

DOGA ULUPINAR

---

# ANCESTRY MAPPING

# BIOLOGICAL BACKGROUND AND RELEVANCE

- ▶ Genetic make up, specifically SNPs within populations are more closely shared than across populations
- ▶ Control for population stratification in genetic association studies
- ▶ Understand how ethnic differences affects disease susceptibility
- ▶ Insight on which genes are more favorable in different populations

# COMPUTATIONAL FORMULATION OF PROBLEM

- ▶ Input: Genotype data (n individuals by m SNPs)
- ▶ Output: Assign global ancestry to each individual
  - ▶ Ancestry = {African, Asian, European, American}
- ▶ Benchmark: Accuracy (F1 score) and Runtime

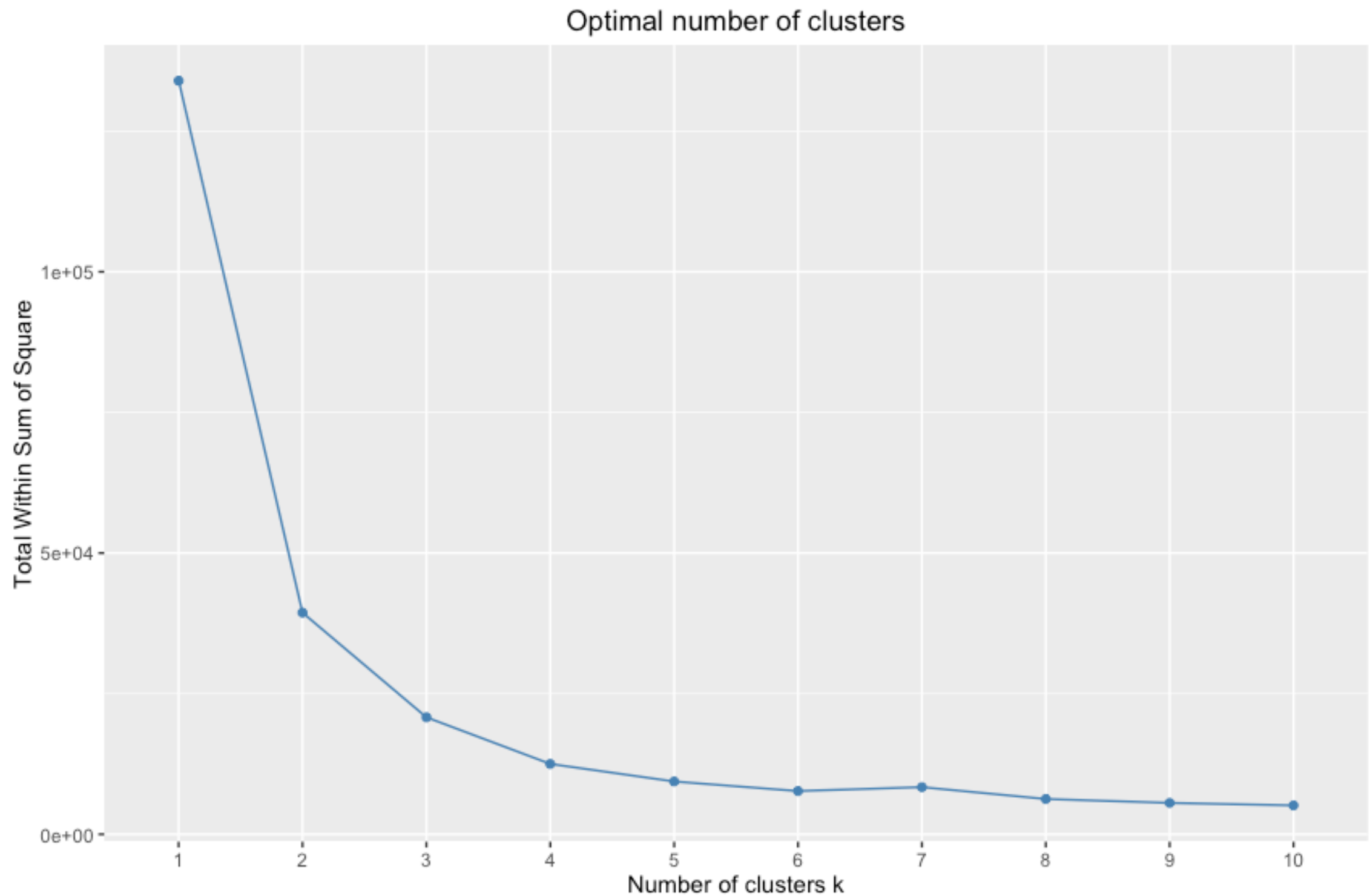
$$F1score = 2 * \frac{precision * recall}{precision + recall}$$

