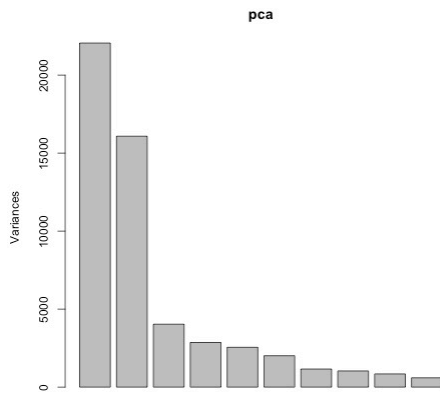
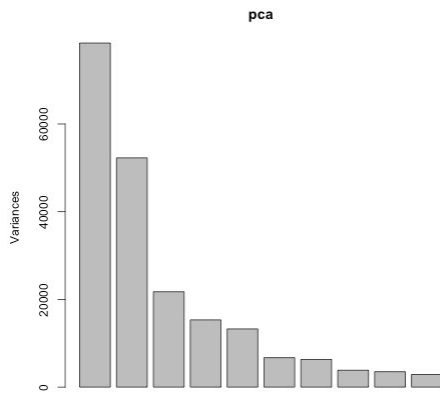
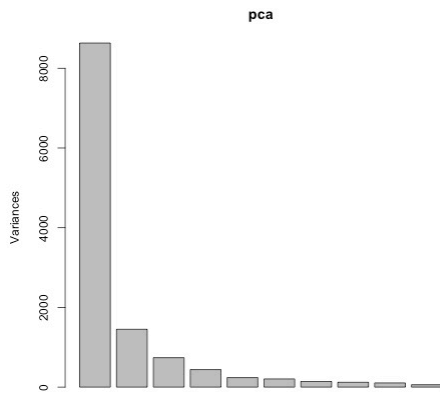
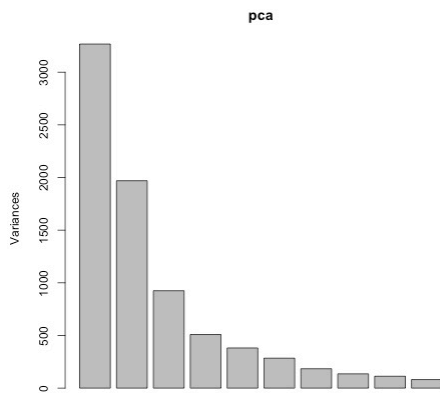
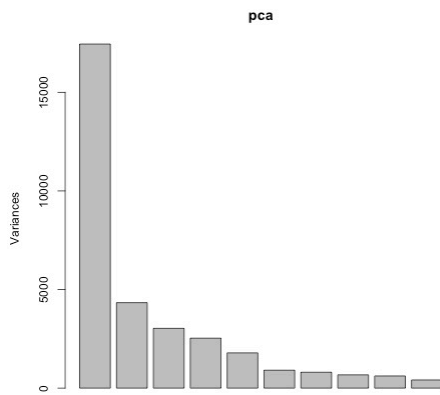
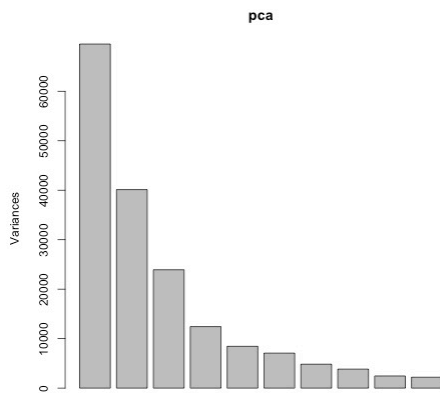
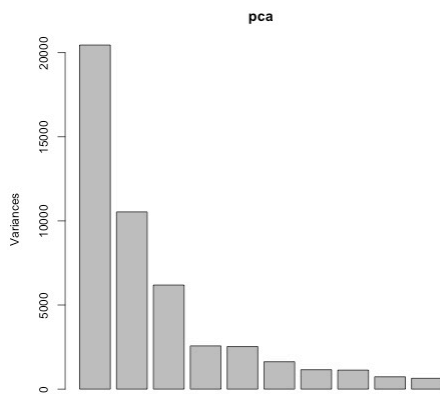
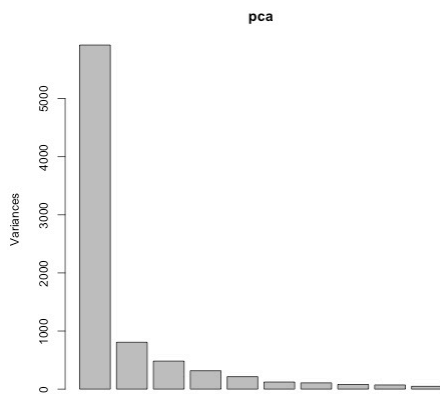
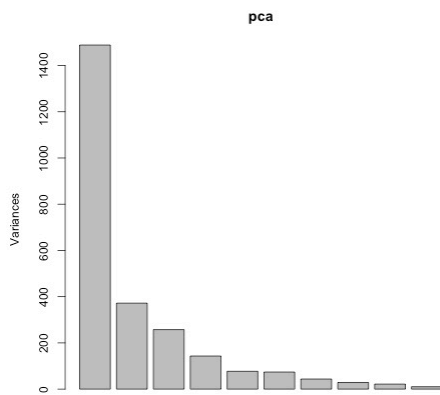


# Demonstrating that Histone Modifications are Cell-Type Specific Through Clustering and Identifying Biologically Relevant Domains

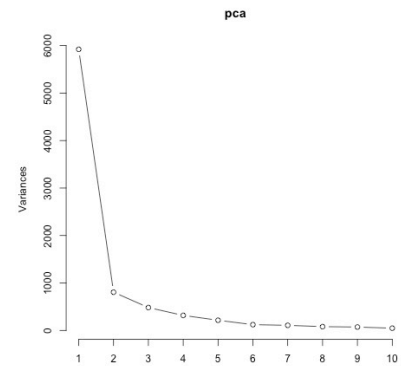
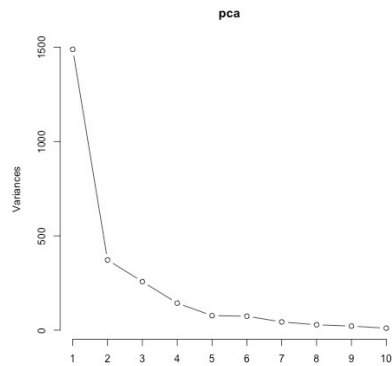
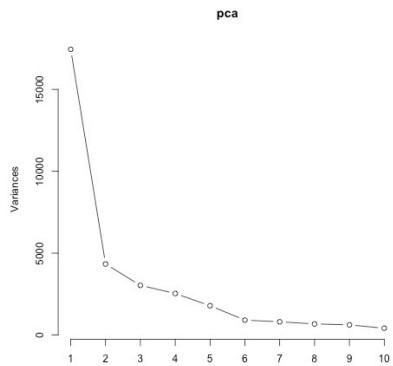
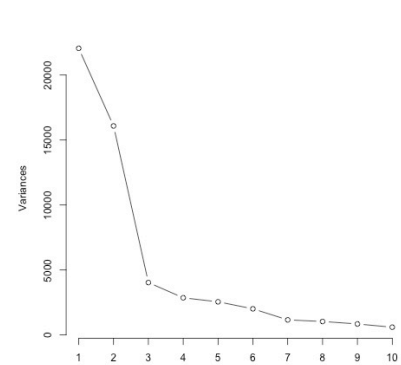
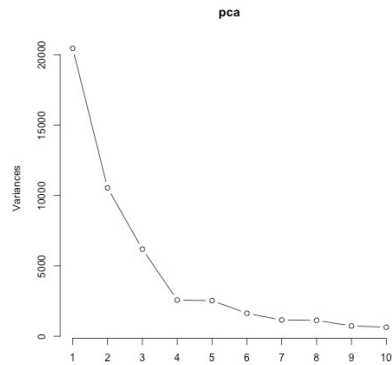
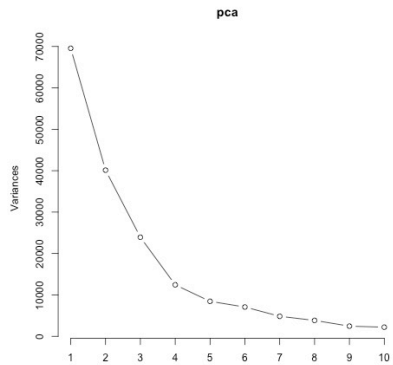
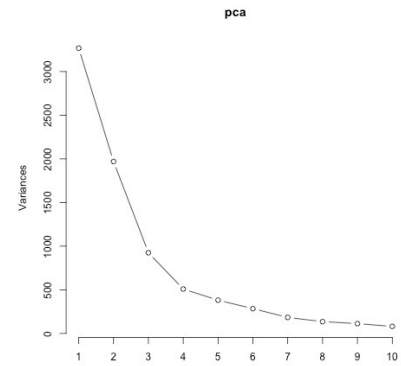
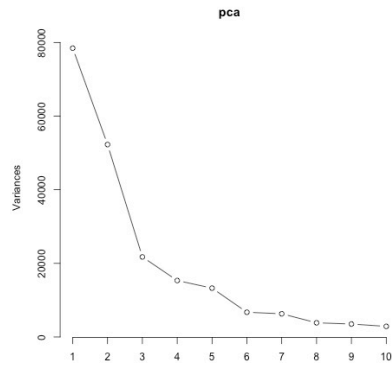
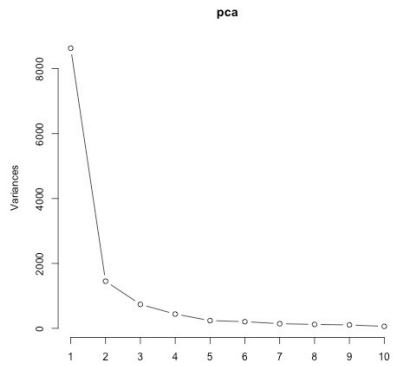
Adriana Sperlea & Douglas Arneson

