

# **Clustering Astronomy Data**

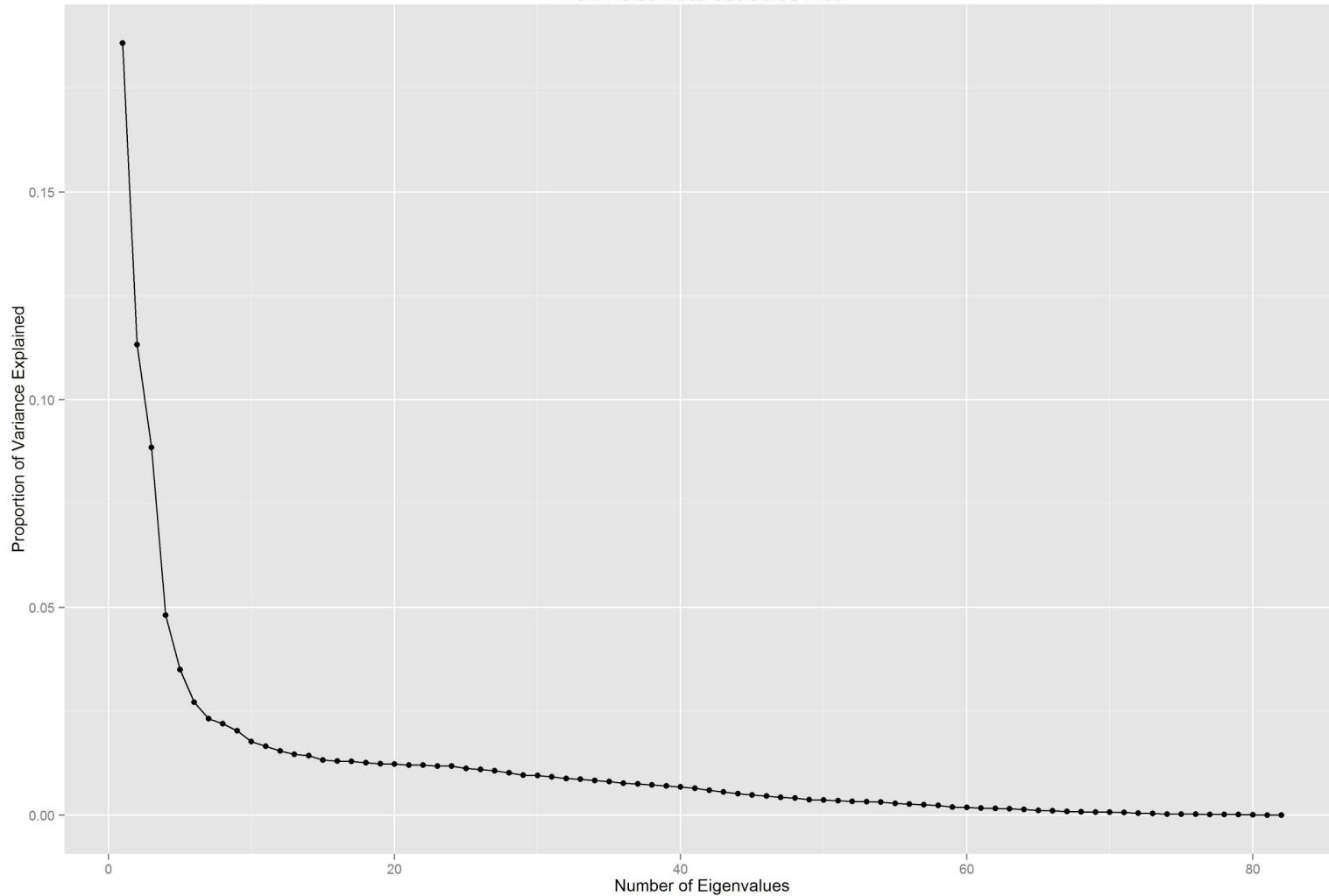
Katherine Eng, Jie Hu, Vanessa Lepe, Kalbi  
Zongo

