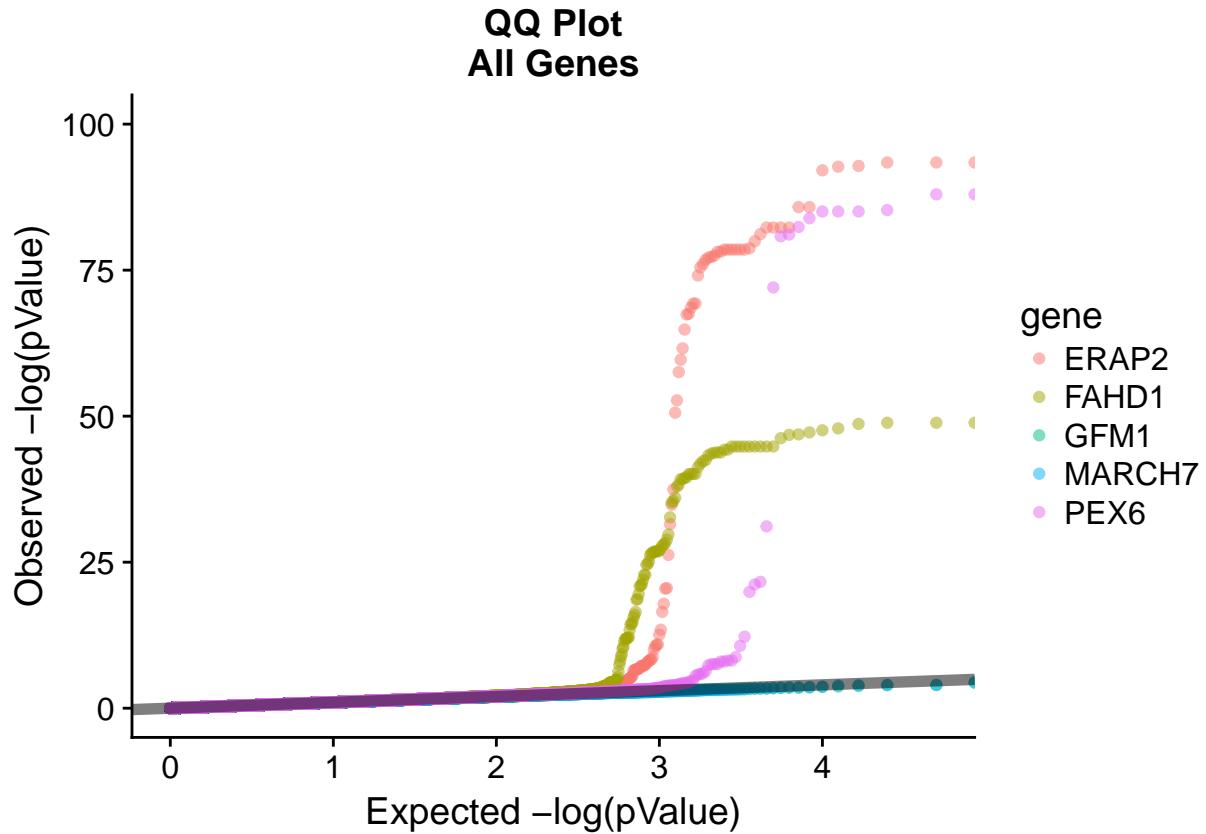
Including Population as covariate:

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805	
16725	16745	16767	16765	16785	16805	



Manhattan Plot All genes after Including Population as Covar

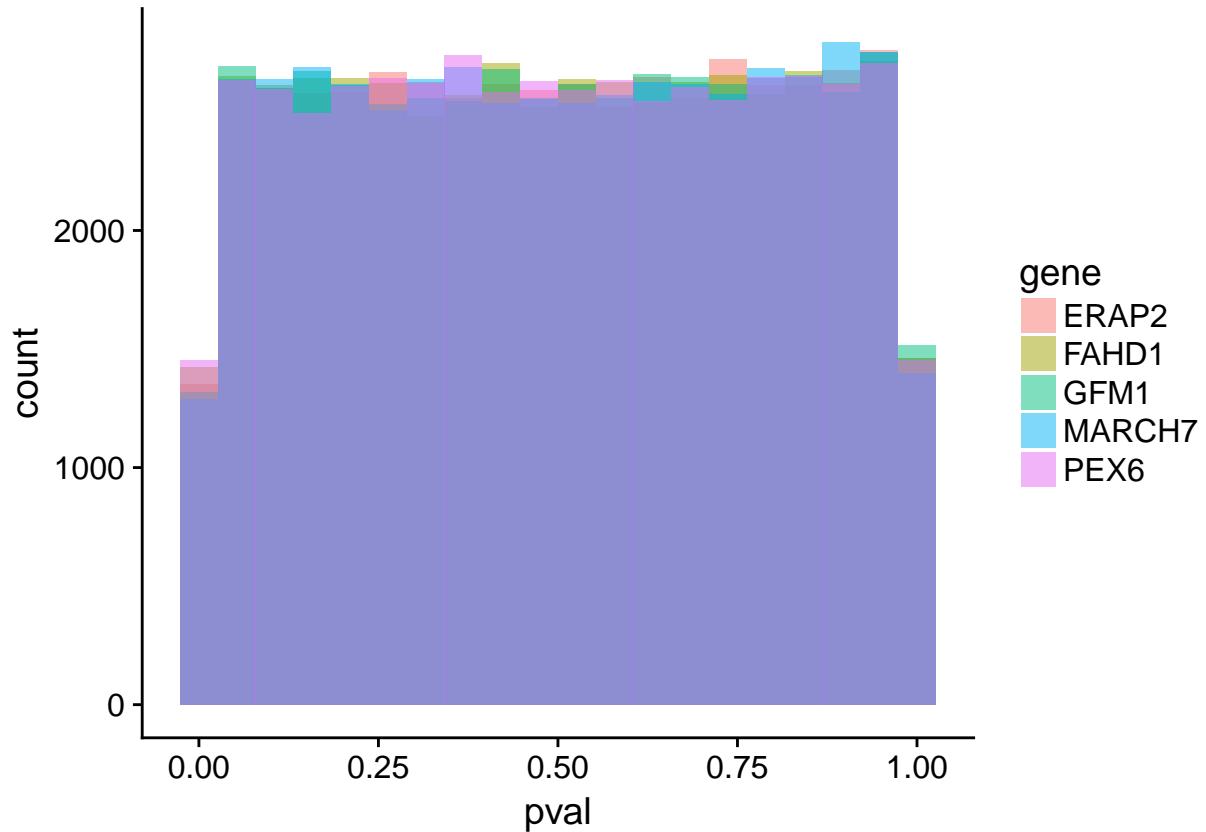




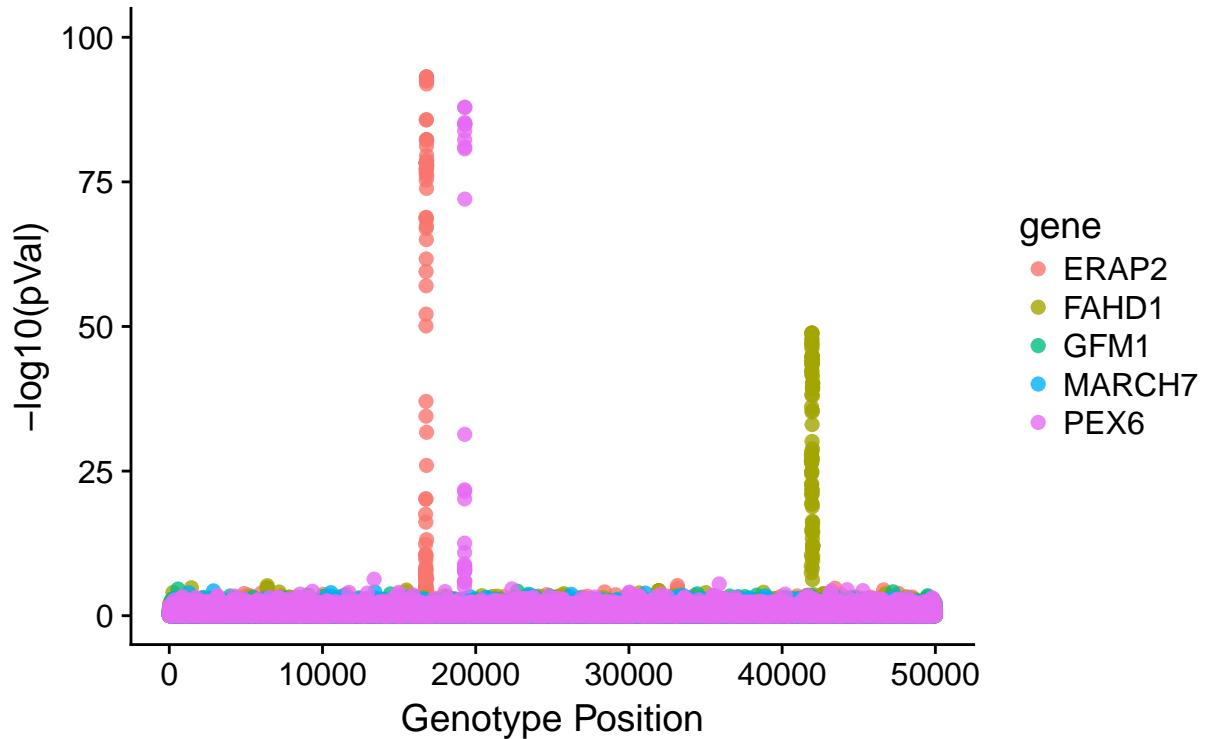
Including Population and Sex as Covariates:

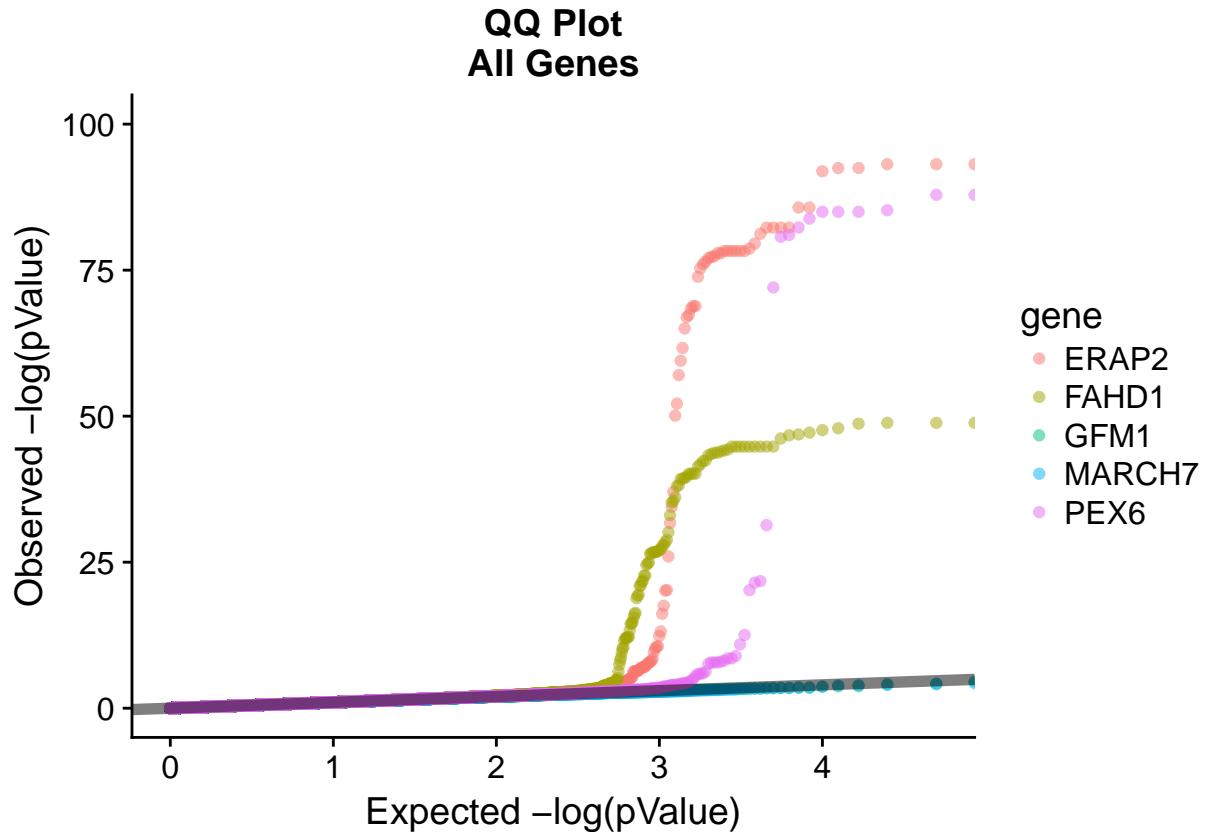
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805