- ▶ Difficulties: High dimensionality of data (large amount of SNPs) and size of data

# THE DATASET AND ASSUMPTIONS MADE

- ▶ Dataset: 1000 Genomes Project
  - ▶ 1092 Individuals
  - ▶ Chromosome 20, 21 and 22 ~ 1.5 million SNPs
- ▶ Assumptions
  - ▶ Number of Populations is known
  - ▶ Population of each individual is not known

# CHECKING ASSUMPTIONS OF KNOWN 4 POPULATIONS



# BASELINE METHOD – KMEANS ACROSS WHOLE CHROMOSOME

- ▶ Objective: Minimize distance between each point and the center of the cluster that this individual is assigned to

$$\sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

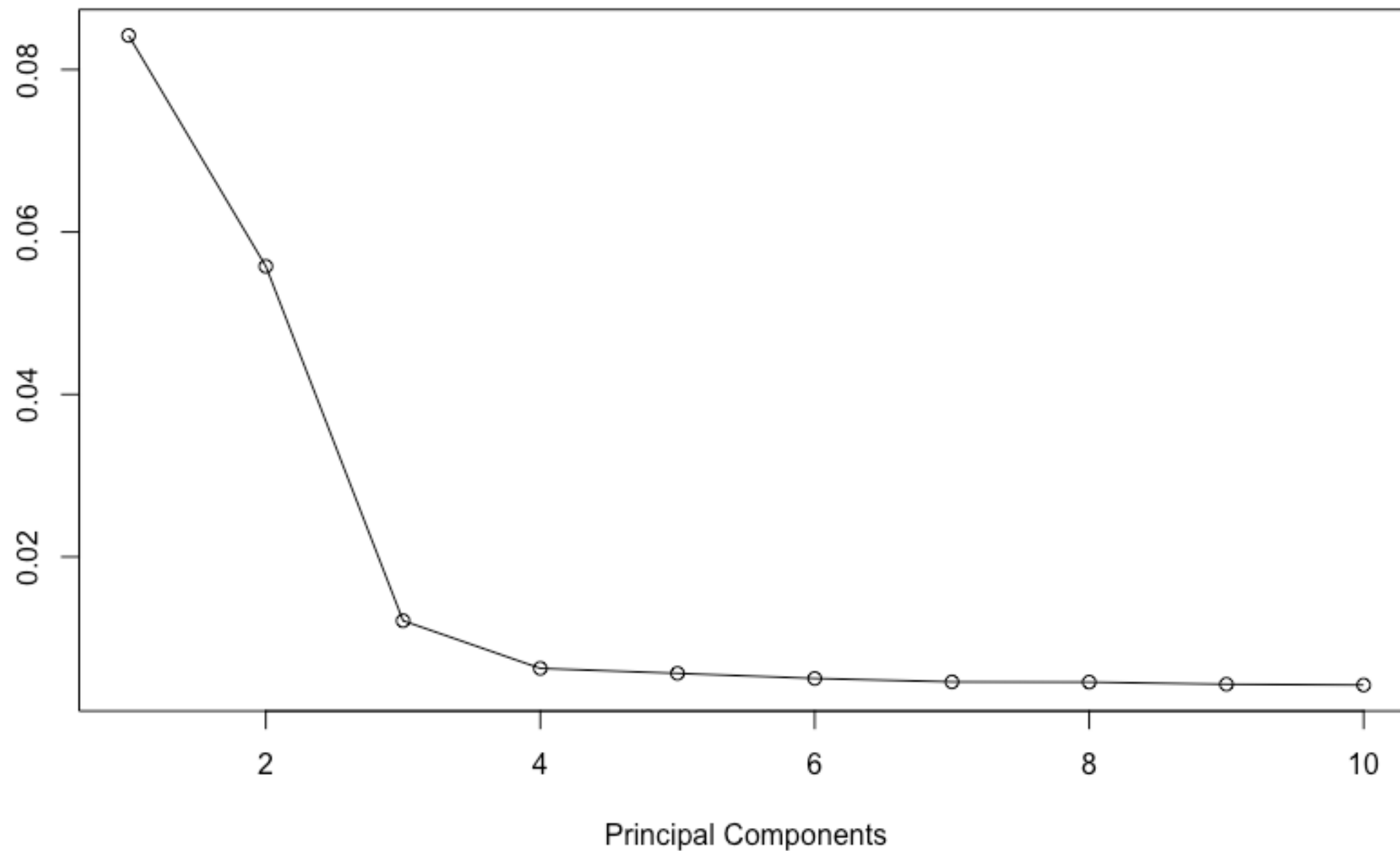
- ▶ Starting points matter, multiple restarts with different start seeds
- ▶ Time Complexity:  
 $O(\text{\#iterations} * \text{\#clusters} * \text{\#individuals} * \text{\#SNPs})$
- ▶ Space Complexity:  $O((\text{\#individuals} + \text{\#clusters}) * \text{\# SNPs})$