# Methods

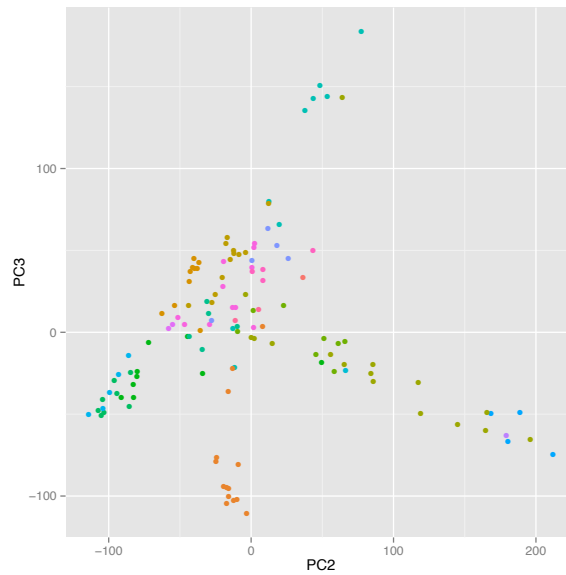
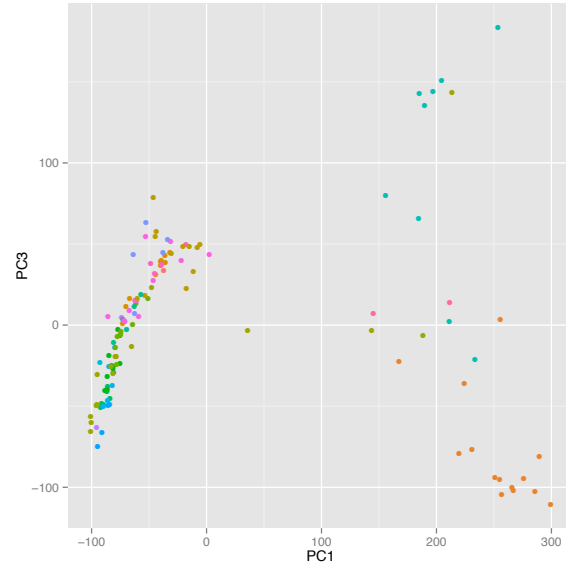
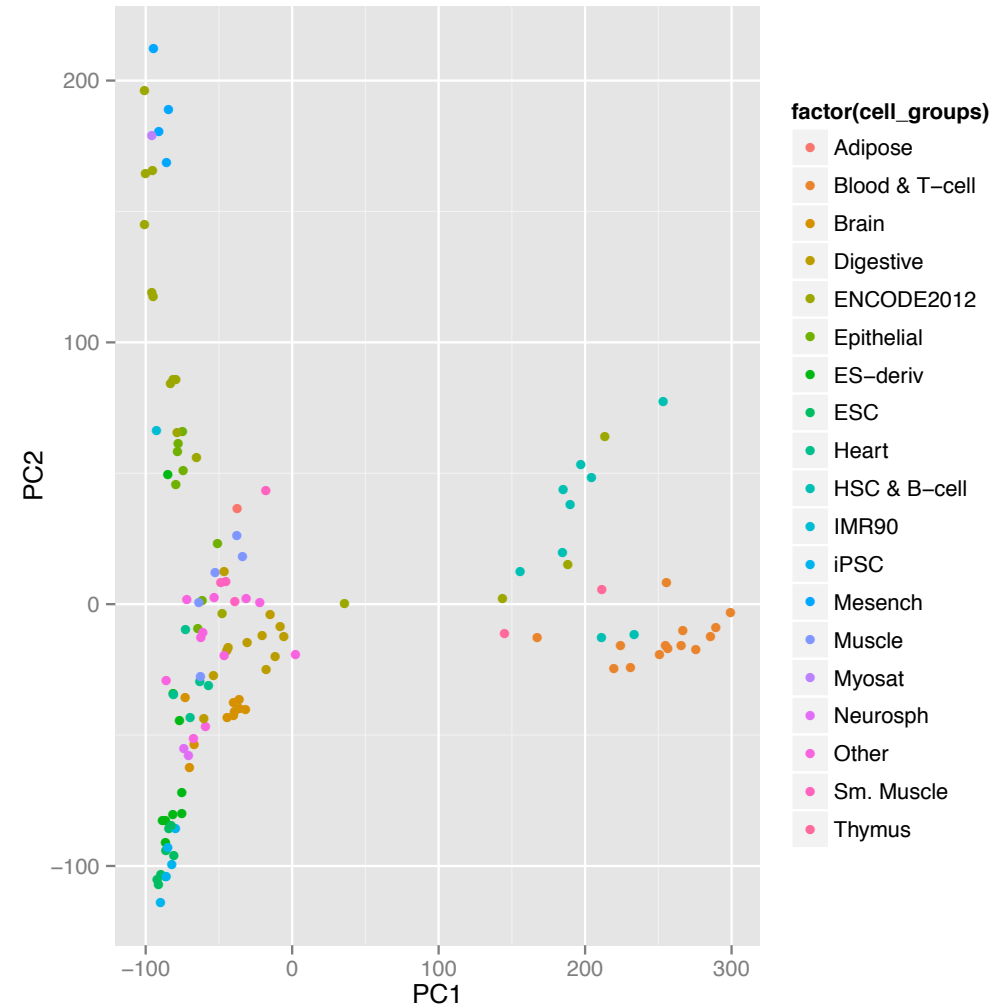
- PCA analysis across cell-types
- K-means clustering to identify cell type specificity of histone modifications
- Analysis of variance to identify important genomic locations



# Scree Plots



# K-Means Clustering



# Subset Data: Remove “Other” and “ENCODE” Cell Types

GROUP					
IMR90	Blood & T-cell	Mesench	Brain	Digestive	ENCODE2012
ESC	Blood & T-cell	Mesench	Brain	Digestive	ENCODE2012
ESC	Blood & T-cell	Mesench	Brain	Digestive	ENCODE2012
ESC	Blood & T-cell	Mesench	Brain	Digestive	ENCODE2012
ESC	Blood & T-cell	Myosat	Adipose	Digestive	ENCODE2012
ESC	Blood & T-cell	Epithelial	Muscle	Digestive	ENCODE2012
ESC	Blood & T-cell	Epithelial	Muscle	Digestive	ENCODE2012
ESC	Blood & T-cell	Epithelial	Muscle	Digestive	ENCODE2012
ESC	Blood & T-cell	Epithelial	Muscle	Other	ENCODE2012
iPSC	Blood & T-cell	Epithelial	Muscle	Other	ENCODE2012
iPSC	Blood & T-cell	Epithelial	Heart	Other	ENCODE2012
iPSC	Blood & T-cell	Epithelial	Heart	Other	ENCODE2012
iPSC	Blood & T-cell	Epithelial	Heart	Other	
iPSC	Blood & T-cell	Neurosph	Heart	Other	
ES-deriv	HSC & B-cell	Neurosph	Heart	Other	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	Other	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	Other	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	Other	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	Other	
ES-deriv	HSC & B-cell	Brain	Digestive	ENCODE2012	
ES-deriv	HSC & B-cell	Brain	Digestive	ENCODE2012	
ES-deriv	HSC & B-cell	Brain	Digestive	ENCODE2012	
ES-deriv	HSC & B-cell	Brain	Digestive	ENCODE2012	

# Subset Data: Remove “Other” and “ENCODE” Cell Types

GROUP				
IMR90	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Myosat	Adipose	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Neurosph	Heart	
ES-deriv	HSC & B-cell	Neurosph	Heart	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	

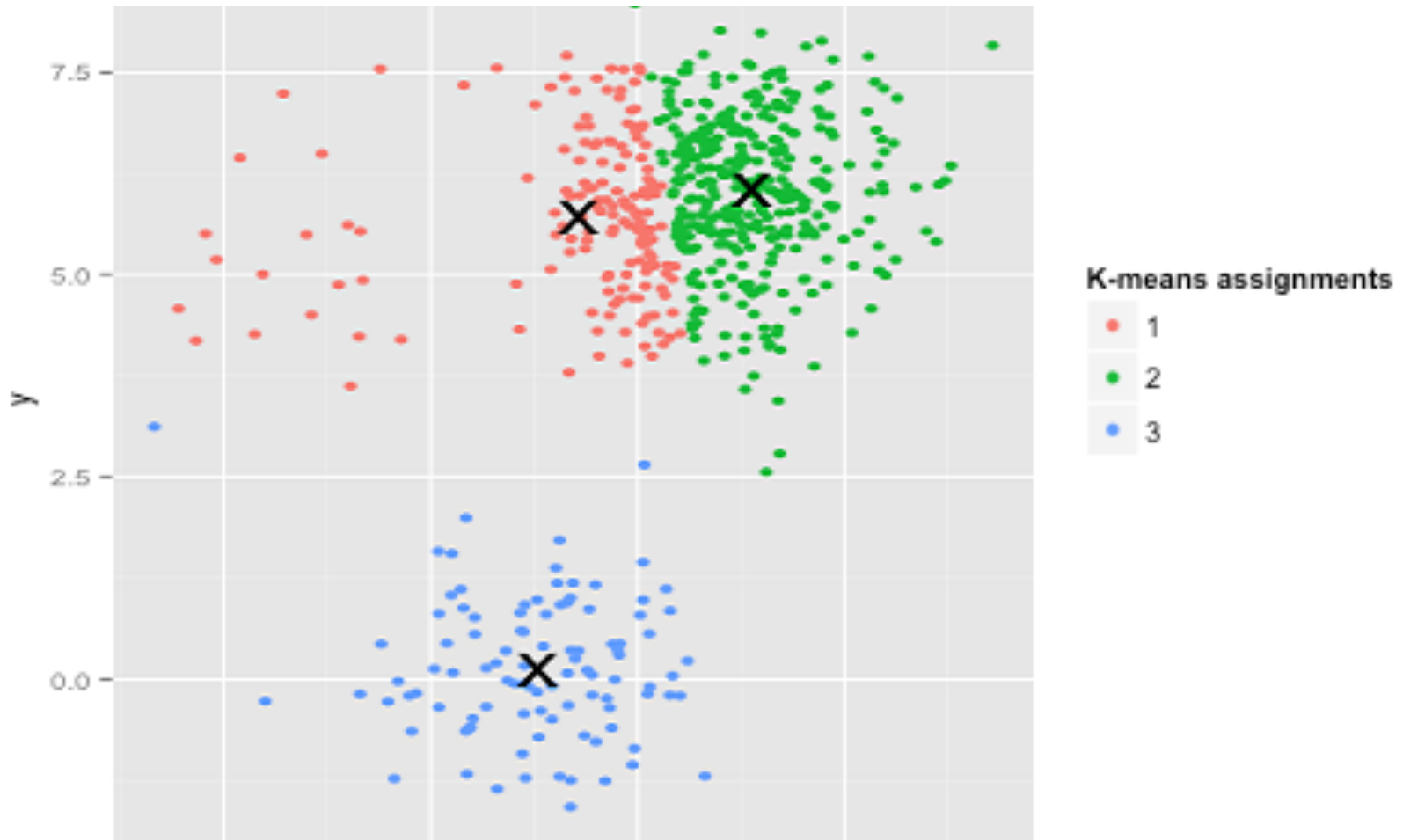
# Problem with K-means Algorithm

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

- K-means attempts to minimize the objective function J
- The term inside the summations gives the distance measure between the chosen point and the cluster center
- An issue with this algorithm is by minimizing the sum of squares within clusters, more weight is given to larger clusters