Manhattan Plot All genes after Including Population and Sex as Covars

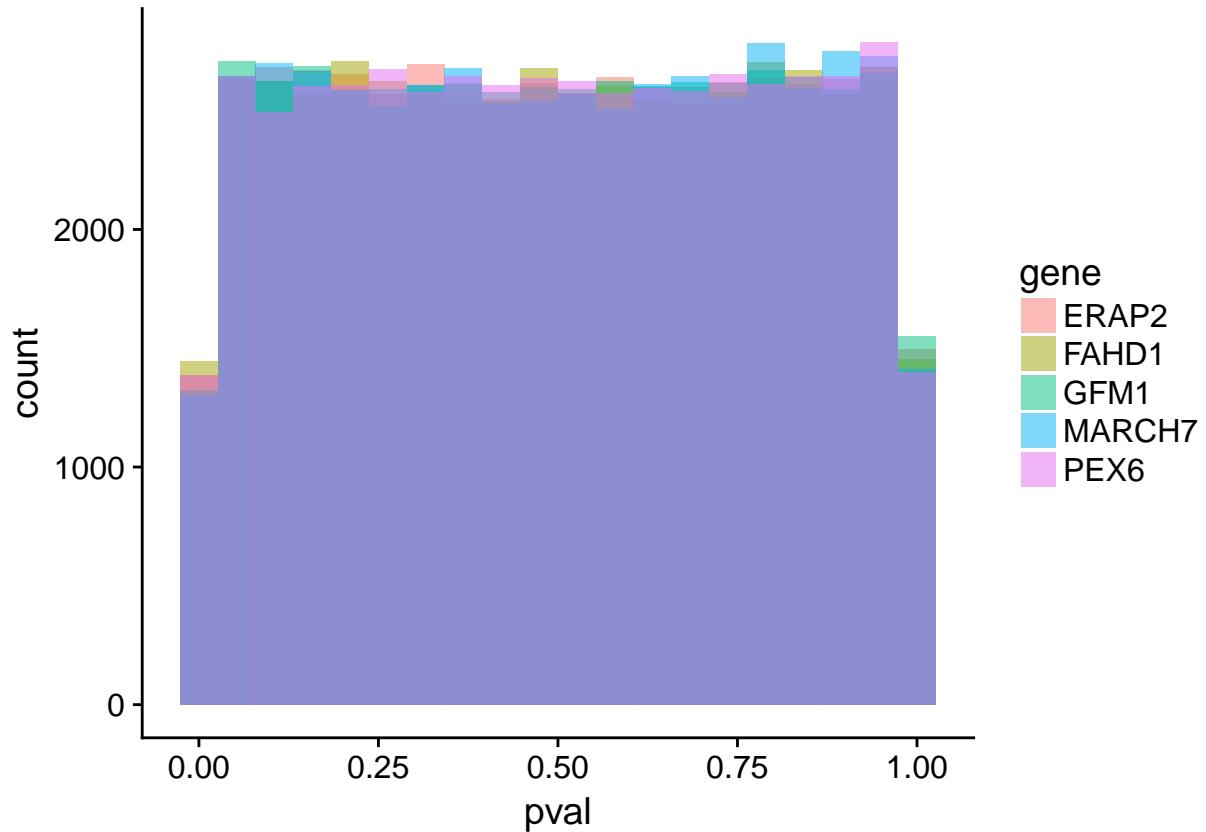


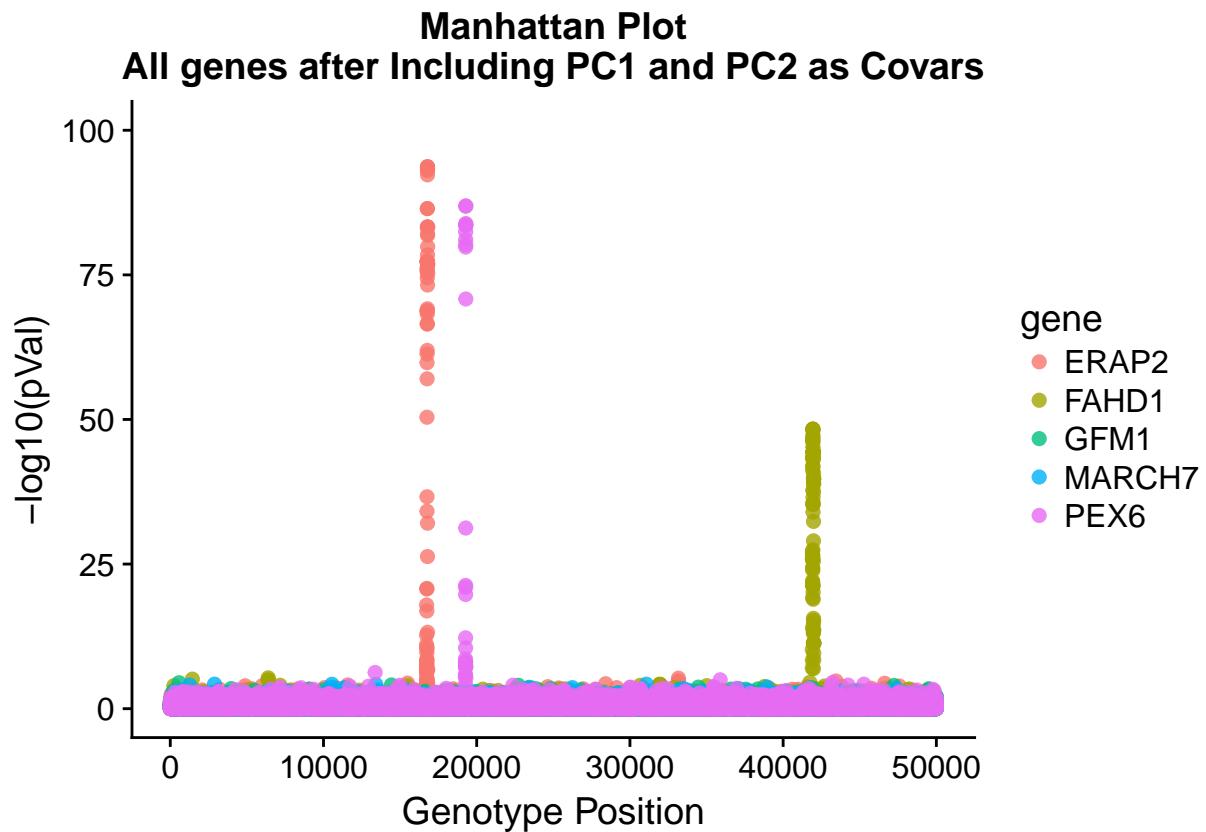


Try PC1 and PC2 as covars:

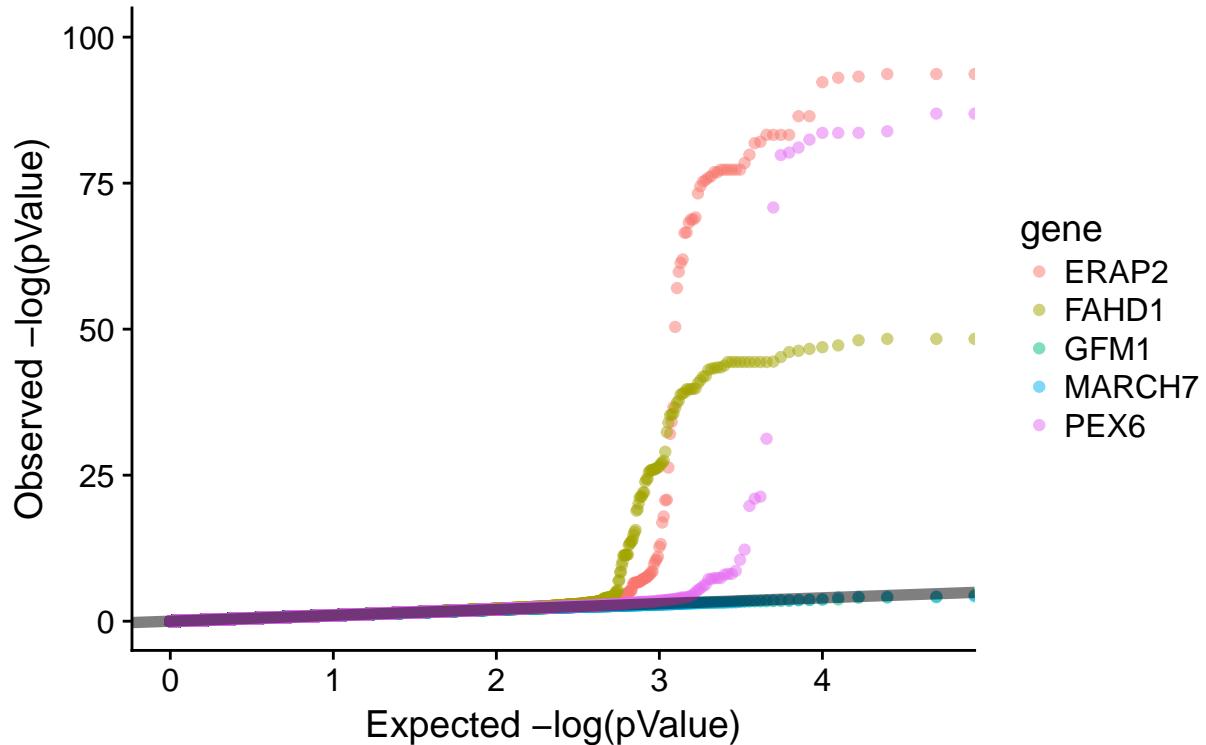
Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805





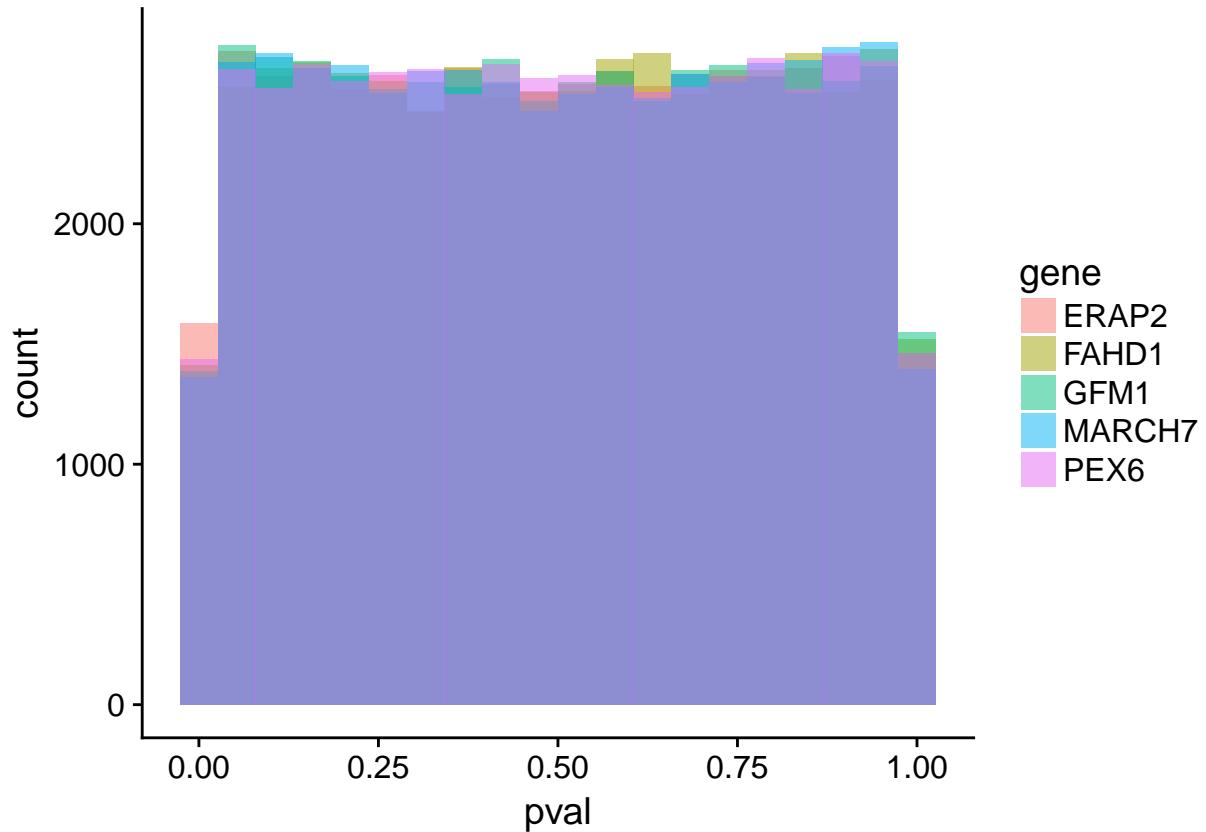
QQ Plot All Genes with PC1 and PC2 as Covars



