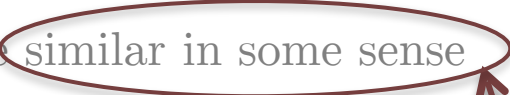


Cluster Analysis

(see also: Segmentation)

Cluster Analysis

- Unsupervised: no target variable for training
- Partition the data into groups (clusters) so that:
 - Observations within a cluster are similar in some sense
 - Observations in different clusters are different in some sense
- There is **no one correct answer**, though there are good and bad cluster solutions
- No method works best all the time



That's not very specific...

(Some) Applications of Clustering

- **Customer segmentation:** groups of customers with similar shopping or buying patterns
- Dimension reduction:
 - Cluster variables together - use one or a linear combination in place of the cluster
 - Cluster individuals together and use cluster variable as proxy for demographic or behavioral variables
- Gather stores with similar characteristics for sales forecasting
- Find topics (clusters of words) in text data
- Find communities in social networks

Methodology:

Hard vs. Fuzzy Clustering

Hard: objects can belong to only one cluster

- k-means (PROC FASTCLUS)
- DBSCAN
- Hierarchical (PROC CLUSTER)

Fuzzy: objects can belong to more than one cluster (usually with some probability)

- Fuzzy C-means (FCM)
- Gaussian Mixture Models / Expectation-Maximization (EM)

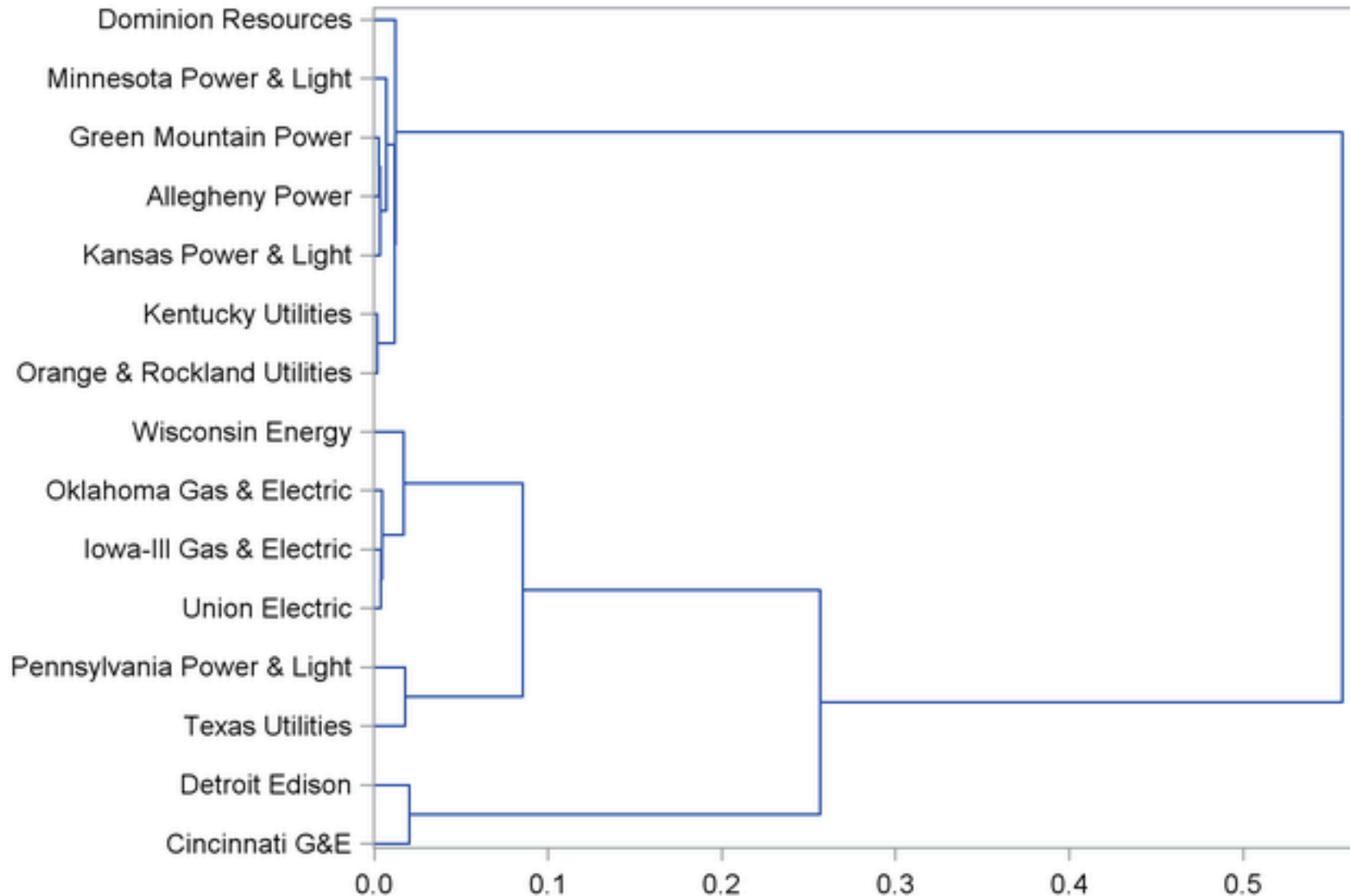
Methodology: Hierarchical vs. Flat

Hierarchical: clusters form a tree so you can visually see which clusters are most similar to each other.

- *Agglomerative*: points start out as individual clusters, and they are combined until everything is in one cluster. (ex: Linkages - coming next class)
- *Divisive*: All points start in same cluster and at each step a cluster is divided into two clusters. (ex: PDDP)

Flat: Clusters are created according to some other process, usually iteratively updating cluster assignments

Dendrogram from Hierarchical Clustering



K-Means Clustering

...

The most popular method (by far)

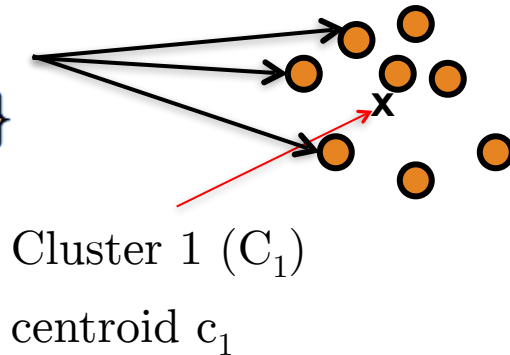
k-Means Clustering

(PROC FASTCLUS in SAS)

- The most popular clustering algorithm

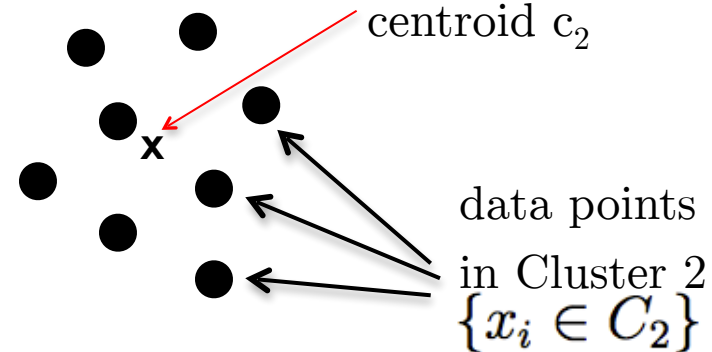
data points

in Cluster 1
 $\{x_i \in C_1\}$



Cluster 2 (C_2)

centroid c_2



- Tries to minimize the sum of squared distances from each point to its cluster centroid. (Global objective function)

$$\sum_{C_k} \sum_{x_i \in C_k} \|x_i - c_k\|^2$$

k-Means Algorithm

- Start with k “seed points”
 - Randomly initialized (most software)
 - Determined ‘methodically’ (SAS PROC FASTCLUS)
- Assign each data point to the closest seed point.
- The seed point then represents a cluster of data
- Reset seed points to be the centroids of the cluster
- Repeat steps 2-4 updating the cluster centroids until they do not change.

k-Means Interactive Demo

<https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>

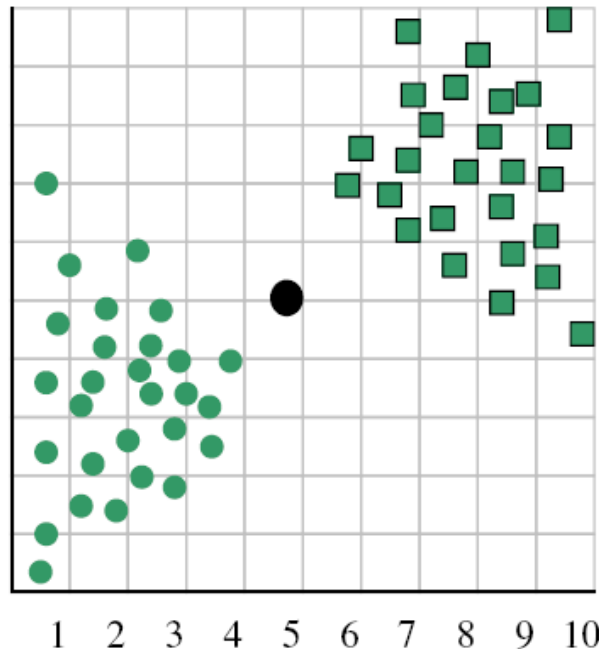
(You may have to add the site to your exceptions list on the Java Control Panel to view.)

Choice of Distance Metric

- Most distances like Euclidean, Manhattan, or Max will provide *similar* answers.
- Use cosine distance (really $1-\cos$ because cosine measures similarity) for text data. This is called **spherical k-means**.
- Using Mahalanobis distance is essentially the Expectation-Maximization (EM method) for Gaussian Mixtures.

Determining Number of Clusters (SSE)

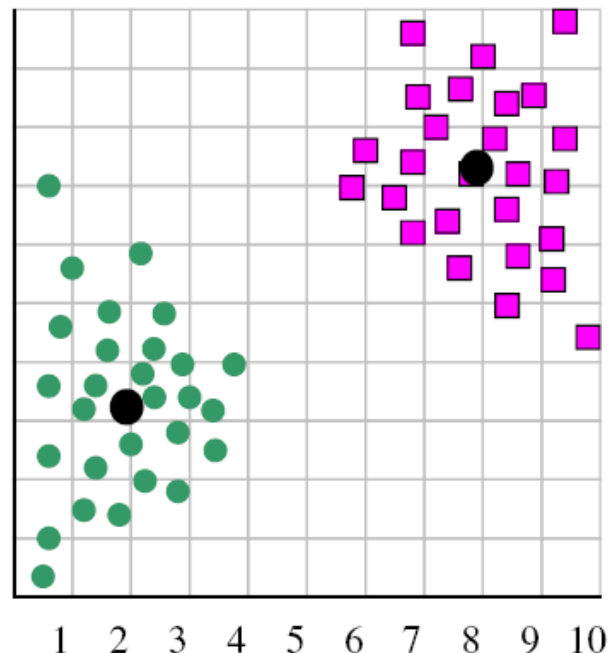
- Try the algorithm with $k=1,2,3,\dots$
- Examine the objective function values
- Look for a place where the marginal benefit to objective function for adding a cluster becomes small



$k=1$ objective
function (SSE) is 902

Determining Number of Clusters (SSE)

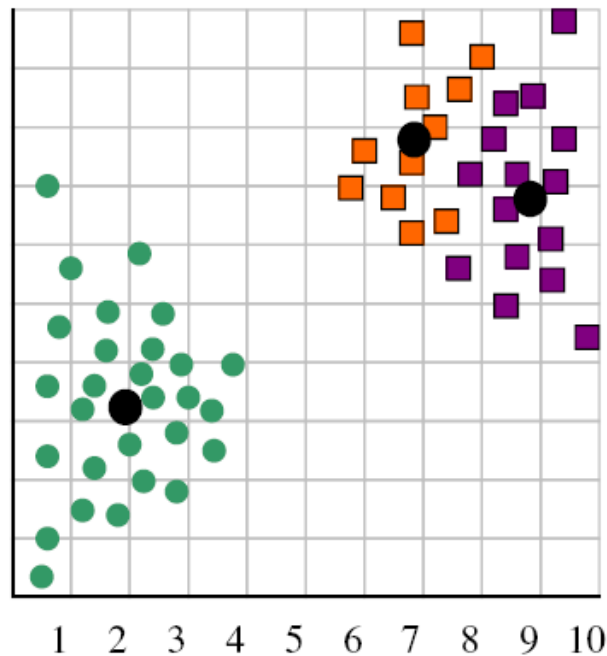
- Try the algorithm with $k=1,2,3,\dots$
- Examine the objective function values
- Look for a place where the marginal benefit to objective function for adding a cluster becomes small



$k=2$ objective
function (SSE) is 213

Determining Number of Clusters (SSE)

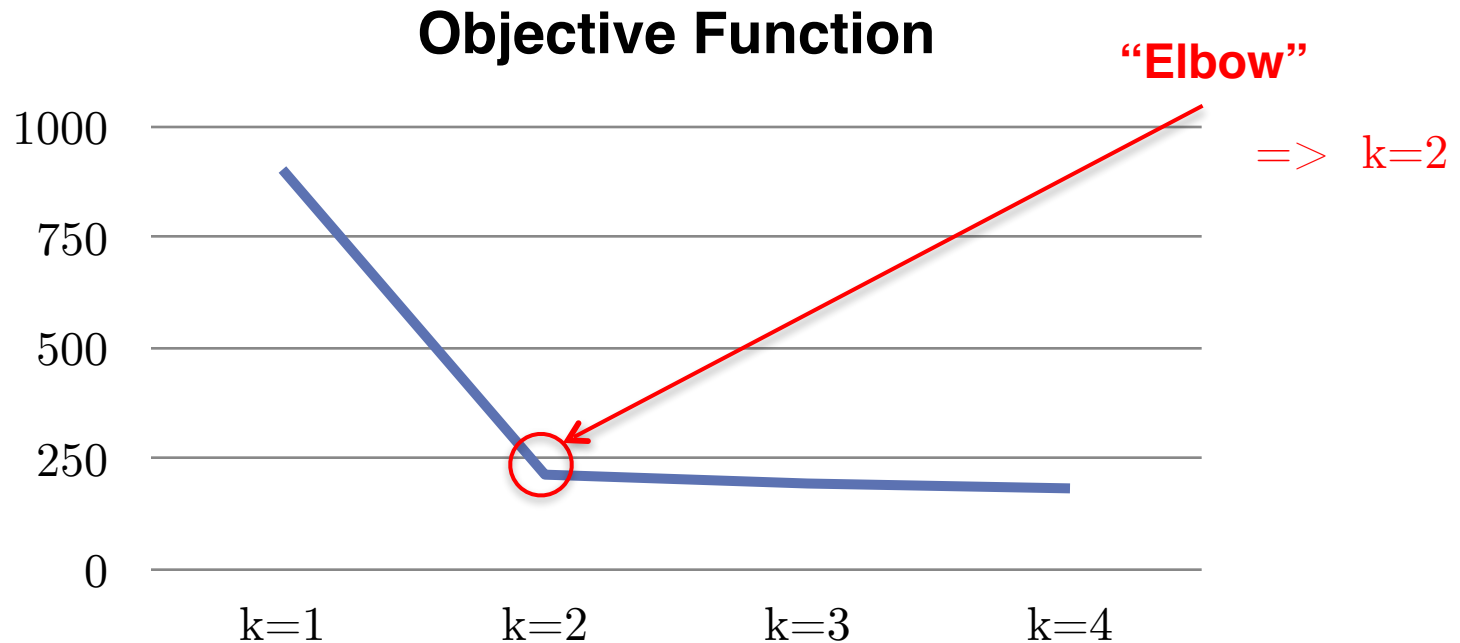
- Try the algorithm with $k=1,2,3,\dots$
- Examine the objective function values
- Look for a place where the marginal benefit to objective function for adding a cluster becomes small



$k=3$ objective
function (SSE) is 193

Determining Number of Clusters (SSE)

- Try the algorithm with $k=1,2,3,\dots$
- Examine the objective function values
- Look for a place where the marginal benefit to objective function for adding a cluster becomes small



k-Means Summary

Disadvantages

- Dependent on initialization (initial seeds)
- Can be sensitive to outliers
 - If problem, should consider k-medoids (uses median not mean)
- Have to input the number of clusters
- Difficulty detecting non-spheroidal (globular) clusters

Advantages

- Modest time/storage requirements.
- Shown you can terminate method after small number of iterations with good results.
- Good for wide variety of data types

R Demo

• • •

Adults dataset

Data Contents

Age

Workclass - Federal-gov, local-gov, never-worked, self-emp-inc, self-emp-not-inc, private, without-pay

Education - Level of education

EducationNumeric - Numeric version of education

MaritalStatus

Occupation - Category of occupation types

Relationship - Relationship to head of household (now called "householder") which is defined as the person (or one of the persons) in whose name the house/apt is owned or rented/maintained.

Race

Sex

capitalGain

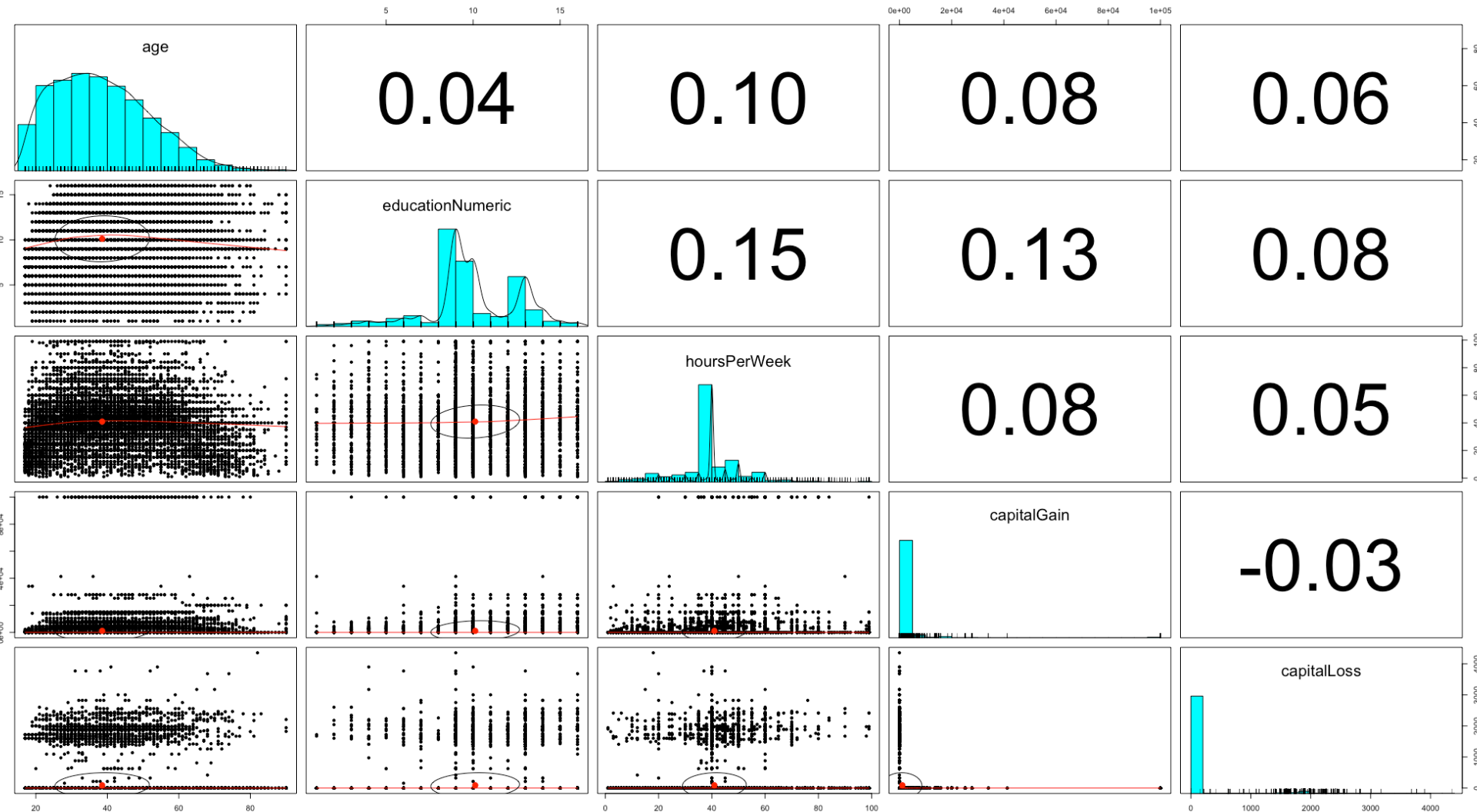
capitalLoss

hoursPerWeek - Average hours worked per week over past year

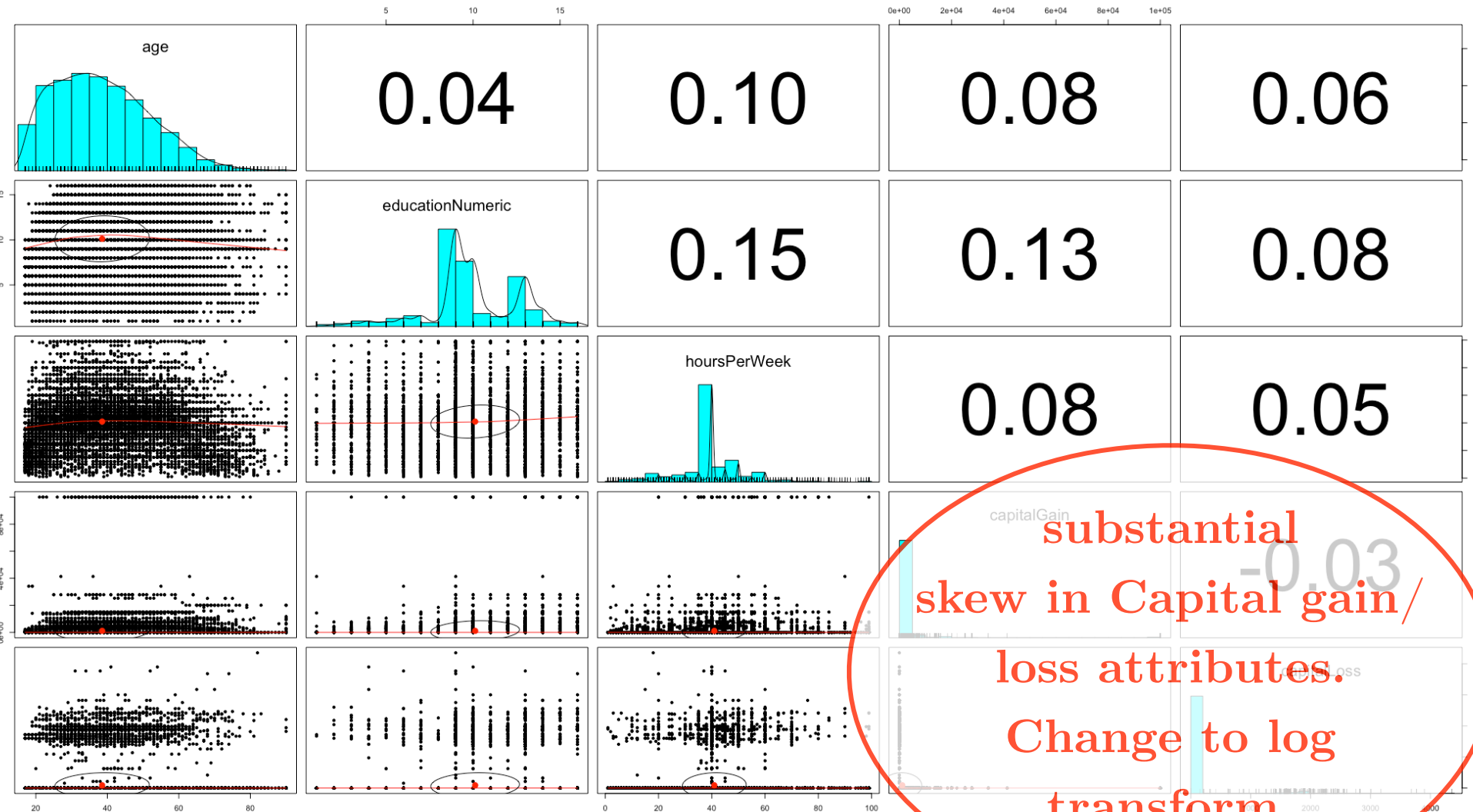
nativeCountry

incomeLevel - Binary value, either $\leq 50K$ or $> 50k$.