# BETTER ANCESTRY MAPPING (B.A.M.)

- ▶ Reduce Dimensionality using Principal Component Analysis (PCA)
  - ▶ Returns orthogonal “dimensions” of highest variance, principal components through eigen decomposition of covariance matrix
  - ▶ Time Complexity:  $O(\min(\text{\#snps}^3, \text{\#individuals}^3))$
- ▶ Benefits of PCA
  - ▶ Easier to visualize and interpret data
  - ▶ Reduce dimensionality of data to decrease runtime and memory usage

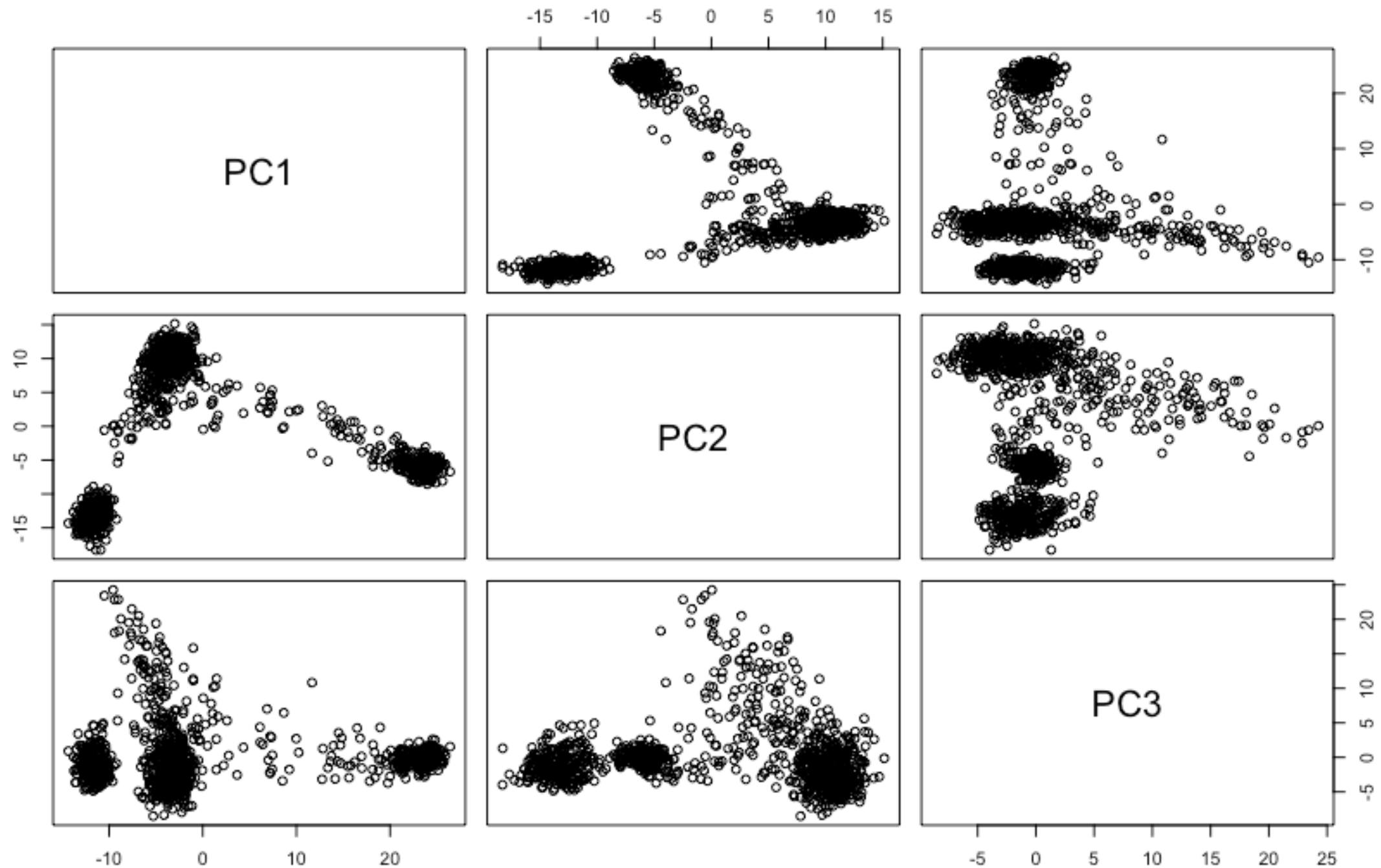
## PCA APPLIED ON CHROMOSOME 22

Scree Plot





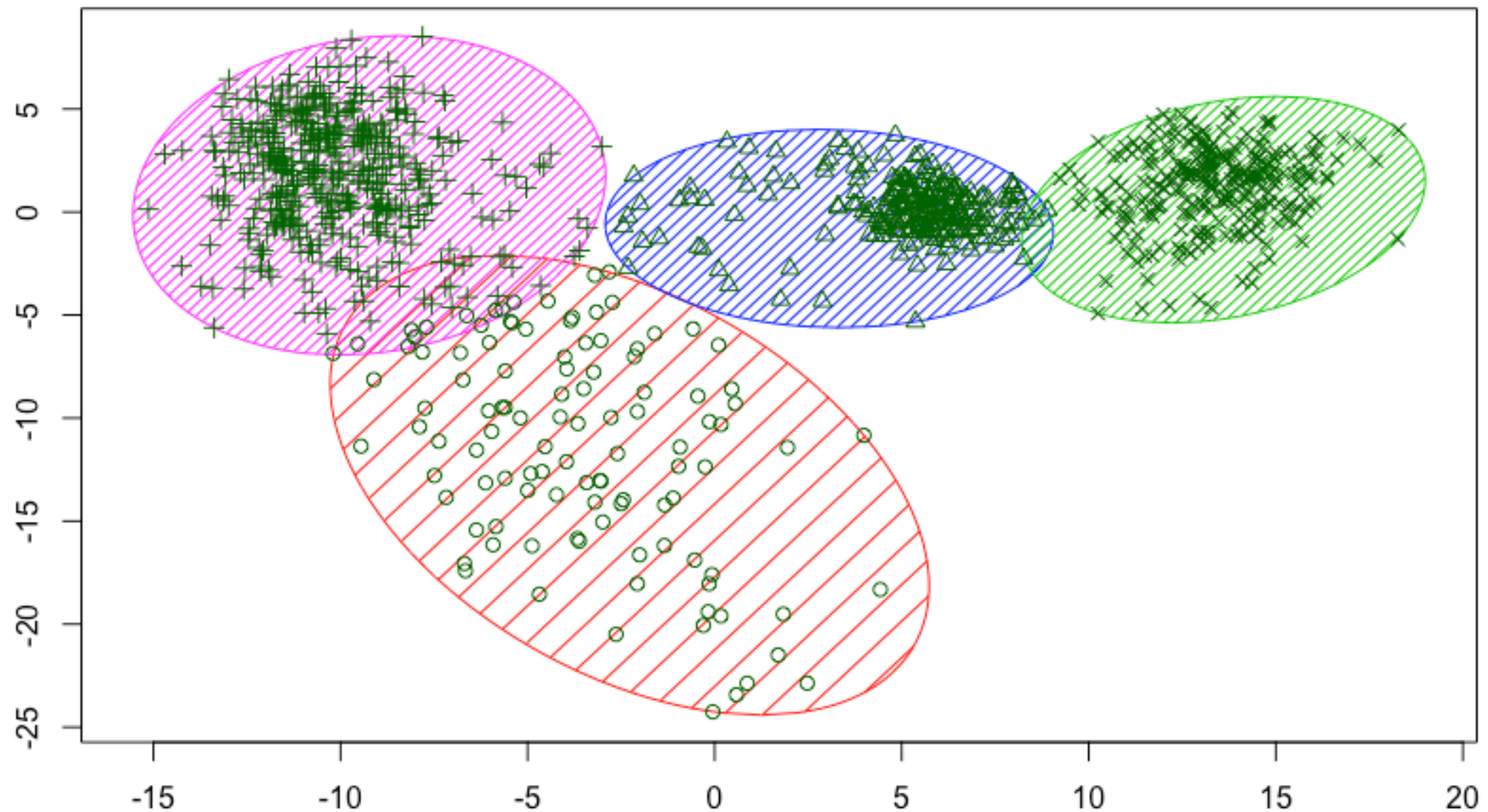
## PCA APPLIED ON CHROMOSOME 22



### BIG PICTURE

- ▶ Compute PCA for each chromosome (can be distributed)