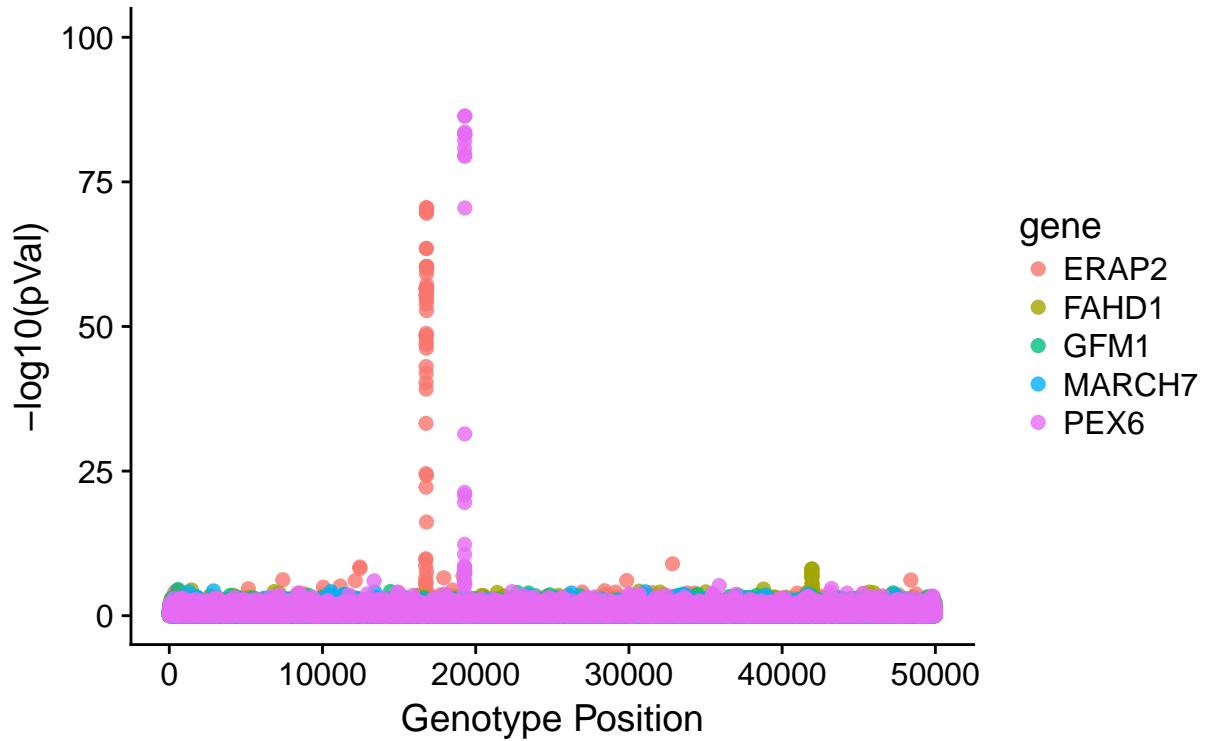
Even more PCs

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805	

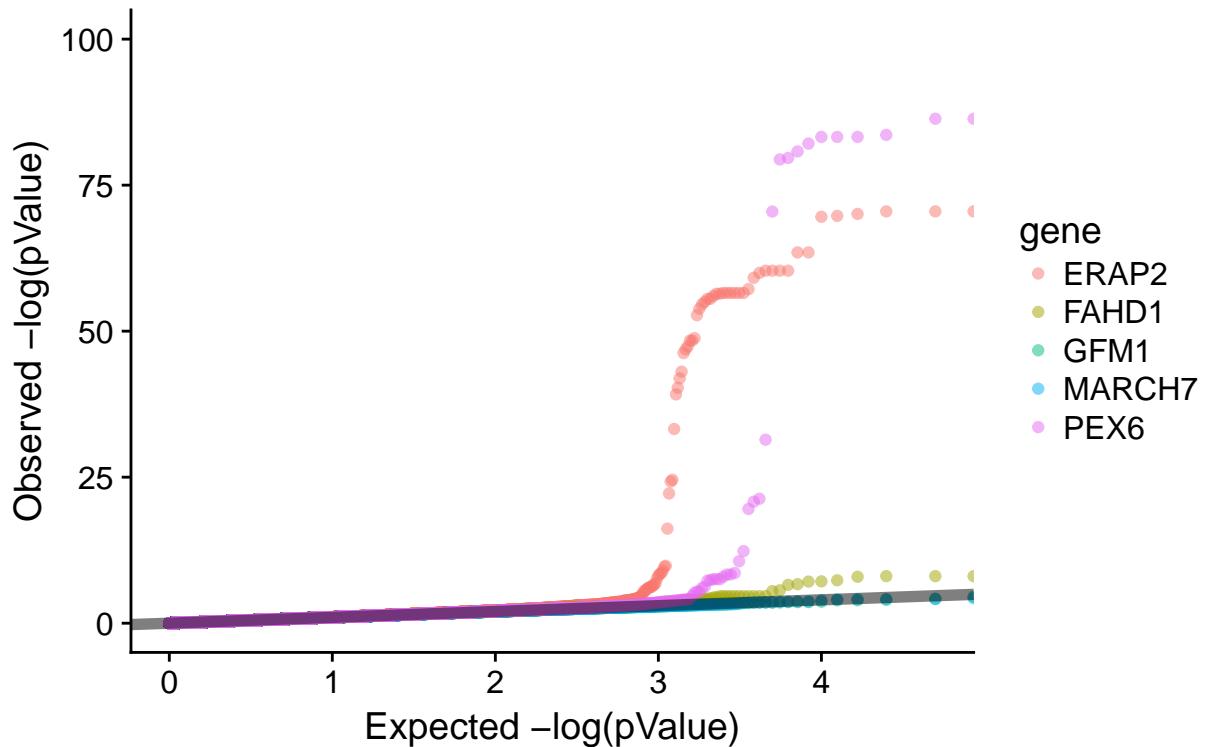
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
16725	16745	16767	16765	16785	16805	



Manhattan Plot All genes after Including PC1 to PC5 as Covars

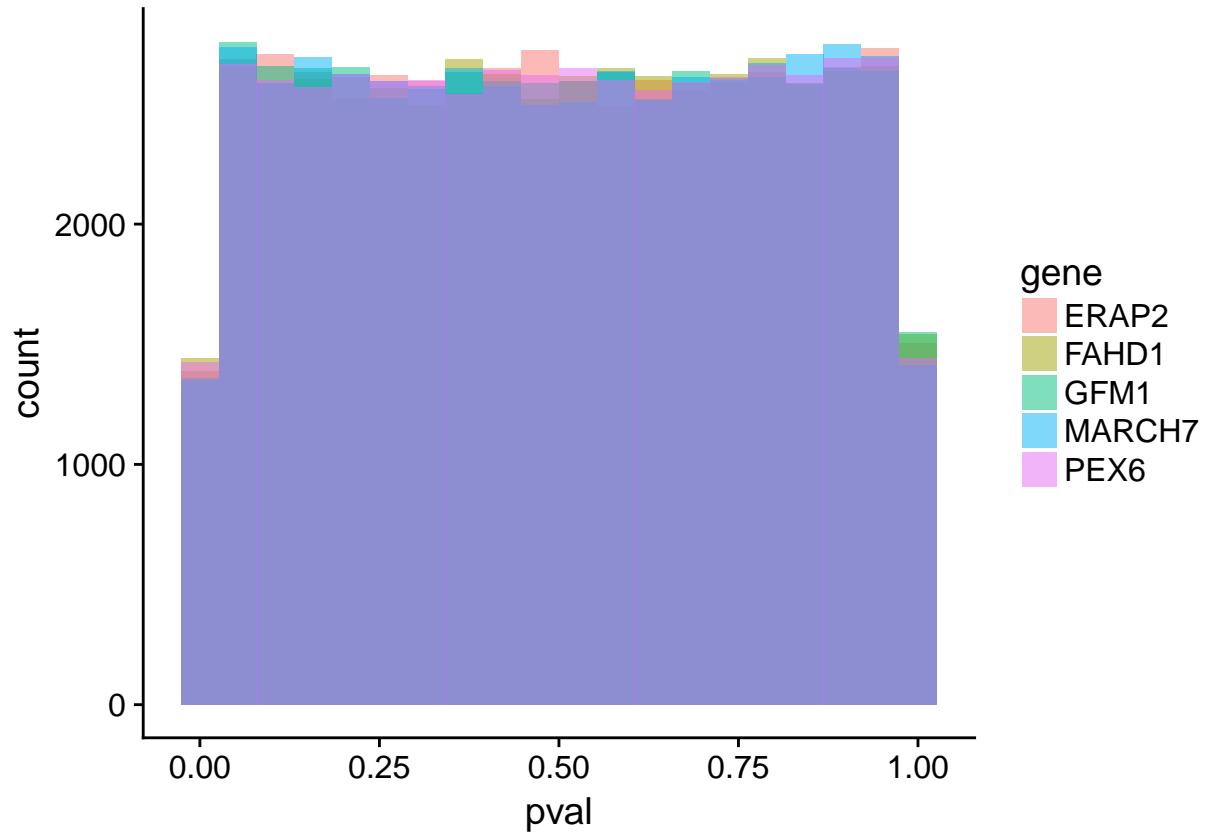


QQ Plot All Genes with PC to PC5 as Covars

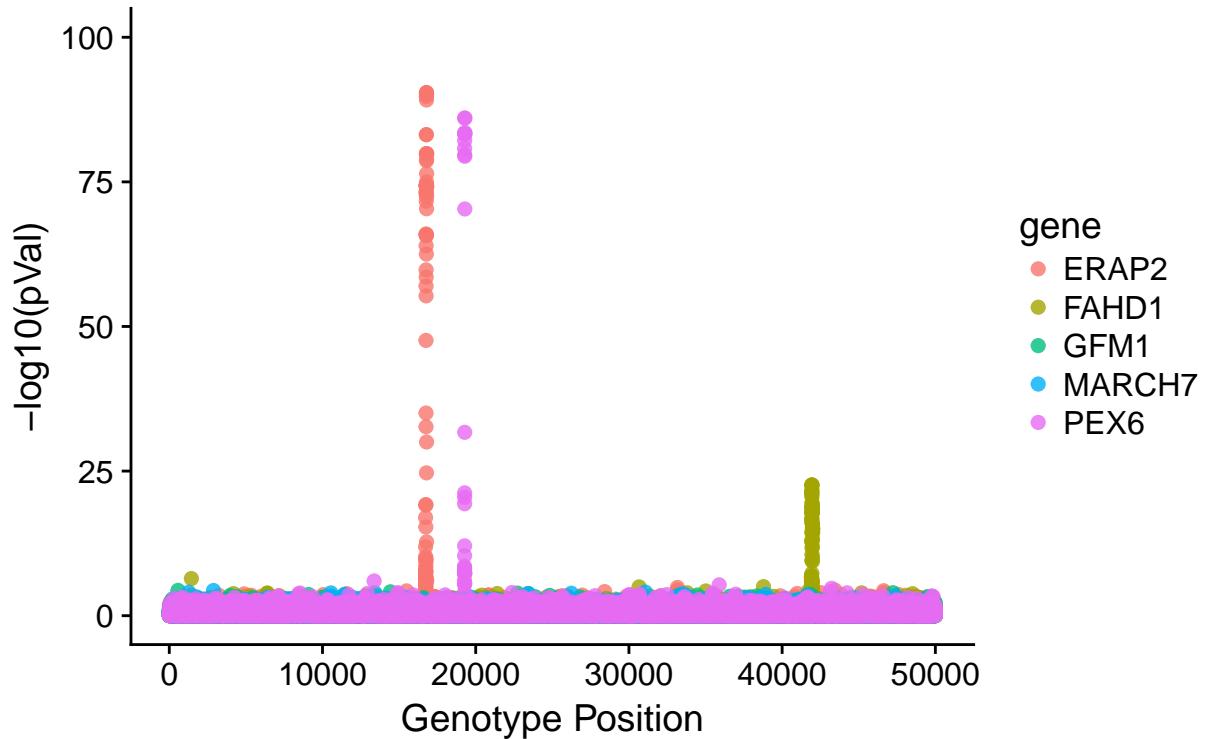


	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
ERAP2	16725	16744	16767	16765	16784	16803

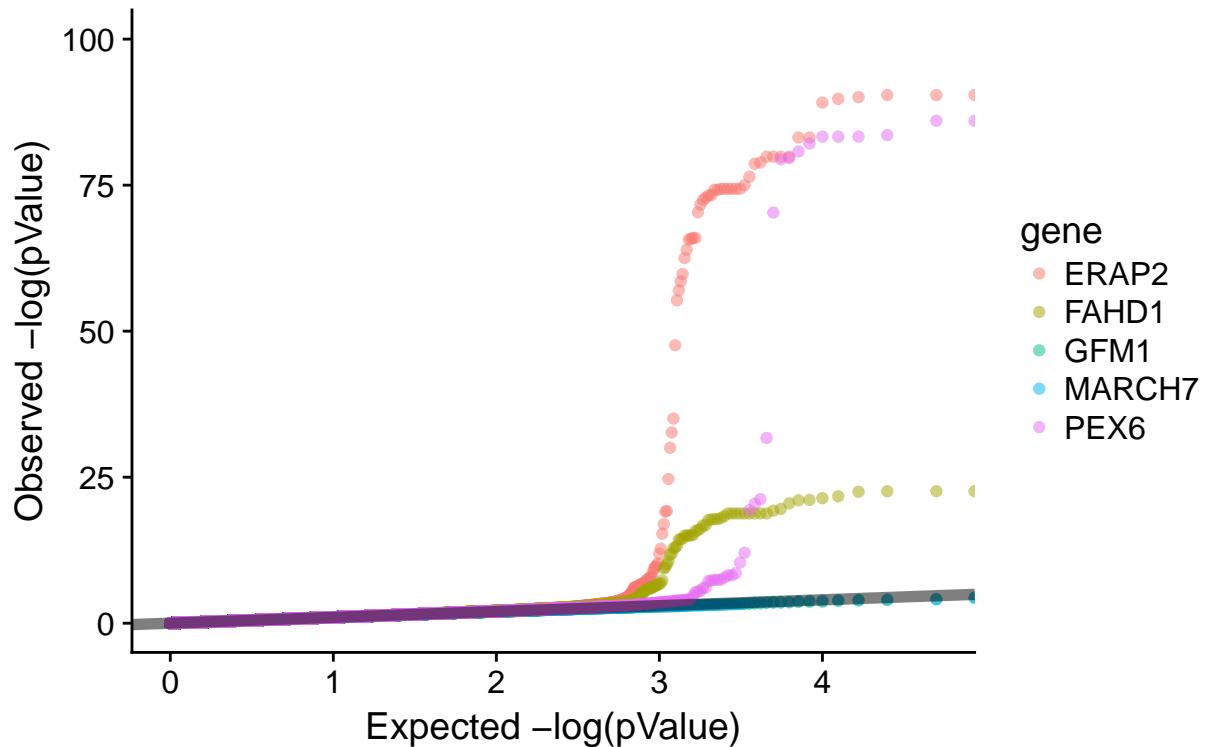
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
FAHD1	16725	16745	16767	16765	16785	16805