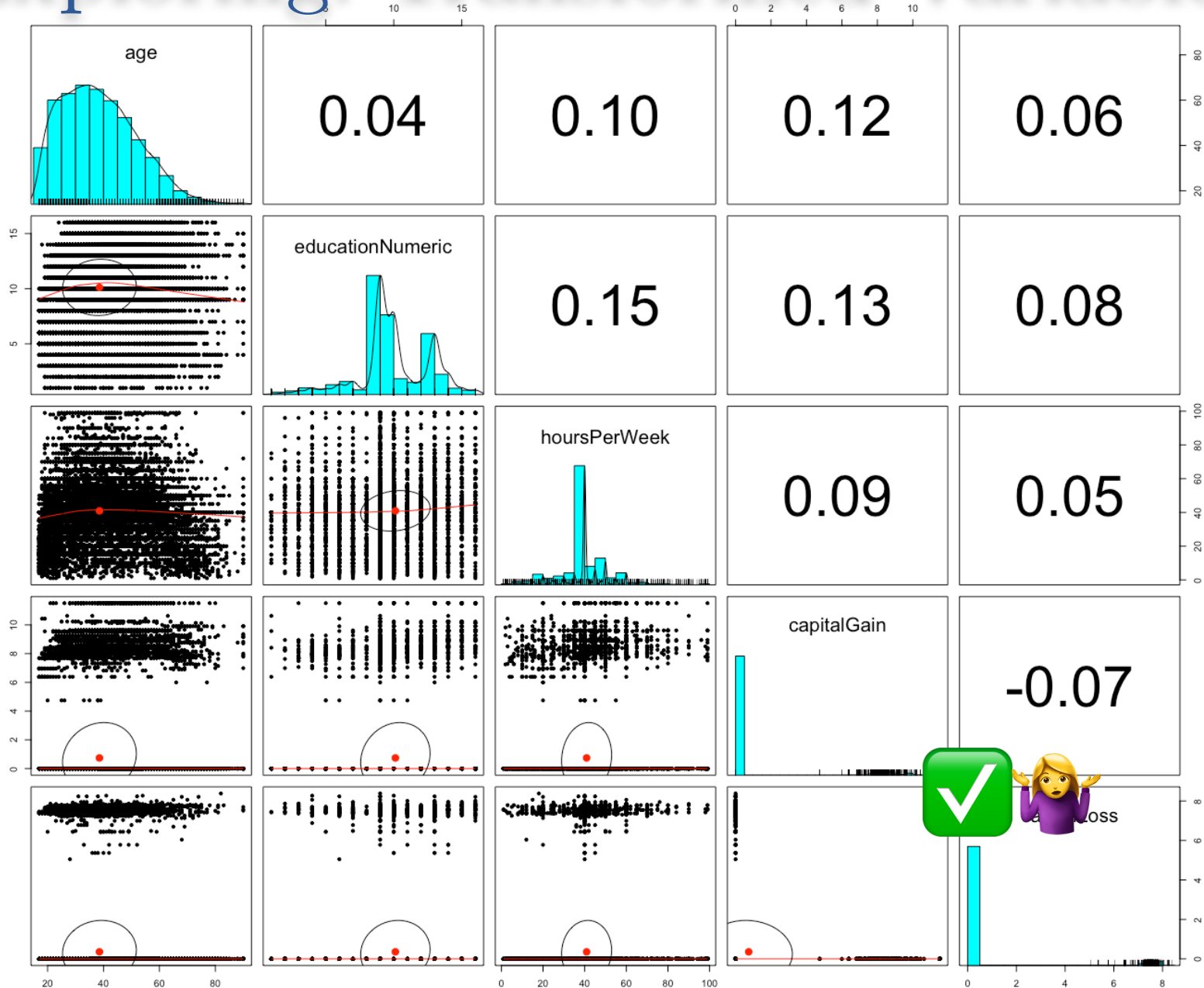
Exploring: Numeric Variables



Exploring: Numeric Variables

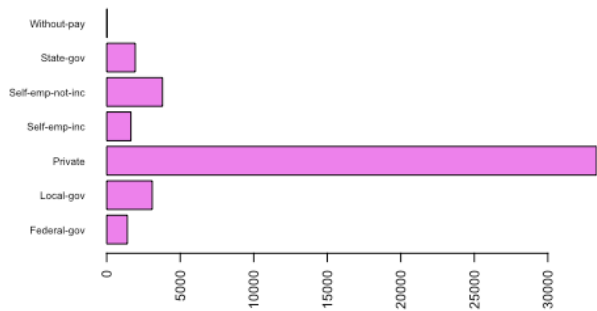


Exploring: Transformed Variables

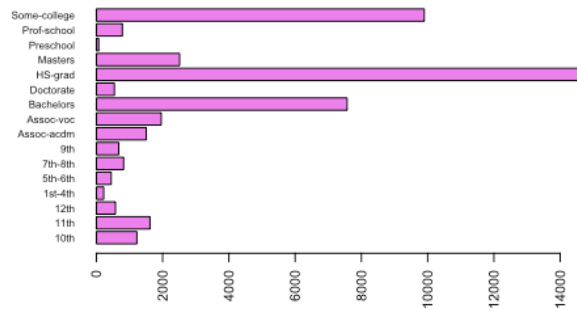


Exploring: Categorical Variables

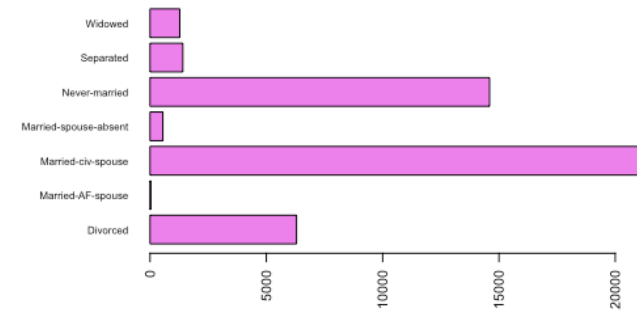
workclass



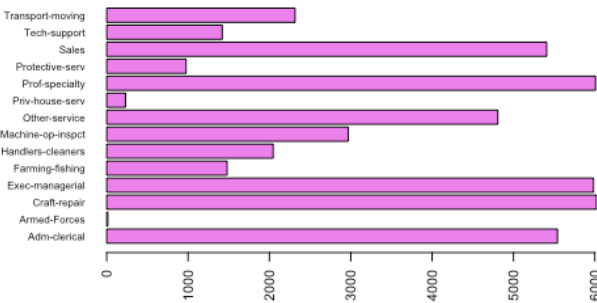
education



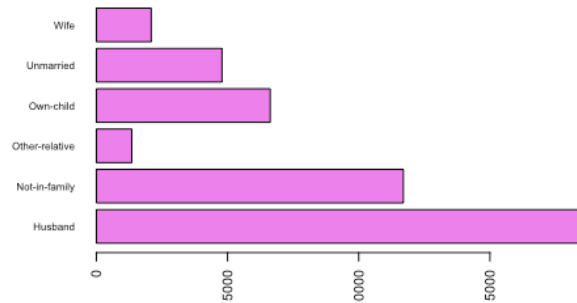
maritalStatus



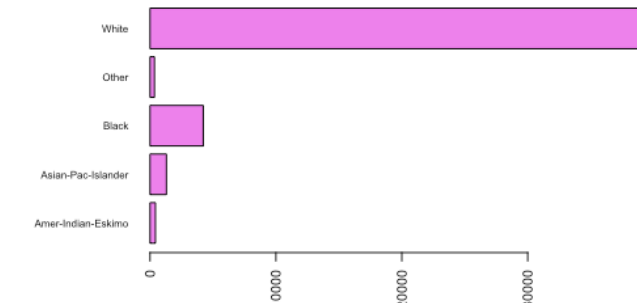
occupation



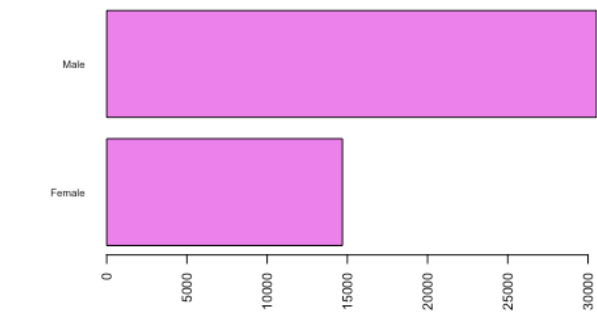
relationship



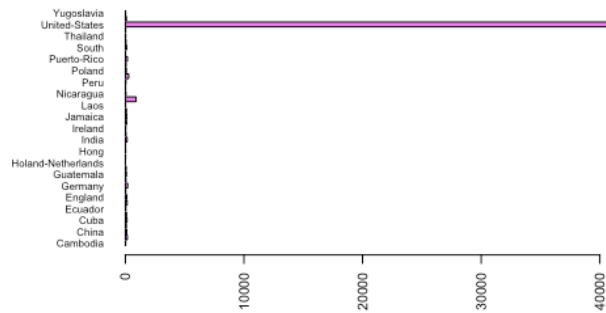
race



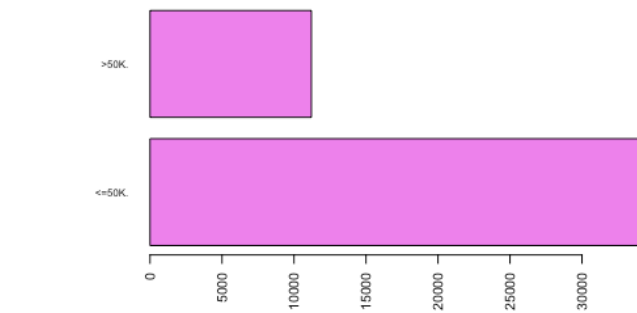
sex



nativeCountry

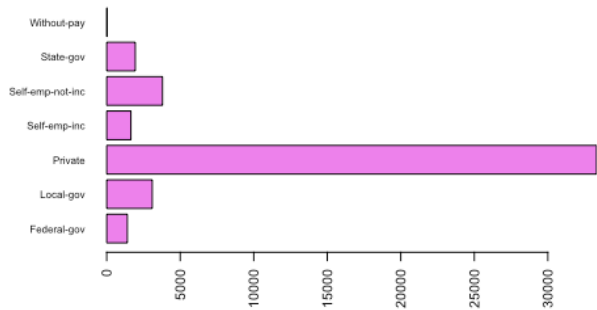


incomeLevel

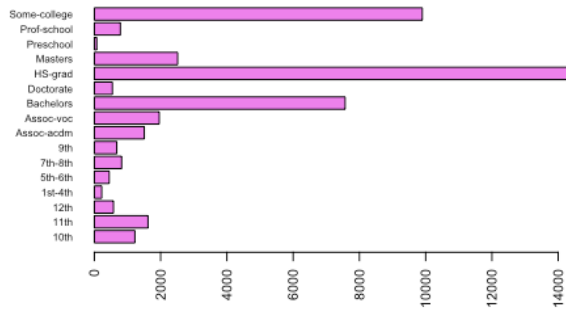


Exploring: Categorical Variables

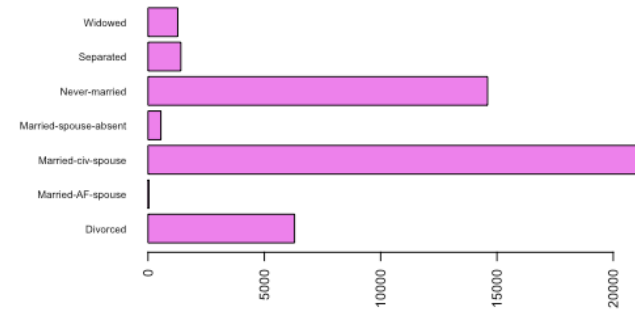
workclass



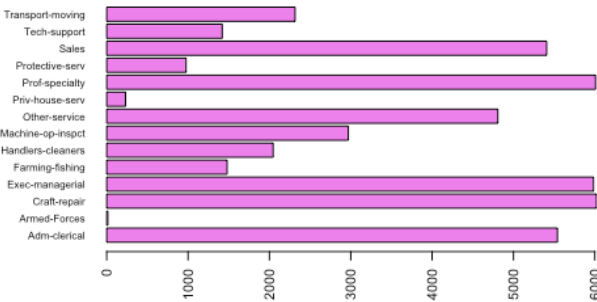
education



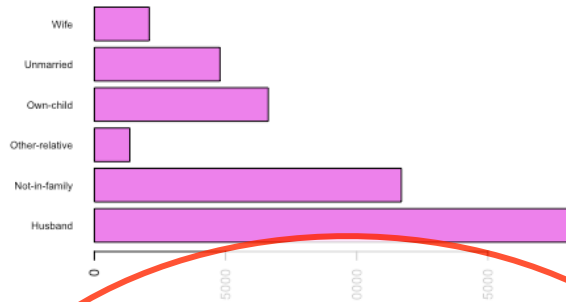
maritalStatus



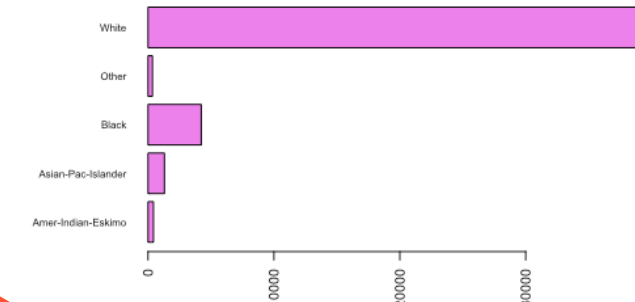
occupation



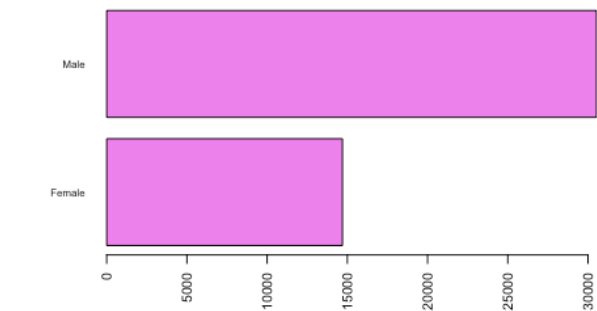
relationship



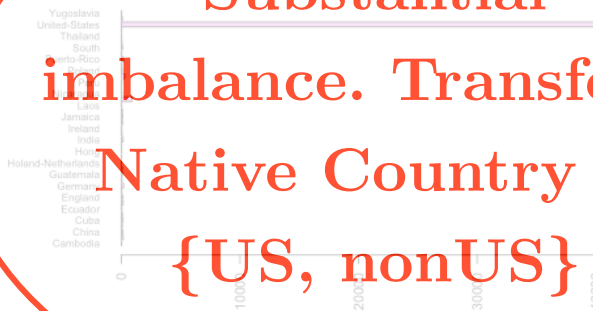
race



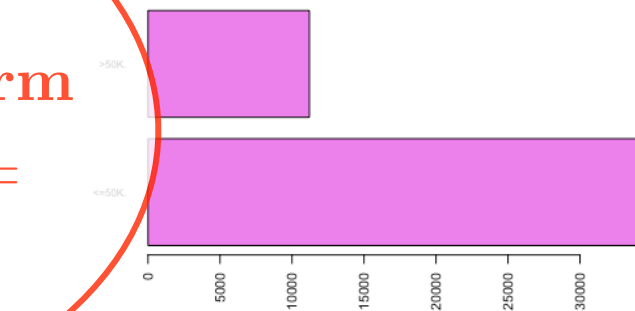
sex



nativeCountry



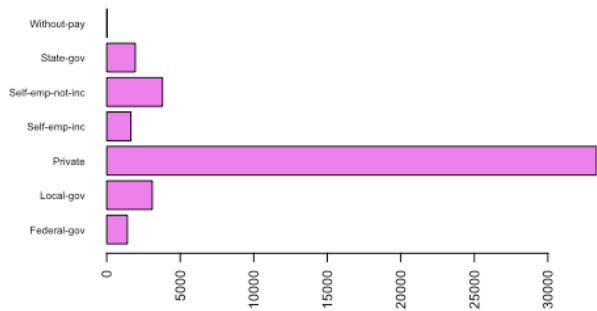
incomeLevel



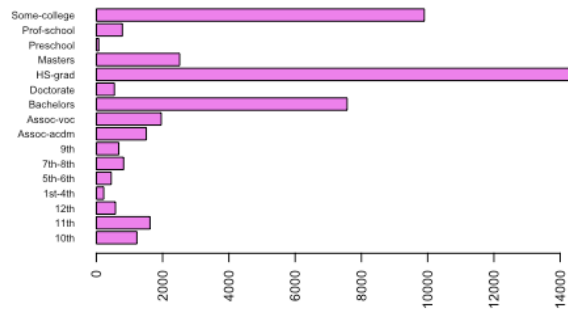
Substantial
imbalance. Transform
Native Country =
{US, nonUS}