- ▶ Determine optimum number of principal component using scree plot, and plot of pairs
- ▶ Run Kmeans on aggregate principal components from each chromosome
- ▶ Train and predict using kfold cross validation

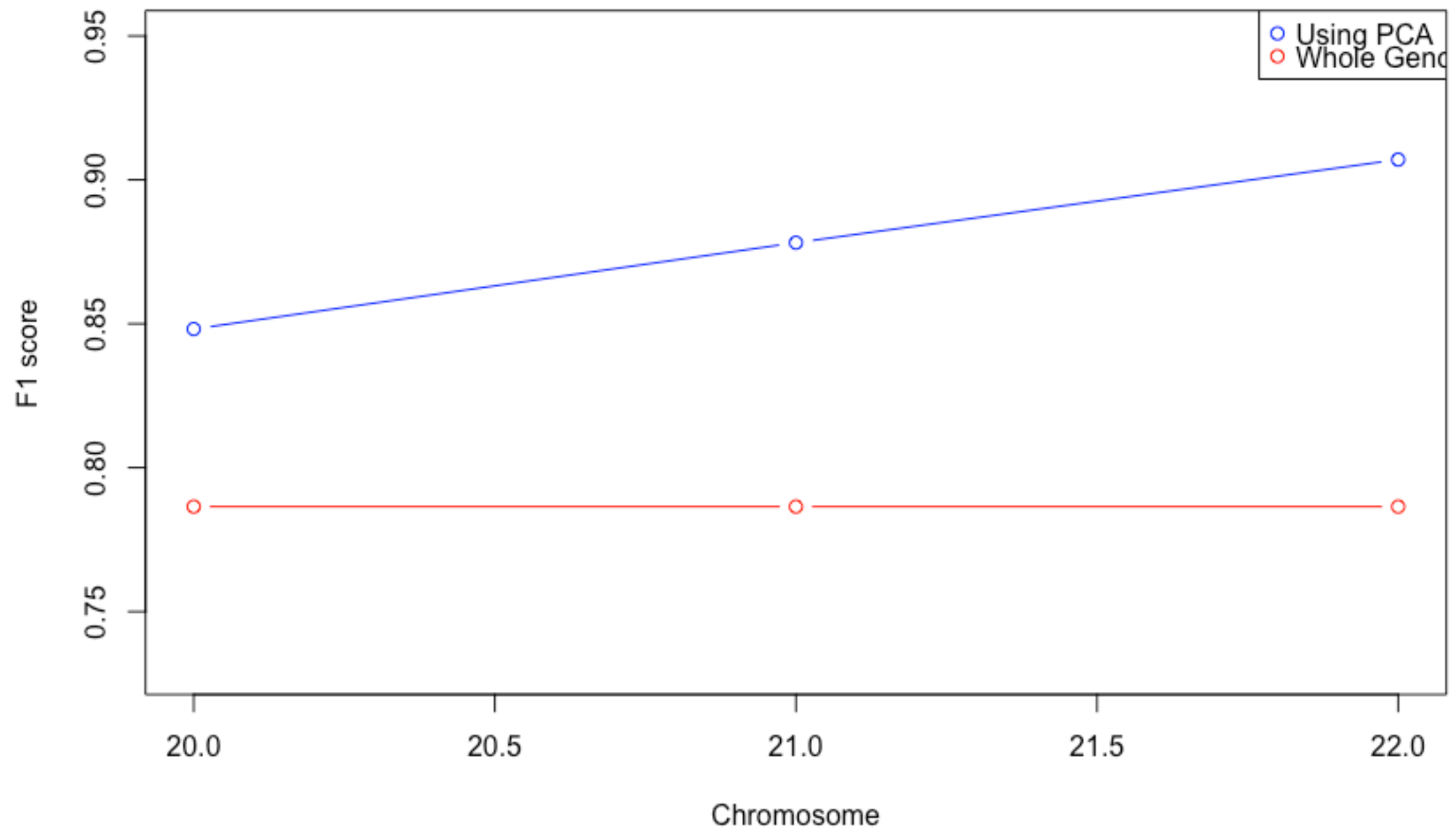
# TWO DIMENSIONAL VISUALIZATION OF KMEANS