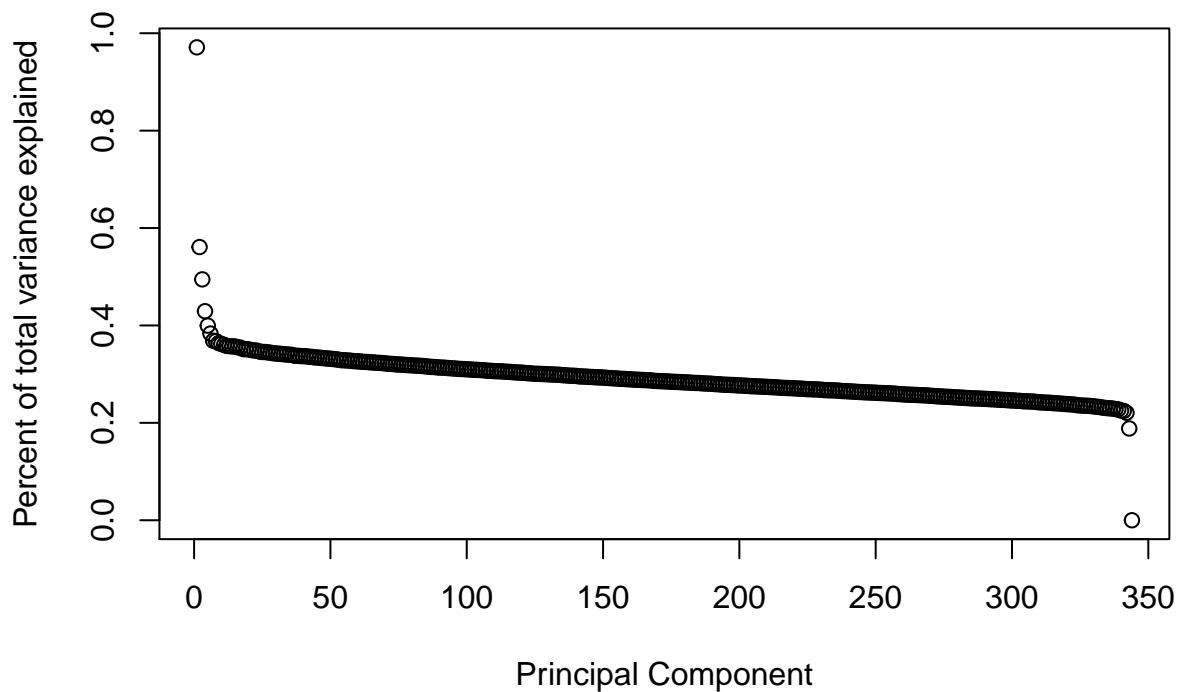


Manhattan Plot All genes after Including PC1 to PC5 as Covars

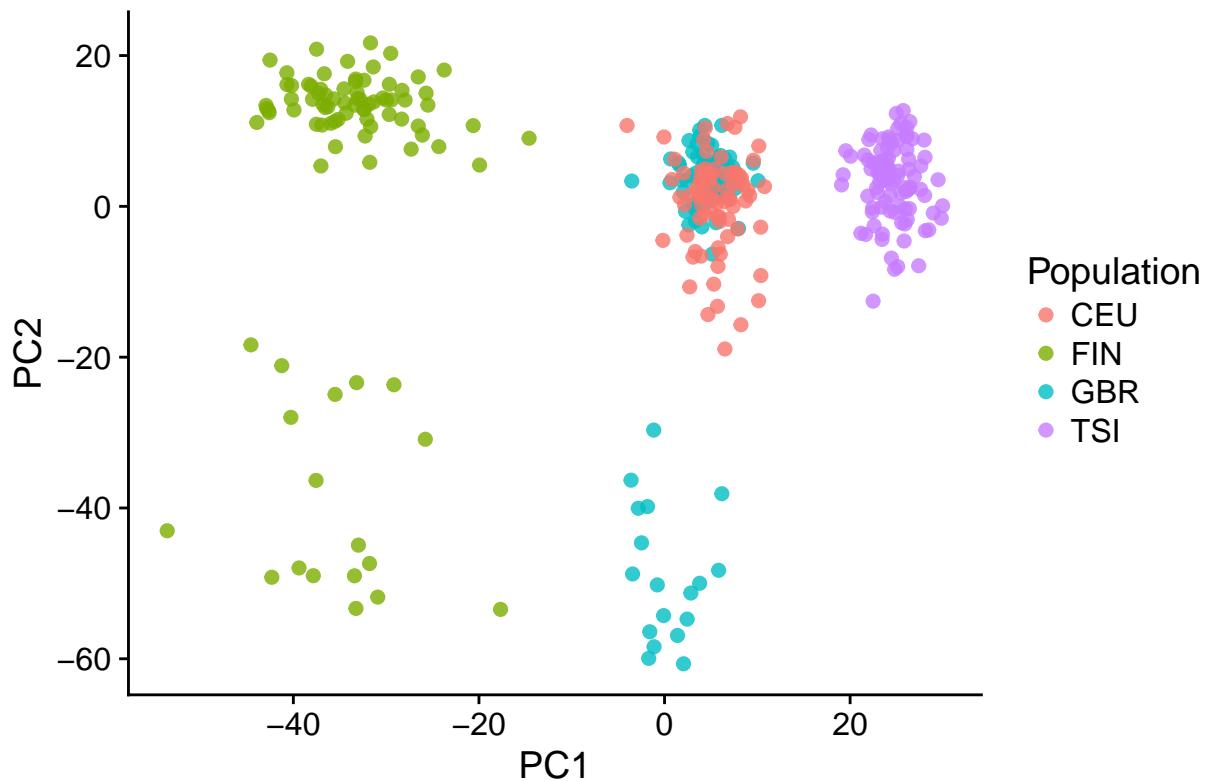


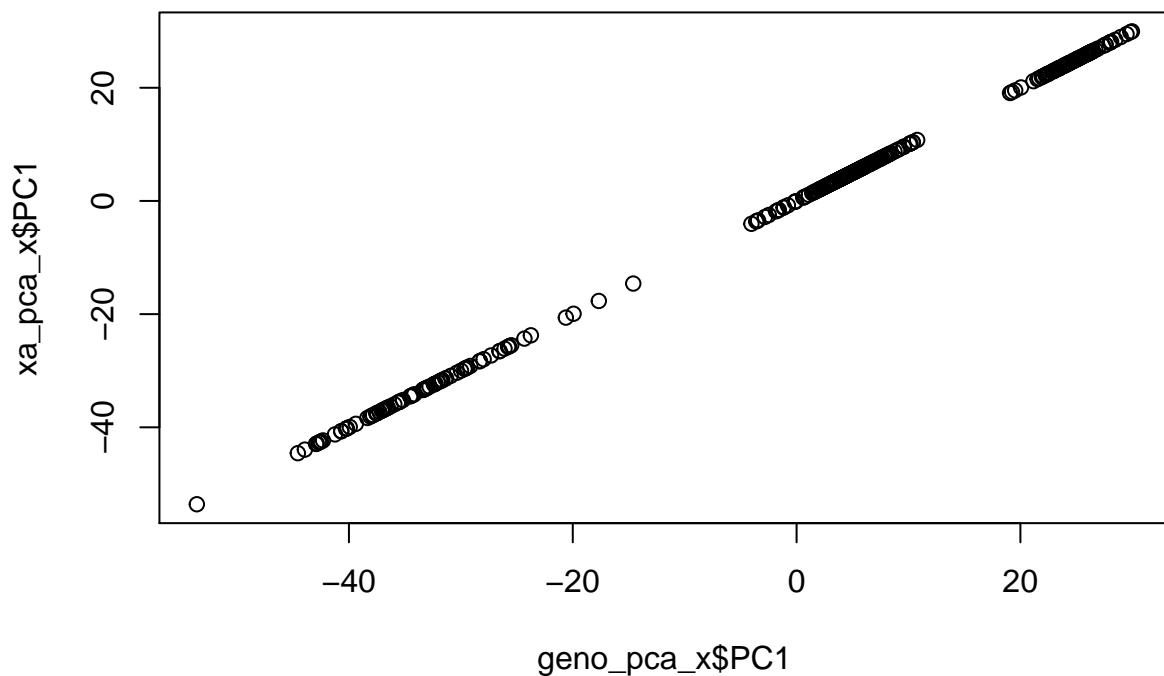
QQ Plot
All Genes with filtered PC1 to PC4 as Covars