# QUESTIONS OF INTEREST

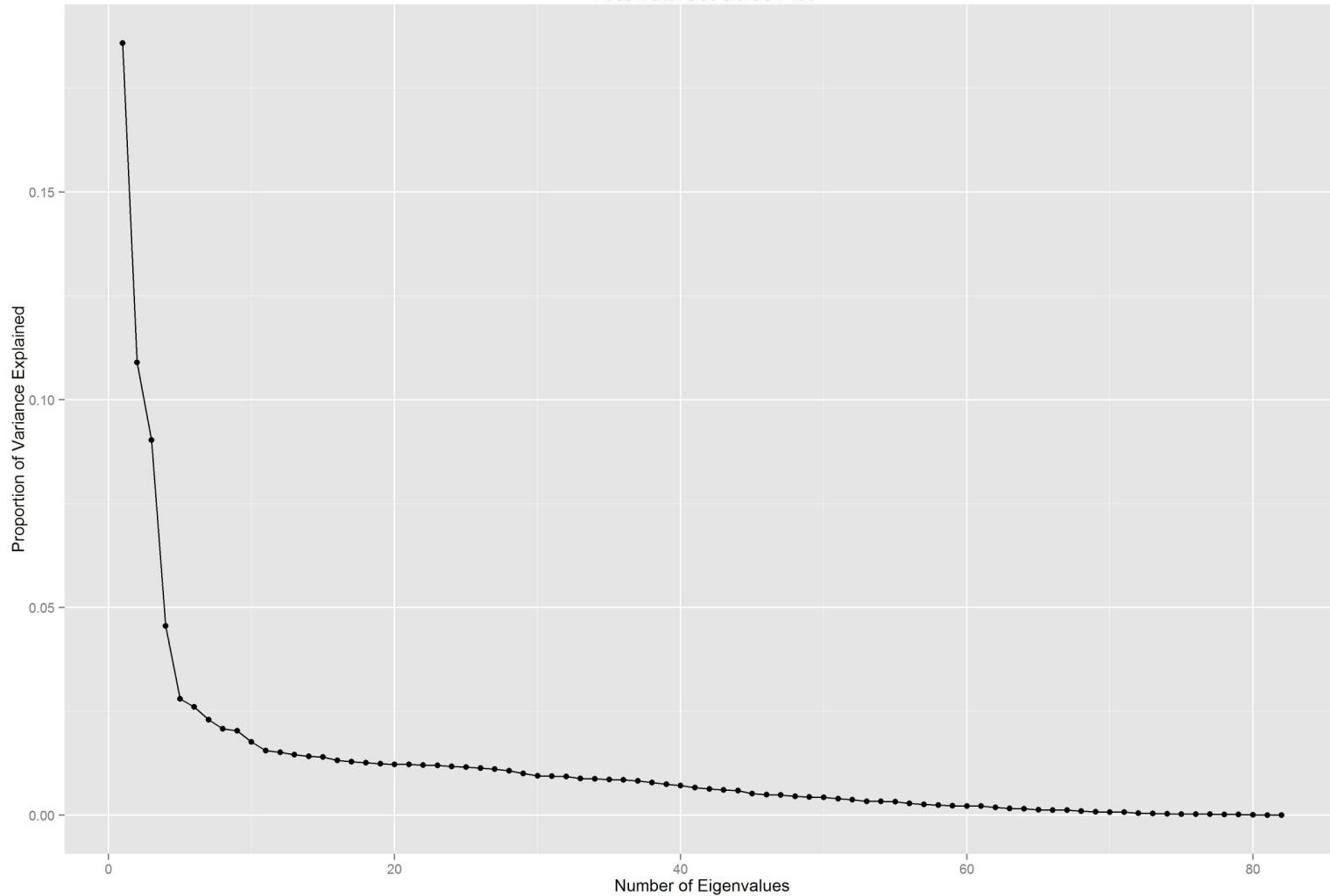
1. How does the principal component analysis (PCA) for low noise data compare to the PCA for the full test data set?
2. How do clustering methods (k-means vs hierarchical) differ in the low noise data set?
3. How do clustering methods (k-means vs hierarchical) differ in the test data set?