Component 1

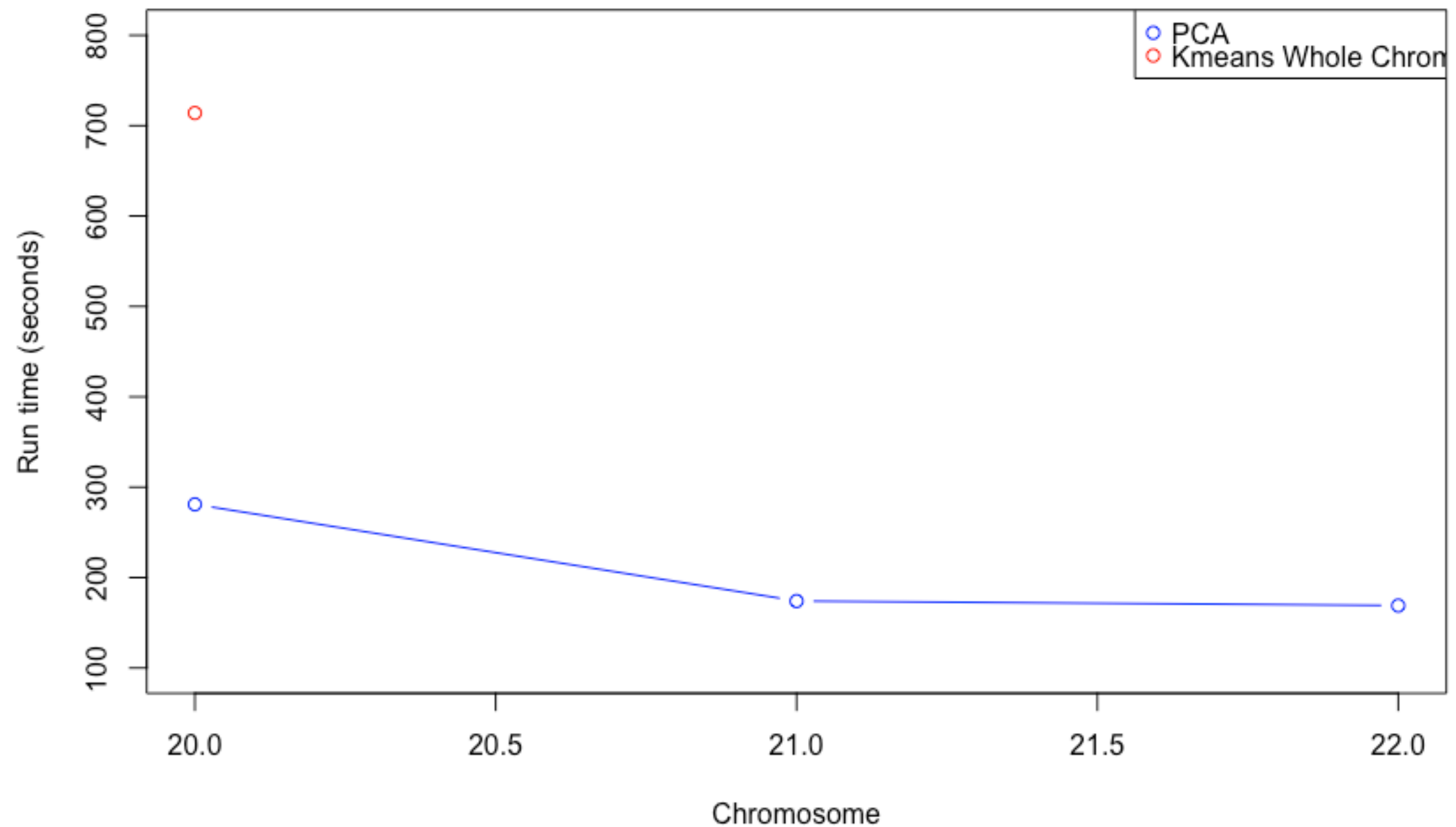
These two components explain 100 % of the point variability.