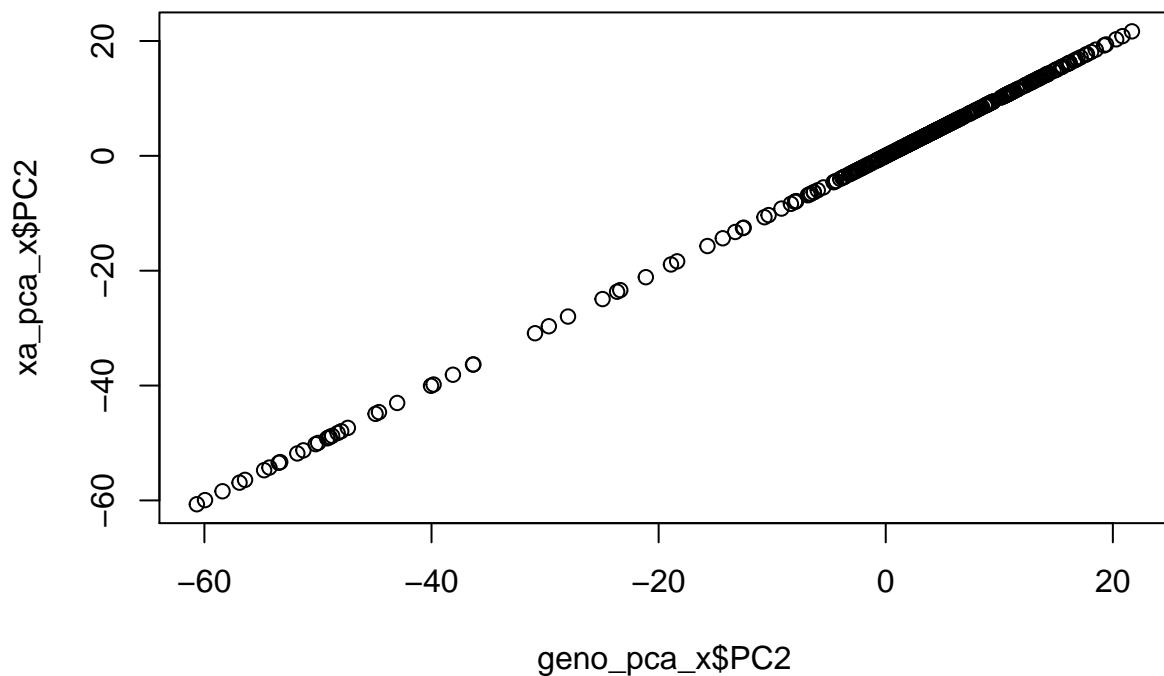


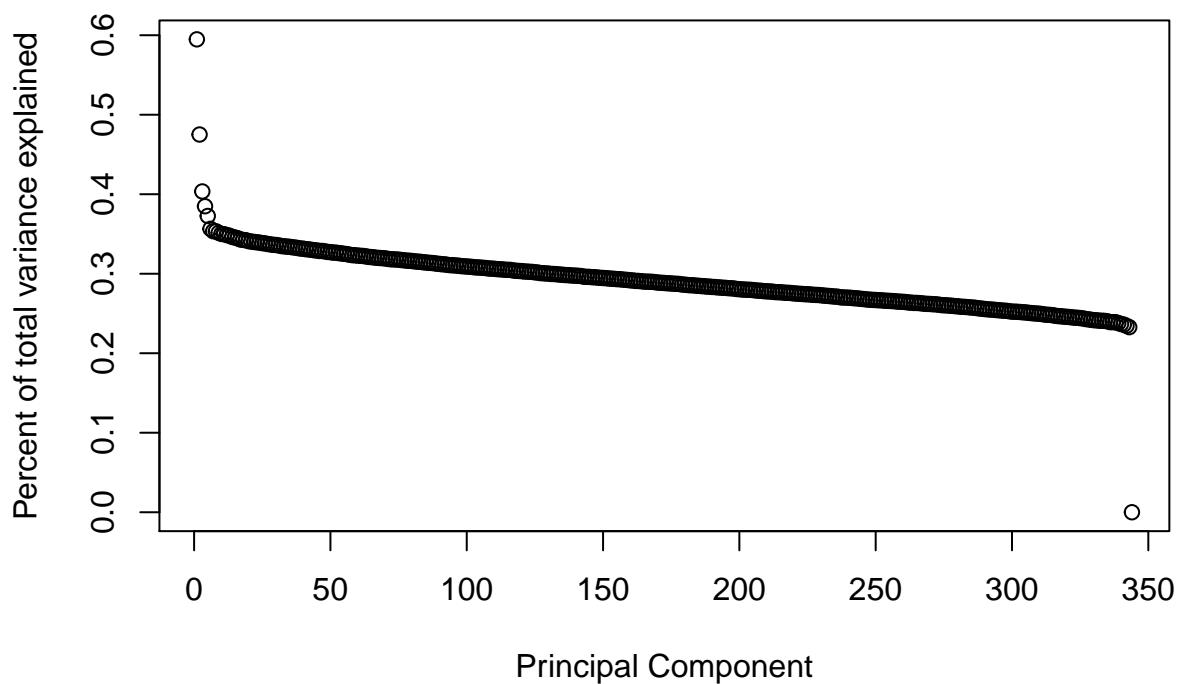


rst two principal components of X_a matrix, colored by Population

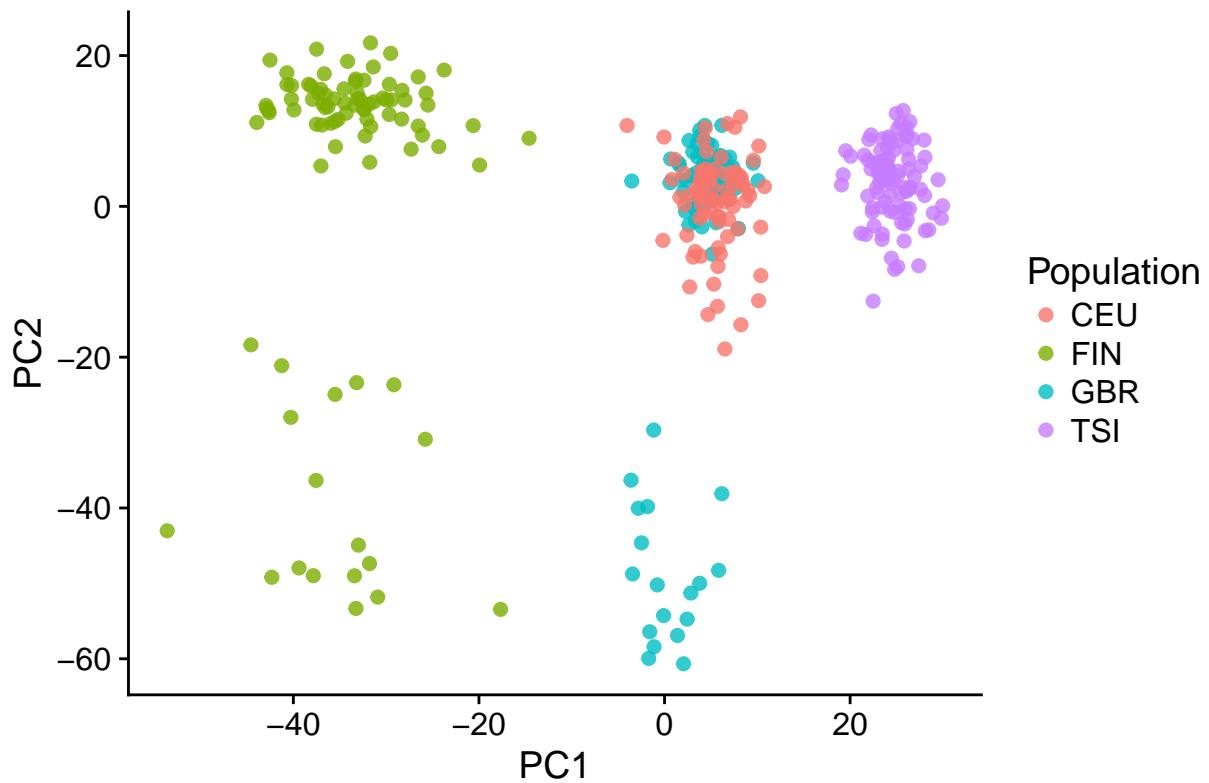


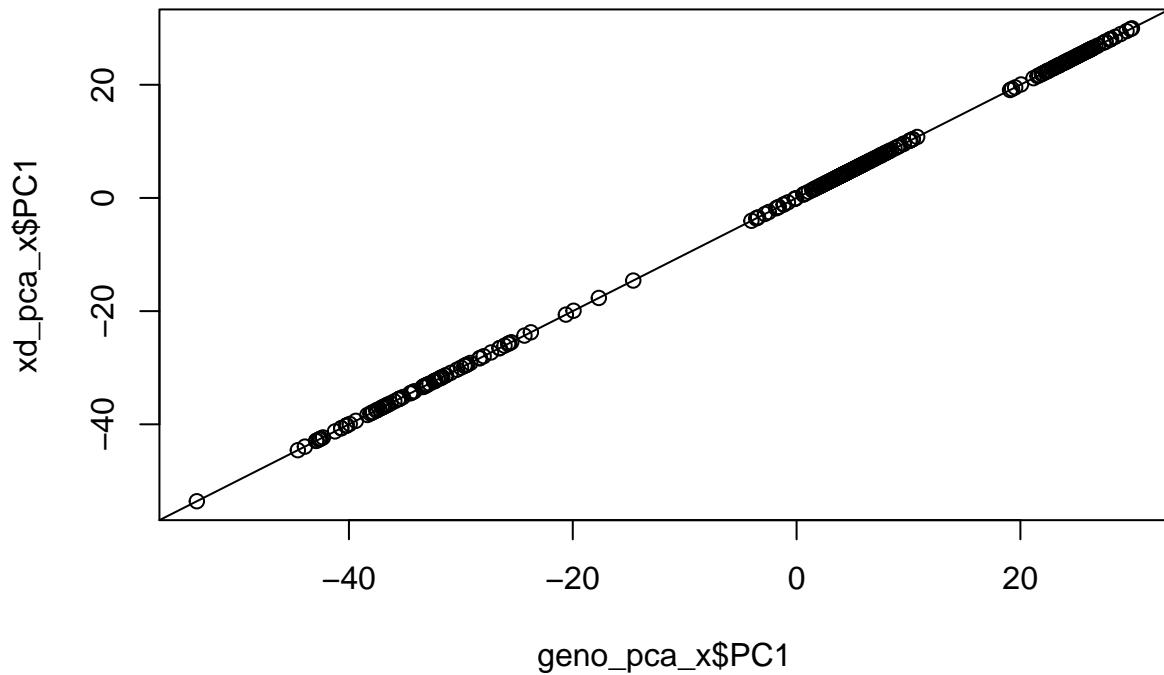




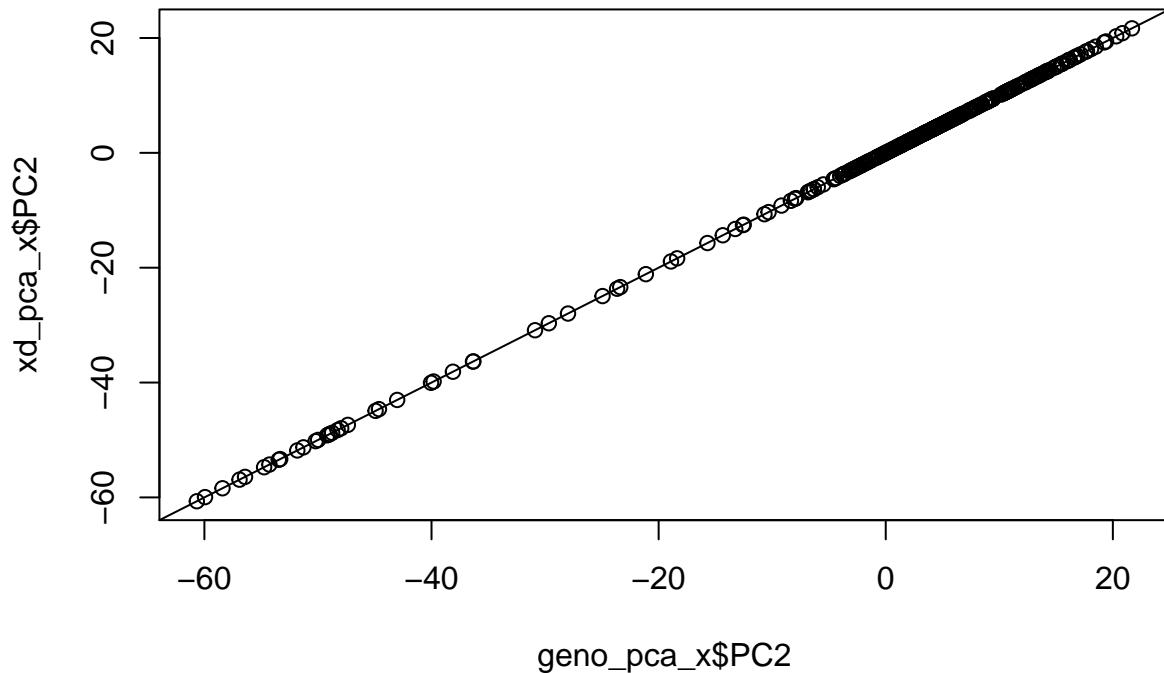


rst two principal components of X_d matrix, colored by Population





```
integer(0)
```