# **RESULTS**

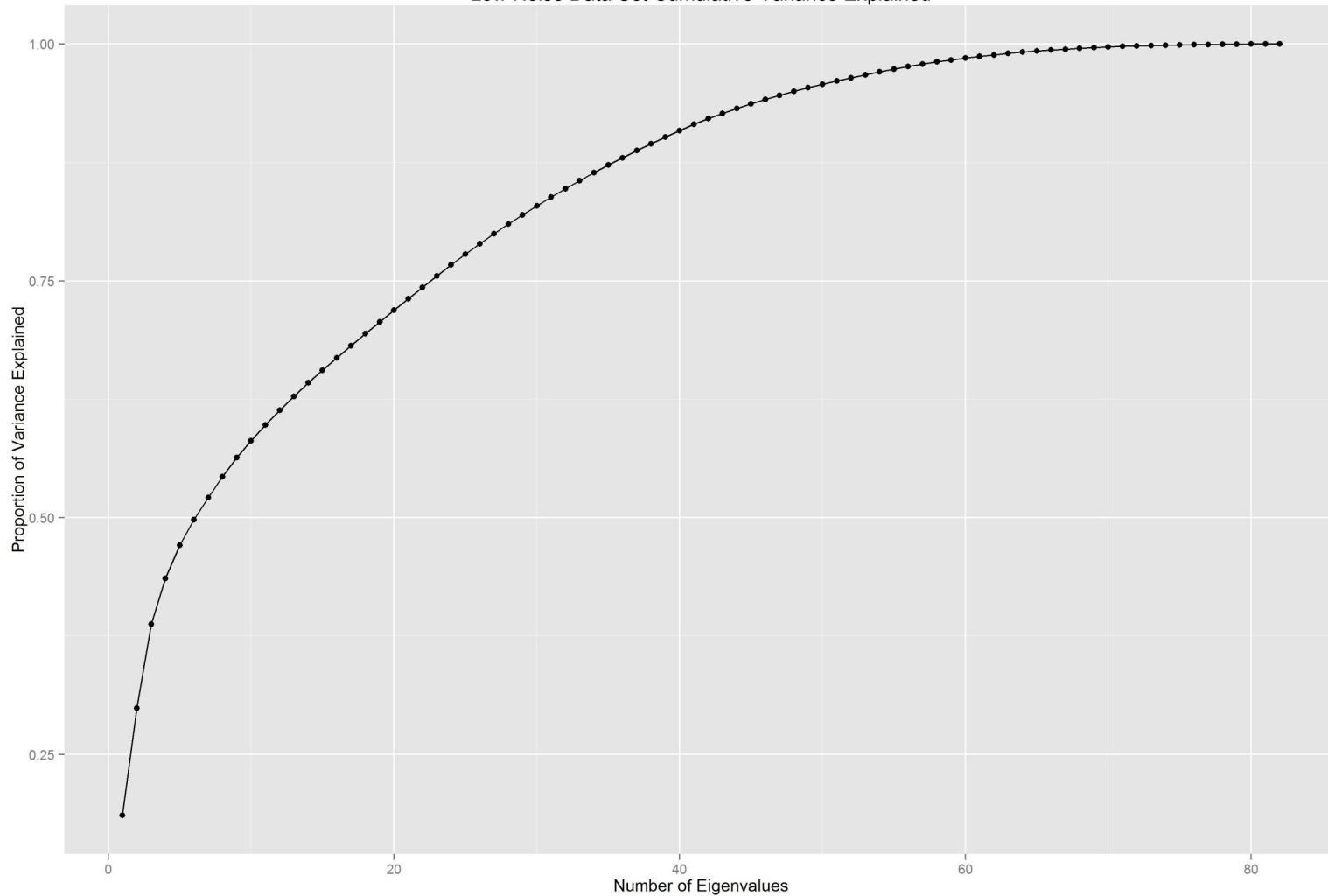
Low Noise Data Set Scree Plot