Exploring: Categorical Variables

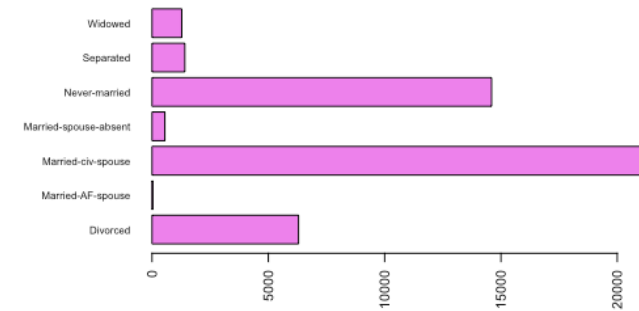
workclass



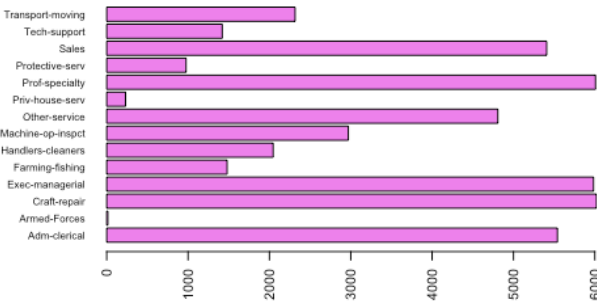
education



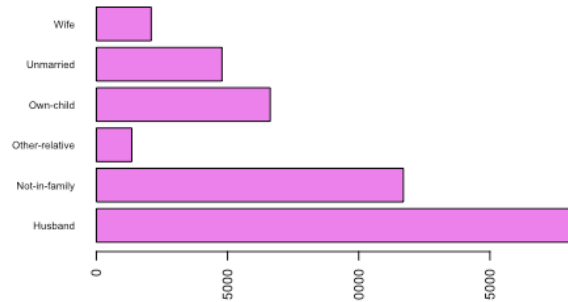
maritalStatus



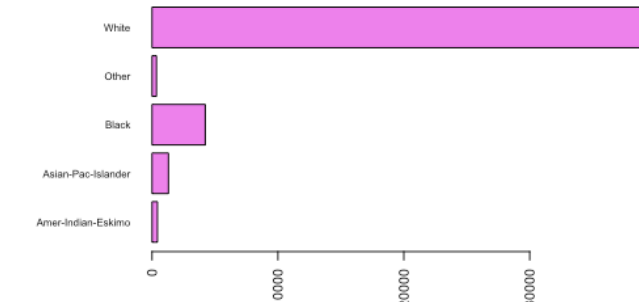
occupation



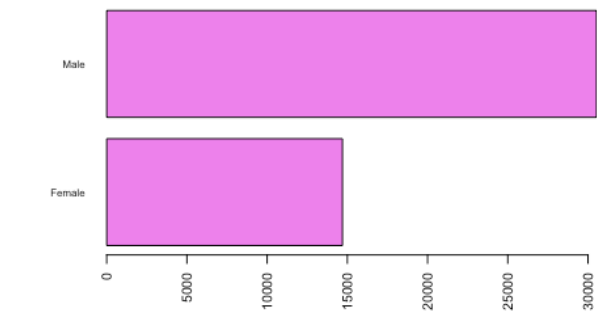
relationship



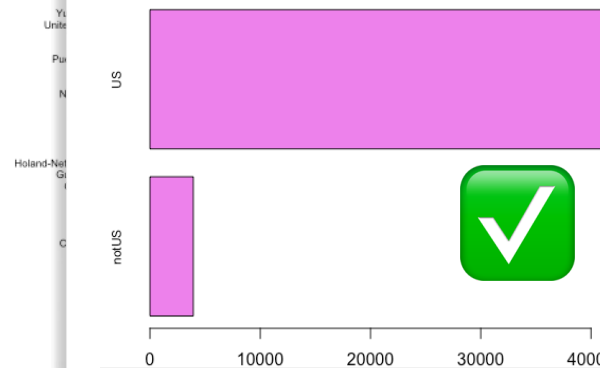
race



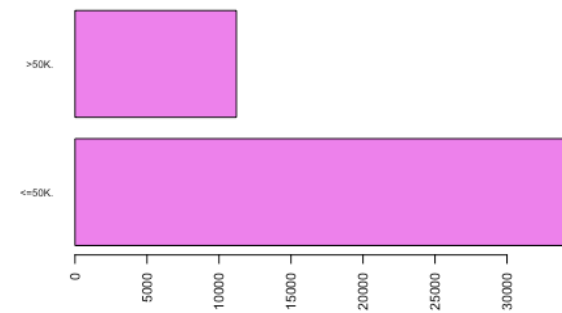
sex



nativeCountry



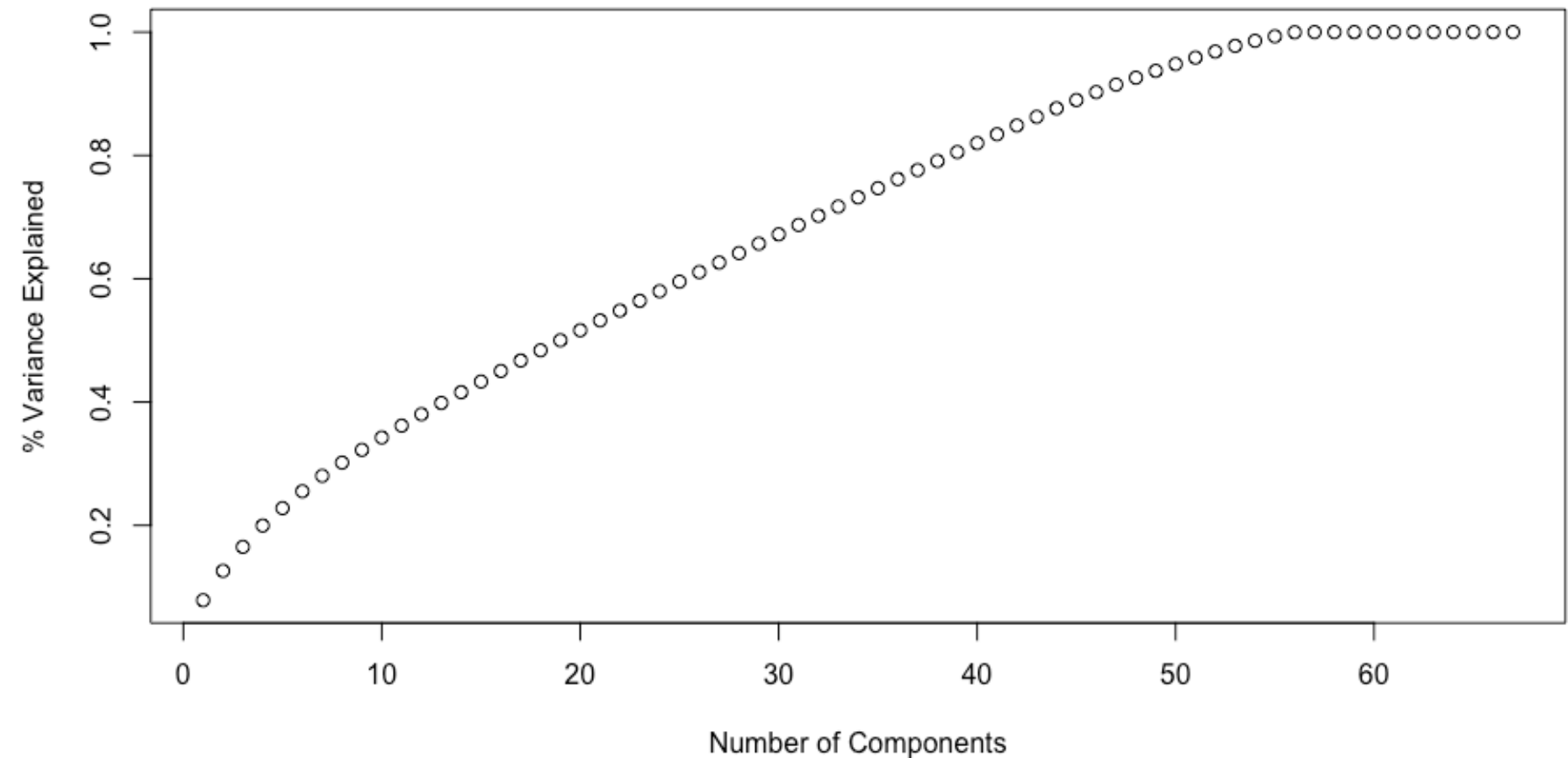
incomeLevel



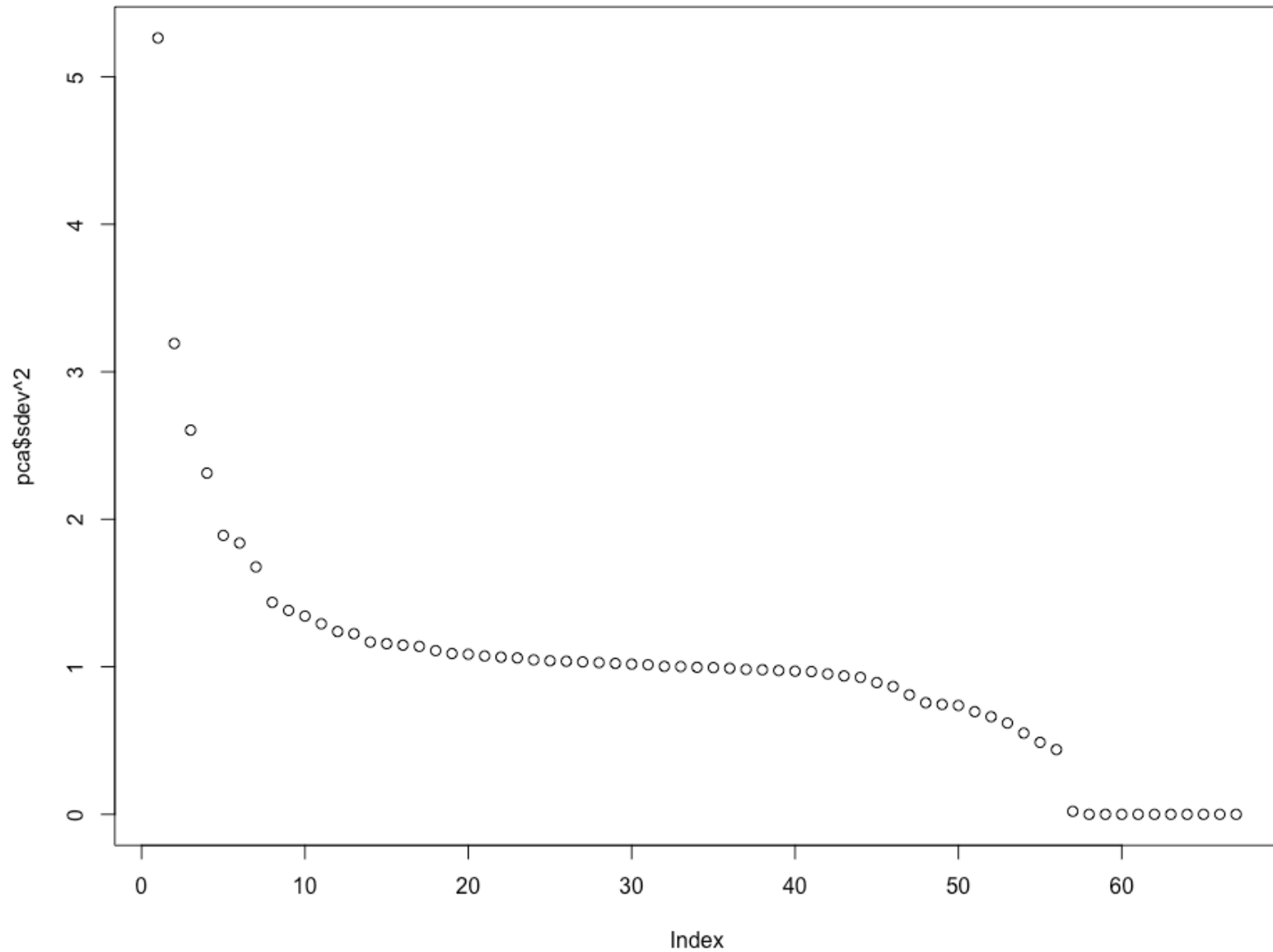
Explore Multivariate Data with PCA

```
#####  
##### Dummy Code the factors  
#####  
adult.x = model.matrix(~. ,contrasts.arg = lapply(adult[,sapply(adult,is.factor) ],  
                                                    contrasts, contrasts=FALSE),  
                      data = adult)  
  
dim(adult.x)  
# Get rid of constant columns  
adult.x = adult.x[,apply(adult.x, 2, sd)>0 ]  
dim(adult.x)  
#####  
##### Go to Principal Components  
#####  
pca = prcomp(adult.x, scale = T)  
#Screeplot  
par(mfrow=c(1,1))  
plot(pca$sdev^2)  
#Cumulative % Variance Explained  
plot(cumsum(pca$sdev^2)/sum(pca$sdev),ylab = '% Variance Explained', xlab='Number of Components')  
#####  
##### Explore Principal Components  
#####  
samplePoints = sample(1:45222, 8000, replace=F)  
plot(pca$x[samplePoints,1:2])  
par(mfrow=c(3,3),mar=c(4,4,1,1))  
plot(pca$x[samplePoints,1:3])  
plot(pca$x[samplePoints,1:4])  
plot(pca$x[samplePoints,1:5])  
plot(pca$x[samplePoints,2:3])  
plot(pca$x[samplePoints,2:4])  
plot(pca$x[samplePoints,2:5])  
plot(pca$x[samplePoints,3:4])
```

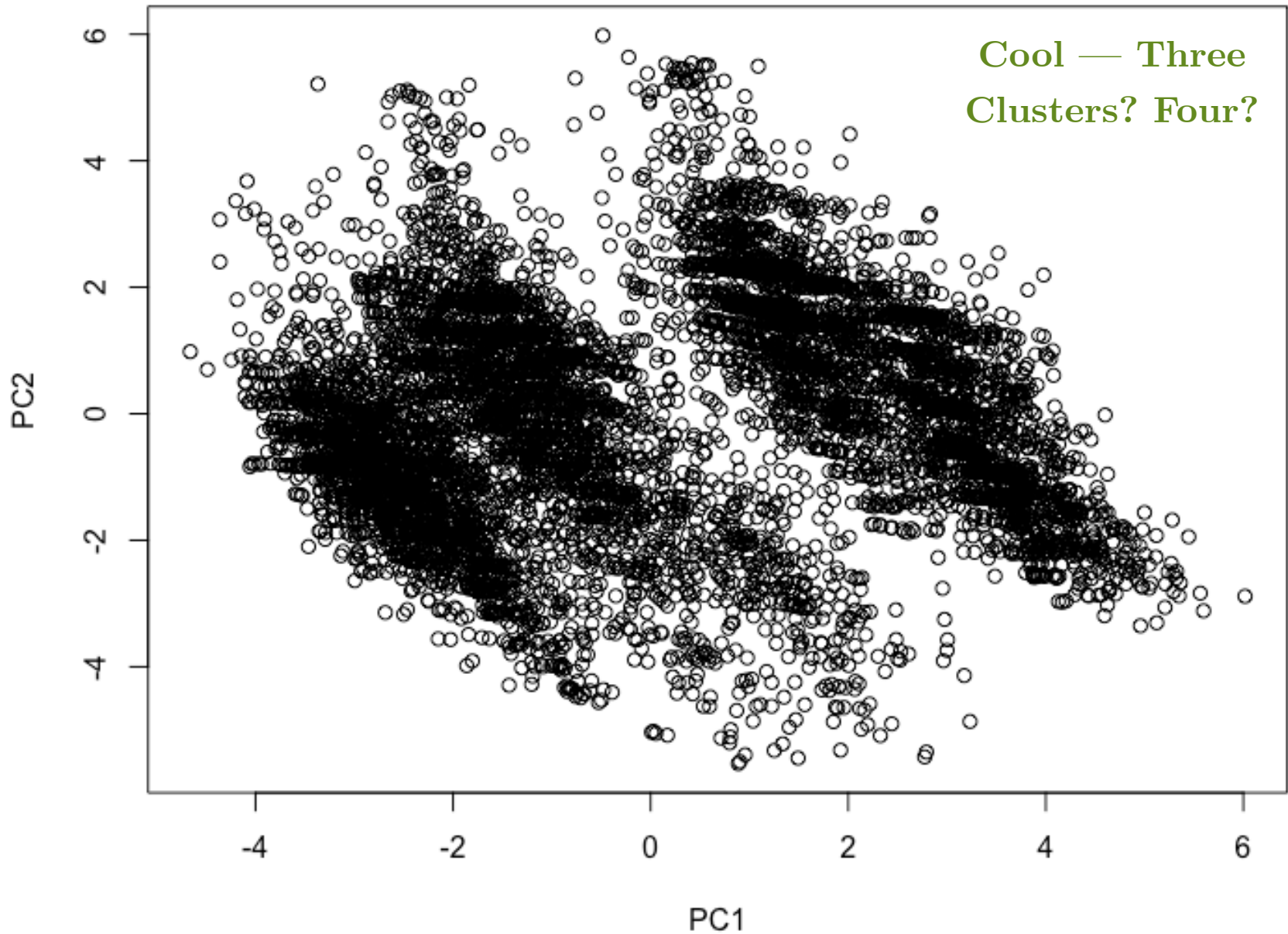
% Variance Explained



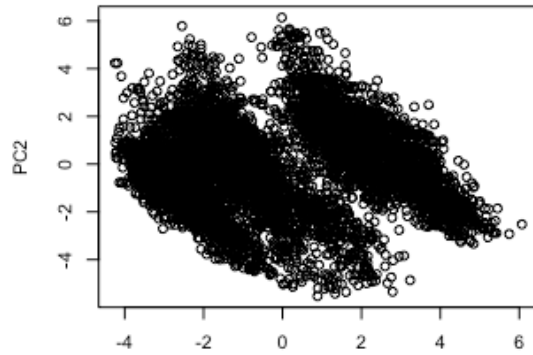
Screeplot



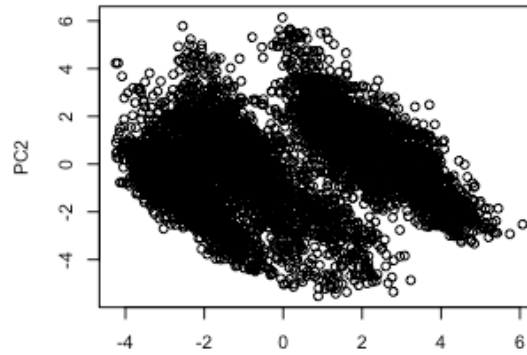
2D Projection



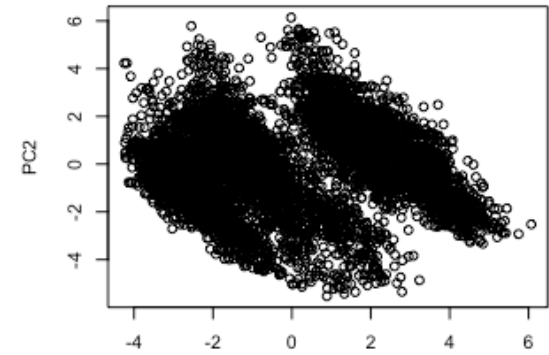
More 2D Projections



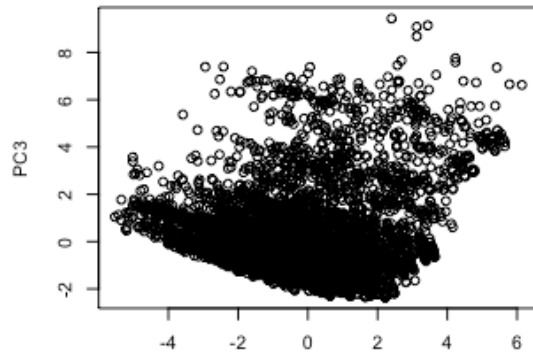
PC1



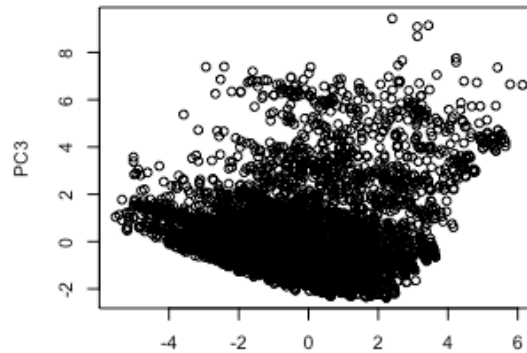
PC1



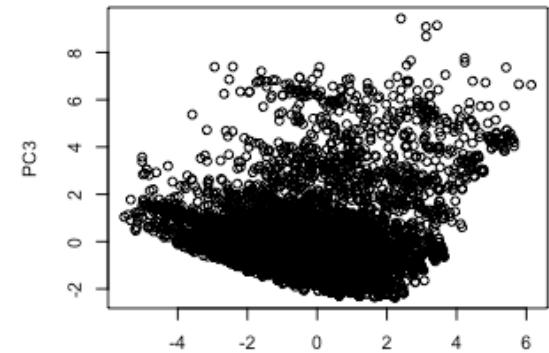
PC1



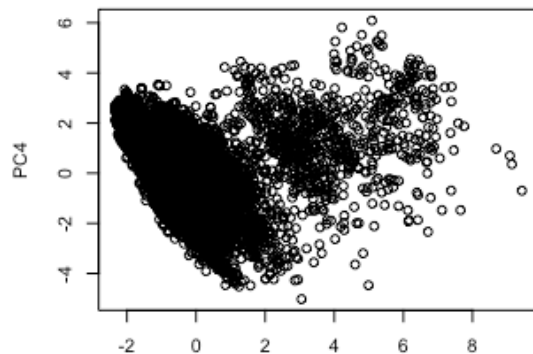
PC2



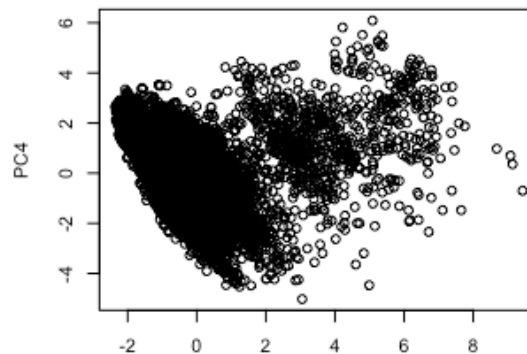
PC2



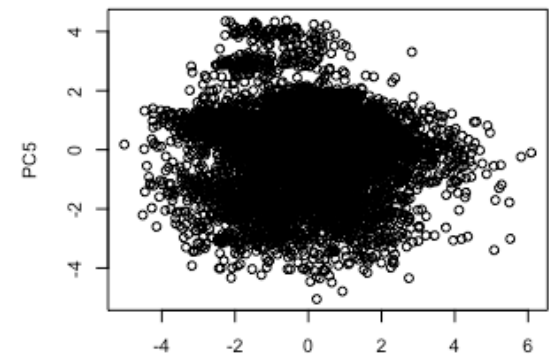
PC2



PC3



PC3



PC4

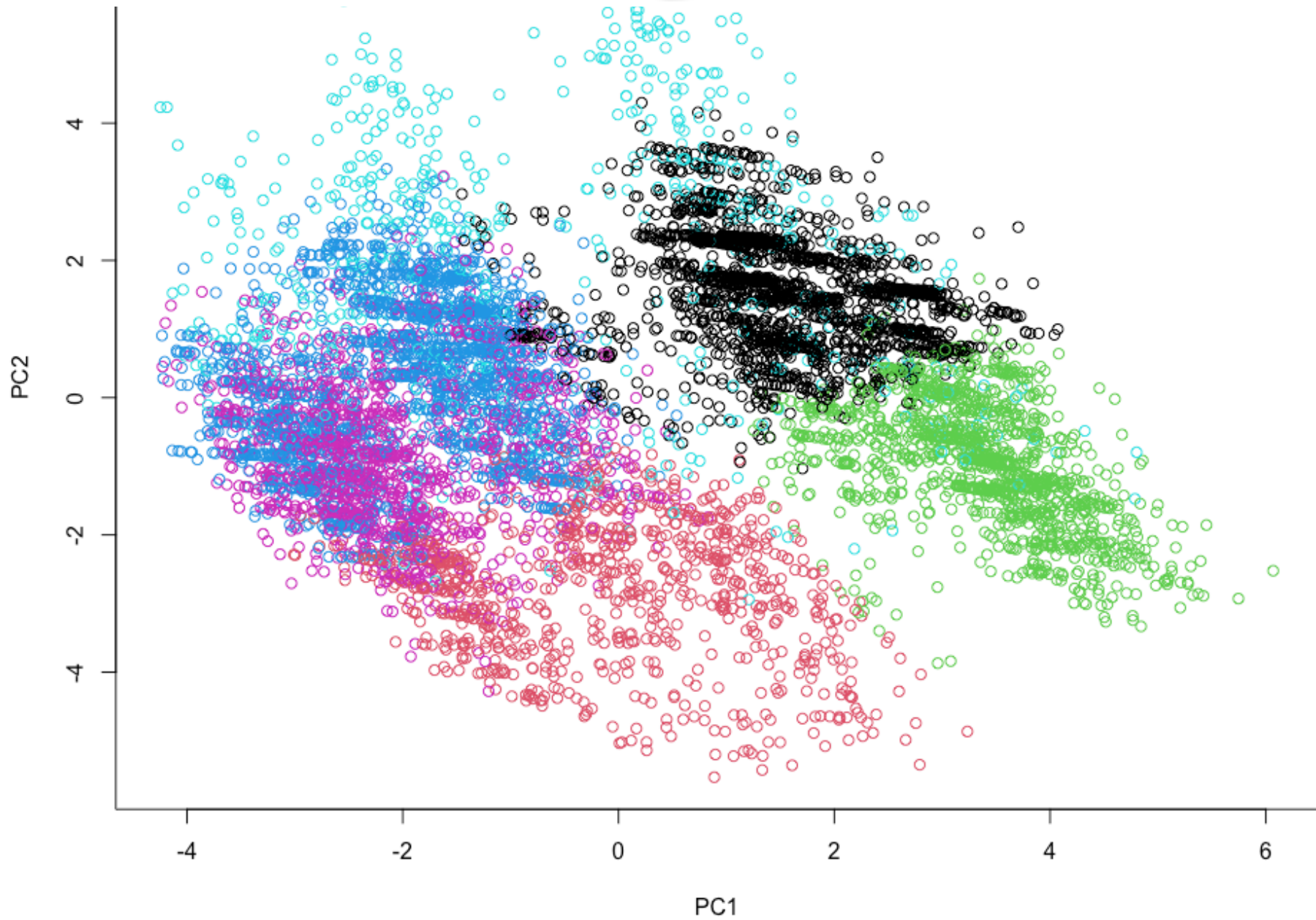
Now Cluster the PCA Scores and Visualize the Clusters in 2D

```
#####  
##### Cluster Principal Components  
#####  
k6v1 = kmeans(pca$x[,1:4],6)
```