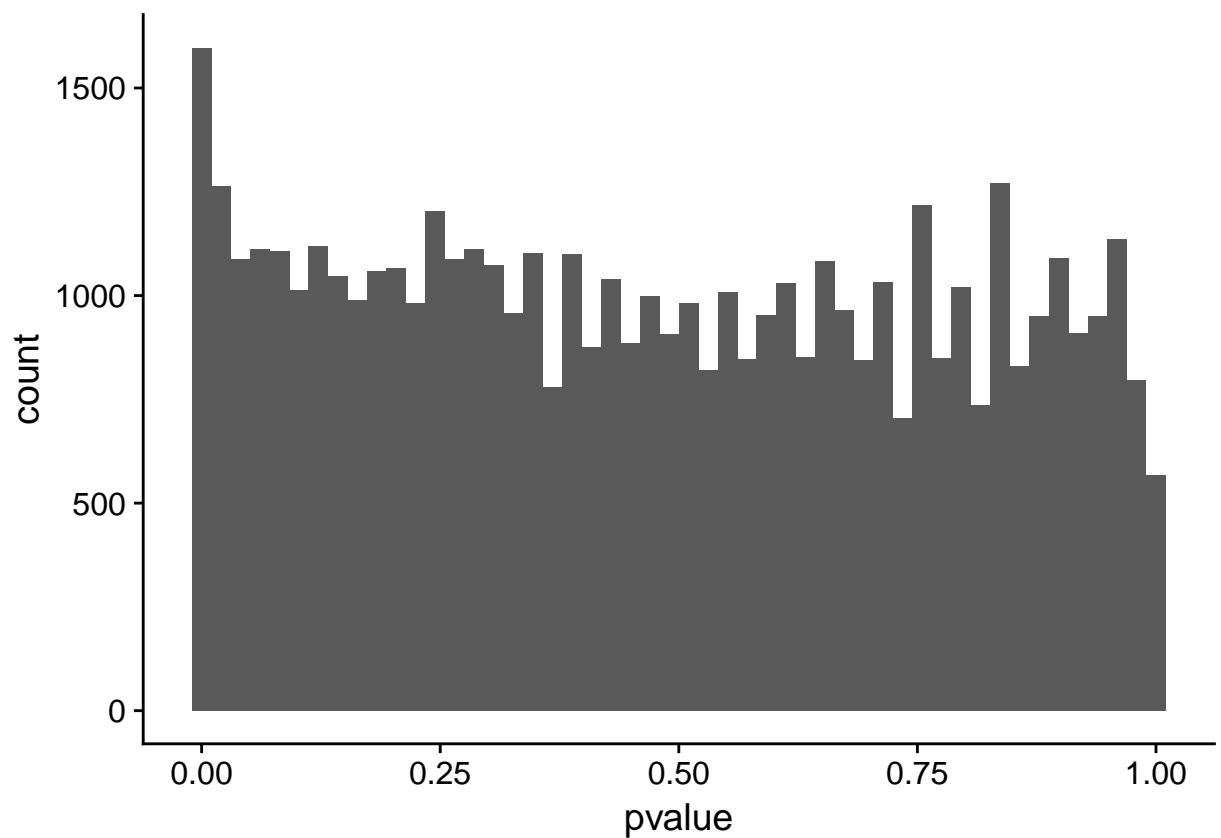
```

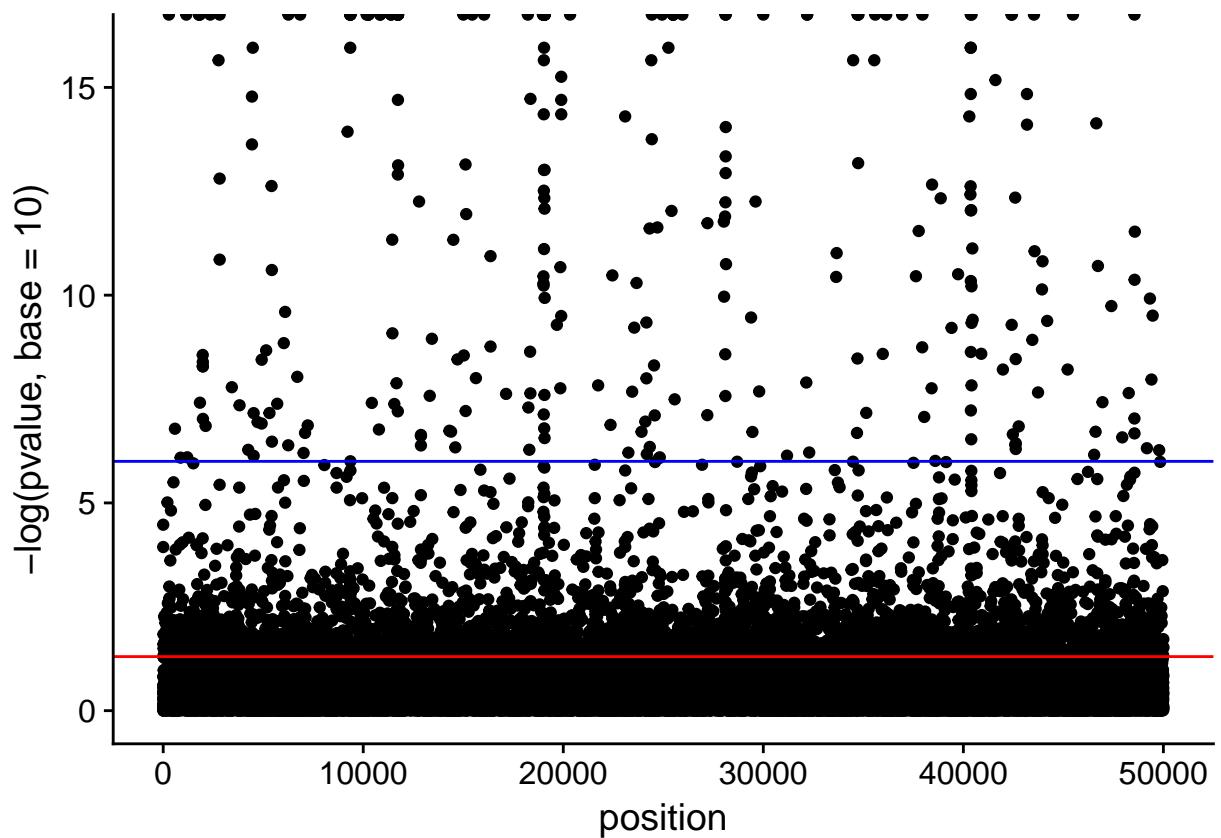
integer(0)

Chi-square test with continuity correction for Hardy-Weinberg equilibrium (autosomal)
Chi2 = 15.45815 DF = 1 p-value = 8.435235e-05 D = -8.499273 f = 0.2236694
[1] 1.381394
[1] 0.2398638

Chi-square test for Hardy-Weinberg equilibrium (autosomal)
Chi2 = 1.381394 DF = 1 p-value = 0.2398638 D = 1.221657 f = -0.0633694

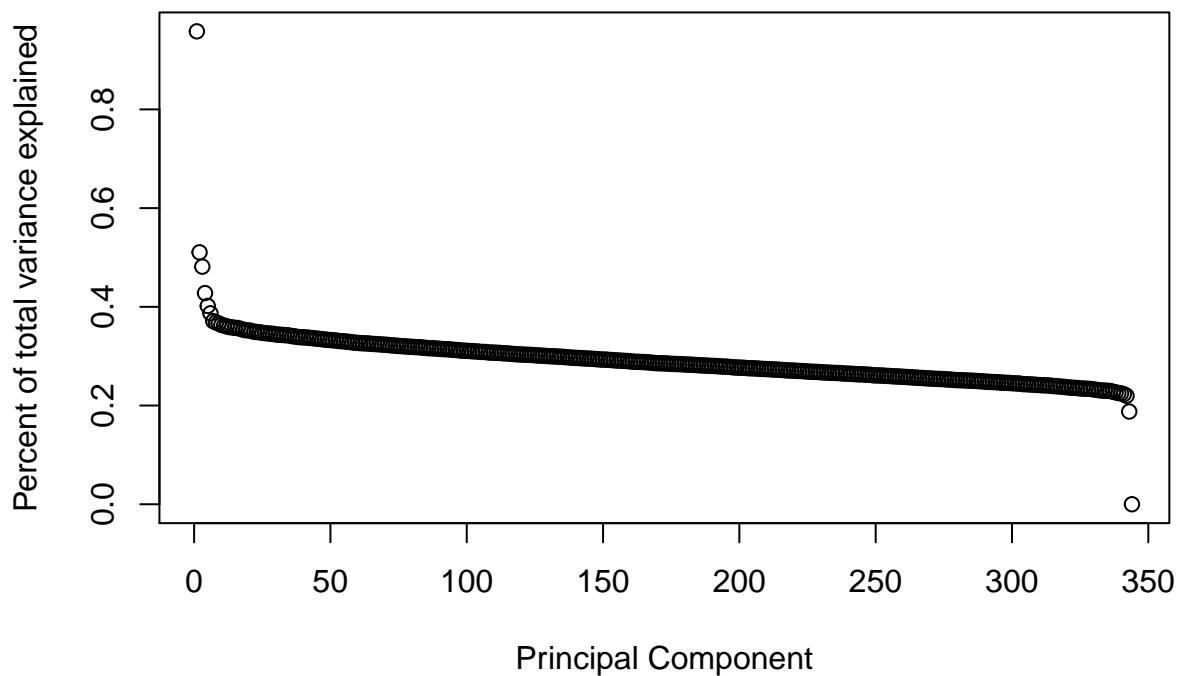
```



