

PGRNseq Sample and Phenotype Processing

Format the phenotype data for downstream processing.

```
pheno.df <- data.frame(id=nomiss.pheno$SM_ID,
                       sex=as.factor(nomiss.pheno$Sex),
                       ANC=as.numeric(nomiss.pheno$ANC_nadir_1000cells.per.mL),
                       site=as.factor(nomiss.pheno$Study.site),
                       race=as.factor(nomiss.pheno$Race),
                       dose=as.factor(nomiss.pheno$Dose..mg.))
```

Empirical Distributions of Phenotypes

```
anc.plt <- ggplot(pheno.df, aes(x=ANC)) + geom_histogram()
print(anc.plt + ggtitle('ANC'))

## stat_bin: binwidth defaulted to range/30. Use 'binwidth = x' to adjust this.

loganc.plt <- ggplot(pheno.df, aes(x=log10(ANC))) + geom_histogram()
print(loganc.plt + ggtitle('log10(ANC)'))

## stat_bin: binwidth defaulted to range/30. Use 'binwidth = x' to adjust this.
```

PCA

Prune SNP set for LD, filter on MAF, and run PCA.

```
ld.snpset <- snpgdsLDpruning(exome.geno)

## SNP pruning based on LD:
## Sliding window: 500000 basepairs, Inf SNPs
## |LD| threshold: 0.2
## Removing 5465 non-autosomal SNPs
## Removing 154888 SNPs (monomorphic, < MAF, or > missing rate)
## Working space: 253 samples, 82543 SNPs
## Chromosome 1: 18.92%, 4668/24674
## Chromosome 2: 20.38%, 3504/17190
## Chromosome 3: 19.84%, 2882/14528
## Chromosome 4: 23.03%, 2349/10198
## Chromosome 5: 21.60%, 2444/11314
## Chromosome 6: 17.92%, 2718/15171
```