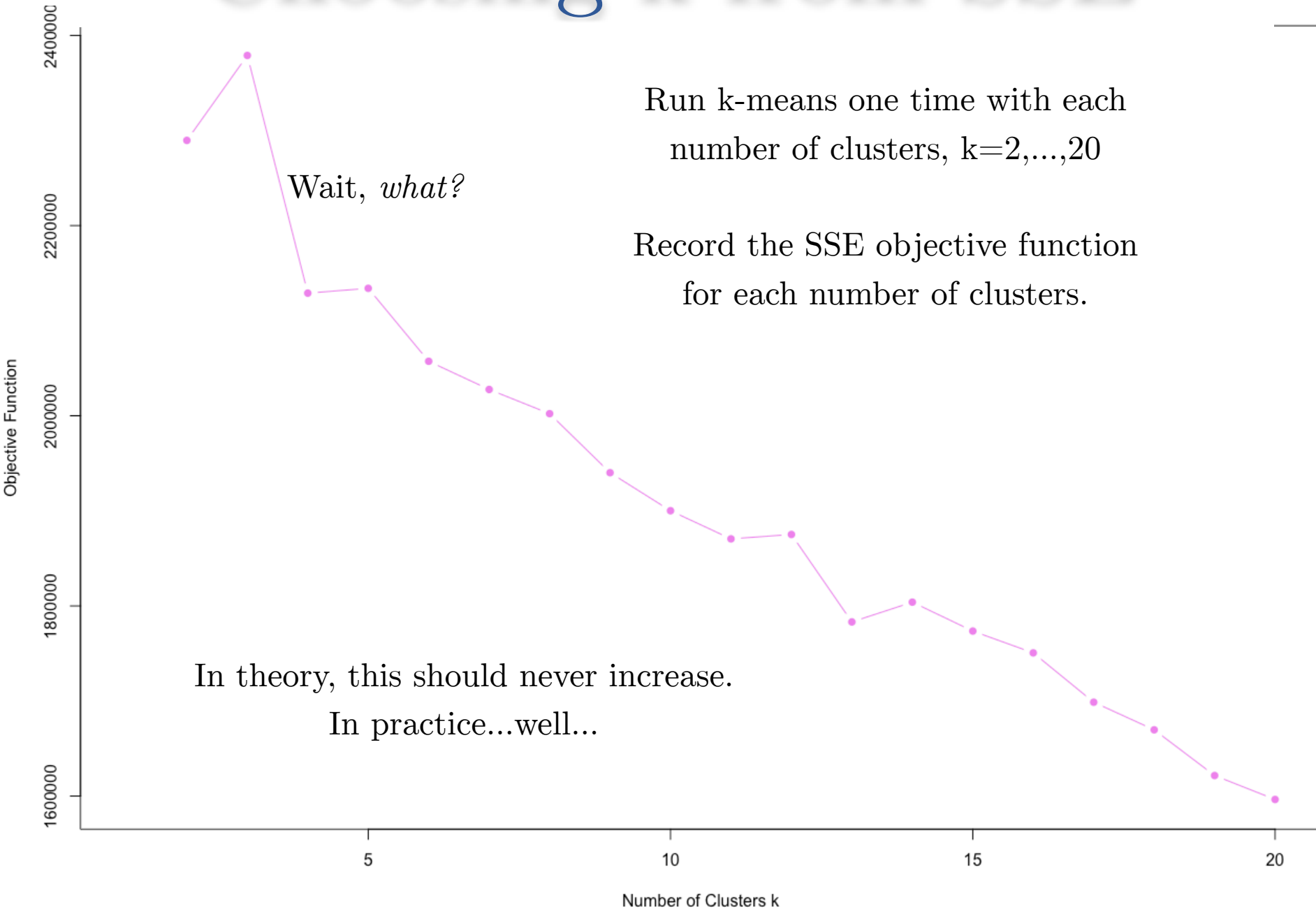
```
#####  
##### Visualize Clusters via PCA  
#####
```

```
adult$k6v1 = k6v1$cluster  
plot(pca$x[samplePoints,1:2], col=adult$k6v1[samplePoints] )
```

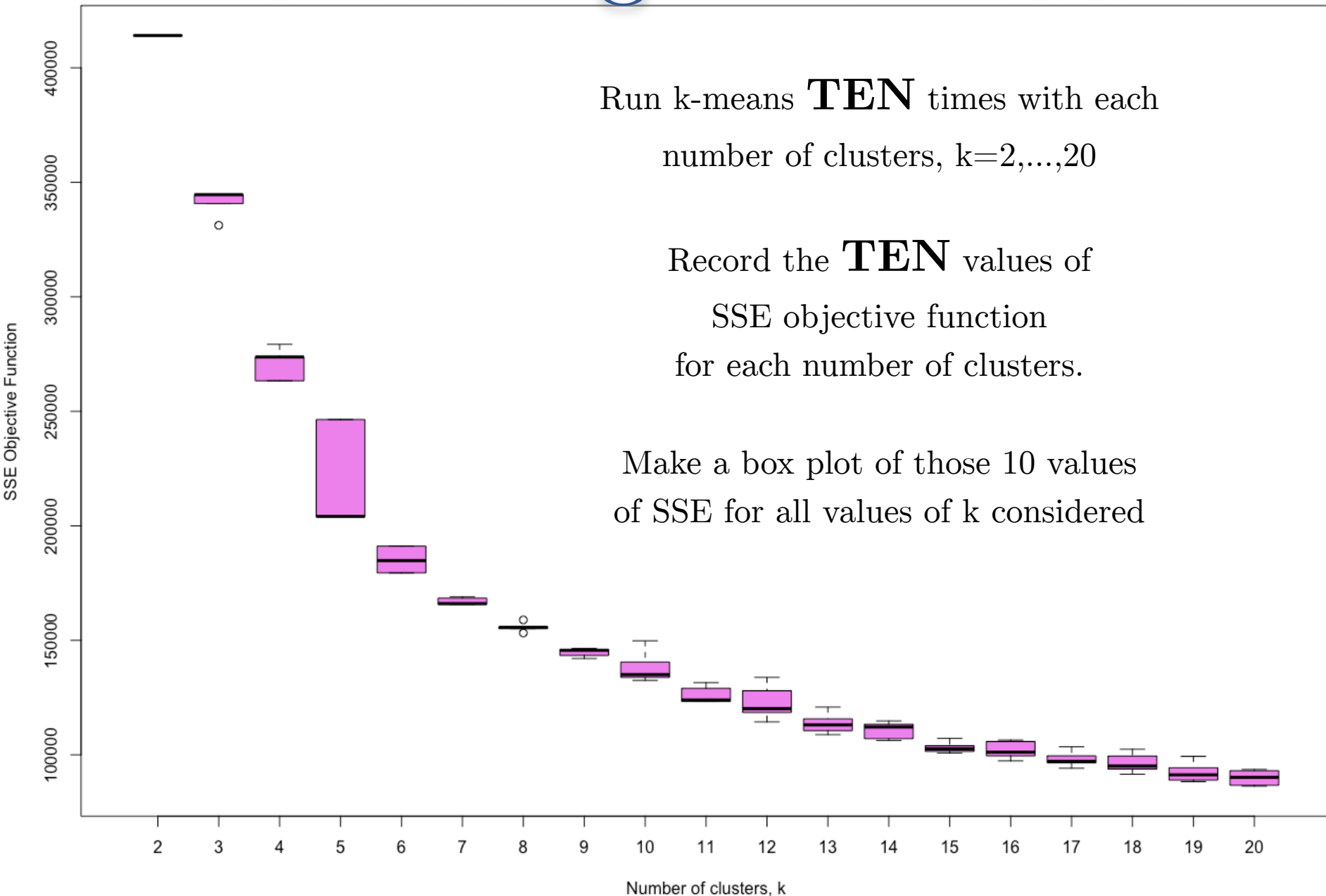

2D Projection



Choosing k from SSE

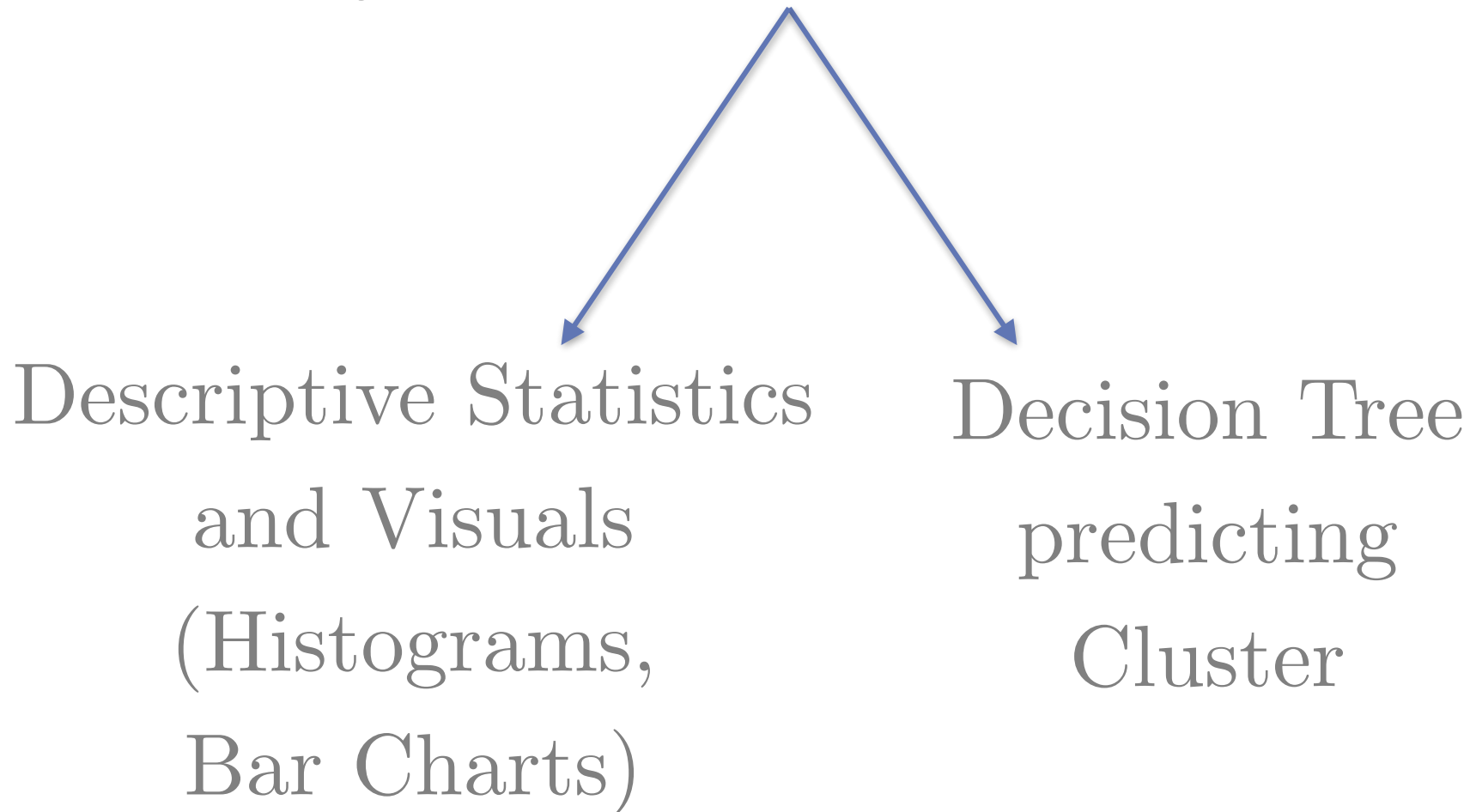


Choosing k from SSE

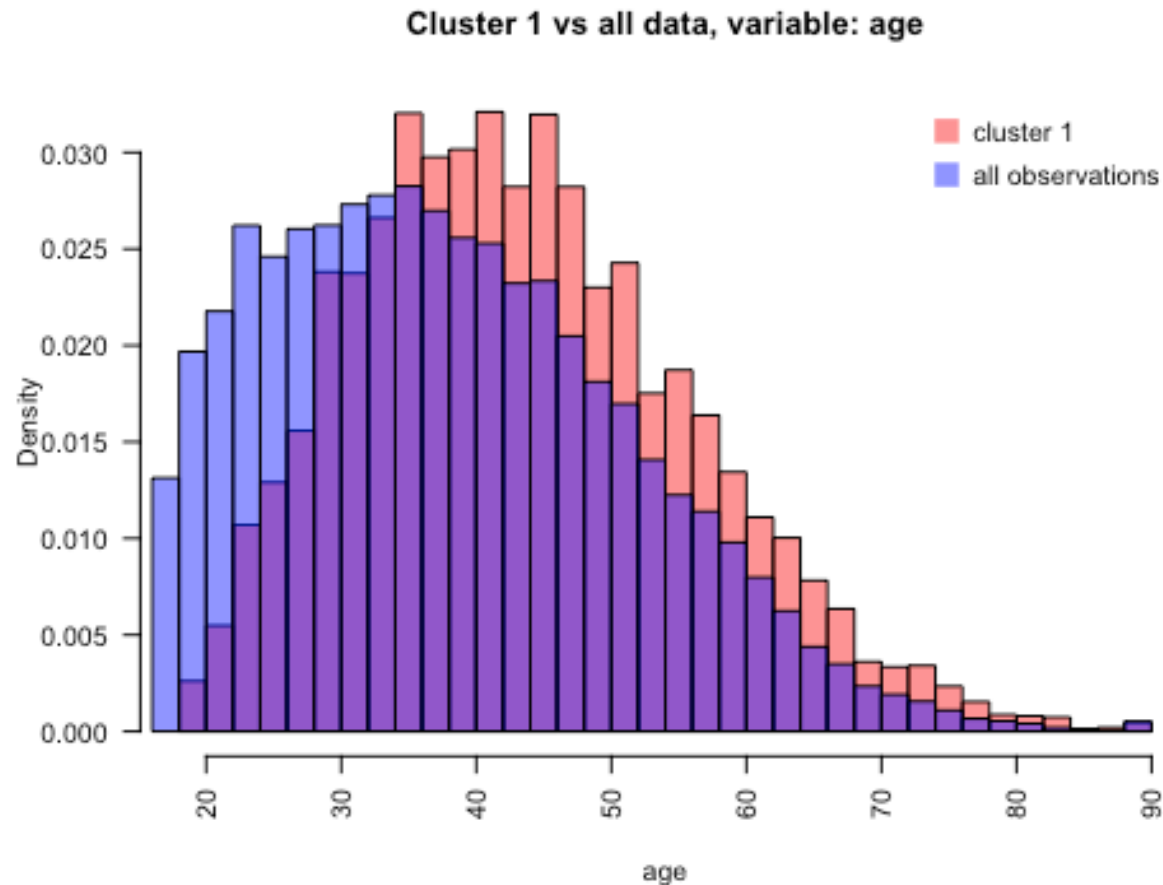


Now Profile the Clusters

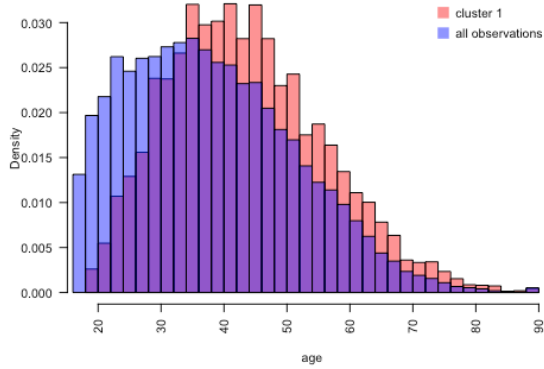
My two favorite approaches



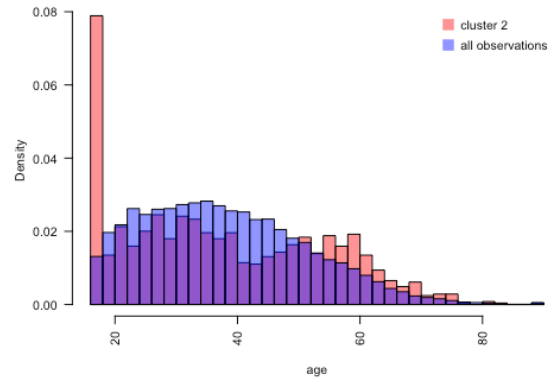
Overlaid Histogram, One cluster vs. all observations