# Problem with K means Algorithm



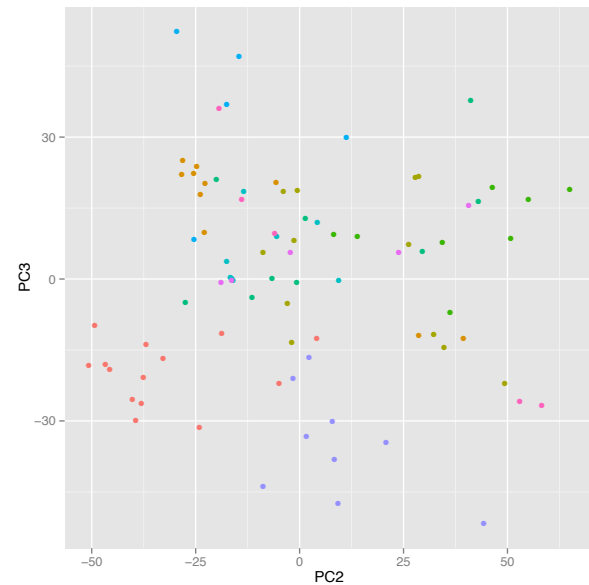
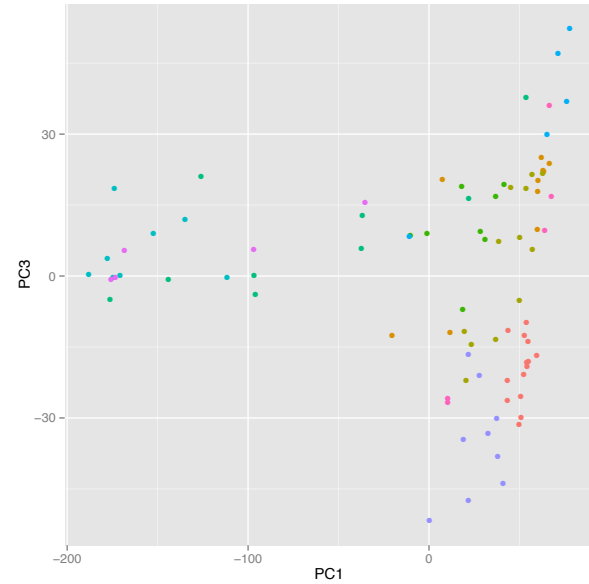
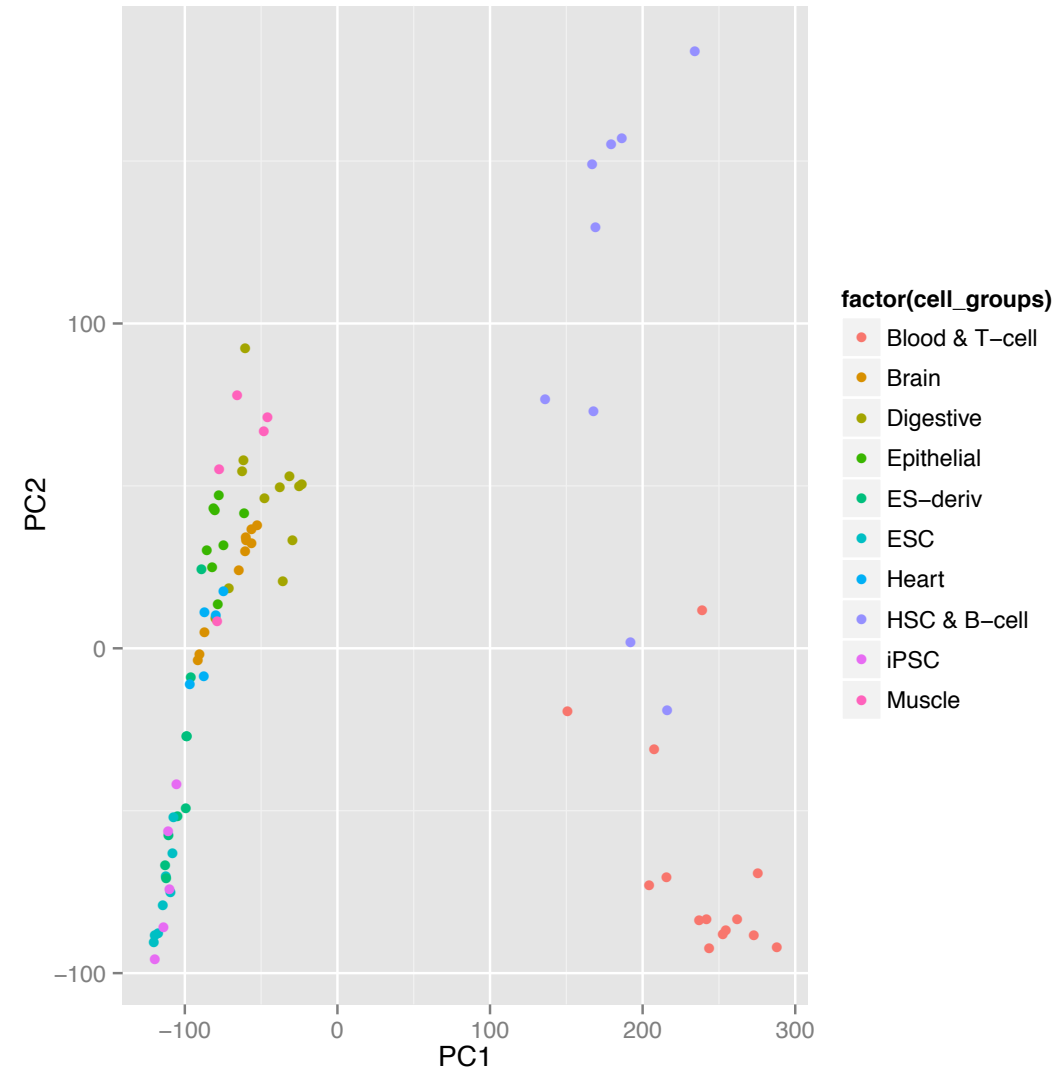
# Further Subset Data: Remove Cell Types with Less than Five Entries

GROUP				
IMR90	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Mesench	Brain	Digestive
ESC	Blood & T-cell	Myosat	Adipose	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Neurosph	Heart	
ES-deriv	HSC & B-cell	Neurosph	Heart	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	
ES-deriv	HSC & B-cell	Thymus	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Sm. Muscle	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	

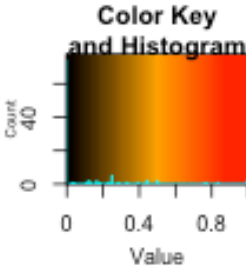
# Further Subset Data: Remove Cell Types with Less than Five Entries

GROUP				
	Blood & T-cell		Brain	Digestive
ESC	Blood & T-cell		Brain	Digestive
ESC	Blood & T-cell		Brain	Digestive
ESC	Blood & T-cell		Brain	Digestive
ESC	Blood & T-cell			Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	Digestive
ESC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Muscle	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell	Epithelial	Heart	
iPSC	Blood & T-cell		Heart	
ES-deriv	HSC & B-cell		Heart	
ES-deriv	HSC & B-cell			
ES-deriv	HSC & B-cell			
ES-deriv	HSC & B-cell	Brain		
ES-deriv	HSC & B-cell	Brain		
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	
ES-deriv	HSC & B-cell	Brain	Digestive	