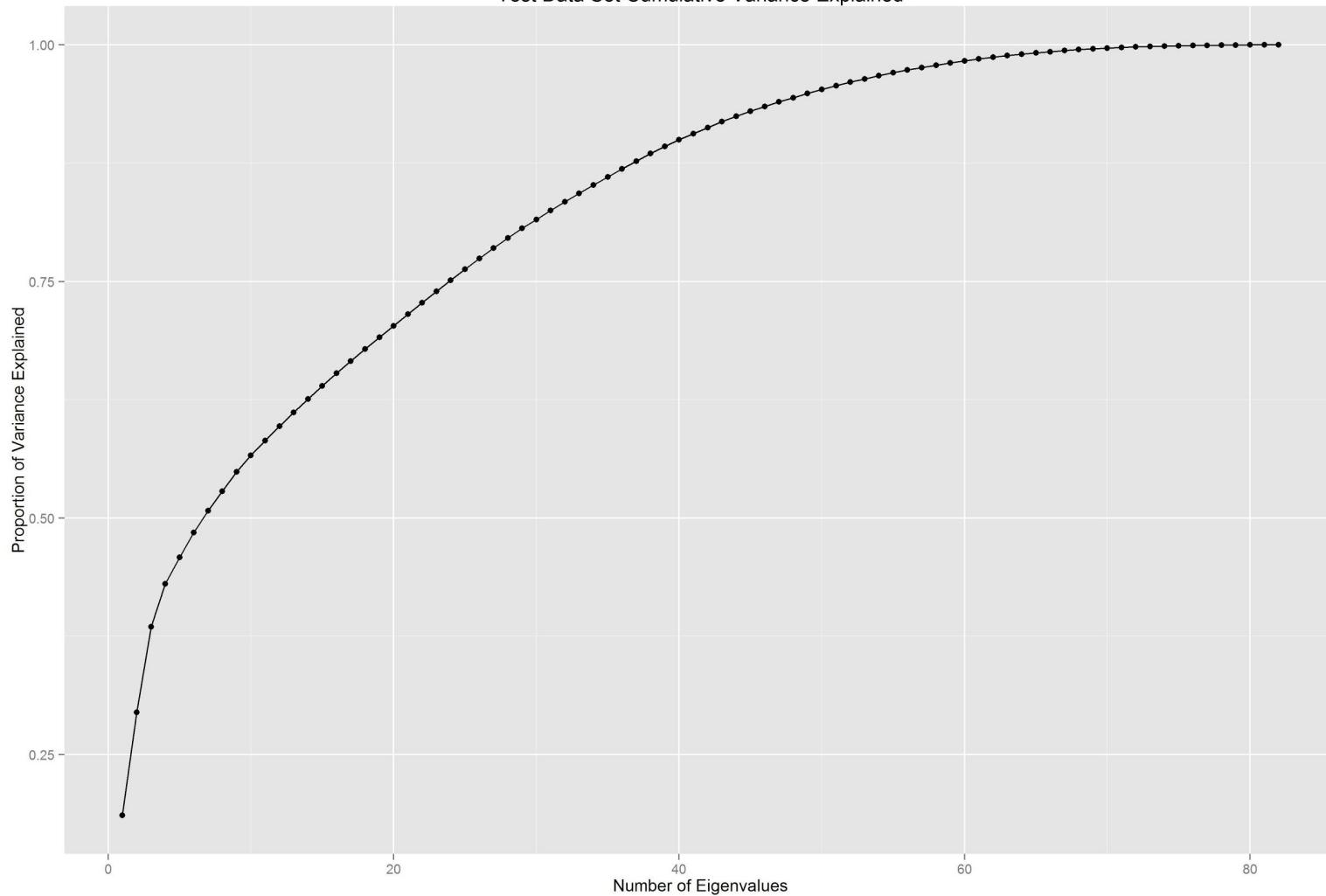
Test Data Set Scree Plot