# BENCHMARK - F1 SCORE



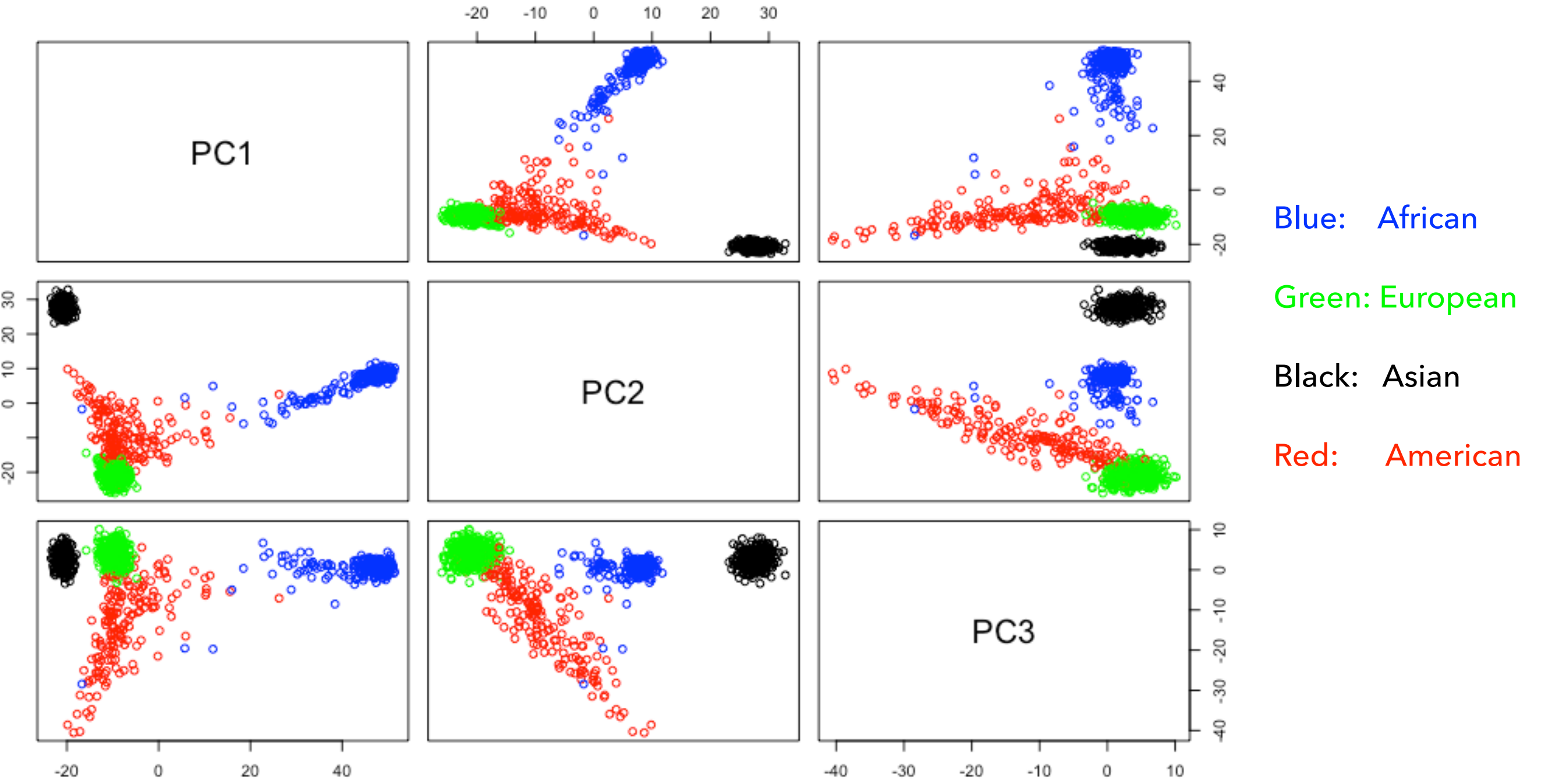
NOTE: Kmeans was only able to run on 1 chromosome