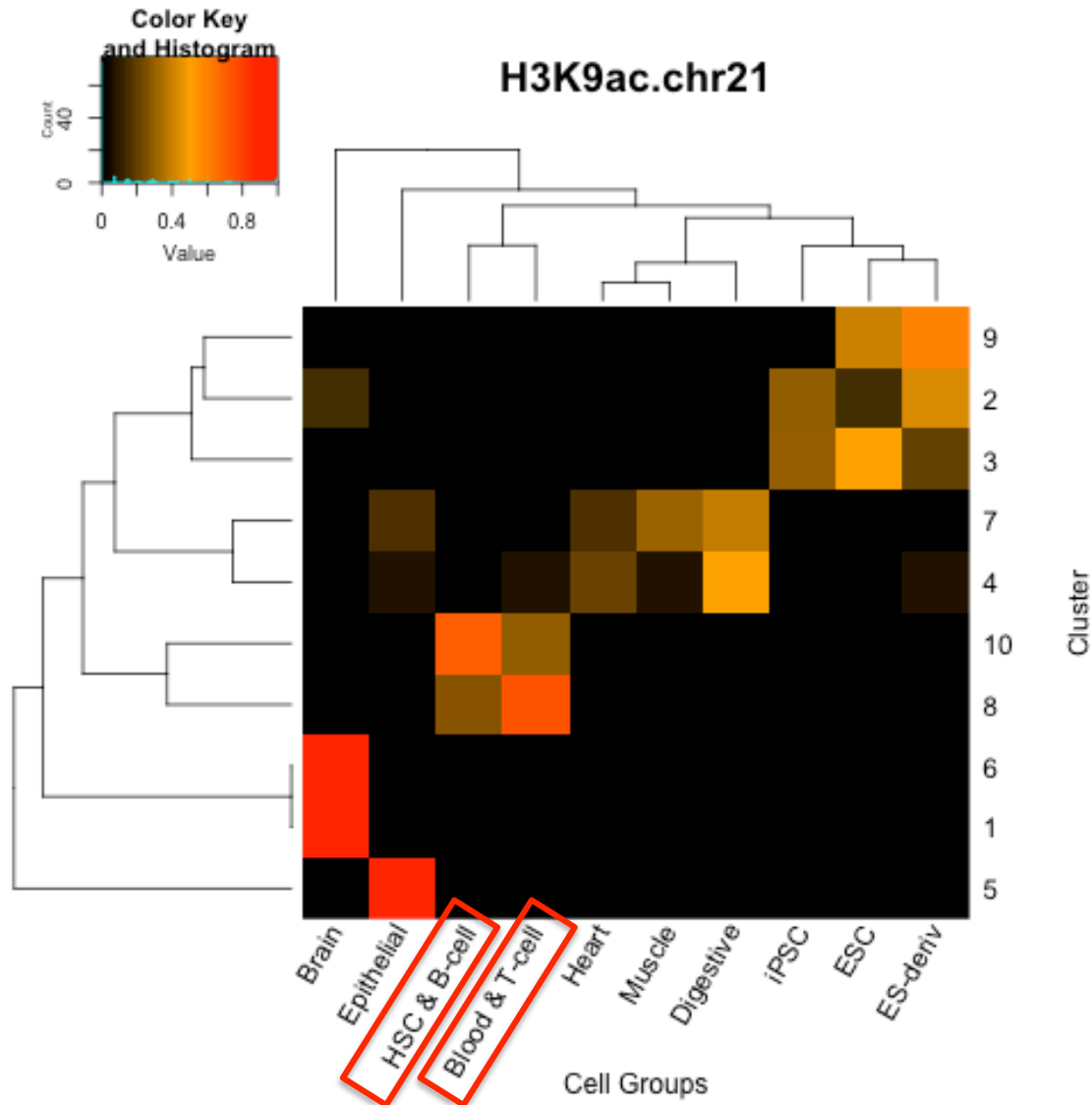
# K-Means Clustering: Refined



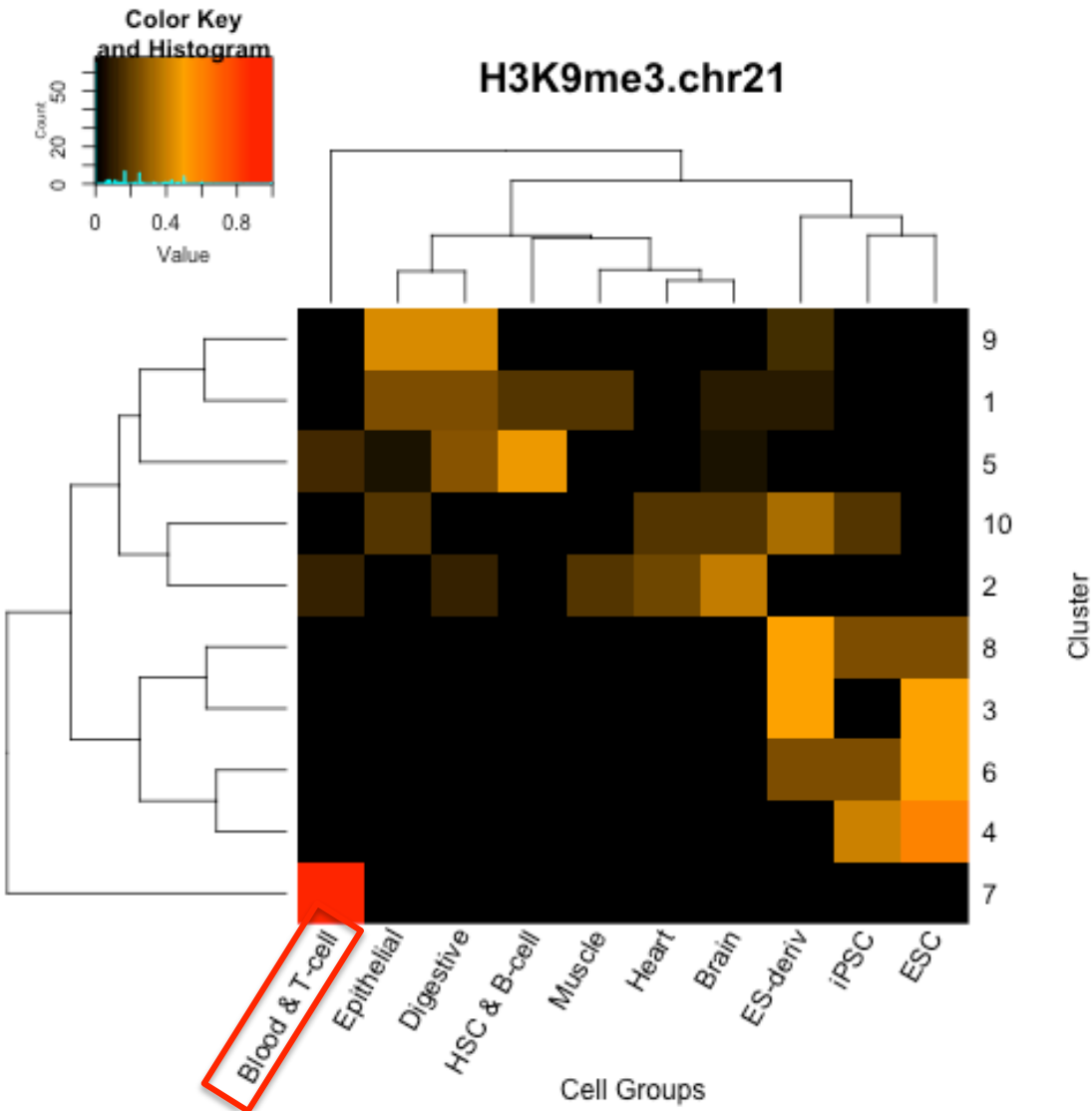
# Heatmaps



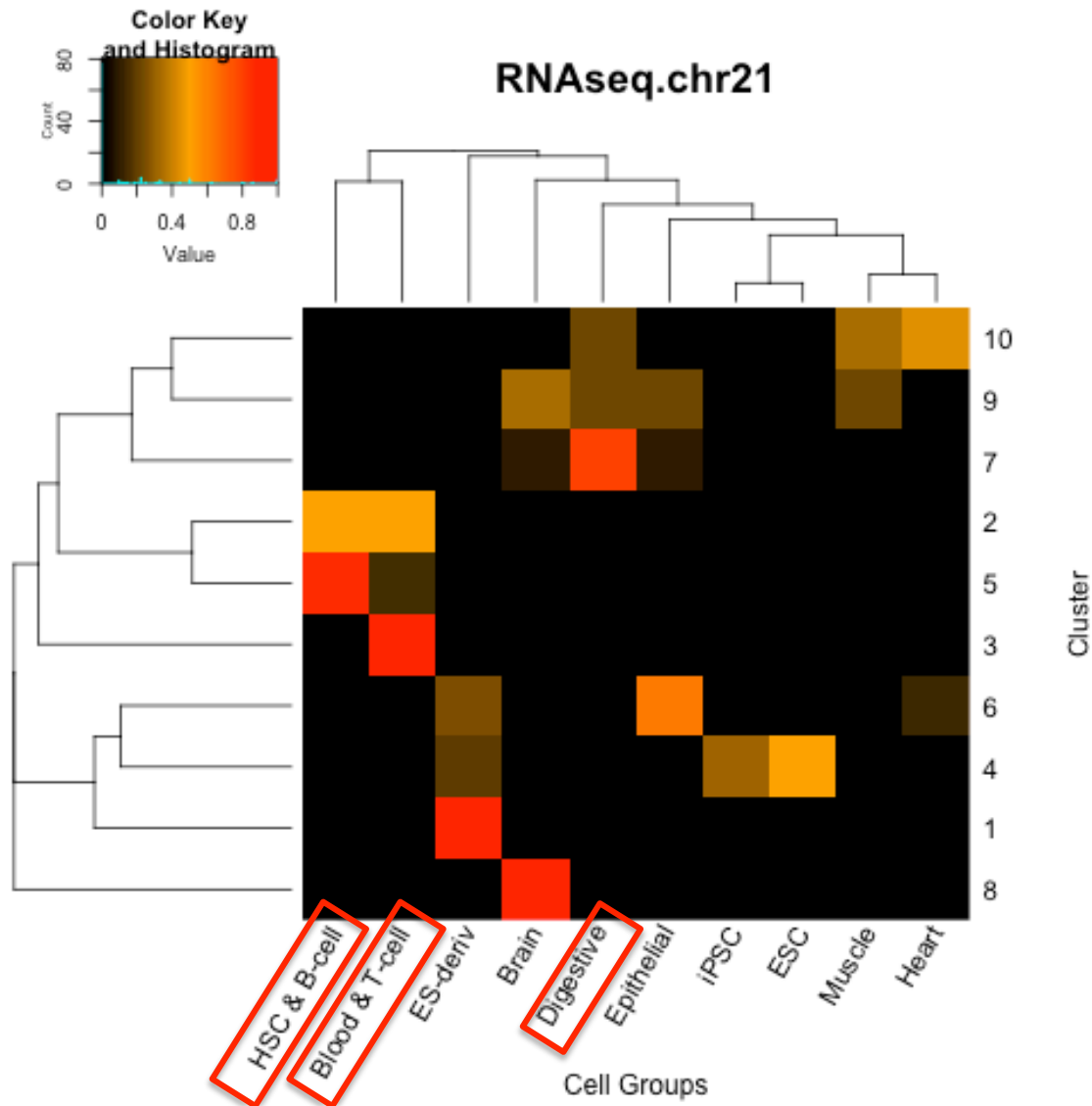
# Visualizing Clustering Efficacy: Heatmaps