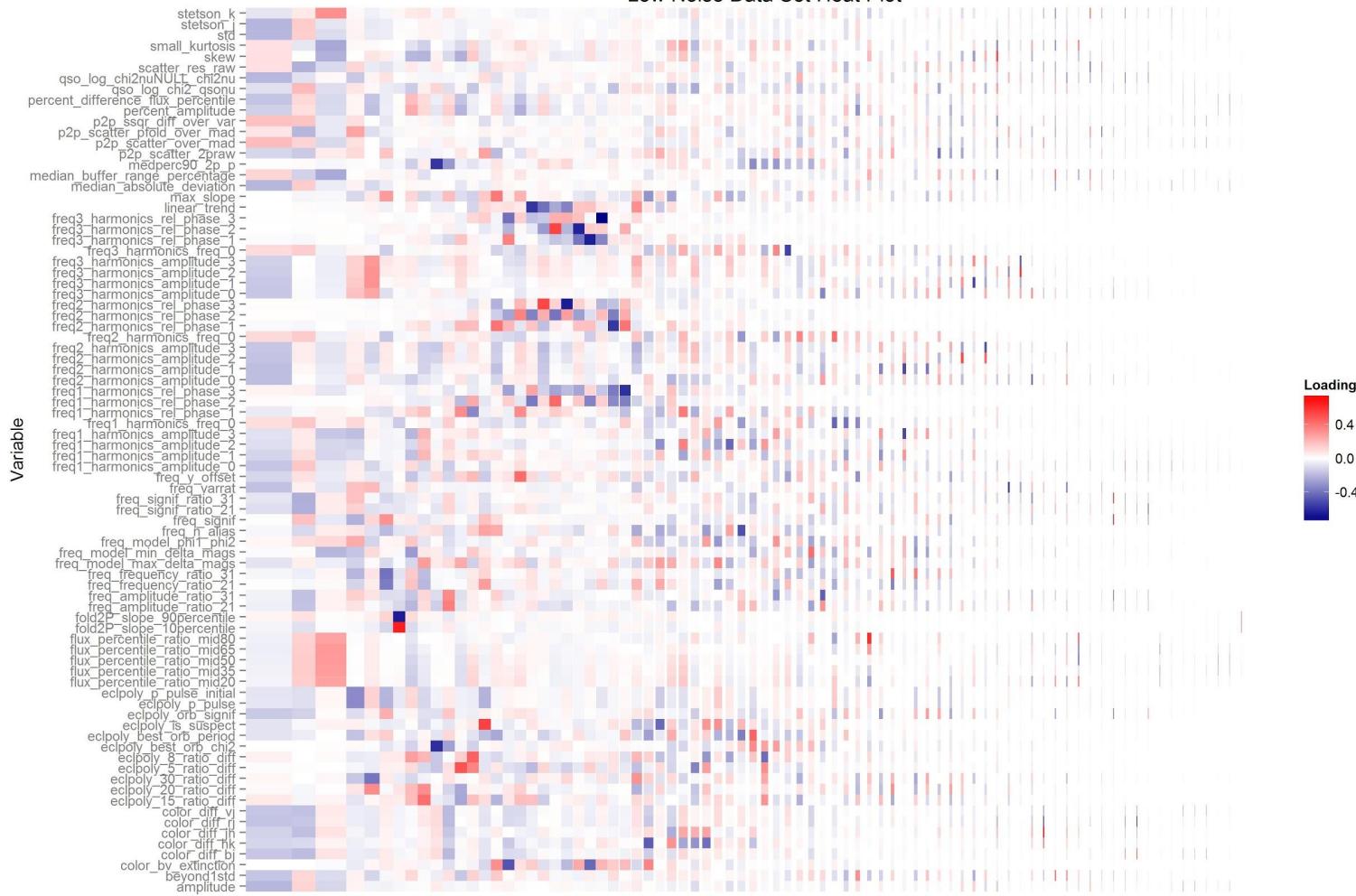
Low Noise Data Set Cumulative Variance Explained