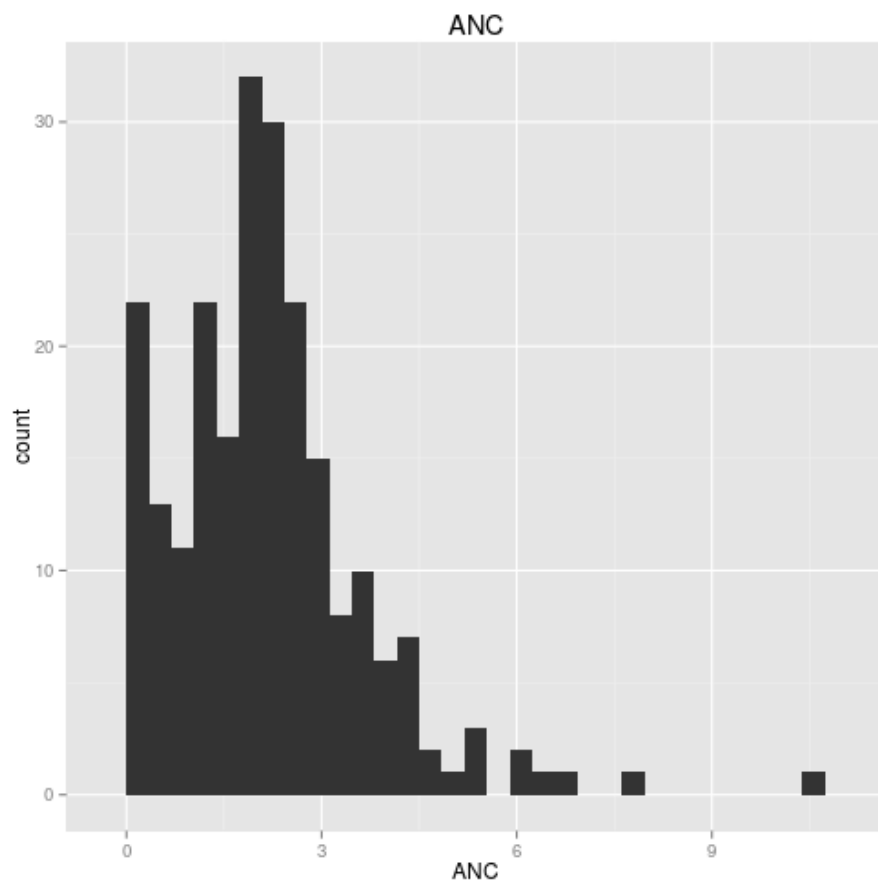


Figure 1: plot of chunk anc_plot

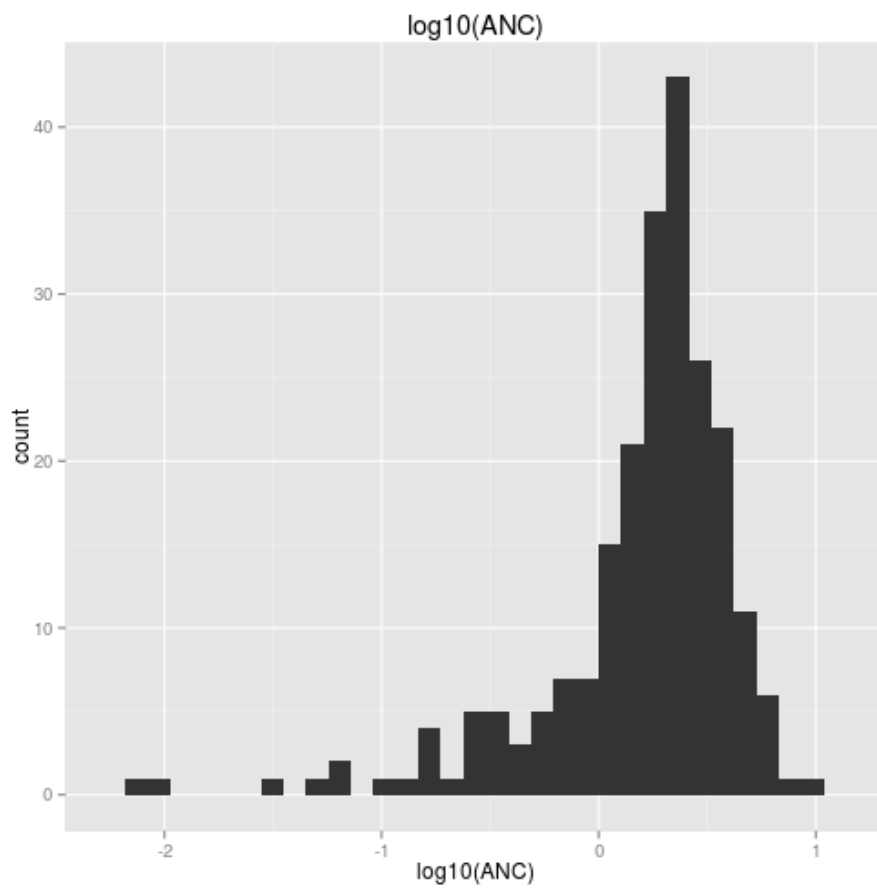


Figure 2: plot of chunk loganc.plt

```

## Chromosome 7: 21.33%, 2366/11092
## Chromosome 8: 21.53%, 1921/8921
## Chromosome 9: 20.35%, 2076/10201
## Chromosome 10: 21.80%, 2073/9511
## Chromosome 11: 18.46%, 2866/15528
## Chromosome 12: 19.82%, 2442/12318
## Chromosome 13: 25.00%, 1095/4380
## Chromosome 14: 20.32%, 1569/7721
## Chromosome 15: 19.46%, 1610/8273
## Chromosome 16: 17.75%, 1843/10384
## Chromosome 17: 16.86%, 2185/12963
## Chromosome 18: 24.11%, 904/3749
## Chromosome 19: 15.31%, 2289/14955
## Chromosome 20: 19.42%, 1248/6428
## Chromosome 21: 21.08%, 595/2823
## Chromosome 22: 18.54%, 947/5109
## 46594 SNPs are selected in total.

pca <- snpgdsPCA(exome.geno, snp.id = unlist(ld.snpset), maf = 0.05)

## Principal Component Analysis (PCA) on SNP genotypes:
## Removing 35251 SNPs (monomorphic, < MAF, or > missing rate)
## Working space: 253 samples, 11343 SNPs
## Using 1 CPU core.
## PCA: the sum of all working genotypes (0, 1 and 2) = 1642261
## PCA: Fri Aug 8 15:36:27 2014 0%
## PCA: Fri Aug 8 15:36:27 2014 100%
## PCA: Fri Aug 8 15:36:27 2014 Begin (eigenvalues and eigenvectors)
## PCA: Fri Aug 8 15:36:27 2014 End (eigenvalues and eigenvectors)

pca.df <- data.frame(sample.id=pca$sample.id, eigenvect=pca$eigenvect)

```

Look at the first two PCs, and color the samples by self-reported ancestry.

```
## Error: object 'eigenvect.1' not found
```