# Heatmaps

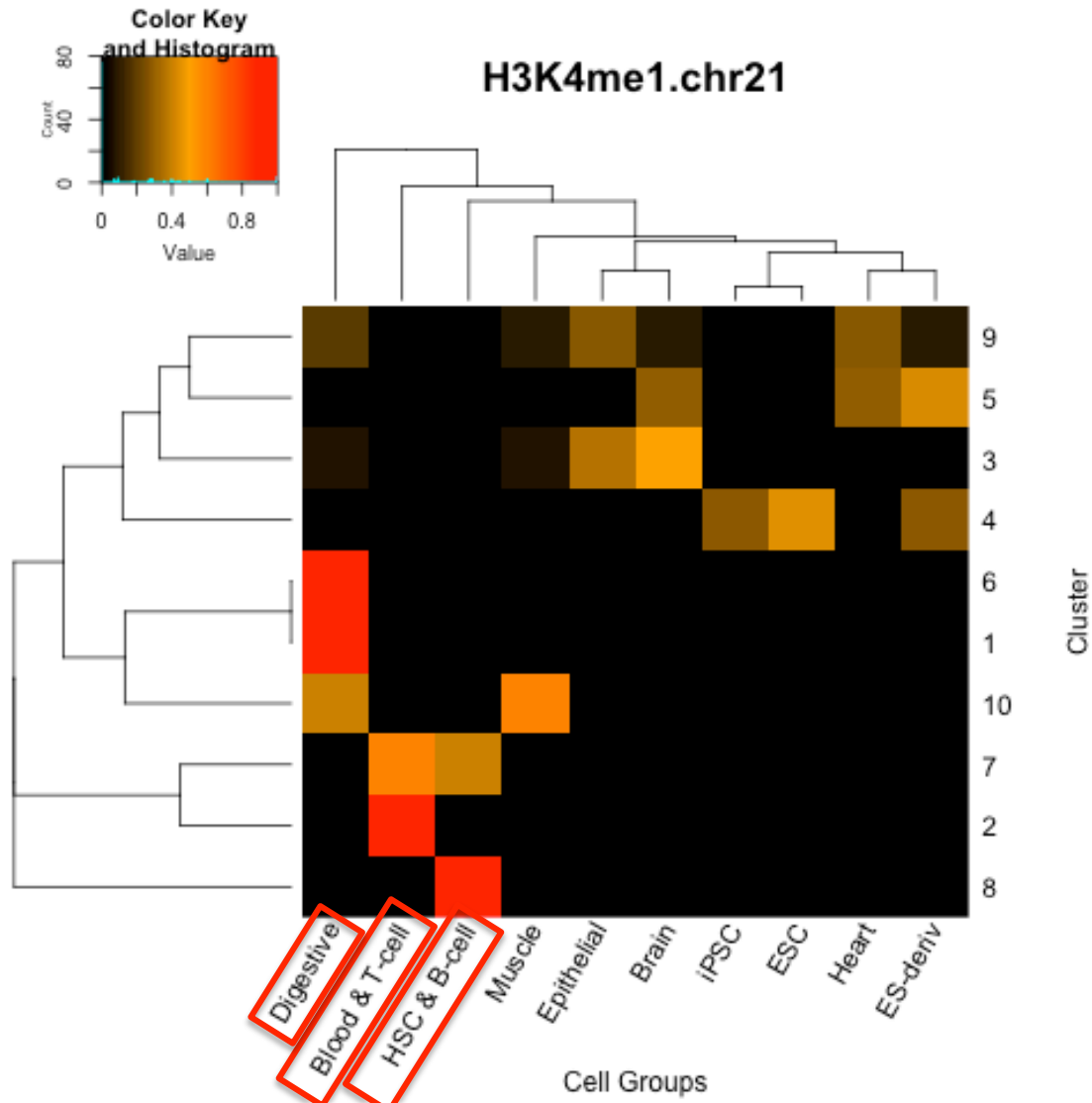


# Heatmaps

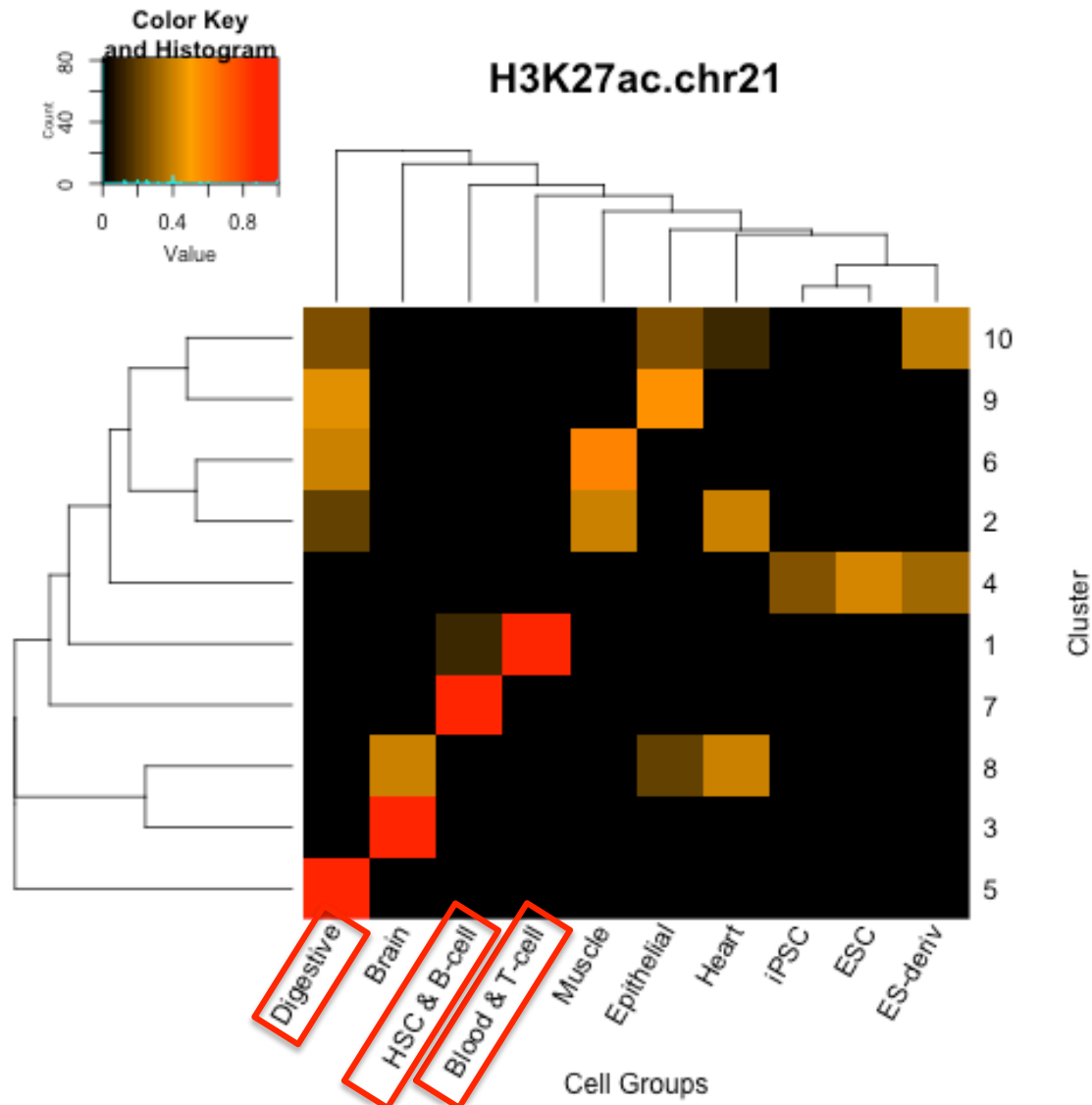




# Visualizing Clustering Efficacy: Heatmaps



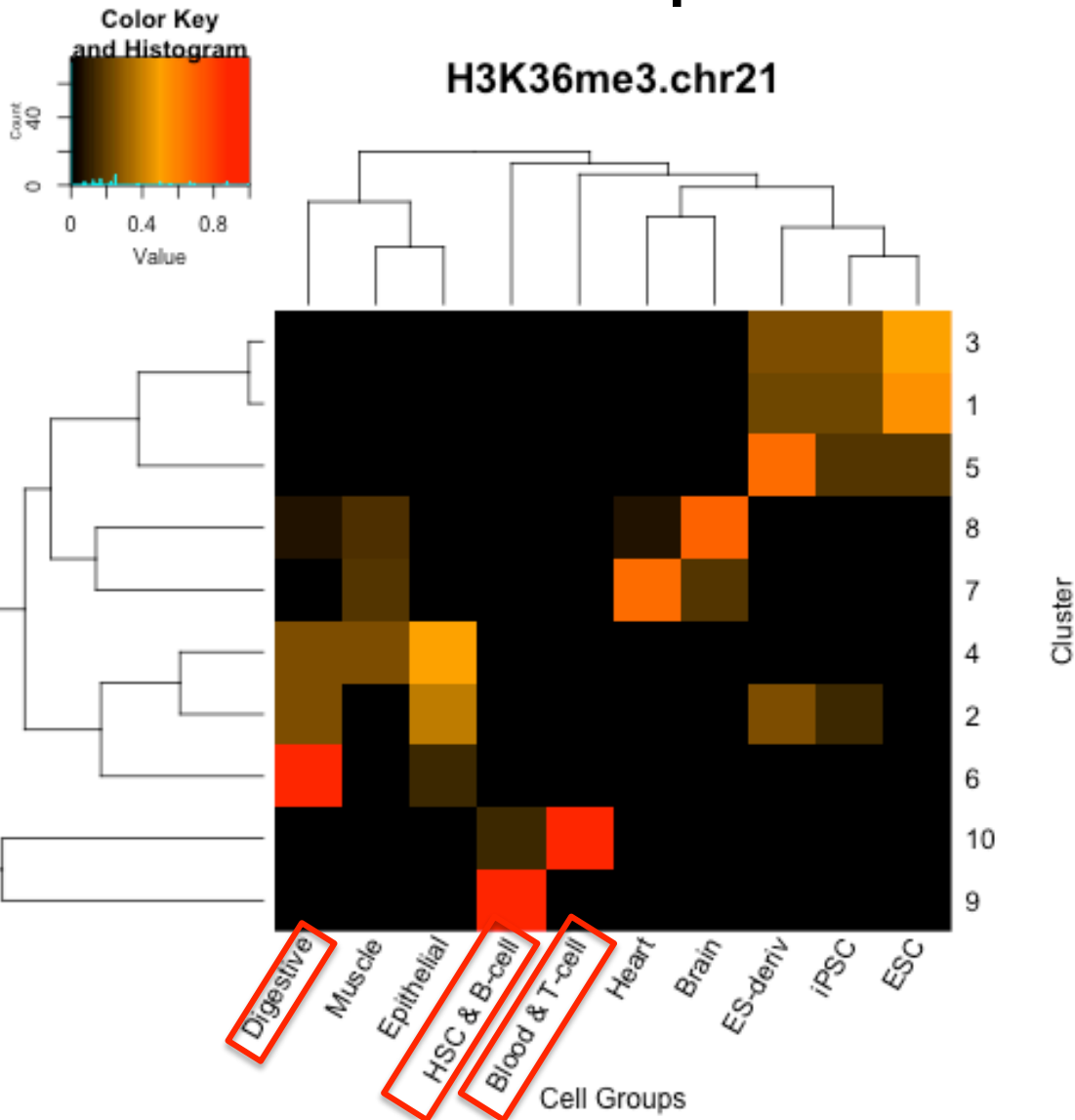
# Visualizing Clustering Efficacy: Heatmaps