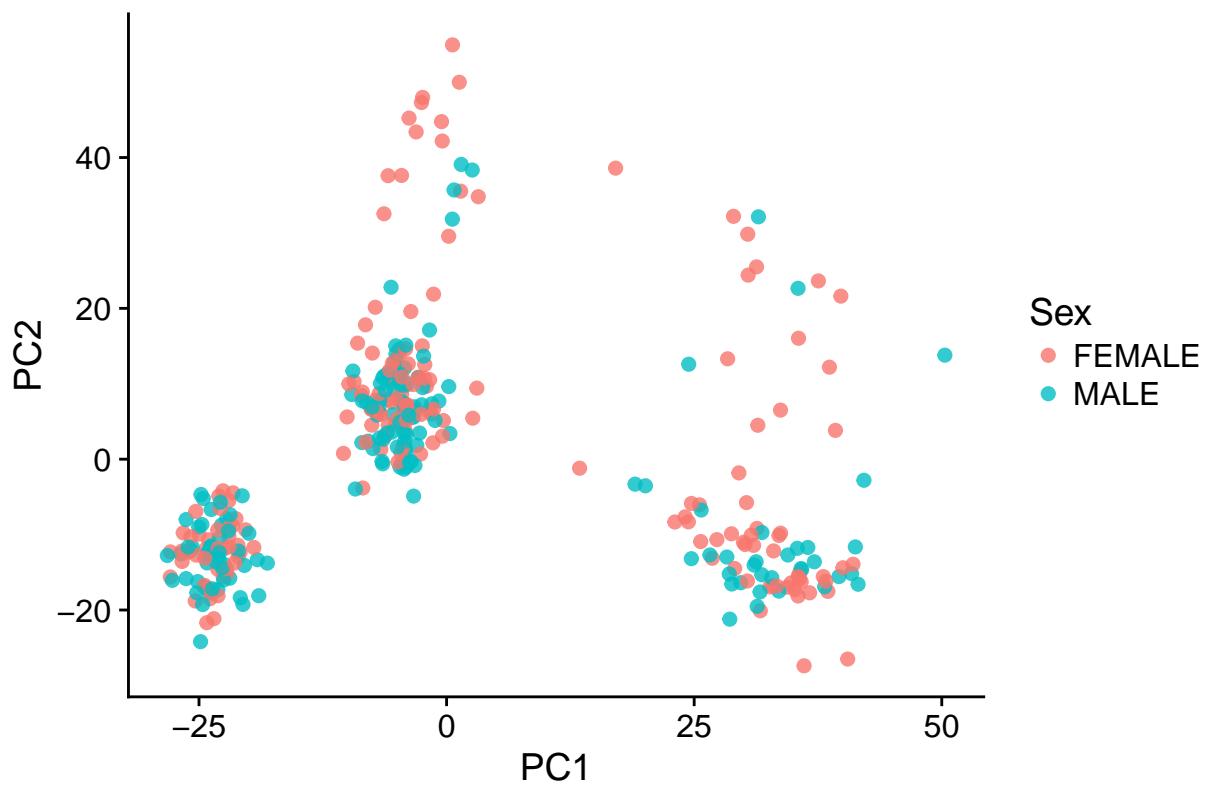


[1] 296

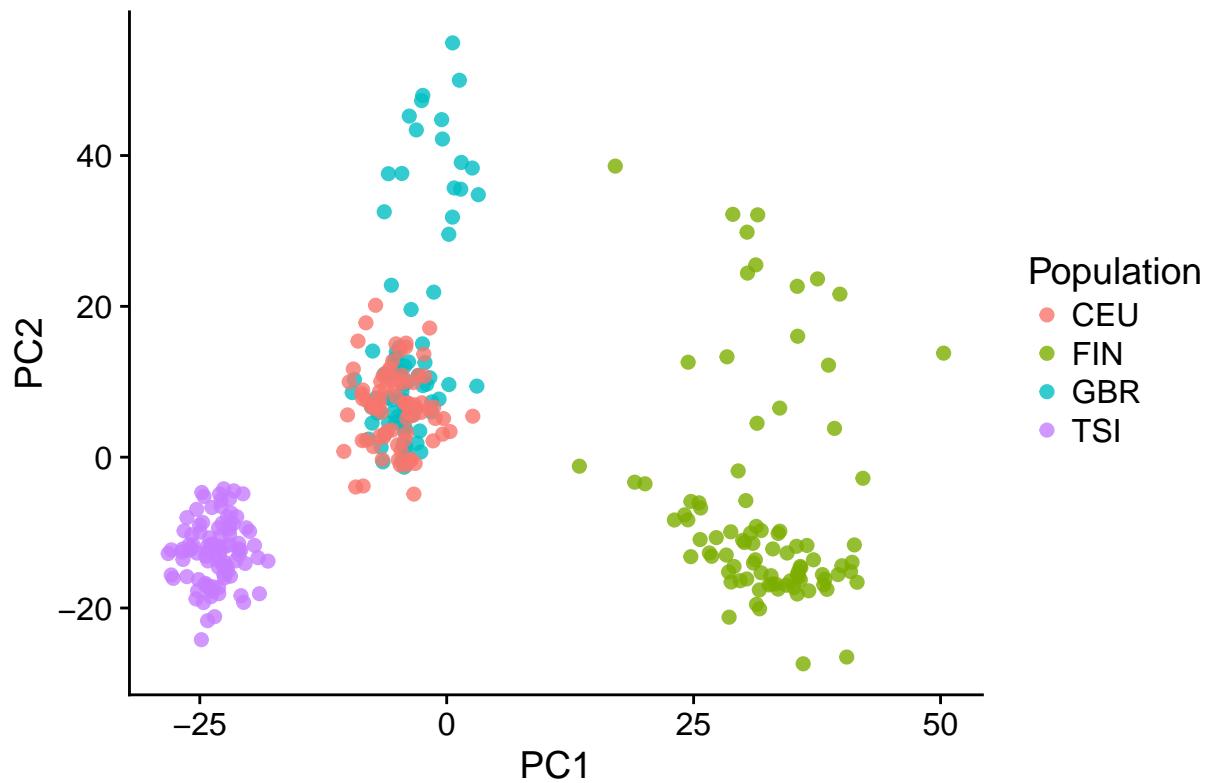
[1] 3908

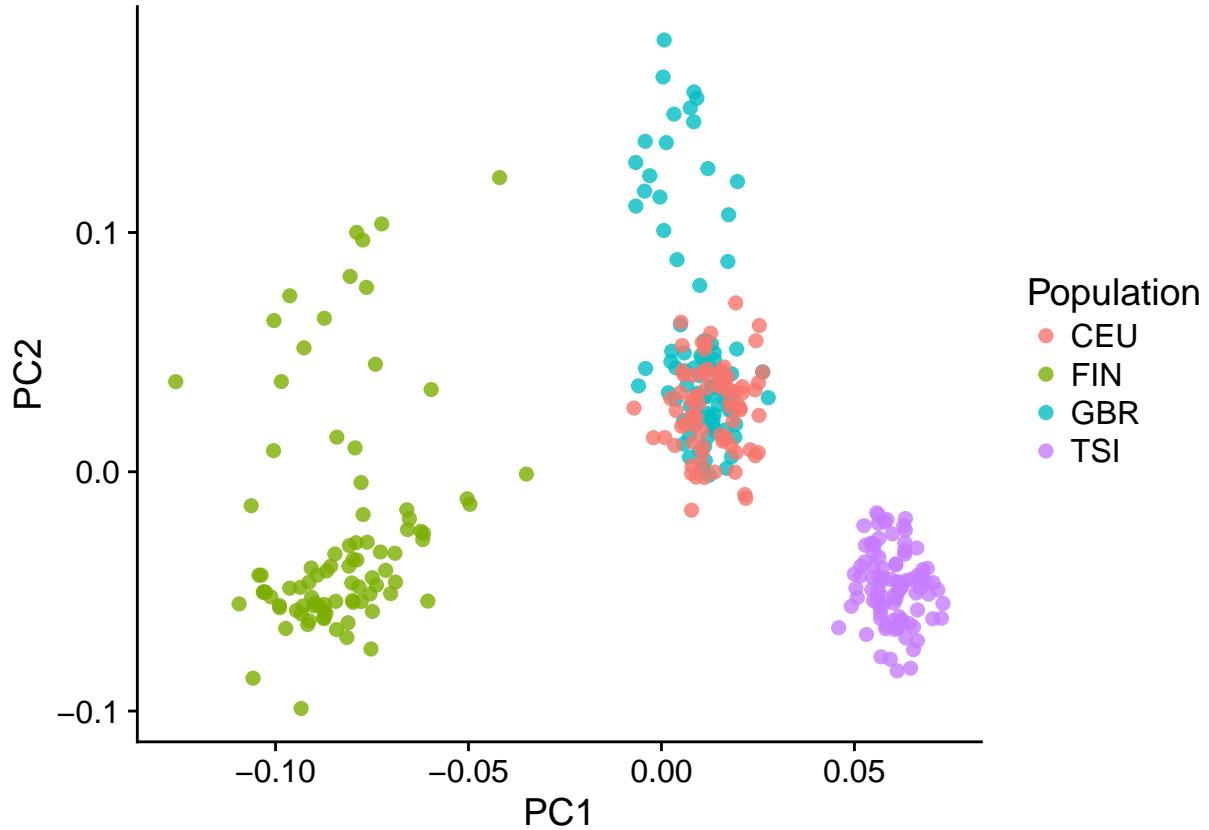


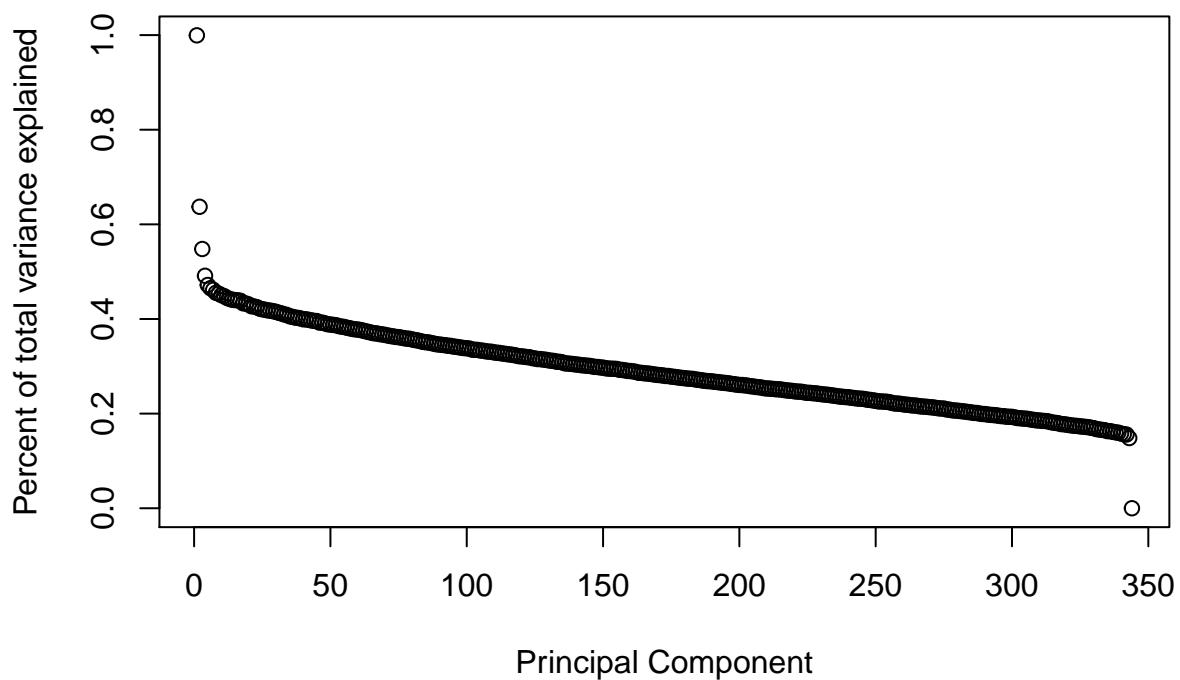
First two principal components, colored by Sex



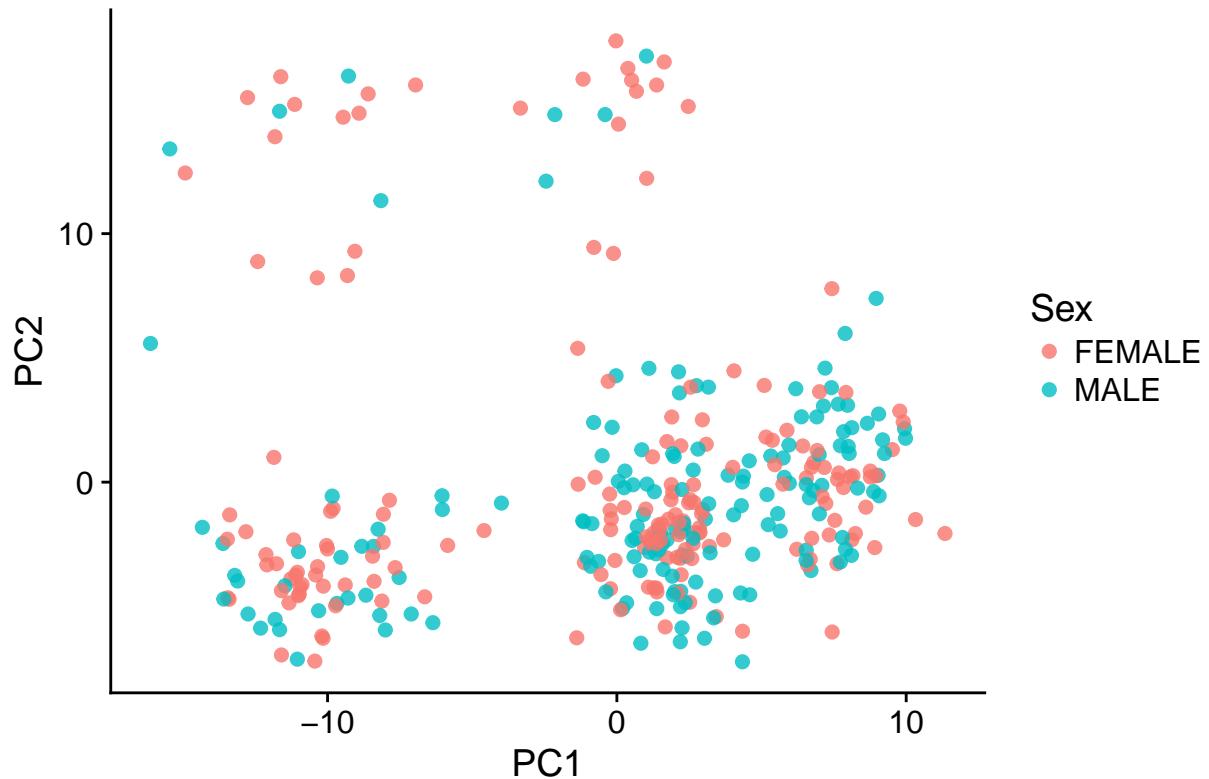
First two principal components, colored by Population



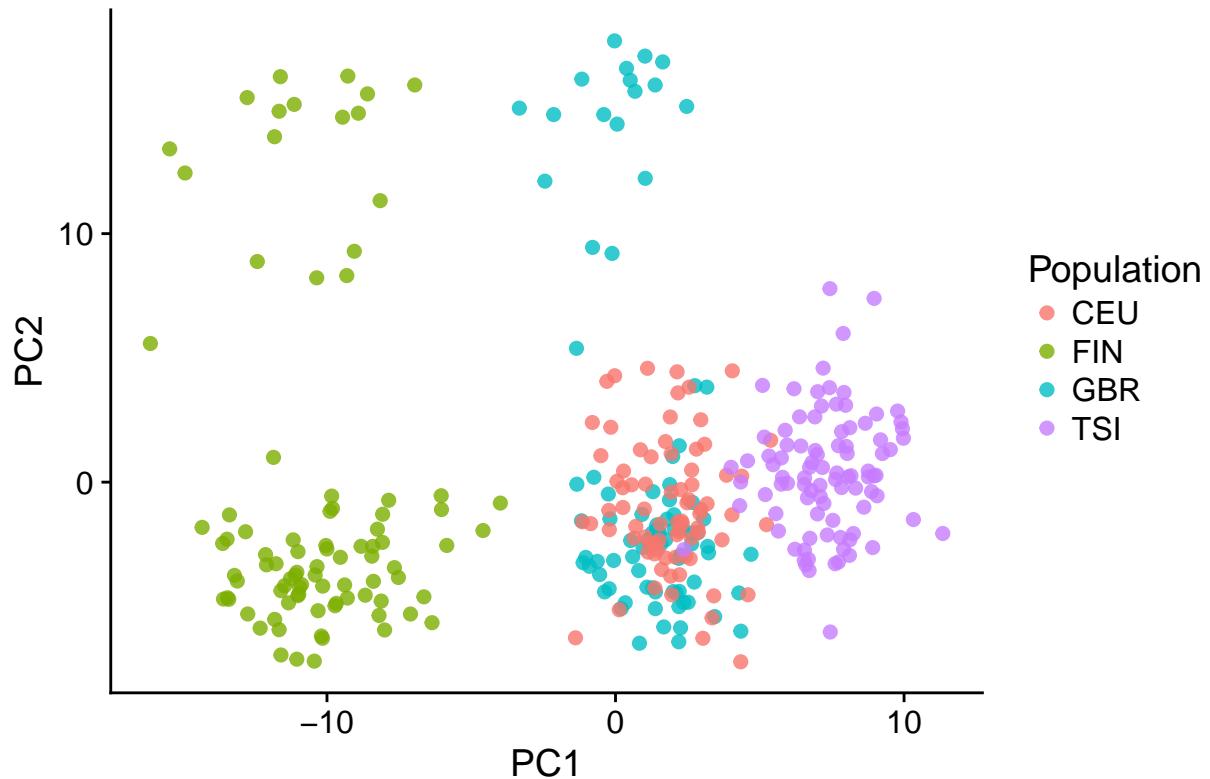




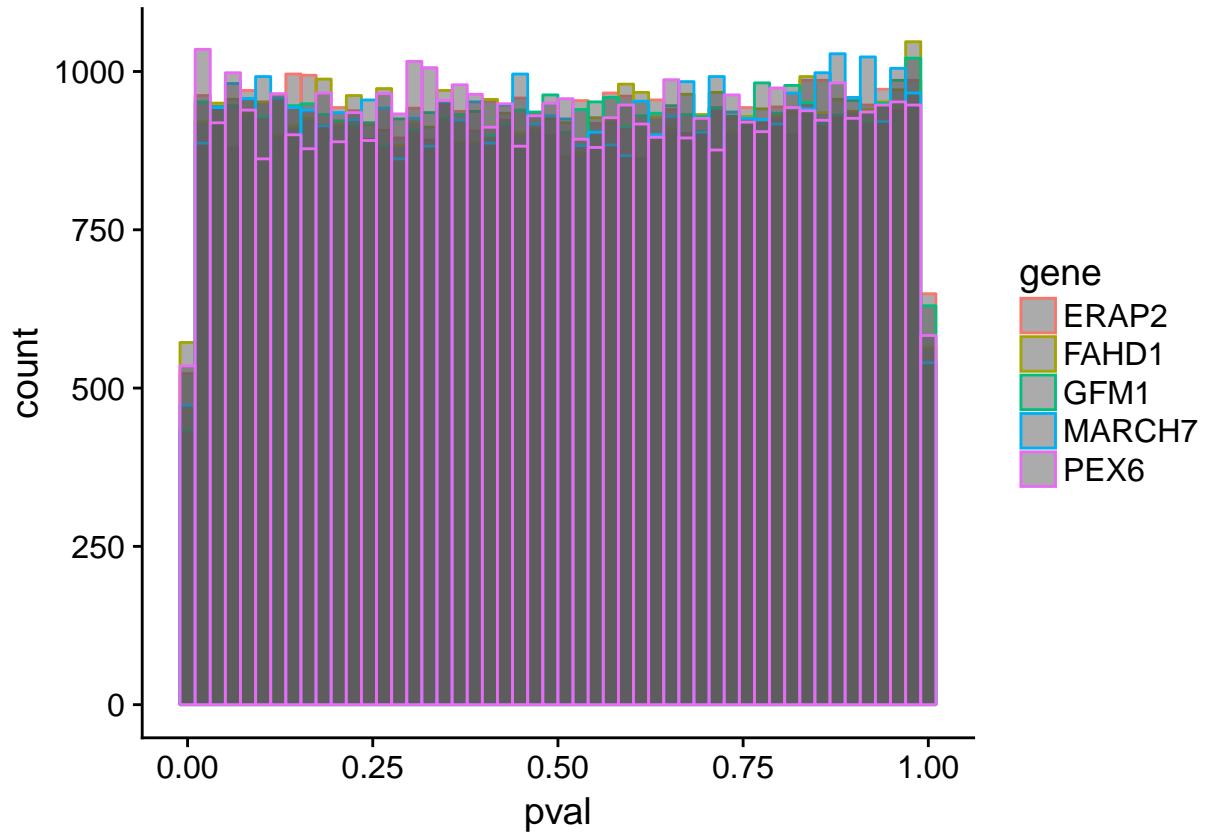
First two principal components, colored by Sex