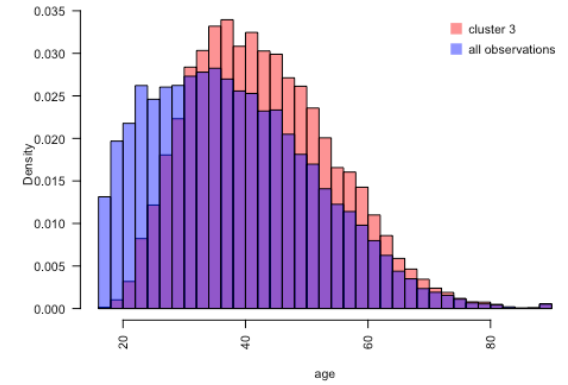
Cluster 1 vs all data, variable: age



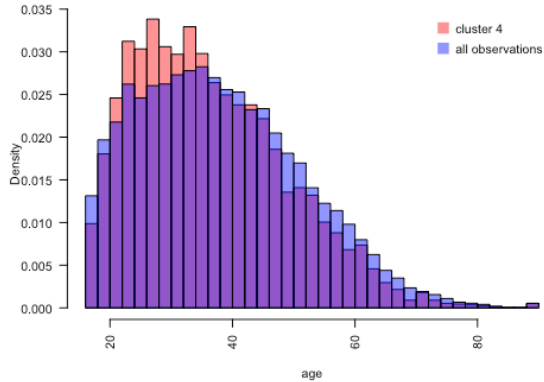
Cluster 2 vs all data, variable: age



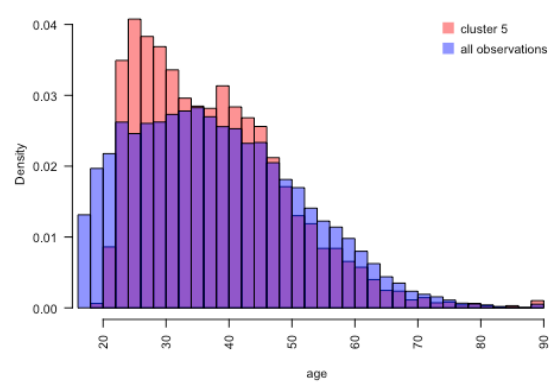
Cluster 3 vs all data, variable: age



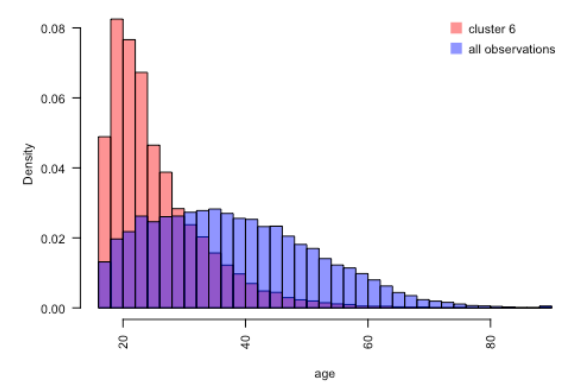
Cluster 4 vs all data, variable: age



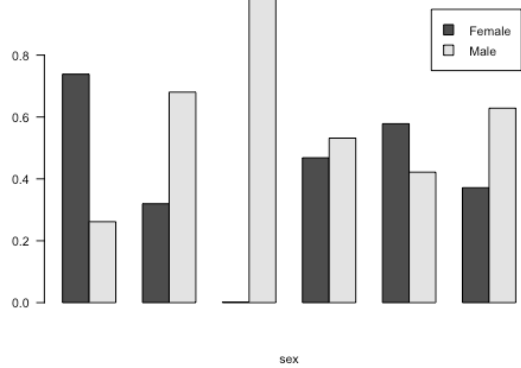
Cluster 5 vs all data, variable: age



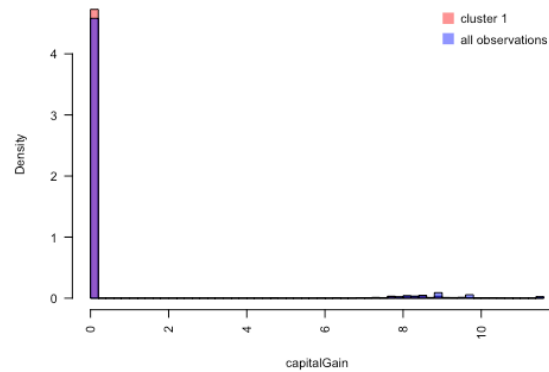
Cluster 6 vs all data, variable: age



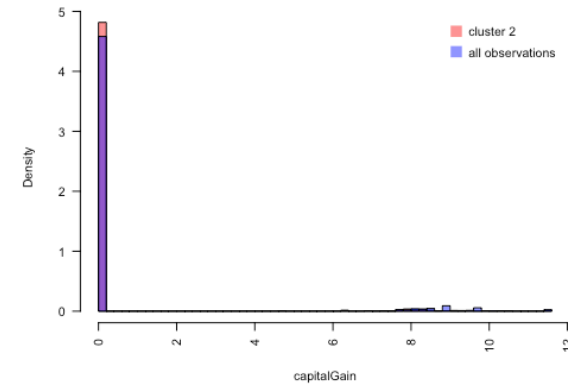
sex by cluster

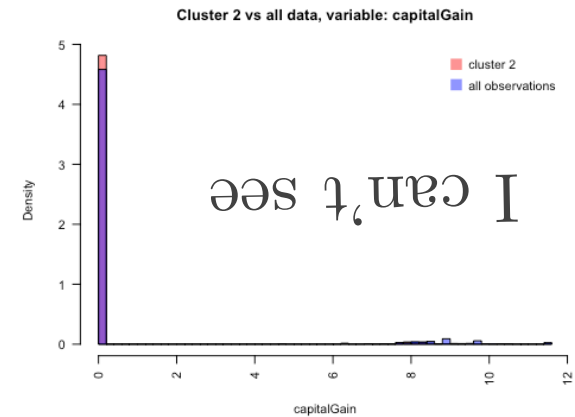
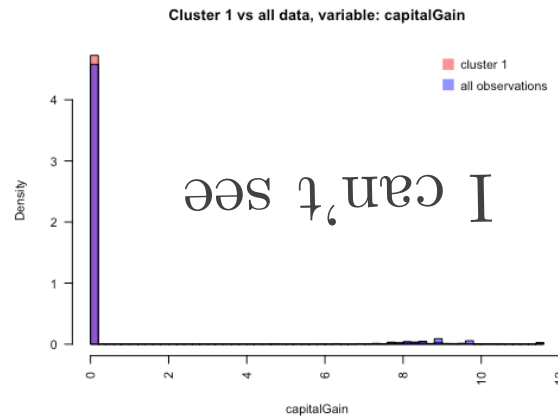
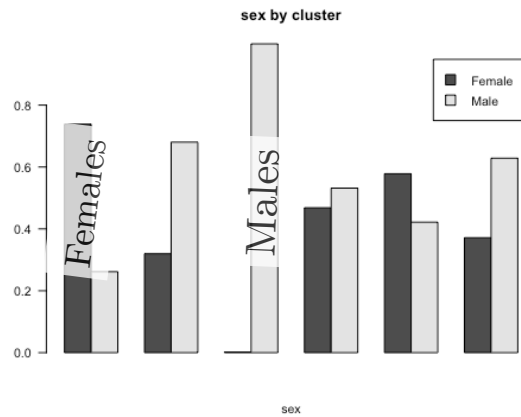
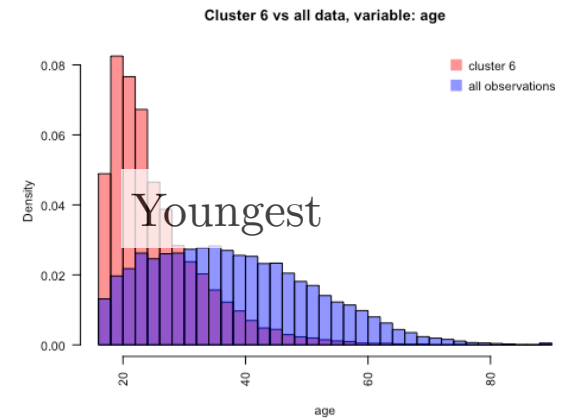
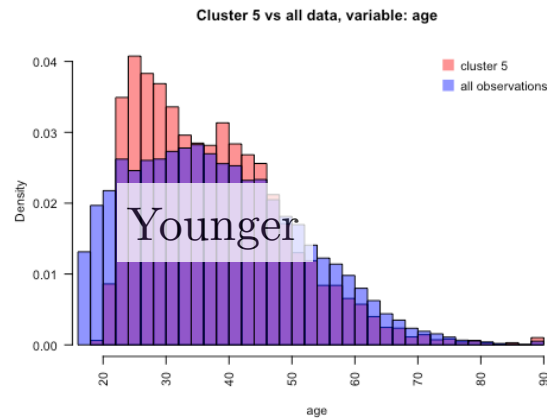
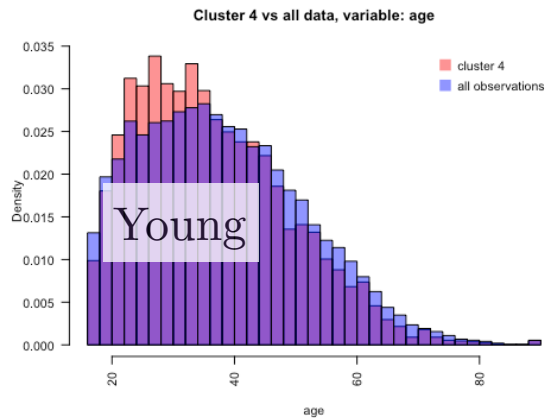
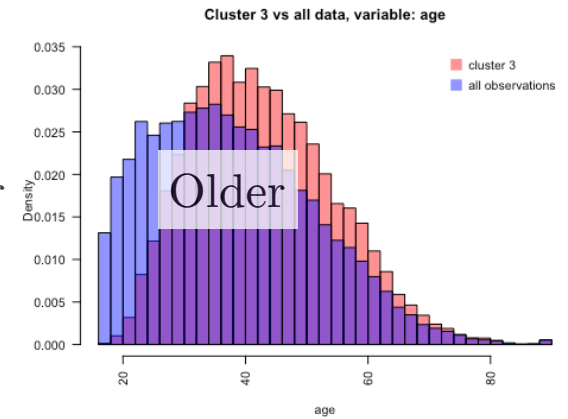
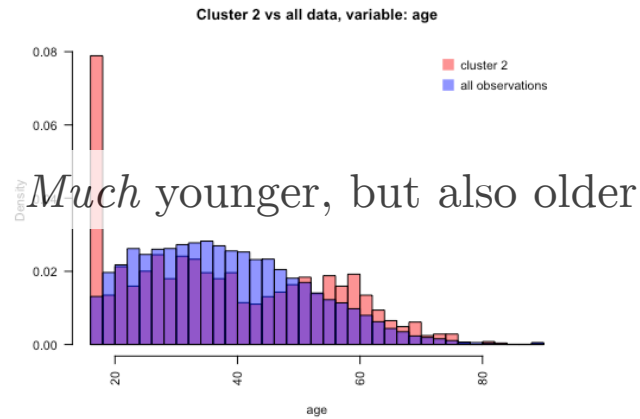
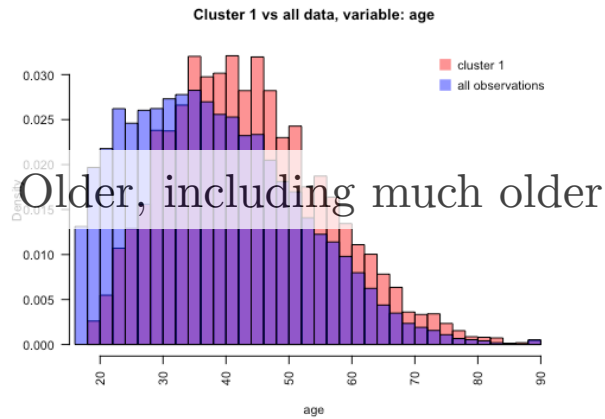


Cluster 1 vs all data, variable: capitalGain

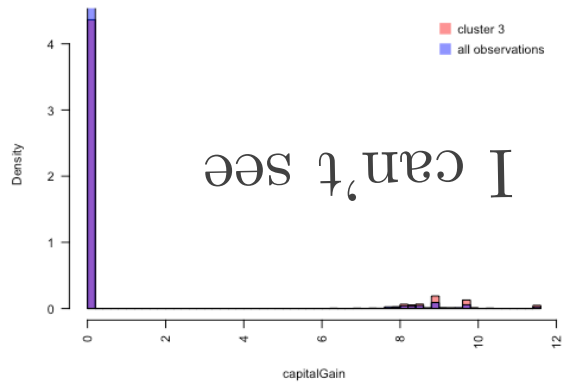


Cluster 2 vs all data, variable: capitalGain

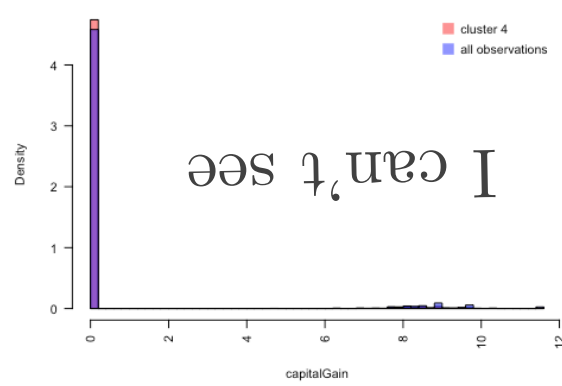




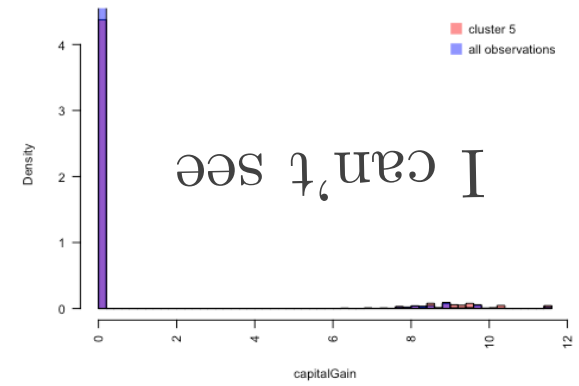
Cluster 3 vs all data, variable: capitalGain



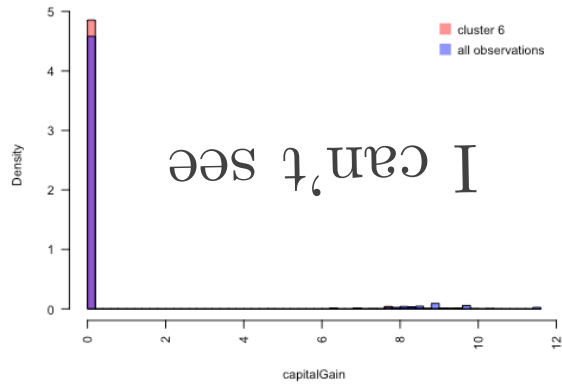
Cluster 4 vs all data, variable: capitalGain



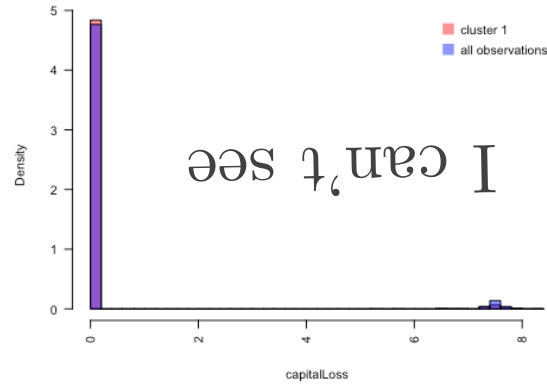
Cluster 5 vs all data, variable: capitalGain



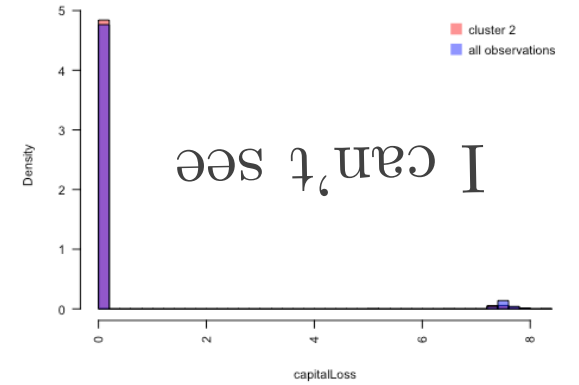
Cluster 6 vs all data, variable: capitalGain



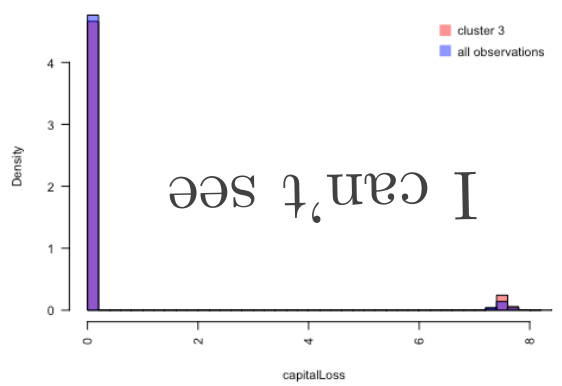
Cluster 1 vs all data, variable: capitalLoss



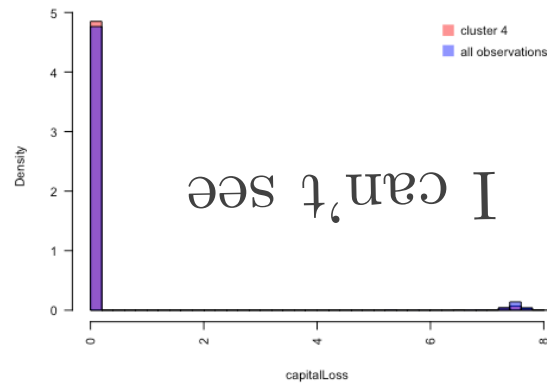
Cluster 2 vs all data, variable: capitalLoss



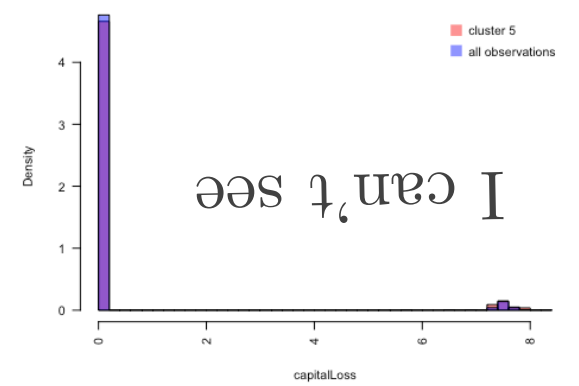
Cluster 3 vs all data, variable: capitalLoss



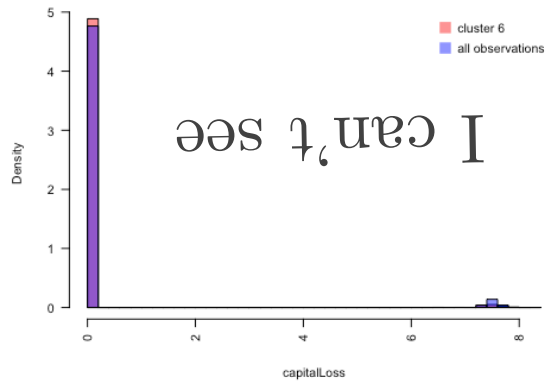
Cluster 4 vs all data, variable: capitalLoss



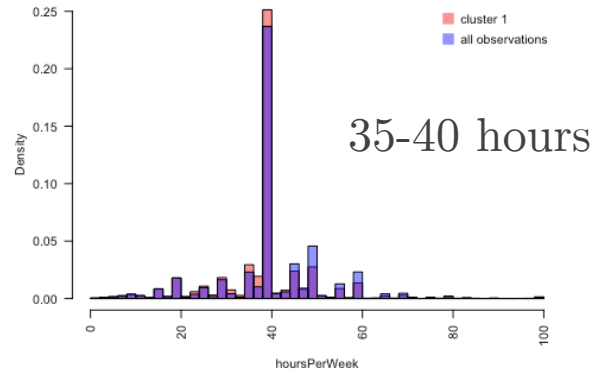
Cluster 5 vs all data, variable: capitalLoss



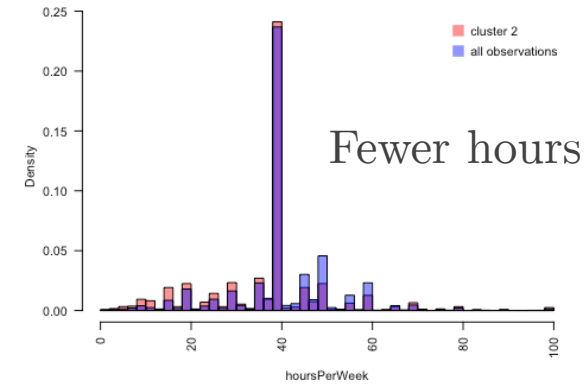
Cluster 6 vs all data, variable: capitalLoss



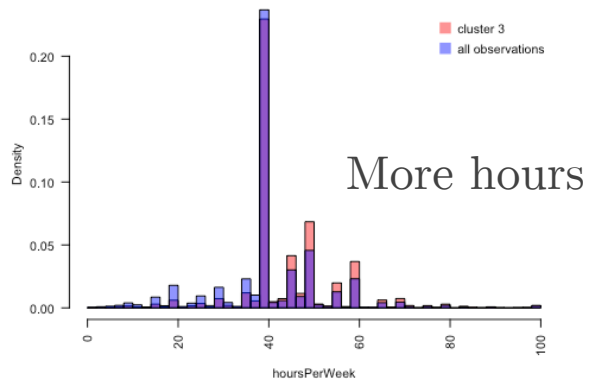
Cluster 1 vs all data, variable: hoursPerWeek



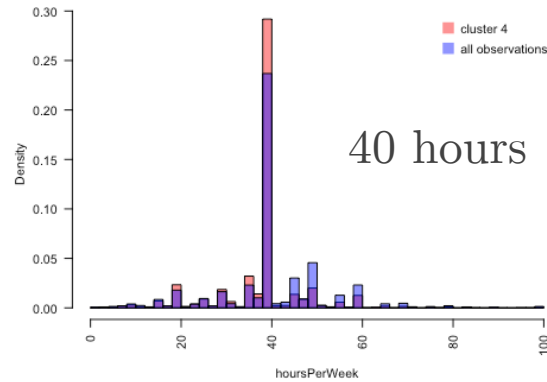
Cluster 2 vs all data, variable: hoursPerWeek



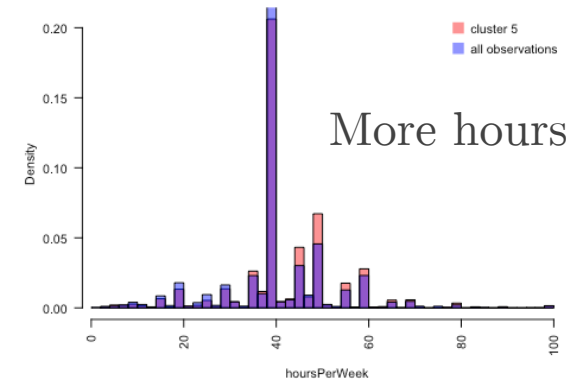
Cluster 3 vs all data, variable: hoursPerWeek



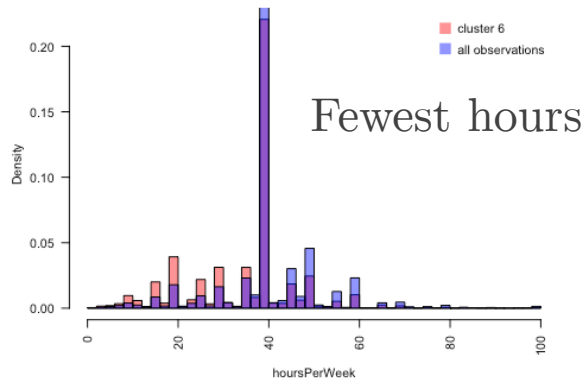
Cluster 4 vs all data, variable: hoursPerWeek



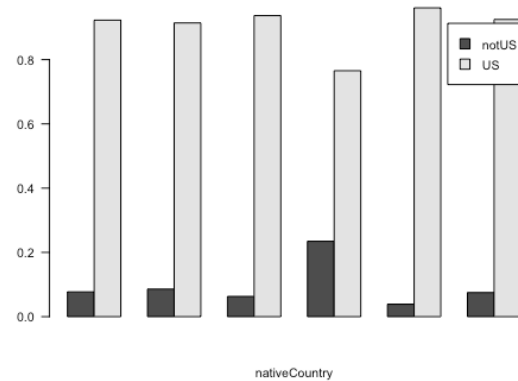
Cluster 5 vs all data, variable: hoursPerWeek



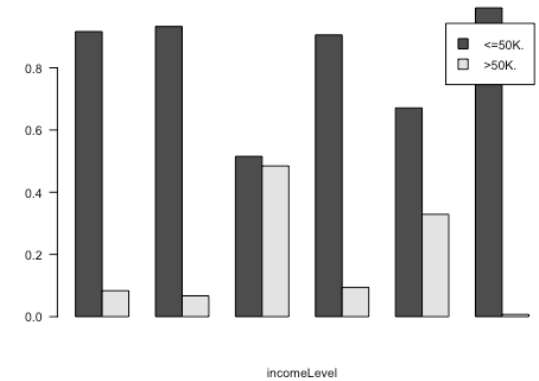
Cluster 6 vs all data, variable: hoursPerWeek



nativeCountry by cluster

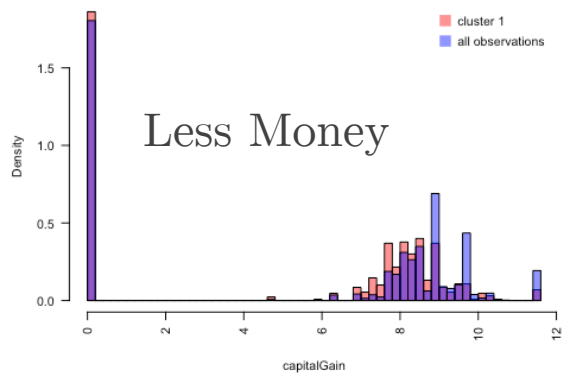


incomeLevel by cluster

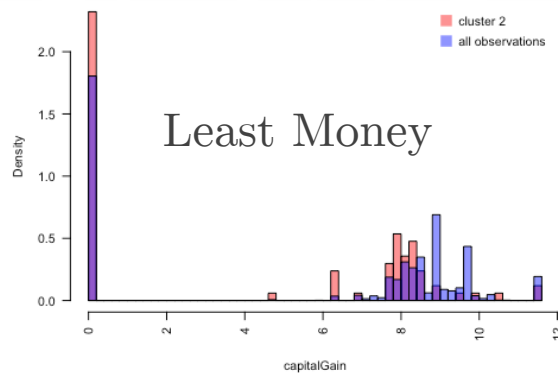


Profiling Capital Gains and Losses by cutting out obs with Zeros on both

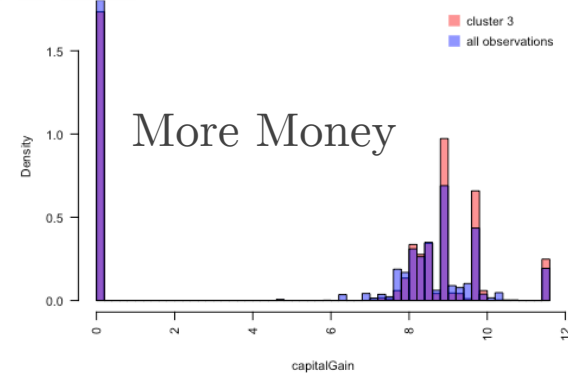
Cluster 1 vs all data, variable: capitalGain