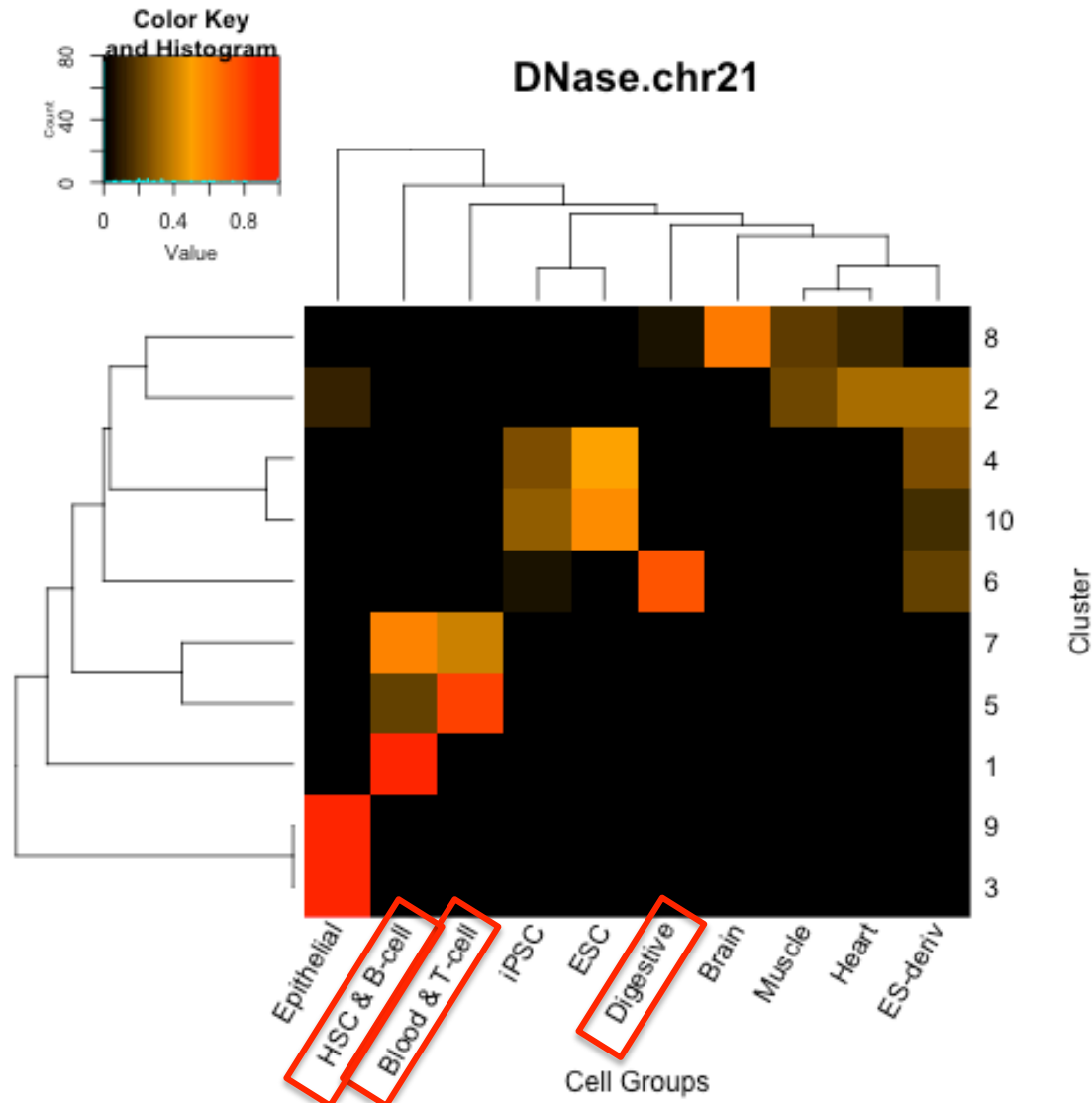
# Heatmaps



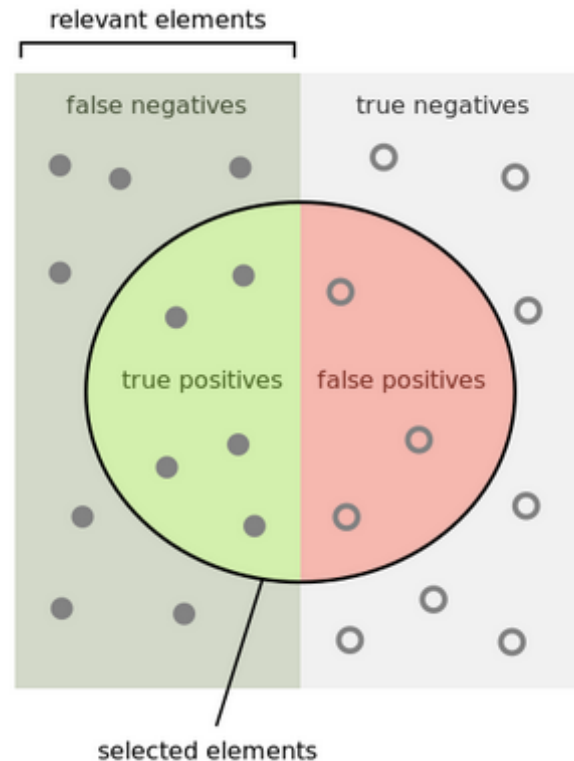
# Heatmaps



# Visualizing Clustering Efficacy: Heatmaps



# Precision and Recall



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

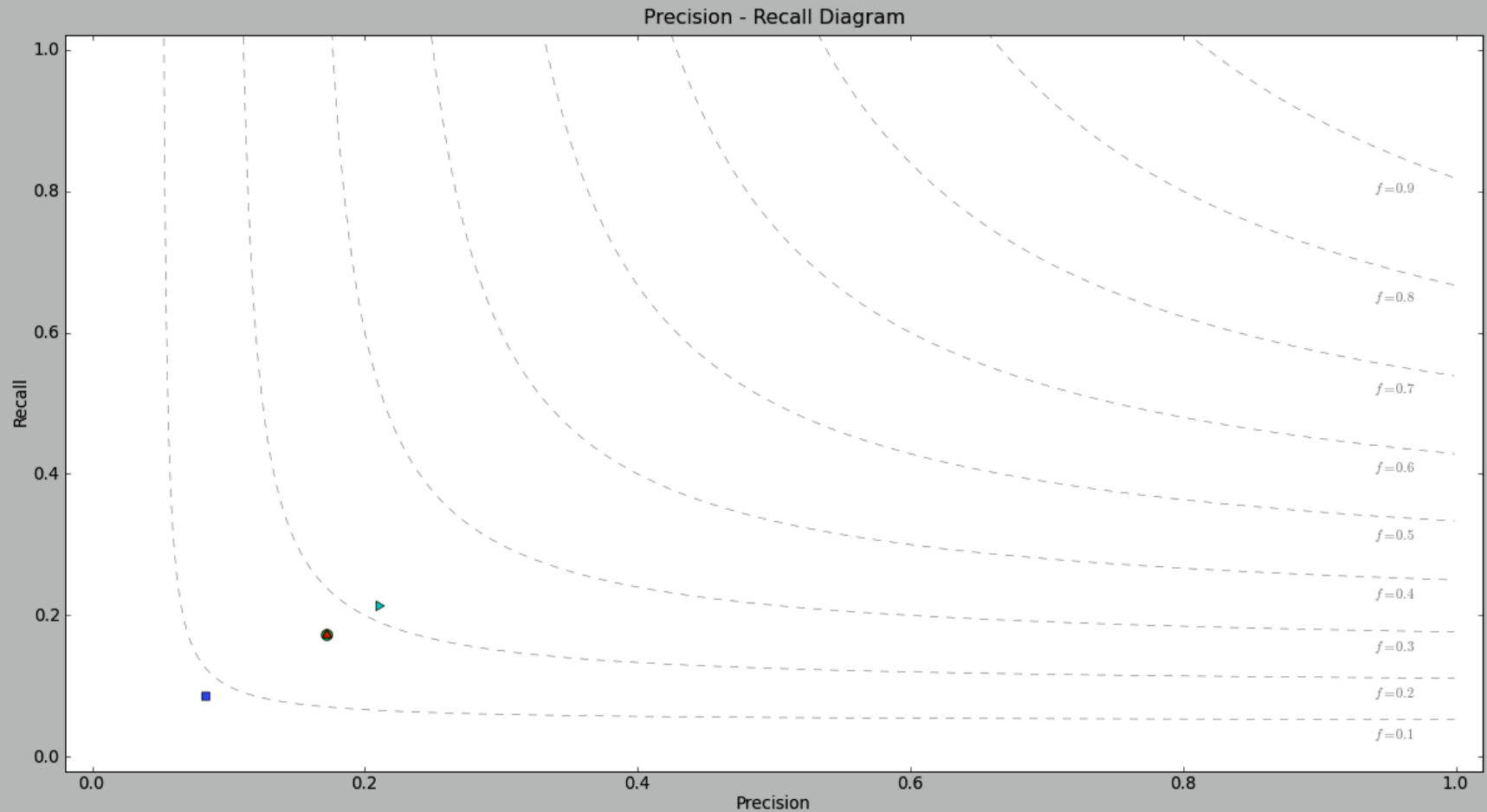
$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

# Evaluating Clustering Efficacy: F-Measure (F-Beta Score)

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- P – precision
- R – recall
- TP – True Positives
- FP – False Positives
- FN – False Negatives
- $\beta$  – weights recall (for values  $>1$ , penalizes false negatives more strongly than false positives)

# Precision & Recall vs F Score





# Evaluating Clustering Efficacy: F-Measure (F-Beta Score)

- **DNase.chr21 - 0.5944444**
- **H3K27ac.chr21 - 0.6277778**
- H3K27me3.chr21 - 0.3944444
- **H3K36me3.chr21 - 0.5805556**
- H3K4me1.chr21 - 0.5166667
- H3K4me3.chr21 - 0.5416667
- H3K9ac.chr21 - 0.4555556
- H3K9me3.chr21 - 0.3555556
- RNAseq.chr21 - 0.5083333

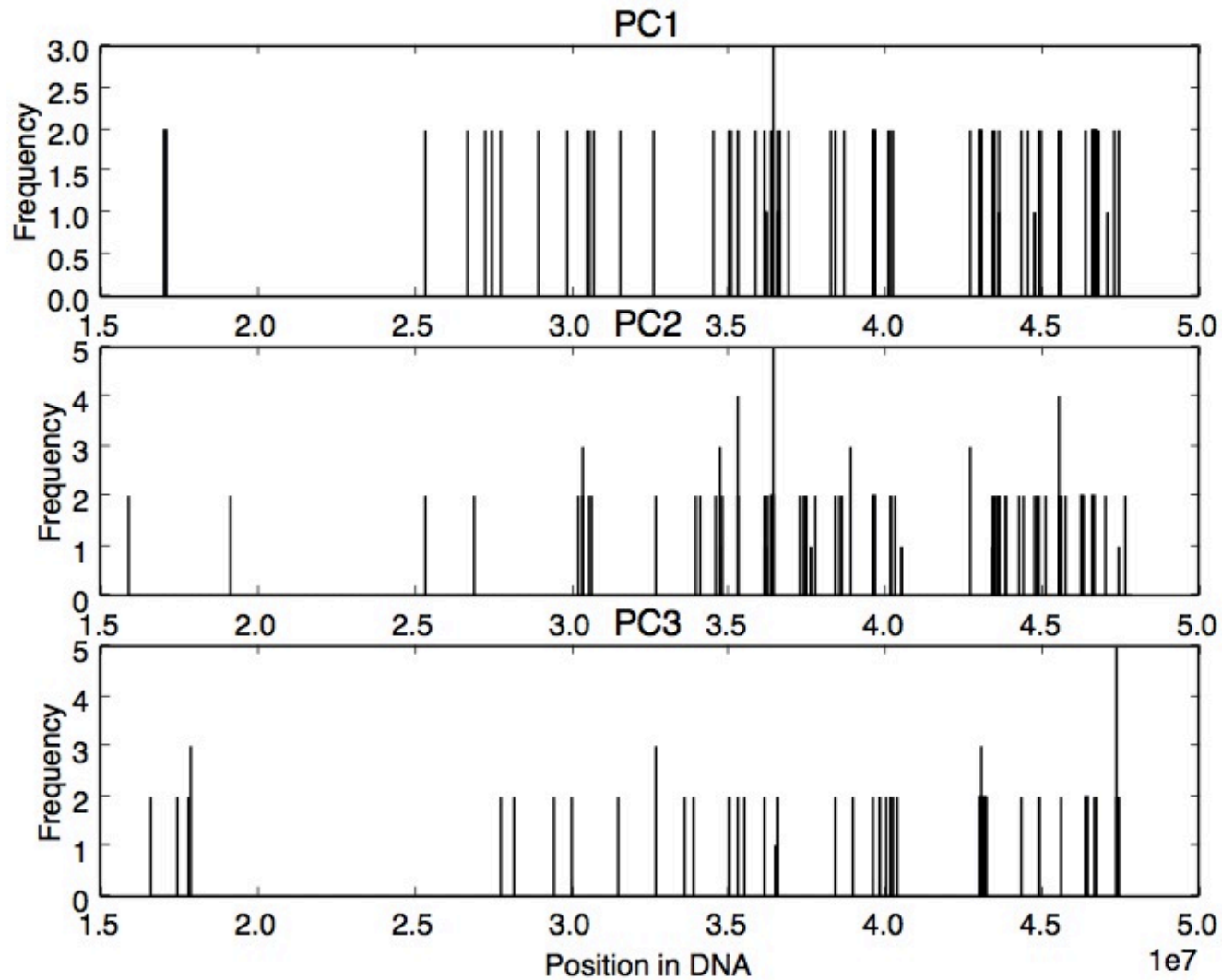
# Analysis of Variance

- Question: Are there certain locations in the genome that explain the PCs and to these locations have an important biological function?

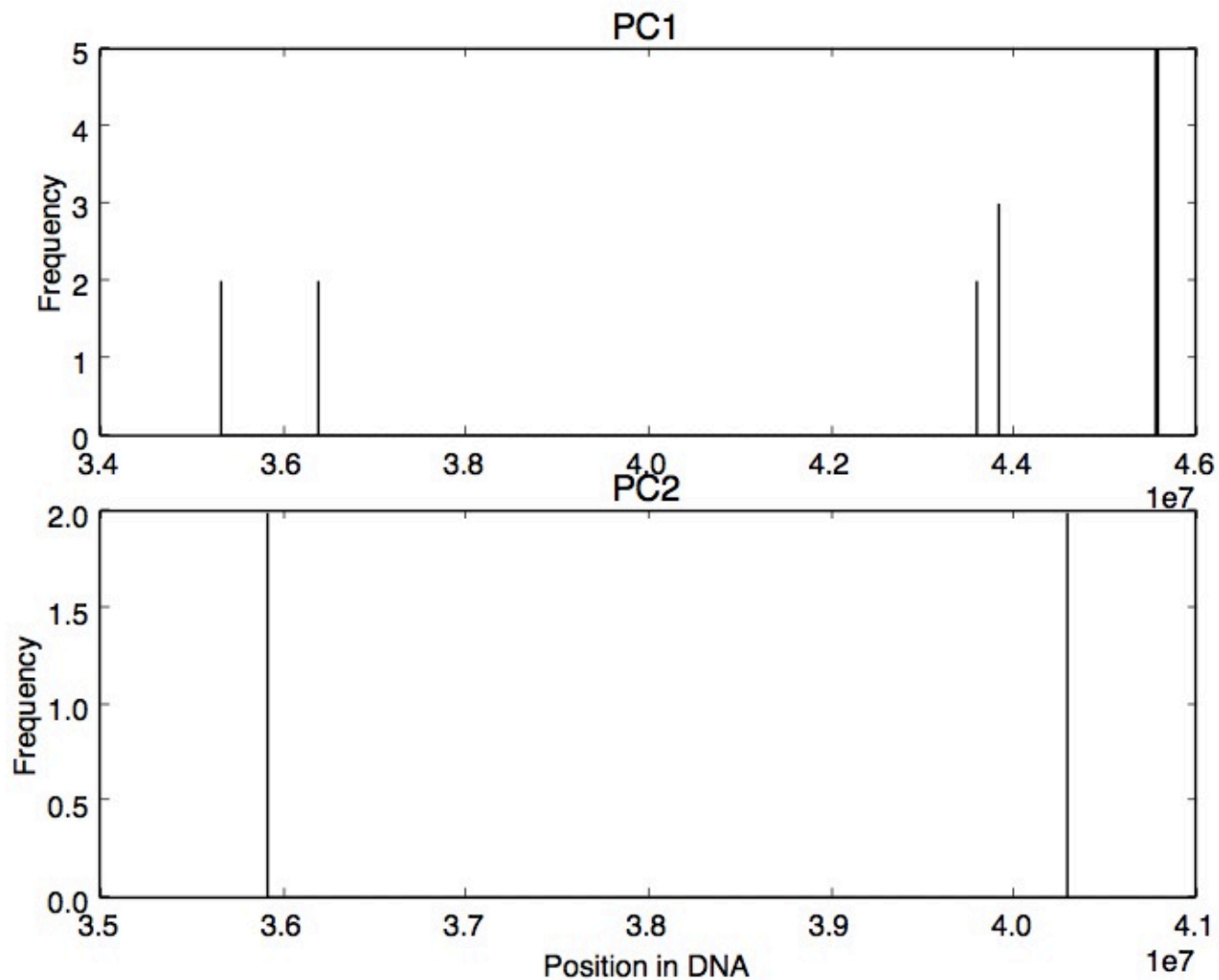
# Analysis of Variance

- Sorted components of each PC in descending order of absolute value
- Looked at patterns in the components that contributed more than 5% to the variance
- Binned regions into 5kb segments