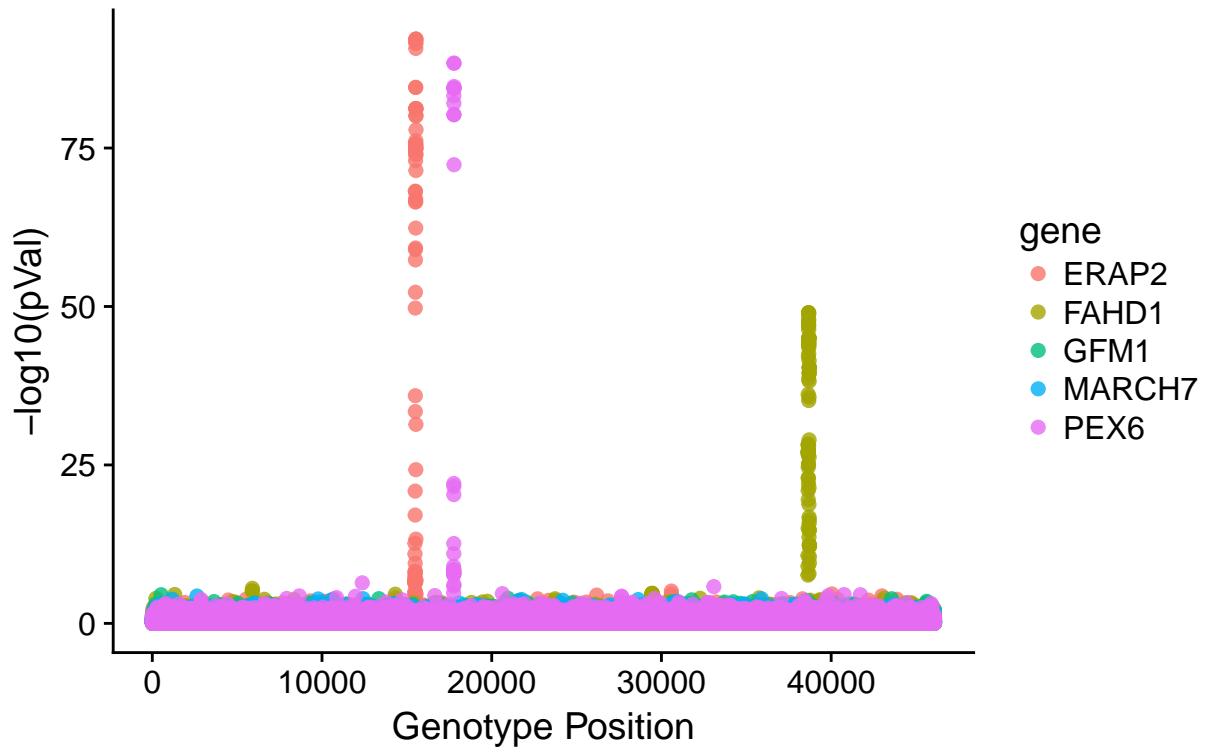
First two principal components, colored by Population

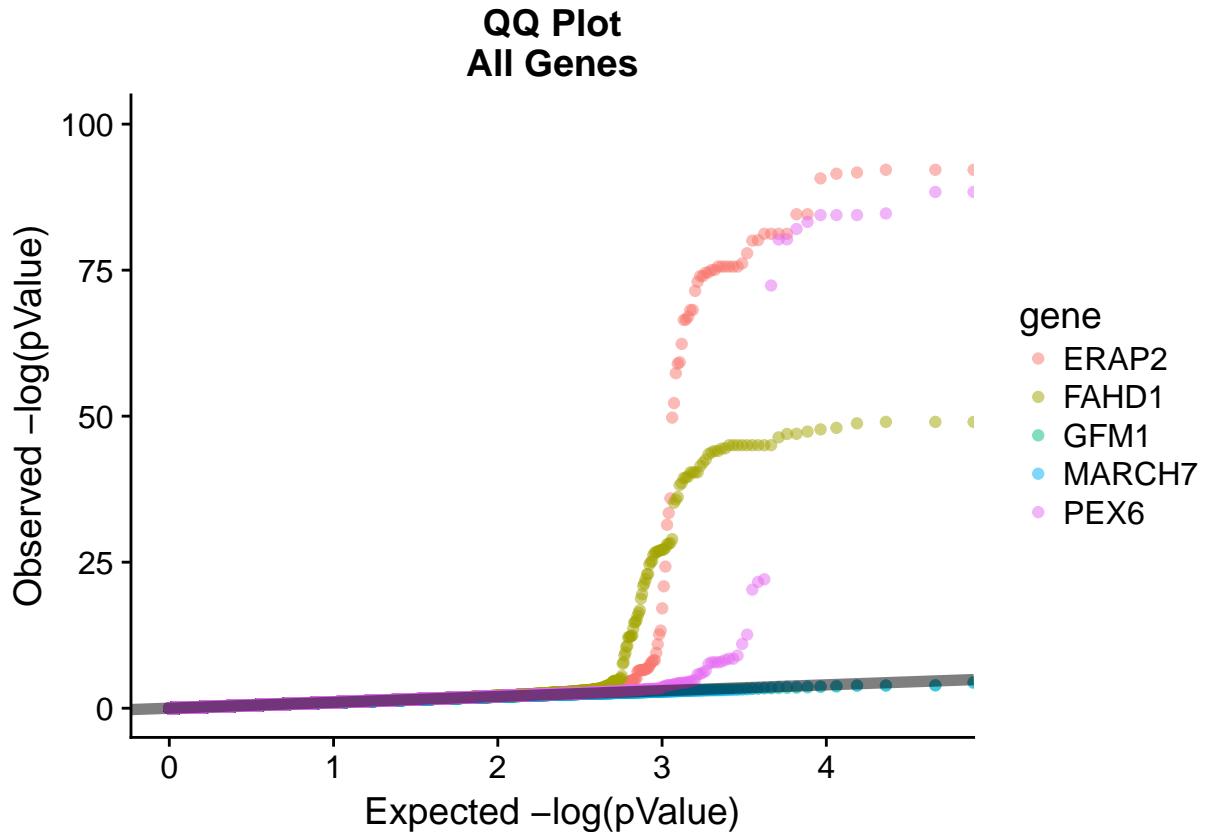


Try just a lin reg and F-test with no covars after filtering.



Manhattan Plot All Genes





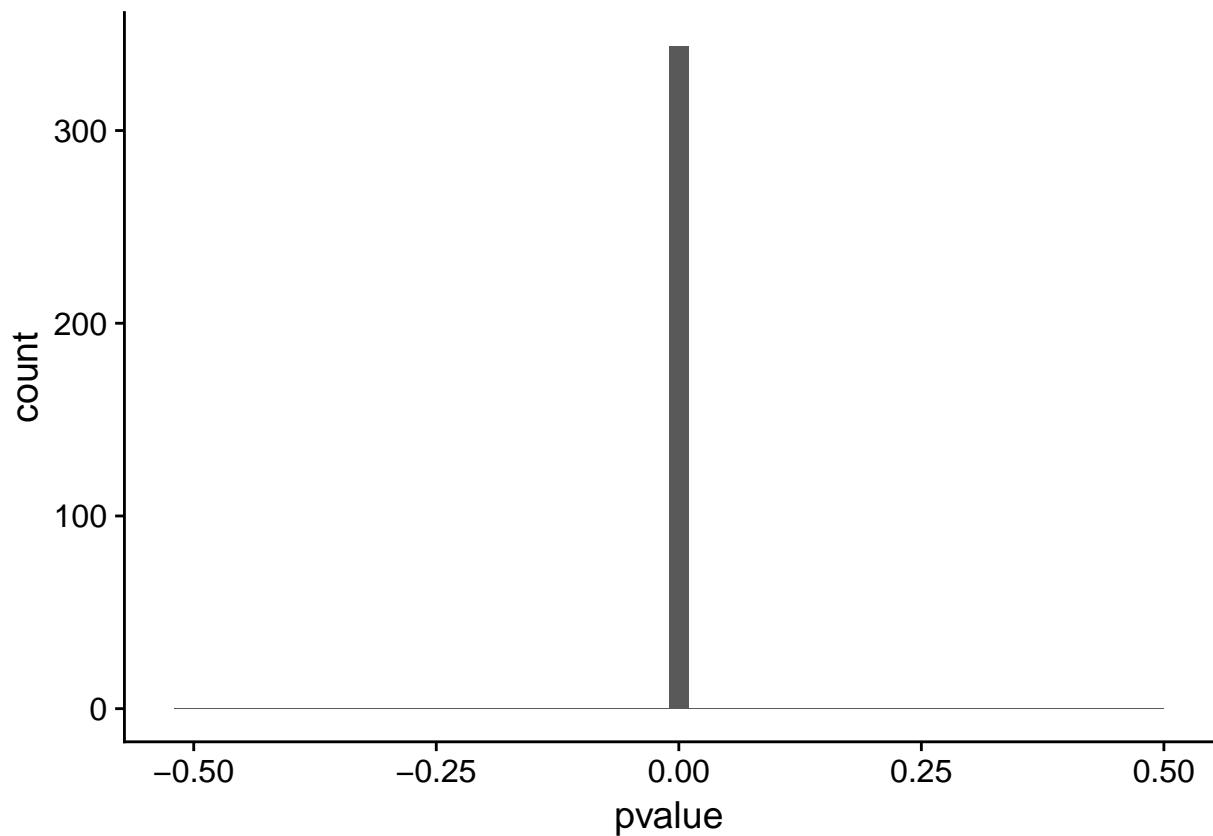
```

indv_HW <- map(data.frame(t(geno)), function(x) HWChisq(table(factor(x, levels = c(0, 1, 2))), cc = 0,
HWChisq(table(factor(t(geno)[ 1,]), levels = c(0, 1, 2))), cc = 0.5)

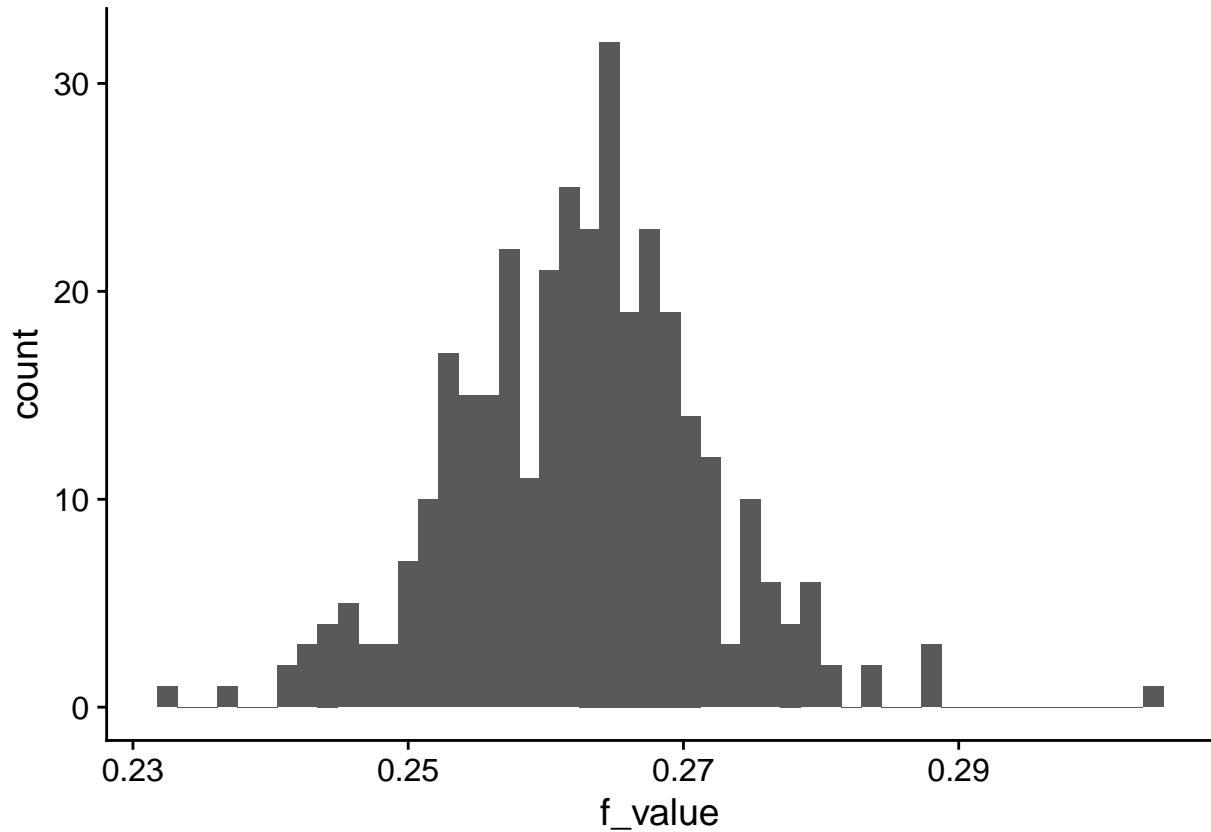
## Chi-square test with continuity correction for Hardy-Weinberg equilibrium (autosomal)
## Chi2 = 3361.967 DF = 1 p-value = 0 D = -2980.787 f = 0.259339

indv_HW2 <- purrr::transpose(indv_HW )
indv_HW_simple <- data.frame(unlist(indv_HW2$pval))
indv_HW_f <- data.frame(unlist(indv_HW2$f))
colnames(indv_HW_simple) <- "pvalue"
colnames(indv_HW_f) <- "f_value"
indv_HW_simple$position <- 1:nrow(indv_HW_simple)
indv_HW_f$position <- 1:nrow(indv_HW_f)
indv_HW_simple %>%
  ggplot(aes(x = pvalue)) +
  geom_histogram(bins = 50)

```



```
indv_HW_f %>%
  ggplot(aes(x = f_value)) +
  geom_histogram(bins = 50)
```



```
lm_test <- lm(pheno$ENSG00000164308.12 ~ x_a[, 1000] + x_d[, 1000] + factor(covars$Population))
lm_tidy <- tidy(lm_test)
fstat <- summary(lm_test)$fstatistic
fstat_2 <- glance(lm_test)$statistic
pf(fstat[1], fstat[2], fstat[3], lower.tail = FALSE)
```