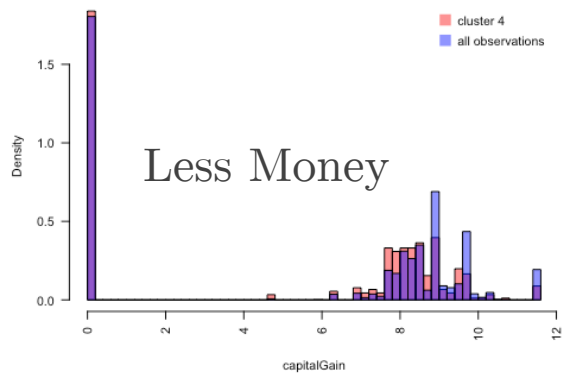
Cluster 2 vs all data, variable: capitalGain



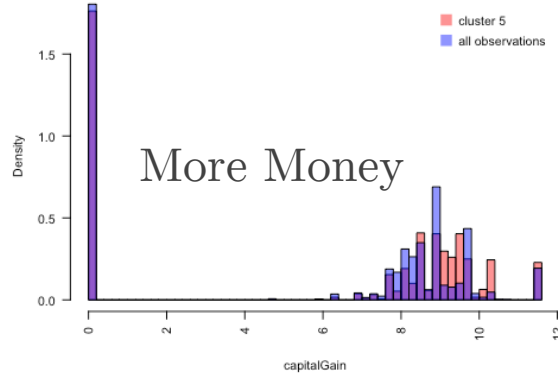
Cluster 3 vs all data, variable: capitalGain



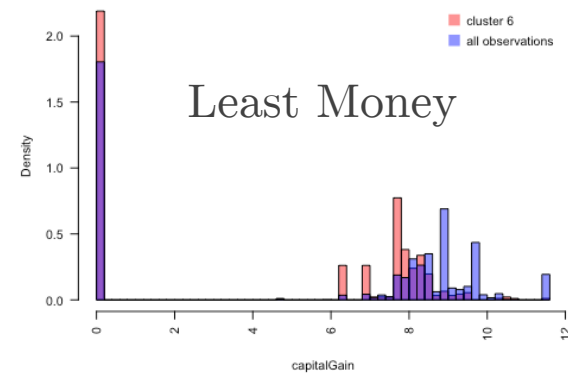
Cluster 4 vs all data, variable: capitalGain



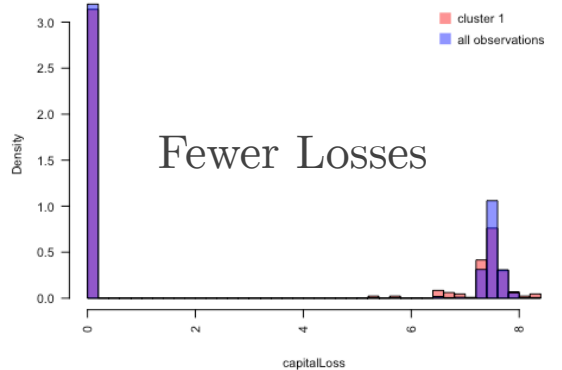
Cluster 5 vs all data, variable: capitalGain



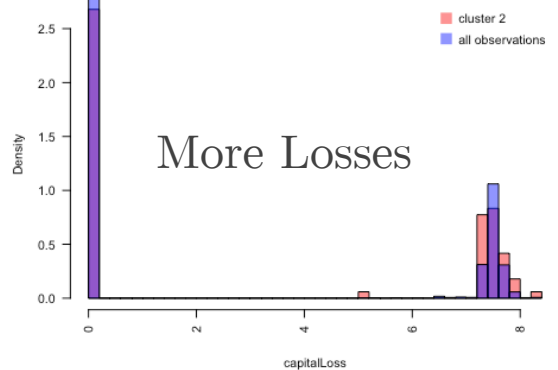
Cluster 6 vs all data, variable: capitalGain



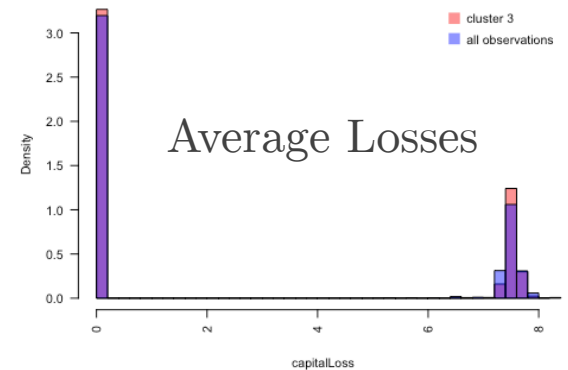
Cluster 1 vs all data, variable: capitalLoss



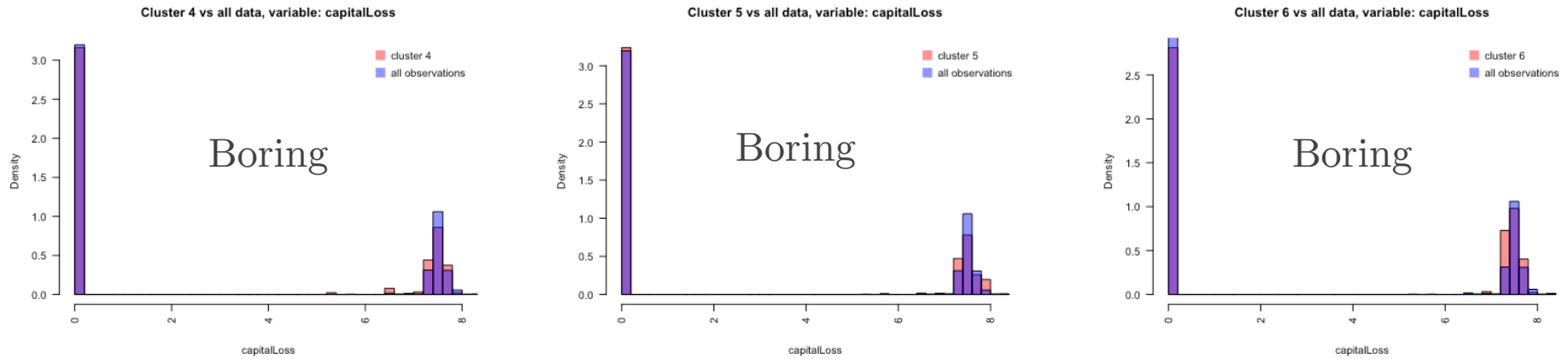
Cluster 2 vs all data, variable: capitalLoss



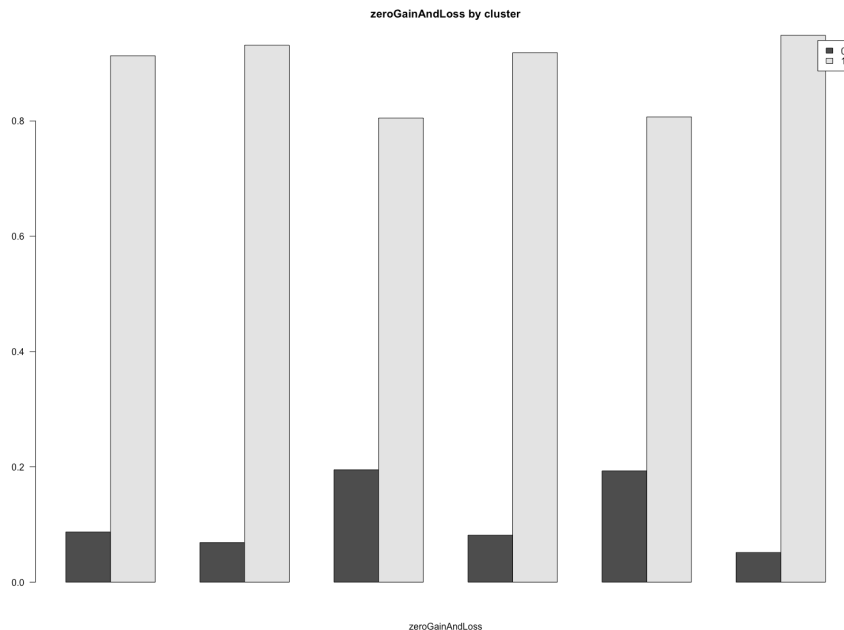
Cluster 3 vs all data, variable: capitalLoss



Profiling Capital Gains and Losses by cutting out obs with Zeros on both



Look also at % of each cluster with Zeros on both:



This variable is 1 (light grey) if **both** **capitalGains** and **capitalLoss** were 0

It is 0 otherwise (i.e. individual had some capital gain OR loss..)

Part 2

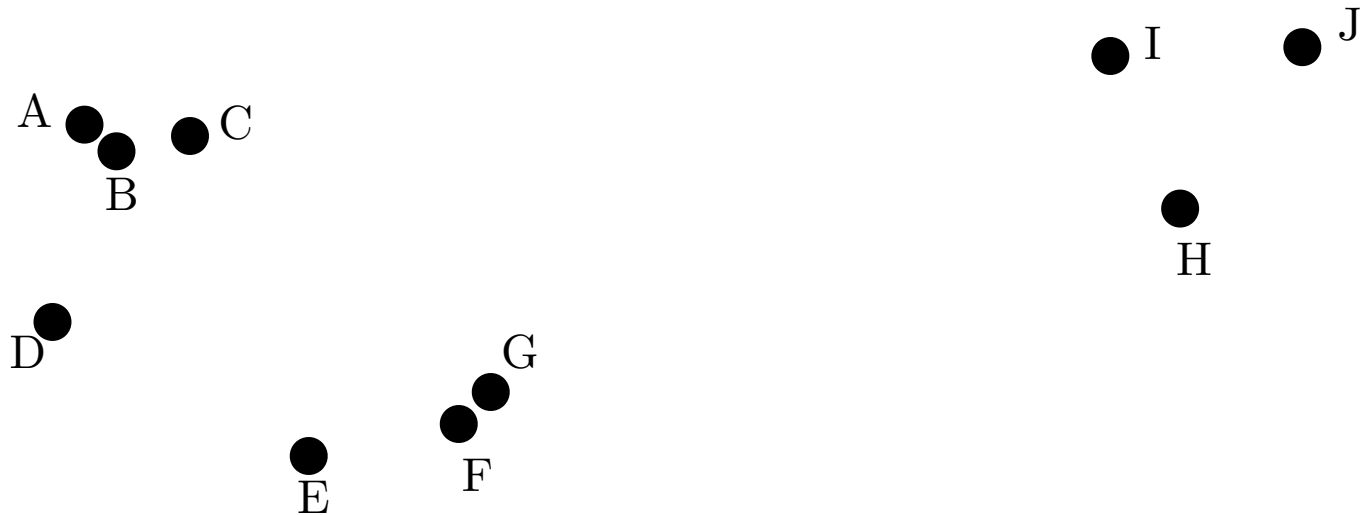
• • •

Hierarchical clustering and other algorithms

Hierarchical Clustering

(Agglomerative)

Some Data



Hierarchical Clustering

(Agglomerative)

Some Data

A ●
B ●
C ●

D ●

E ●

F ●
G ●

Starts with each
point in its own
cluster.

I ●
J ●
H ●

At each step of the algorithm,
computes the distance between
all pairs of clusters and combines
the two closest clusters.

Hierarchical Clustering

(Agglomerative)

First Step



Hierarchical Clustering

(Agglomerative)

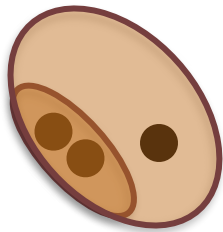
Second Step



Hierarchical Clustering

(Agglomerative)

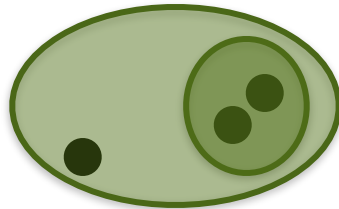
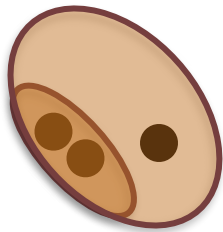
Third Step



Hierarchical Clustering

(Agglomerative)

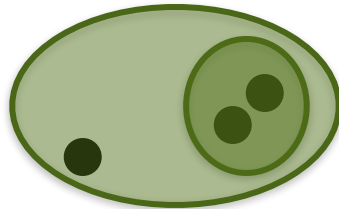
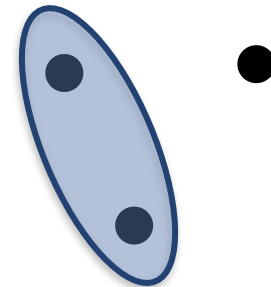
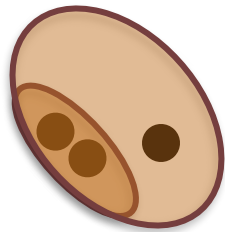
Forth Step



Hierarchical Clustering

(Agglomerative)

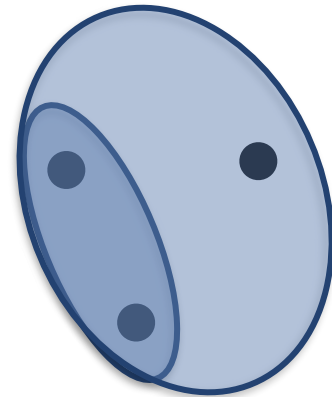
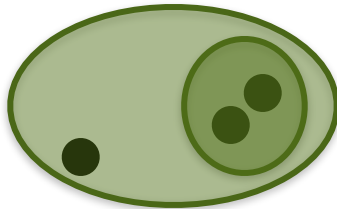
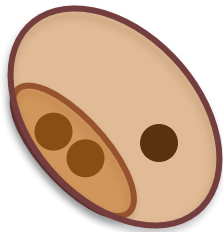
Fifth Step



Hierarchical Clustering

(Agglomerative)

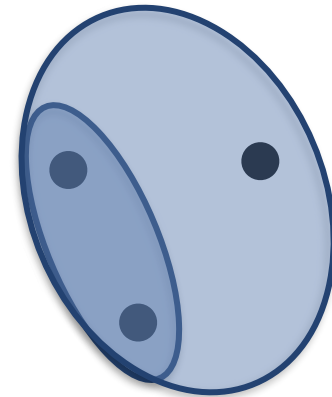
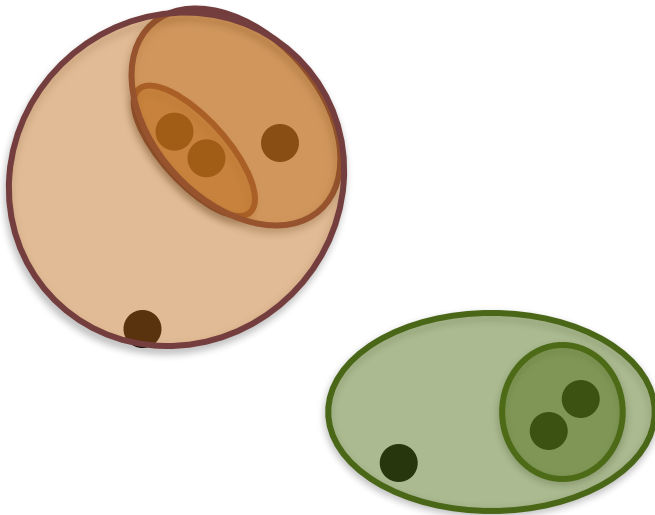
Sixth Step



Hierarchical Clustering

(Agglomerative)

Seventh Step

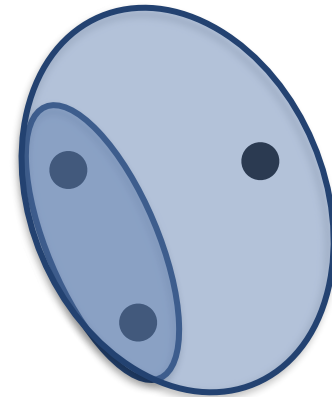
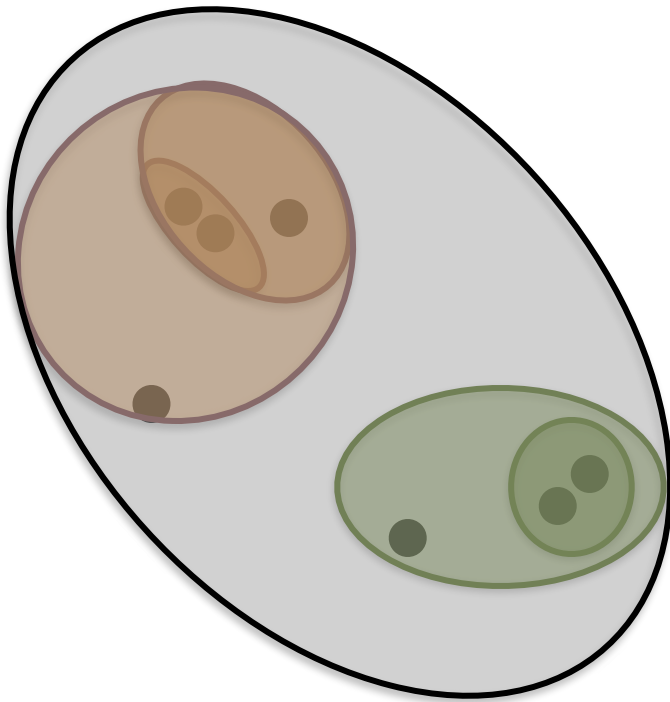


We might have known that we only wanted 3 clusters, in which case we'd stop once we had 3.

Hierarchical Clustering

(Agglomerative)

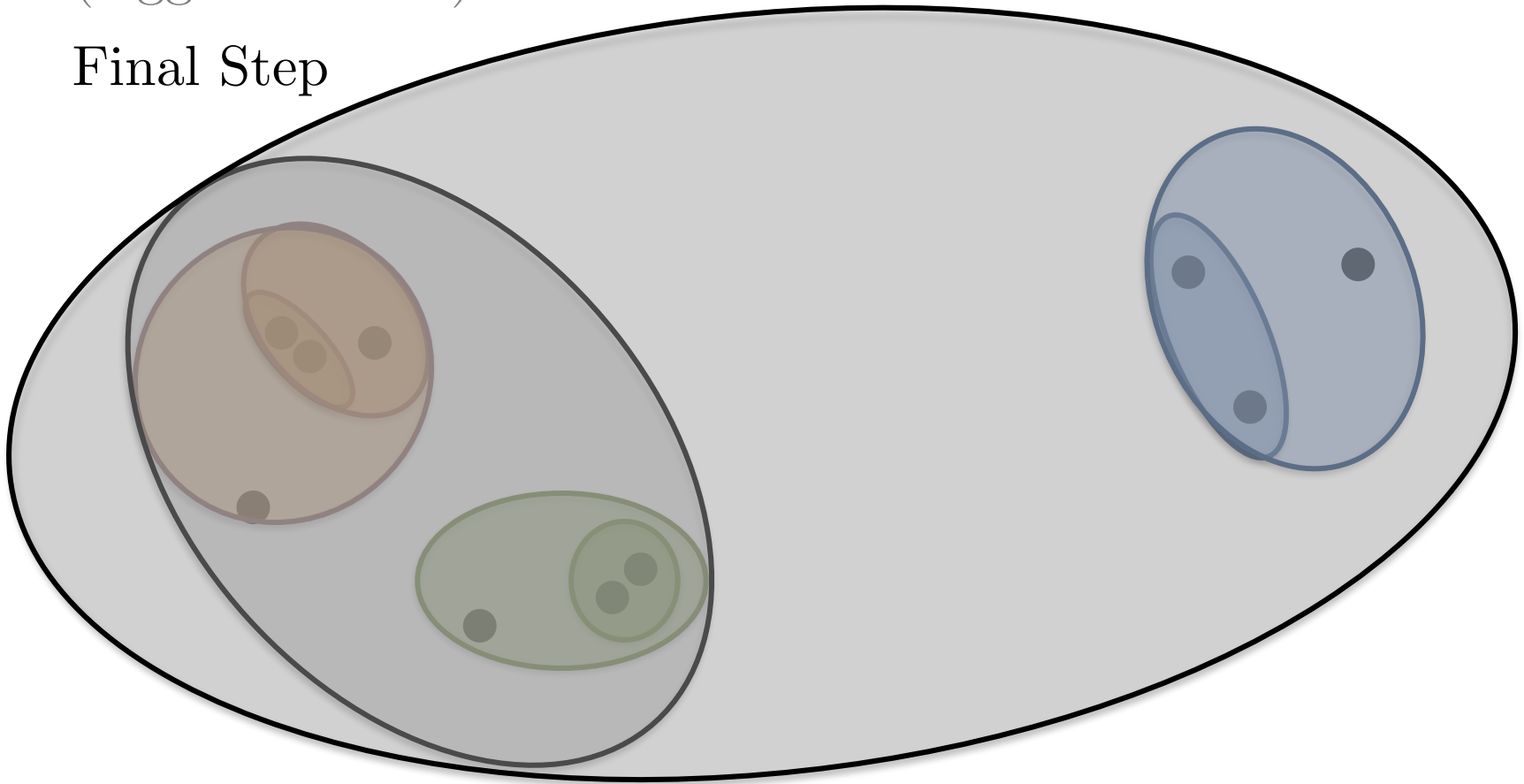
Eighth Step



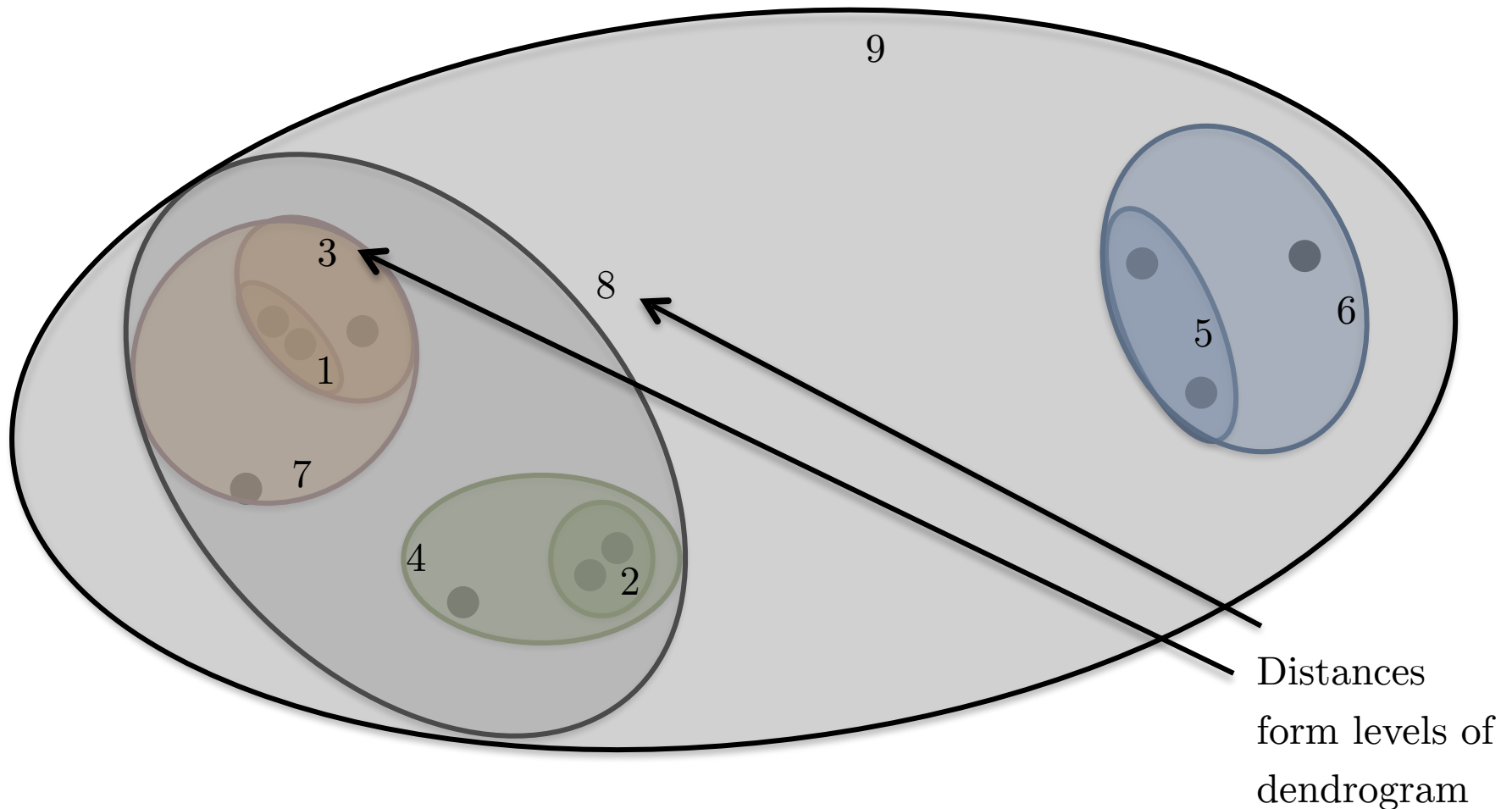
Hierarchical Clustering

(Agglomerative)