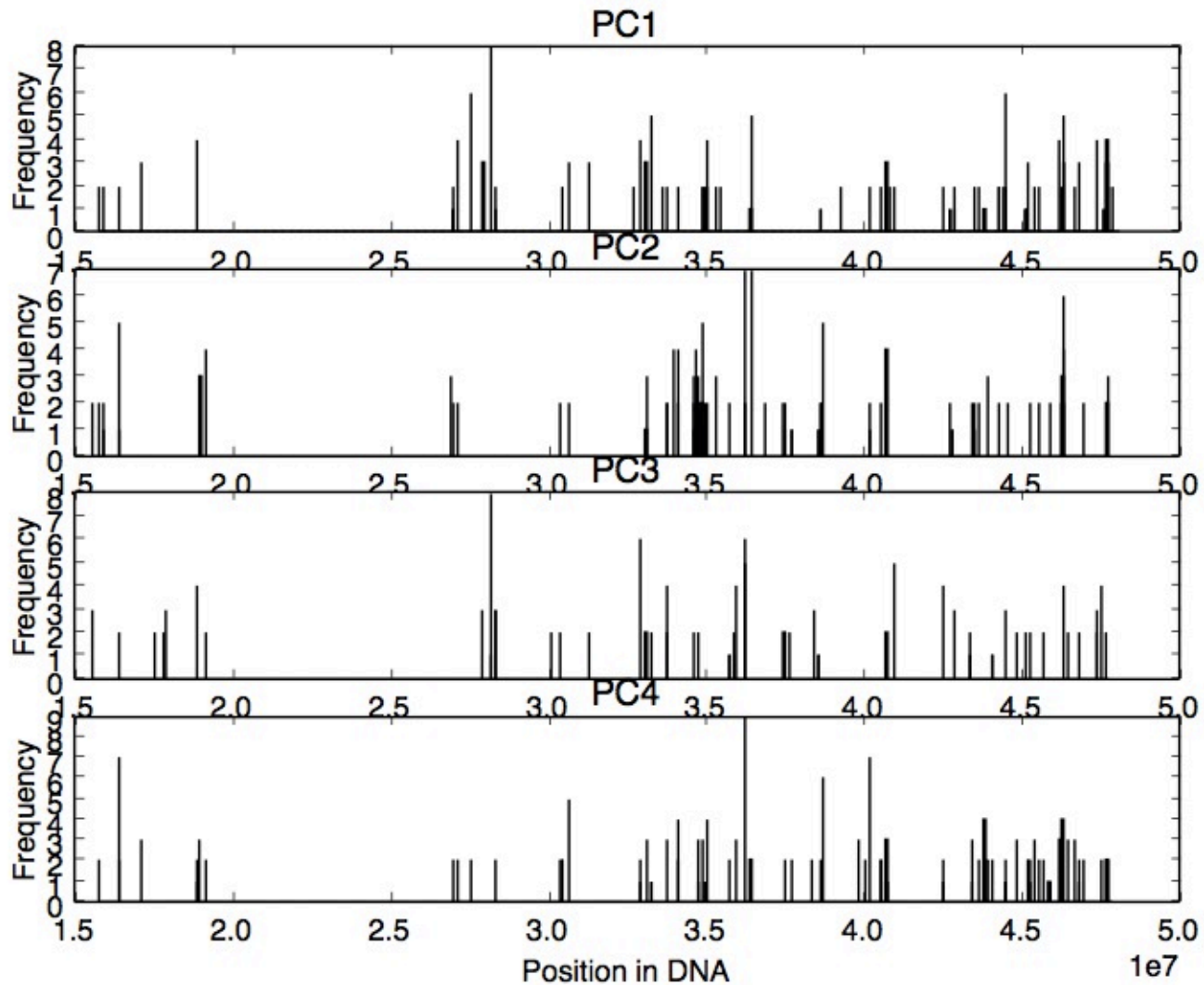
# DNase



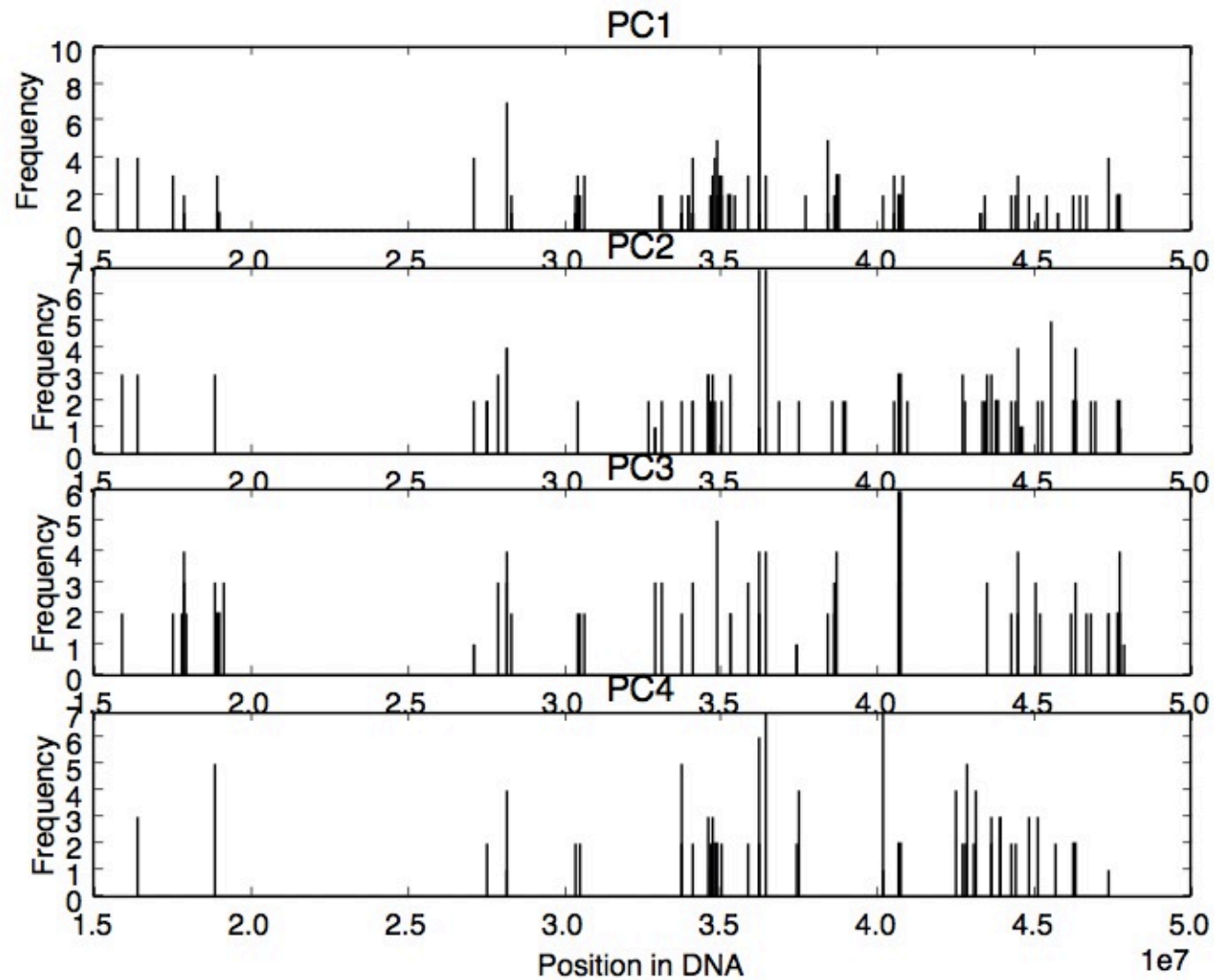
# H3K4me1



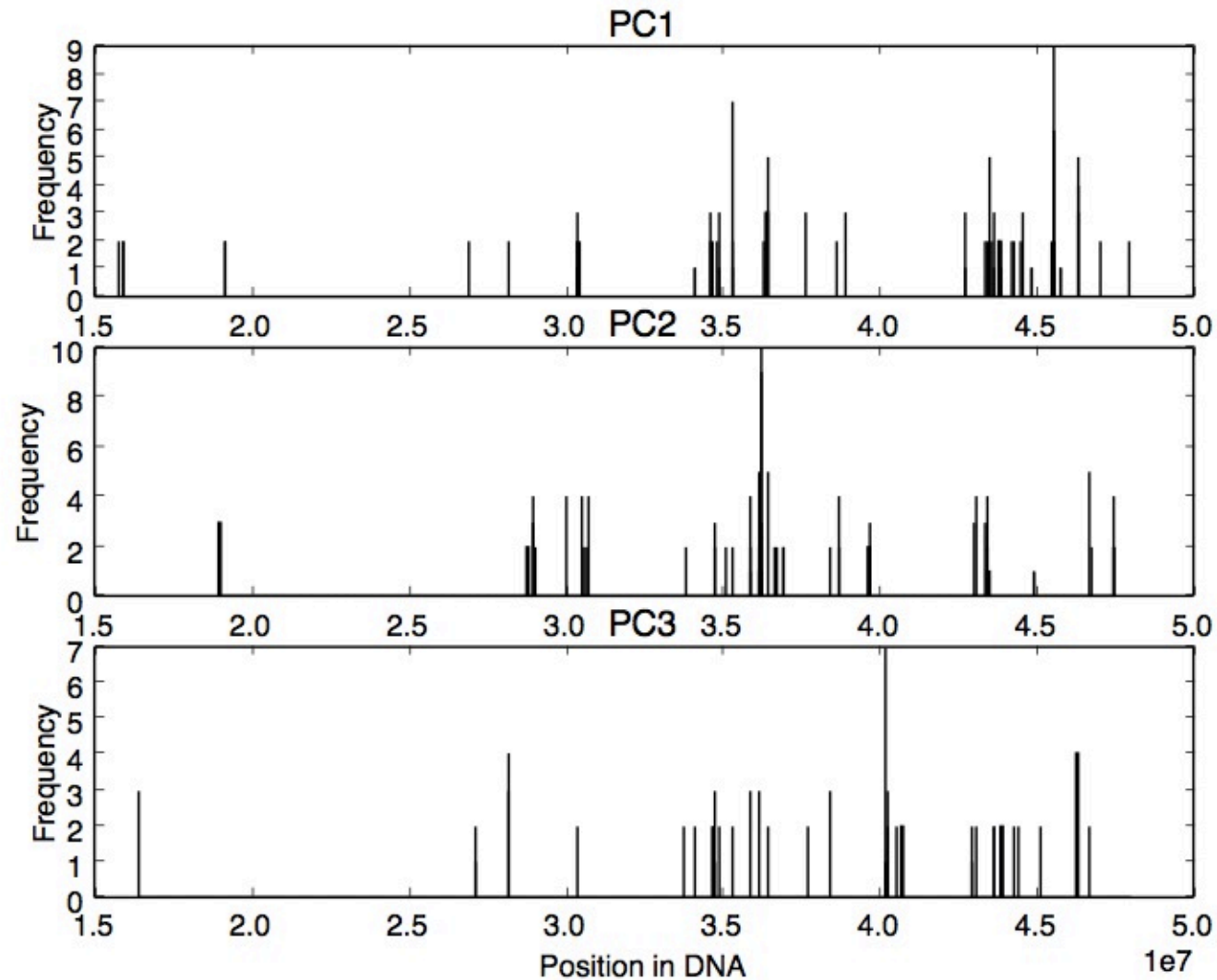
# H3K4me3



# H3K9ac

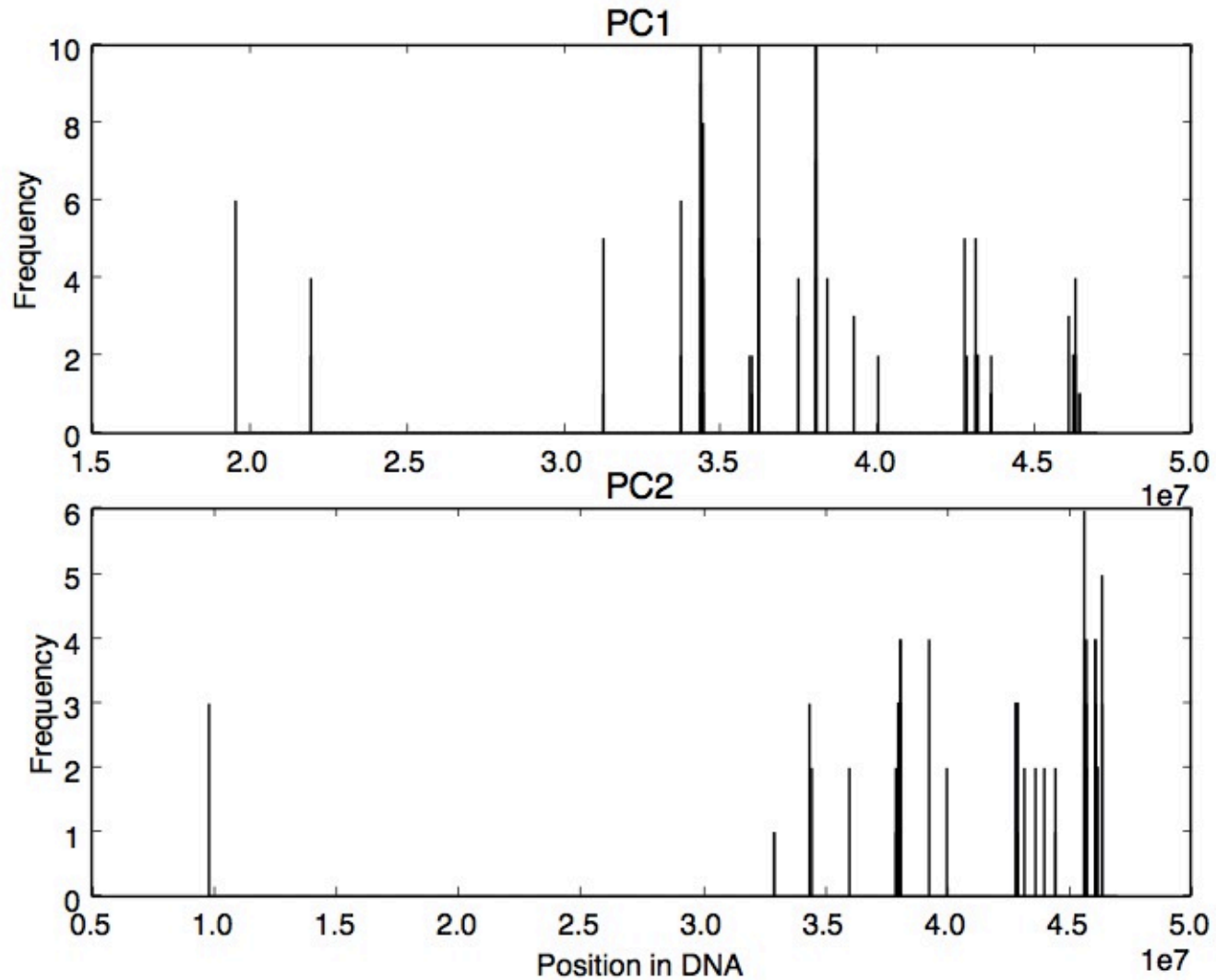


# H3K27ac

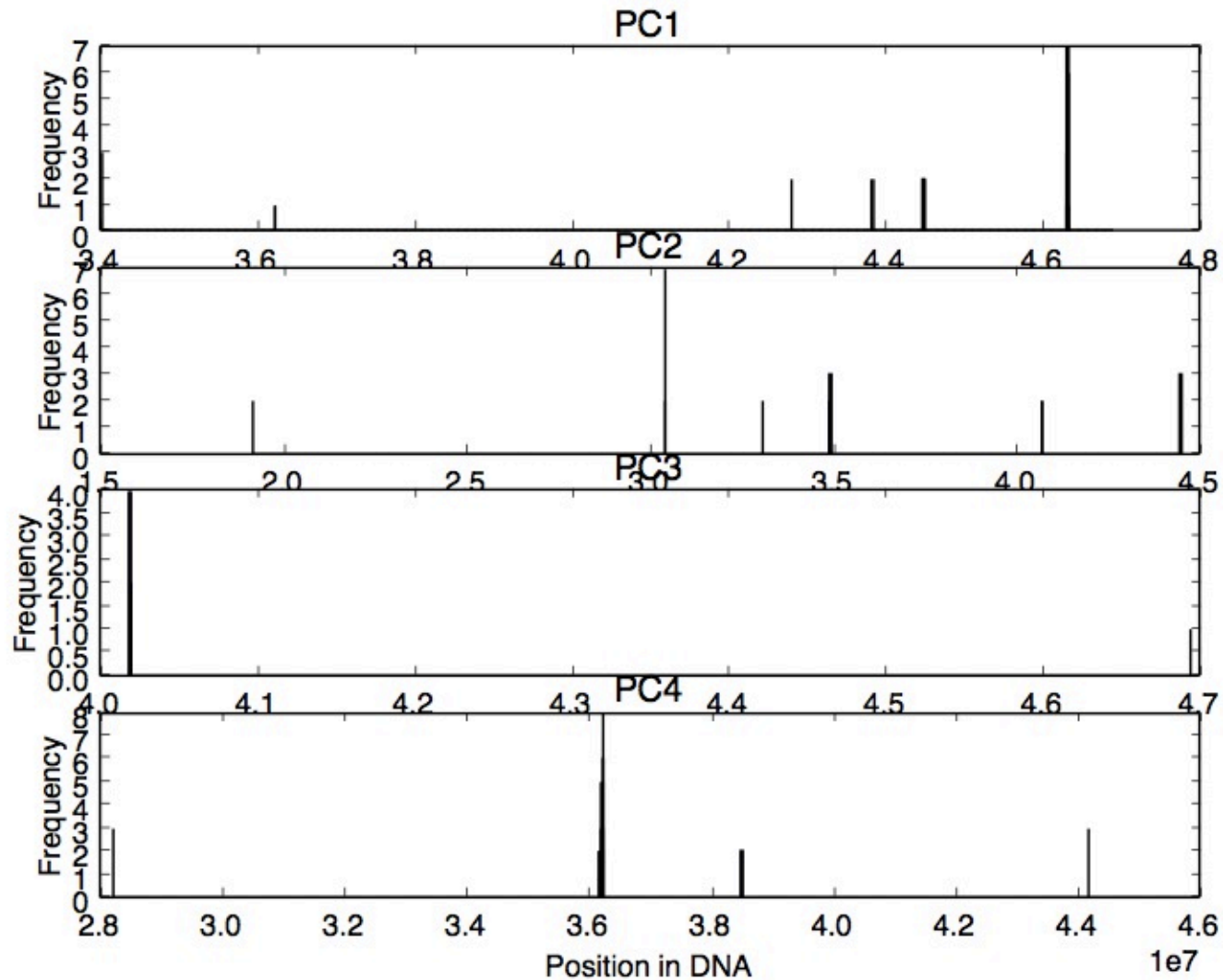




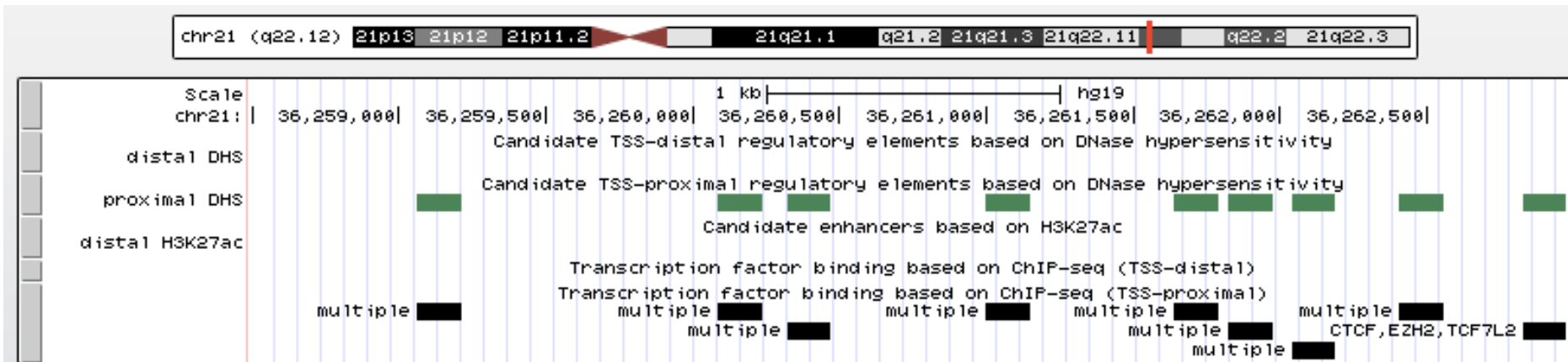
# H3K27me3



# H3K36me3



# Comparison to annotation



- Conclusions:
  - Histone modifications showed the most cell-type specificity in Blood & T-cell, HSC & B-cell and Digestive cells
  - Certain genomic locations contribute more to the variance in cell types
- Future directions:
  - Use a Gaussian mixture model for clustering
  - Use higher resolution data to predict the biological function of important locations in the genome