# BENCHMARK – RUN TIME

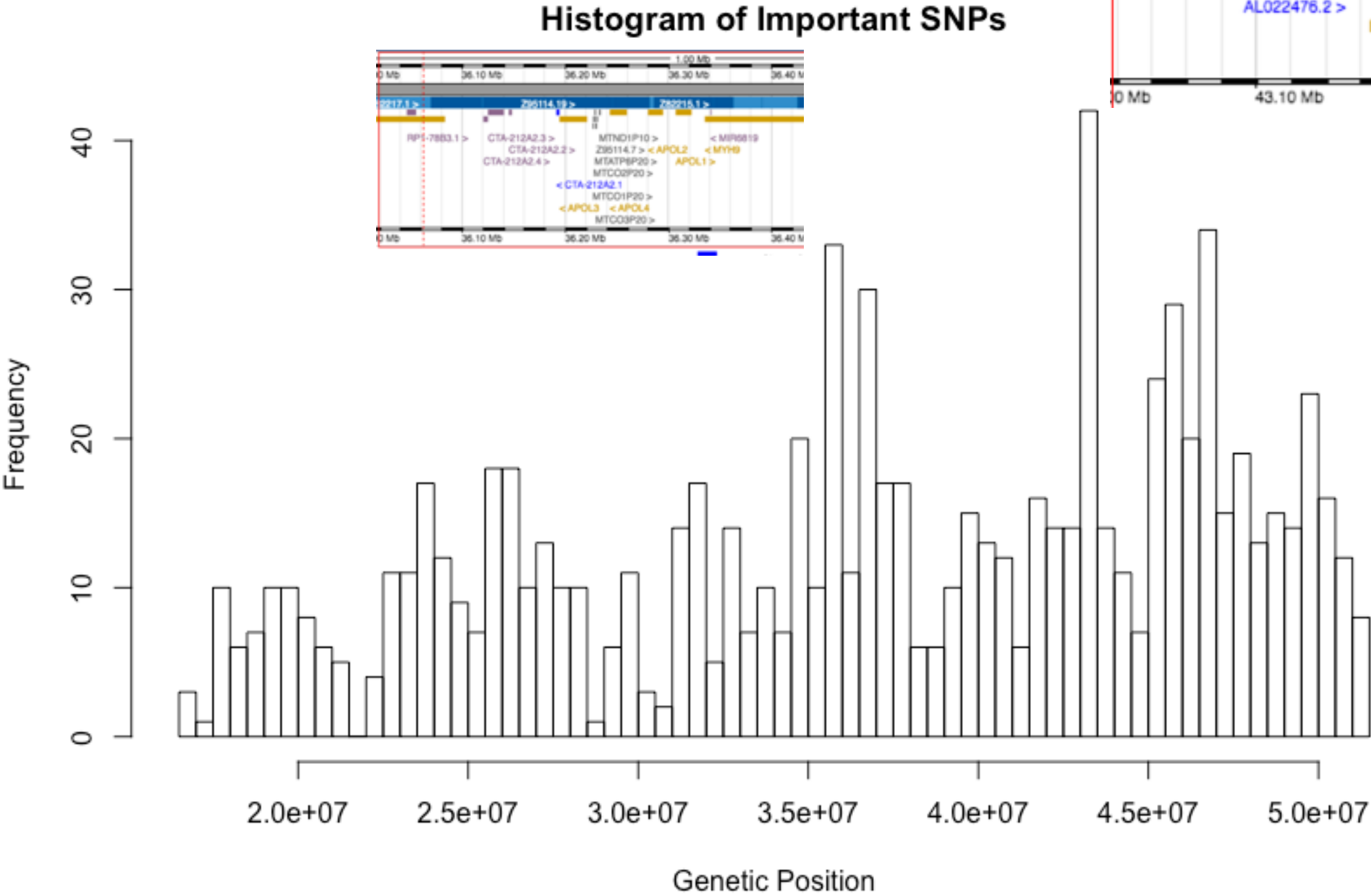


NOTE: Kmeans was only able to run on 1 chromosome

# PCAS PERFORMANCE CLASSIFYING DATA



SPARSE PCA





### DISCUSSION AND FUTURE WORK

- ▶ Kmeans clustering on the entire chromosome seems to suffer from the high dimensionality of the data
- ▶ Reducing each chromosome to its principal components and aggregating based increases F1 score and can be done concurrently
- ▶ Europeans and Americans populations are genetically very similar, a more sophisticated probabilistic approach like admixture can provide for better results