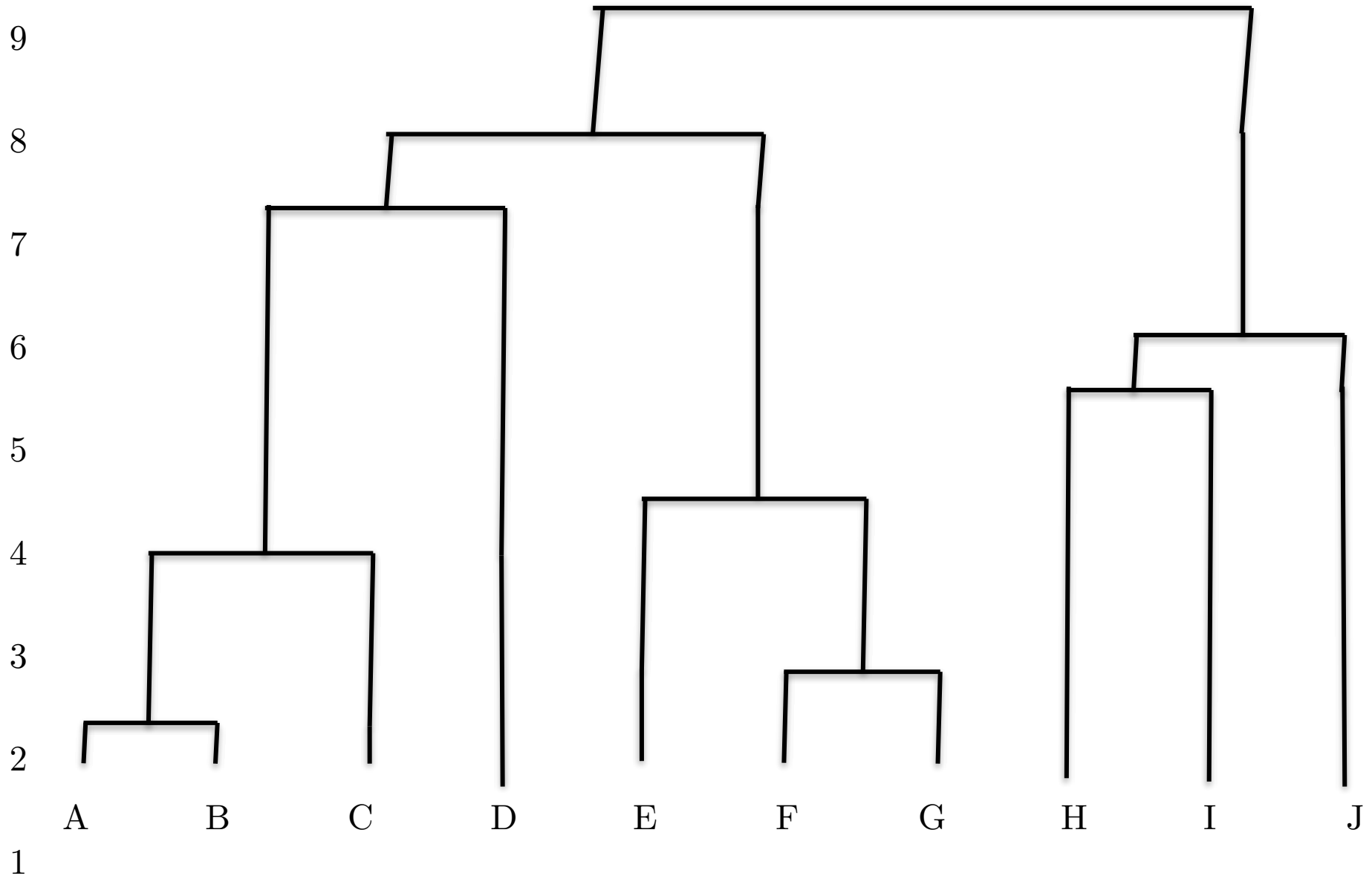
Final Step



Hierarchical Clustering



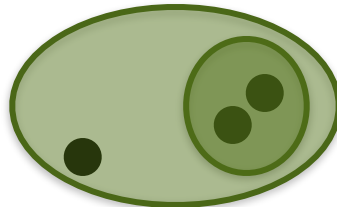
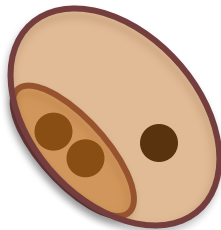
Resulting Dendrogram (not to scale)



Linkages

Which clusters/points are closest to each other?

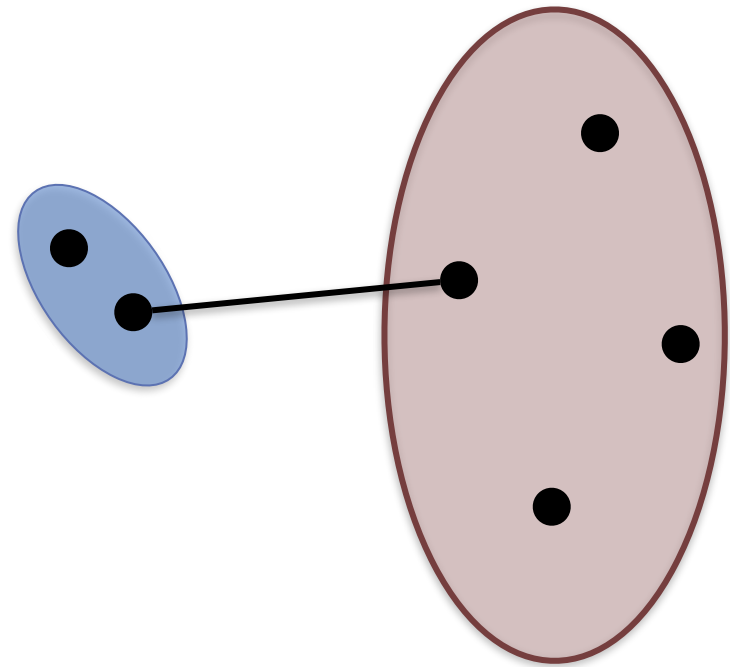
How do I measure the distance between a point/cluster and a cluster?



Linkages

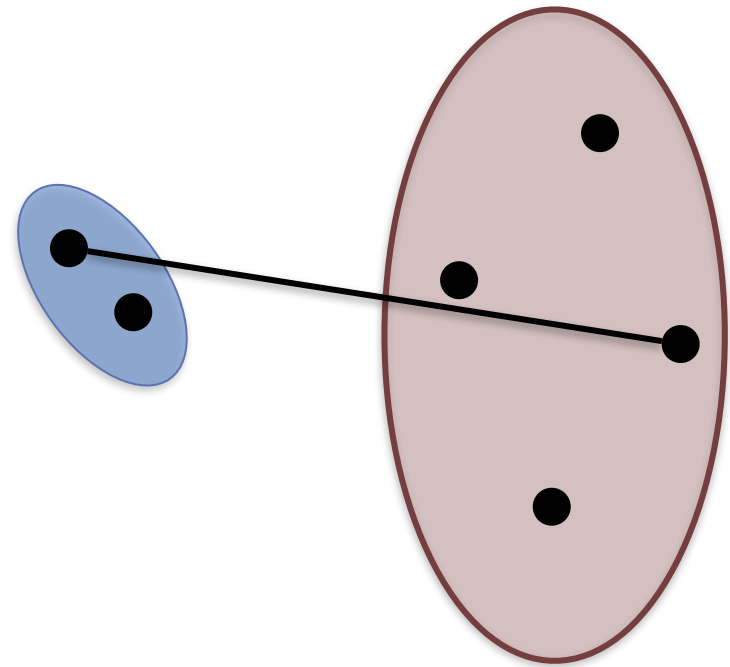
Single Linkage: Distance between the closest points in the clusters.

(Minimum Spanning Tree)



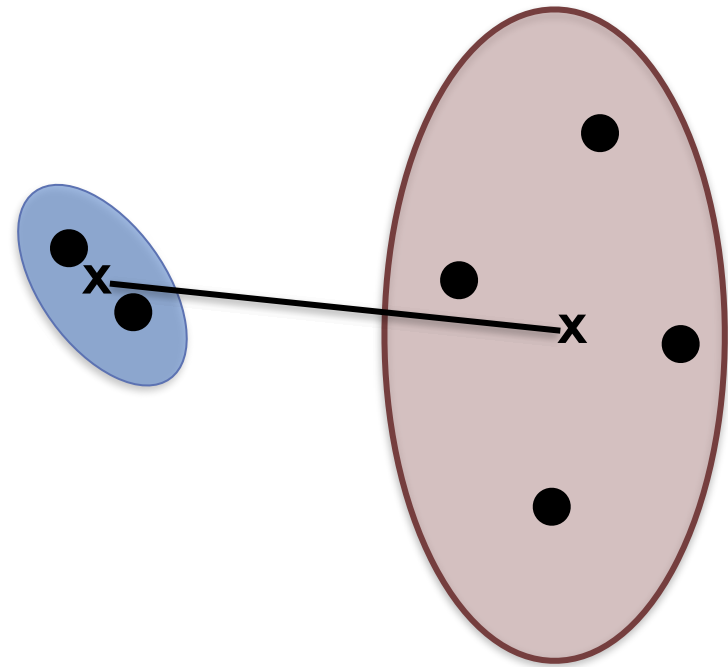
Linkages

Complete Linkage: Distance between the farthest points in the clusters.



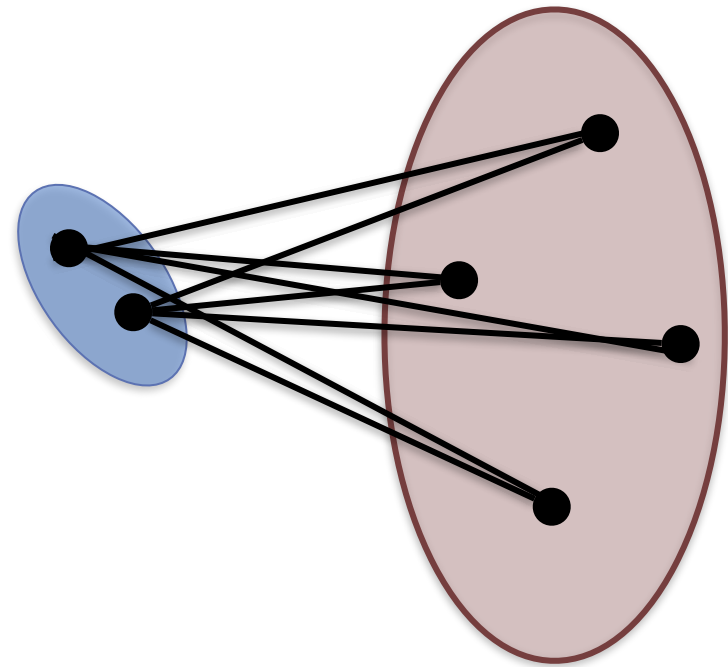
Linkages

Centroid Linkage: Distance between the centroids (means) of each cluster.



Linkages

Average Linkage: Average distance between all points in the clusters.



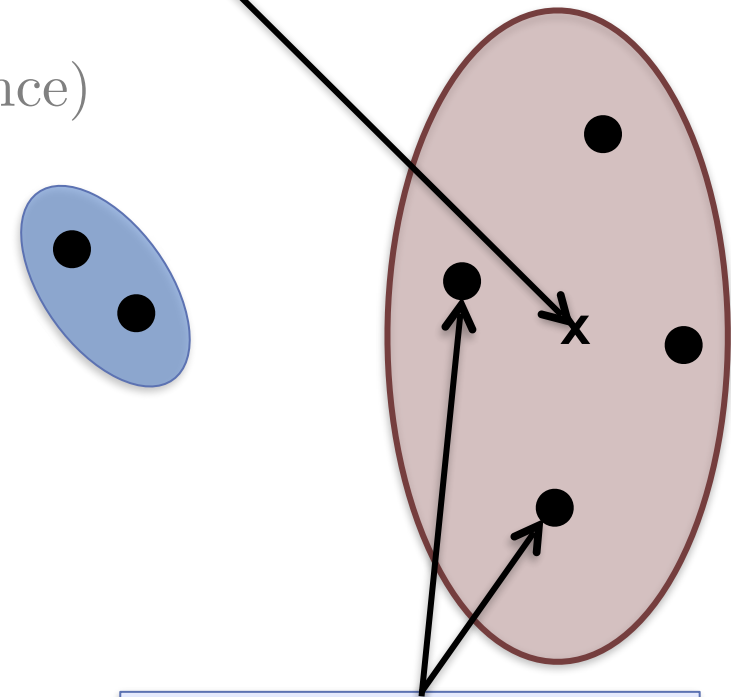
Linkages

Ward's Method: Increase in SSE (variance) when clusters are combined.

$$\sum_{j=1}^{N_i} \|\mathbf{x}_j - \mathbf{c}_i\|_2$$

- Default in SAS **PROC CLUSTER**
- Shown mathematically similar to centroid linkage

centroid for cluster i, \mathbf{c}_i

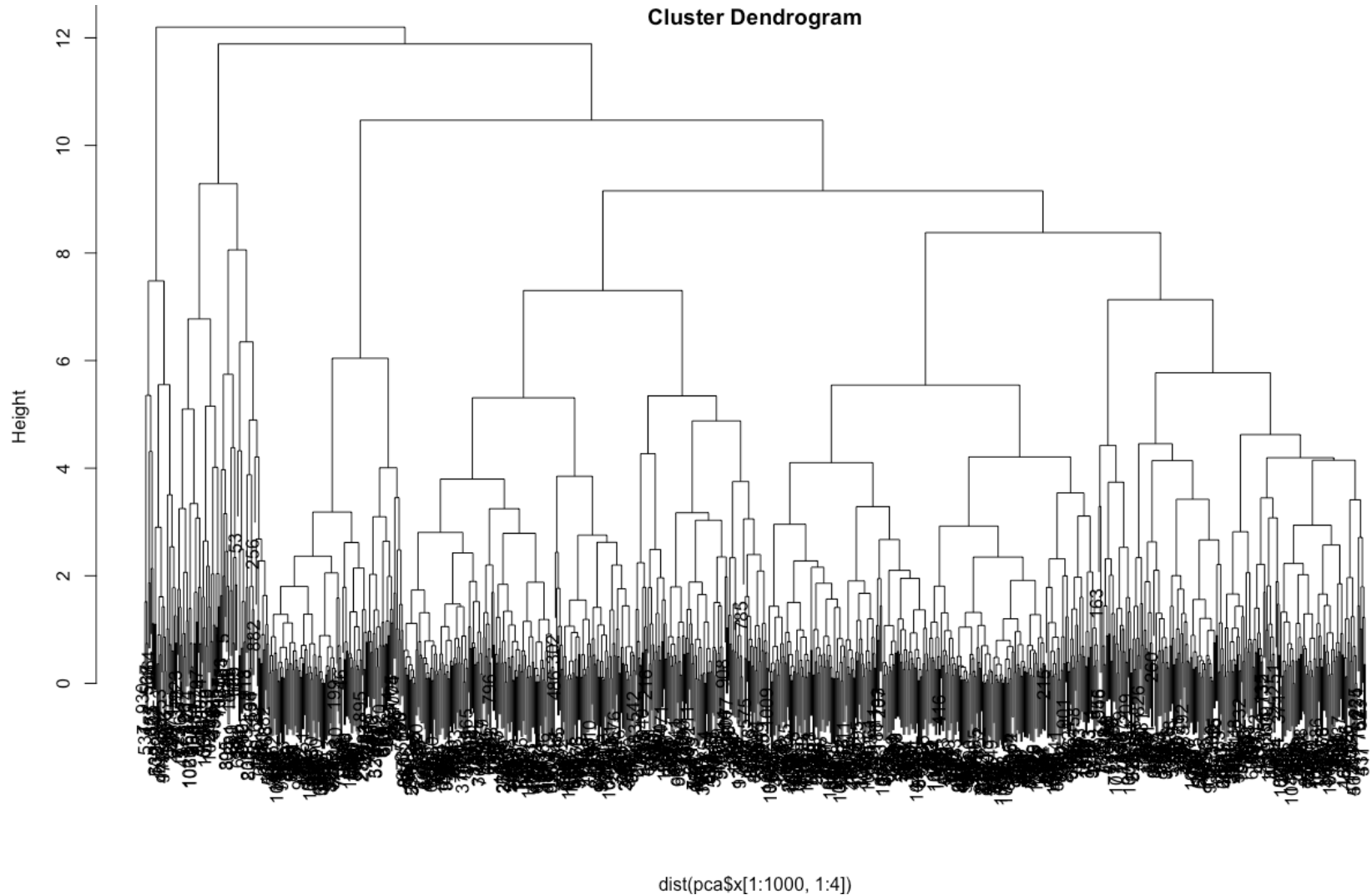


data points in cluster i:

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_i}$

Hierarchical Clustering Summary

Disadvantages



Hierarchical Clustering Summary

Disadvantages

- Computationally intensive, large storage requirements, **not good for large datasets**
- **Lacks global objective** function: only makes decision based on local criteria.
- **Merging decisions are final.** Once a point is assigned to a cluster, it stays there.
- **Poor performance on noisy or high-dimensional data** like text.