Model-based Clustering

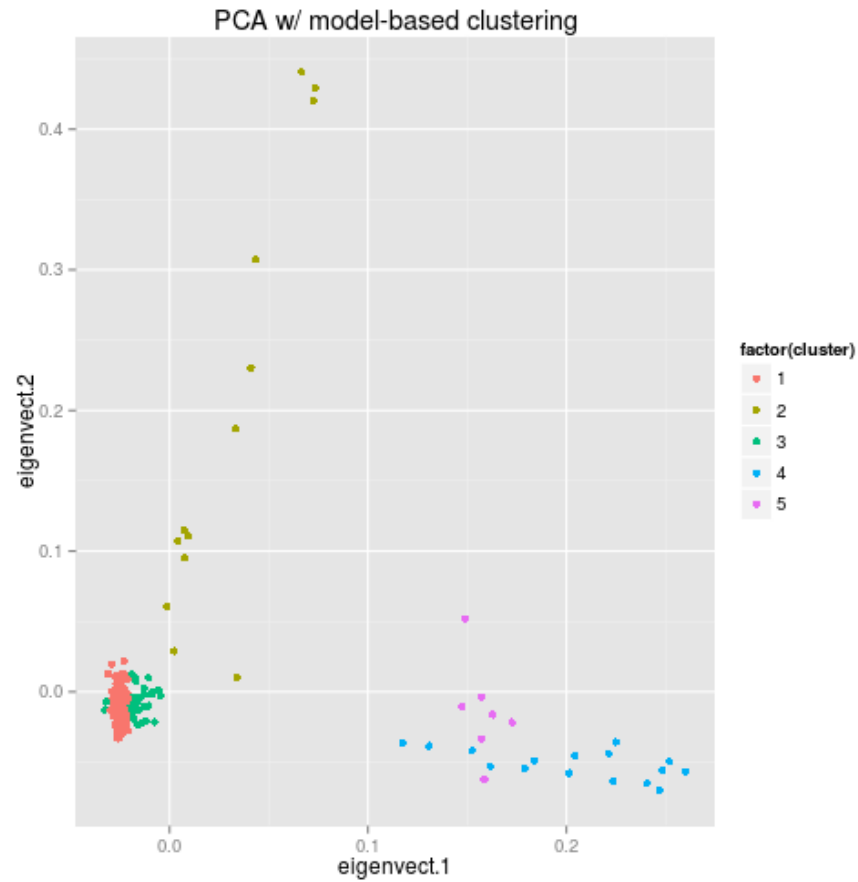
Add principal components to the phenotype data frame, then cluster samples in the first two PCs.

```

pheno.df <- merge(pca.df, pheno.df, by.x = "sample.id", by.y = "id")
names(pheno.df)[which(names(pheno.df) == "sample.id")] <- "id"
cluster <- Mclust(pheno.df[,c("eigenvect.1", "eigenvect.2")])
pheno.df$cluster <- cluster$classification

```

Examine the results of our PCA and clustering the first two PCs.



Taking a closer look at the clustering of the Europeans. We may be able to explain some genetic variation just by site of ascertainment.

```
## Warning: Removed 35 rows containing missing values (geom_point).
```

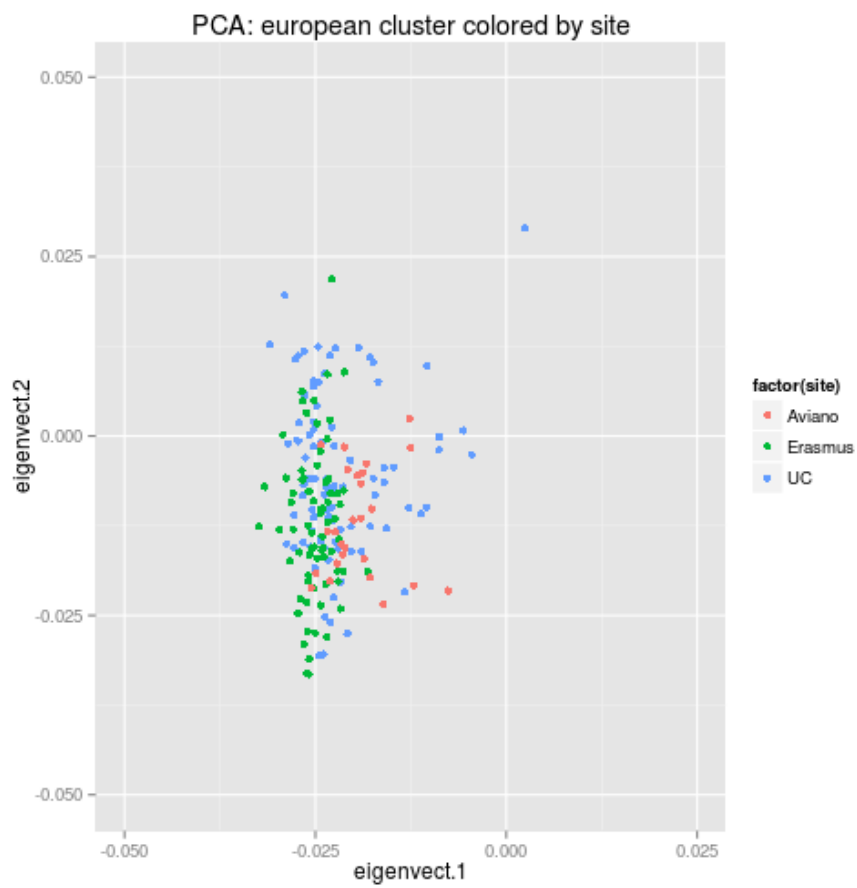


Figure 3: plot of chunk euro_zoom_plot