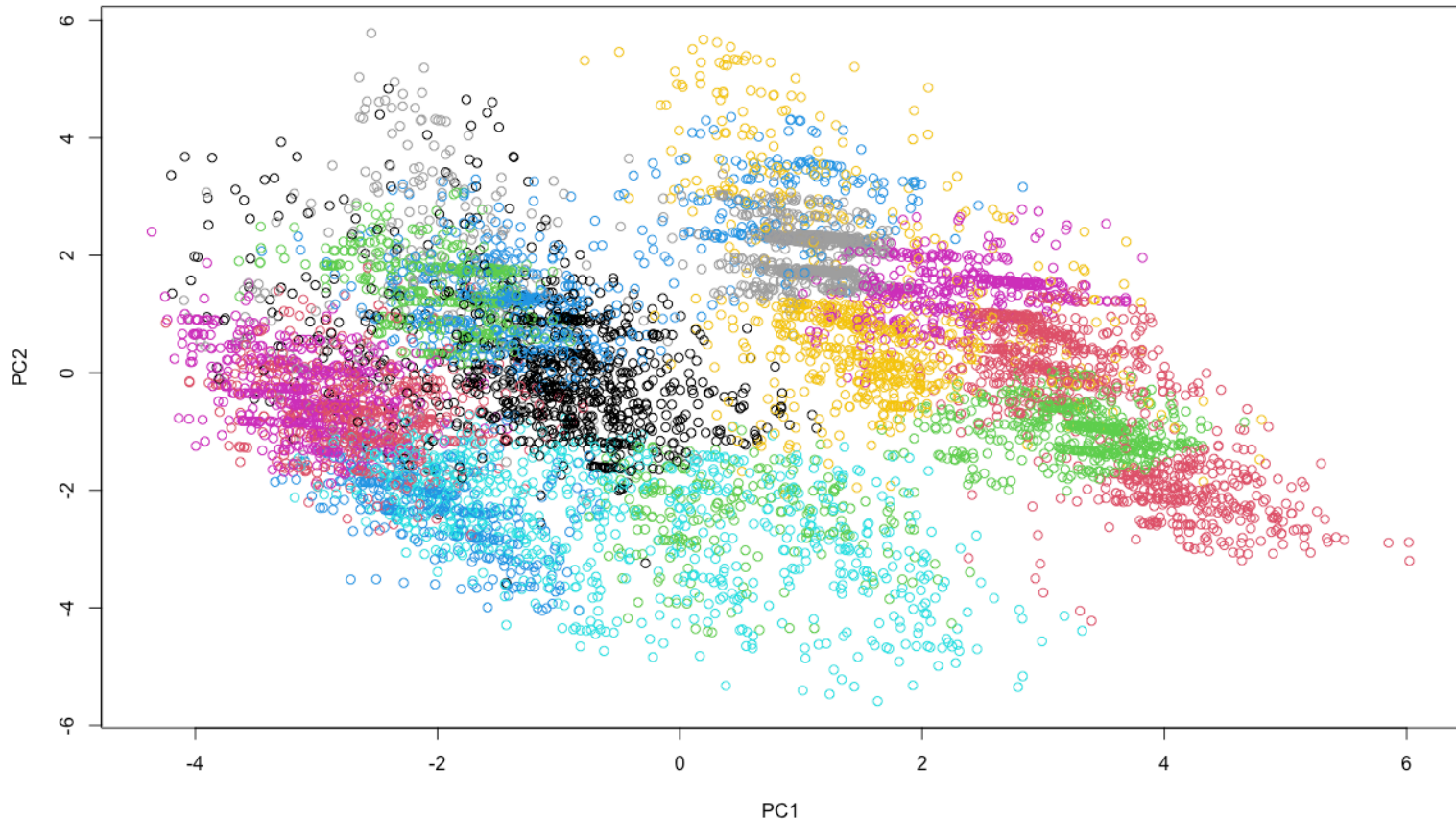
Advantages

- Lacks global objective function: **no complicated algorithm** or problem with local minima
- **Creates hierarchy (dendrogram)** that can help choose the number of clusters and examine how those clusters relate to each other.
- **Can be used in conjunction with other faster methods** (input the centroids from a k-means result with $k=100$) - *next few slides*

Another Approach

Cluster into a large number of clusters using k-means (here, $k=20$)

```
#####  
k6v1 = kmeans(pca$x[,1:4],20)  
#####
```



Another Approach

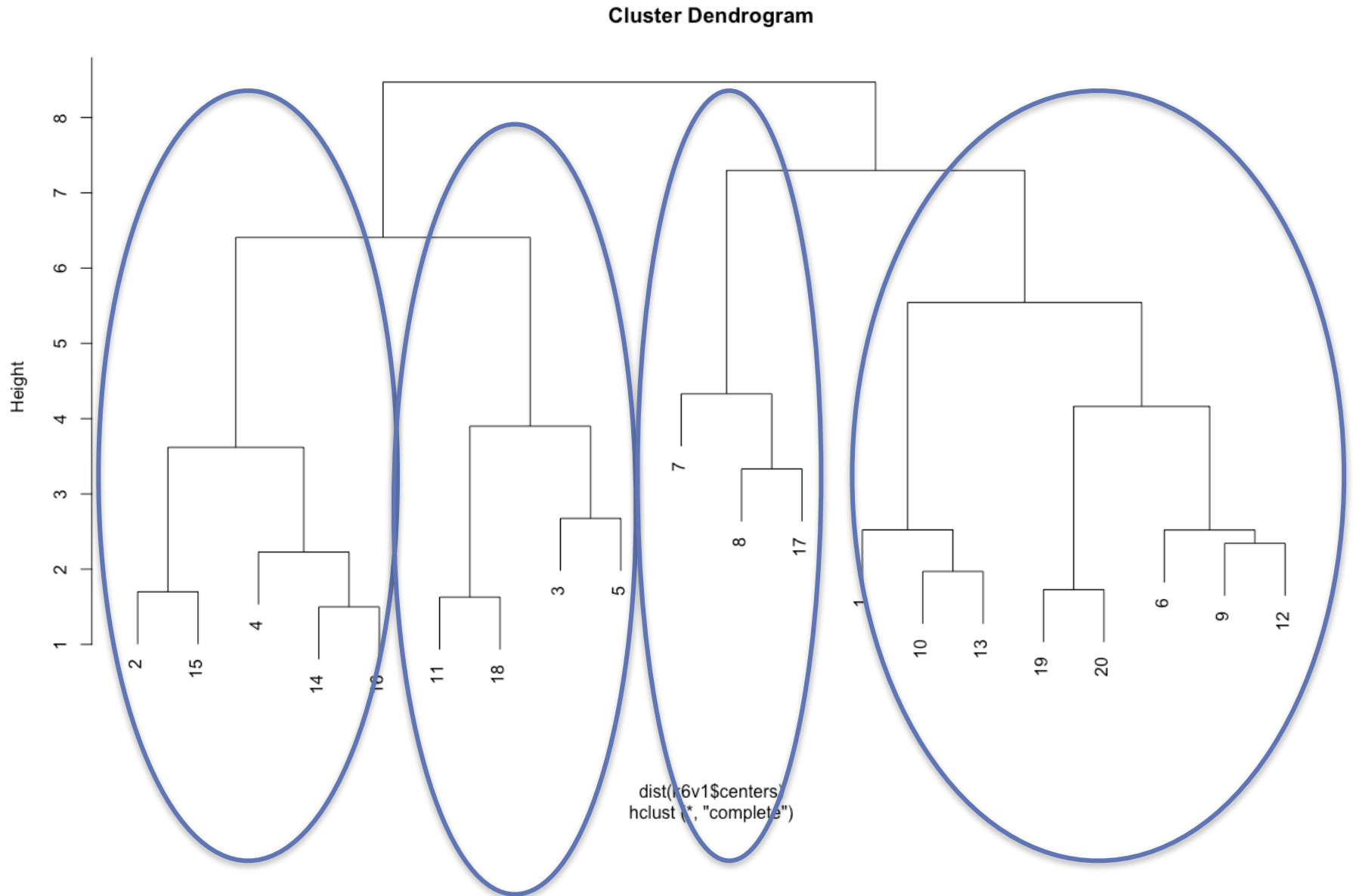
Create a distance matrix using those 20 centroids and input to Hierarchical Clustering.

```
#####  
  hc = hclust(dist(k6v1$centers))  
  plot(hc)  
#####
```

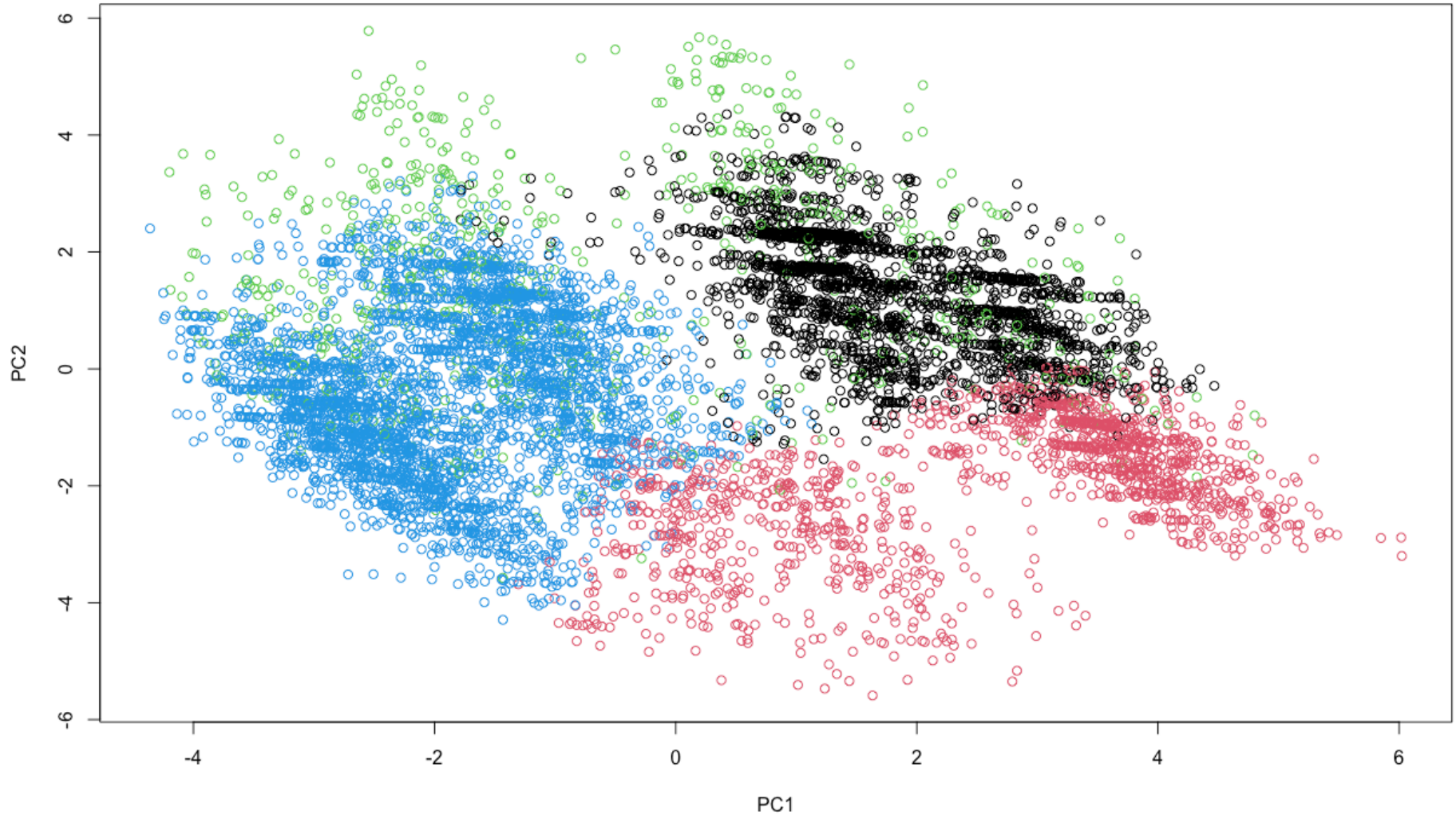
Dendrogram Linking 20 Cluster Centroids



Dendrogram Linking 20 Cluster Centroids



Corresponding 4 cluster solution



Case Study: Breast Cancer Dataset

...

Practice Makes Perfect!

Hierarchical Clustering in Viya

```
cas;  
caslib _all_ assign;
```

```
proc stdize data=PUBLIC.BREAST_CANCER stdonly outsdz=BCdist;  
    var interval(CT Size Shape Margin Epithelial Bare Chromatin  
Normal Mitoses /  
    std=std);  
run;
```

```
proc cluster data = BCdist  
    method = centroid  
    plots(only maxpoints=700)=(dendrogram ccc pseudo)  
    outtree=work.Cluster_tree;  
var bare chromatin ct epithelial margin mitoses normal shape size  
;  
run;
```

Additional Slides (Self Study)

...

Practice Case Studies at the End

Other types of Clustering (self-study)

- **DBSCAN** – Density based algorithm designed to find dense areas of points. Capable of identifying ‘noise’ points which do not belong to any clusters.
 - interactive demo:
<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>
- **Graph/Network Clustering** – Spectral clustering and modularity maximization. Covered in Social Network Analysis in Spring.

Cluster Validation and Profiling

• • •

Additional Notes for Self-Study

Cluster Validation

How do I know that my clusters are actually clusters?

- Lots of techniques/metrics have been proposed
 - Measure **separation** between clusters
 - Measure **cohesion** within clusters
 - Validation metrics usually “algorithm agnostic”
 - NbClust package in R has [many](#)
- All have merit, most are difficult to interpret in the context of statistical significance.

Cluster Validation

To establish statistical significance:

- Show that you can't do just as well with randomized data (i.e. assume the null hypothesis of no clusters)
- Simulate ~ 1000 random data sets choosing from the distributions or ranges of your variables. Cluster them with the same number of clusters. Record the SSE (k-means objective function) or validity metric of choice. Use this to show that your actual SSE is far better than you could expect to achieve if no clusters exist.

Aligned Box Criterion

- SAS proprietary method of validating number of clusters.
- Compares the compactness (SSE) of the observed clusters with the theoretical distribution of compactness from random groups of points (the “reference distribution”)
- Want this gap between the observed and expected compactness to be as great as possible.
- 15 minute video:
<https://video.sas.com/detail/video/4572850292001/the-abcs-of-selecting-clusters>

Profiling Clusters

Now that we have clusters, how do we describe them?

- Use basic descriptives and hypothesis tests to show differences between clusters
- Use a decision tree to predict cluster

Some Explanation of SAS's Clustering Output (SELF-STUDY)

...

Because it's not exceedingly easy to figure out online!

Cubic Clustering Criterion (CCC)

- Only available in SAS (to my knowledge)
- $CCC > 2$ means that clustering is good
- $0 > CCC > 2$ means clustering requires examination
- If slightly negative, risk of outliers is low
- If $\sim < -30$ then risk of outliers is high
- Should not be used with single or complete linkage, but with centroid or ward's method.
- Each cluster must have >10 observations.

Determining Number of Clusters with the Cubic Clustering Criterion (CCC)

- A partition into k clusters is good when we see a dip in CCC for $k-1$ clusters and a peak for k clusters.
- After k clusters, the CCC should either gradually decrease or a gradual rise (the latter event happens when more isolated groups or points are present)¹

Determining Number of Clusters with the Cubic Clustering Criterion (CCC)

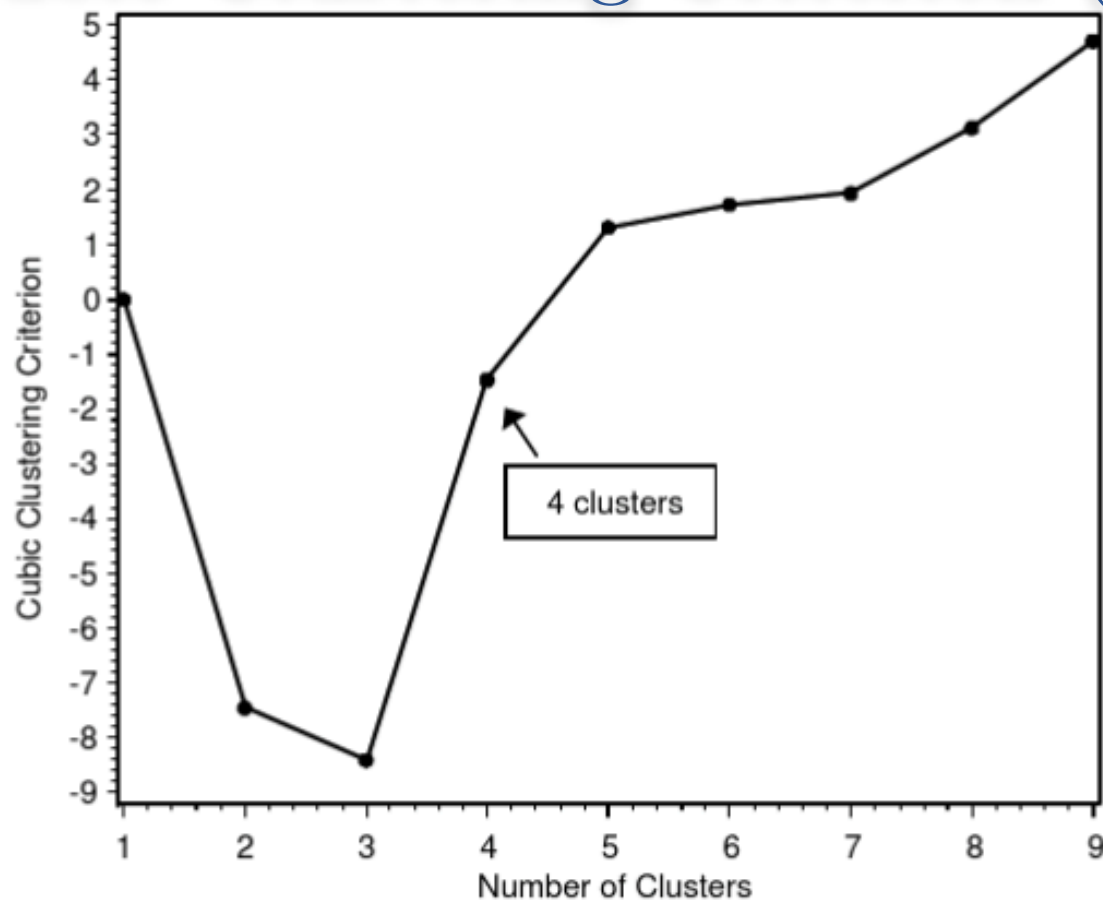


Figure 9.6 CCC as a function of the number of clusters.

Determining Number of Clusters with the Cubic Clustering Criterion (CCC)

WARNING: Do not expect the CCC to be common knowledge outside of the SAS domain.

[CITATION] **SAS technical report A-108**

WS Sarle - The Cubic Clustering Criterion. Cary, NC: SAS Institute, 1983

[Cited by 15](#) [Related articles](#) [Cite](#) [Save](#)

‘Overall R-Squared’ and ‘Pseudo-F’

These statistics draw connections between a final clustering and ANOVA.

- Total Sum of Squares (SST)
- Between Group Sum of Squares (SSB)
- Within Group Sum of Squares (SSW)
 - This is the k-means objective previously referred to as SSE.
 - Minimizing SSW \Rightarrow Maximizing SSB
- $SST = SSB + SSW$.
- ‘Overall R^2 ’ $= SSB/SST$
- $$pseudo\ F = \frac{R^2(k-1)}{(1-R^2)/(n-k)}$$

Case Study using Handwritten Digits (SELF-STUDY)

...

Practice makes perfect!

Example: PenDigit Data

- Goal: Automatic recognition of handwritten digits
- Digit database of 250 samples from 44 writers
- Subjects wrote digits in random order inside boxes of 500 by 500 tablet pixel resolution
- Spatial resampling to obtain a constant number of regularly spaced points on the trajectory
 - (x_1, x_2) give the first point coordinate
 - (x_3, x_4) give the second point coordinate
 - etc.

Example: PenDigit Data

```
proc fastclus   data=datasets.pendigittest  
                maxclusters=10  
                out = clus;  
    var x1--x16;  
run;
```

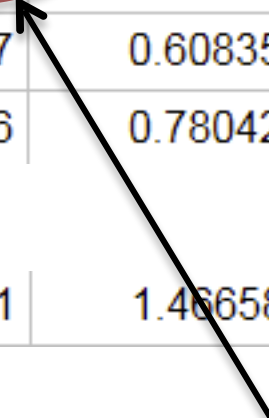
Example: PenDigit Data

Cluster Summary						
Cluster	Frequency	RMS Std Deviation	Maximum Distance from Seed to Observation	Radius Exceeded	Nearest Cluster	Distance Between Cluster Centroids
1	192	17.9838	131.7		7	108.9
2	880	20.9814	127.8		8	83.1740
3	115	16.3779	115.2		10	110.4
4	133	15.6579	118.4		5	101.5
5	575	19.9900	121.8		6	100.6
6	586	20.2352	129.9		5	100.6
7	217	15.5165	115.1		1	108.9
8	357	20.3639	124.4		2	83.1740
9	138	17.8981	119.2		2	100.3
10	305	15.4317	132.1		3	110.4

The first step to creating your own
hierarchical dendrogram.

Example: PenDigit Data

Statistics for Variables				
Variable	Total STD	Within STD	R-Square	RSQ/(1-RSQ)
X1	35.99529	23.35133	0.580229	1.382250
X2	14.73708	11.63537	0.378247	0.608354
X3	26.48636	19.87561	0.438336	0.780422
OVER-ALL	30.24017	19.27952	0.594581	1.466581

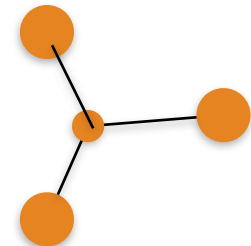
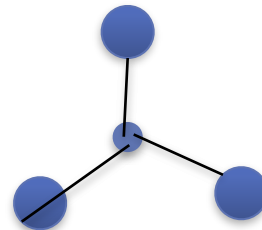


```
proc glm data= clus;  
class cluster;  
model x1 = cluster;  
run; quit;
```


Example: PenDigit Data

Statistics for Variables				
Variable	Total STD	Within STD	R-Square	RSQ/(1-RSQ)
X1	35.99529	23.35133	0.580229	1.382250
X2	14.73708	11.63537	0.378247	0.608354
X3	26.48636	19.87561	0.438336	0.780422
OVER-ALL	30.24017	19.27952	0.594581	1.466581

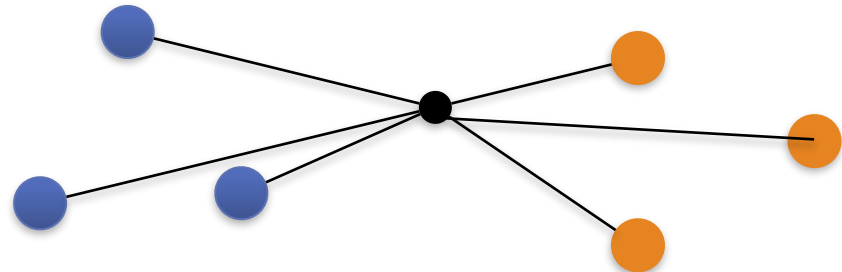
$$\frac{\sum_{C_j} \sum_{x_i \in C_j} \|\mathbf{x}_i - \mathbf{c}_j\|^2}{\sum_{x_k} \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2}$$



Example: PenDigit Data

Statistics for Variables				
Variable	Total STD	Within STD	R-Square	RSQ/(1-RSQ)
X1	35.99529	23.35133	0.580229	1.382250
X2	14.73708	11.63537	0.378247	0.608354
X3	26.48636	19.87561	0.438336	0.780422
OVER-ALL	30.24017	19.27952	0.594581	1.466581

$$\frac{\sum_{C_j} \sum_{x_i \in C_j} \|\mathbf{x}_i - \mathbf{c}_j\|^2}{\sum_{x_k} \|\mathbf{x}_k - \bar{\mathbf{x}}\|^2}$$



Example: PenDigit Data

Statistics for Variables				
Variable	Total STD	Within STD	R-Square	RSQ/(1-RSQ)
X1	35.99529	23.35133	0.580229	1.382250
X2	14.73708	11.63537	0.378247	0.608354
X3	26.48636	19.87561	0.438336	0.780422
OVER-ALL	30.24017	19.27952	0.594581	1.466581

Essentially the % variance explained using the centroids as prototypes for each cluster. (Similar to PCA % variance explained!)