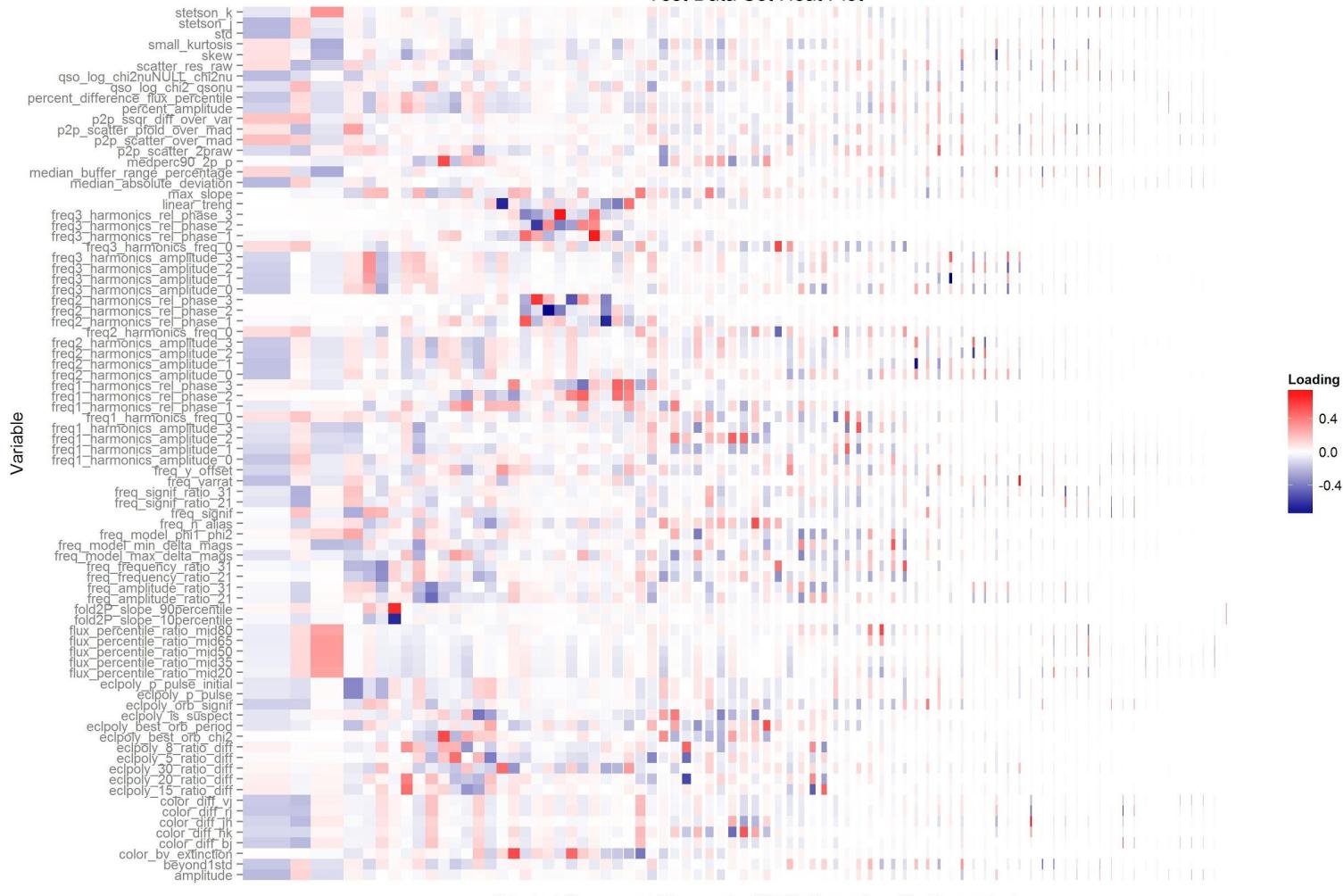
Test Data Set Cumulative Variance Explained



Low Noise Data Set Heat Plot



Test Data Set Heat Plot



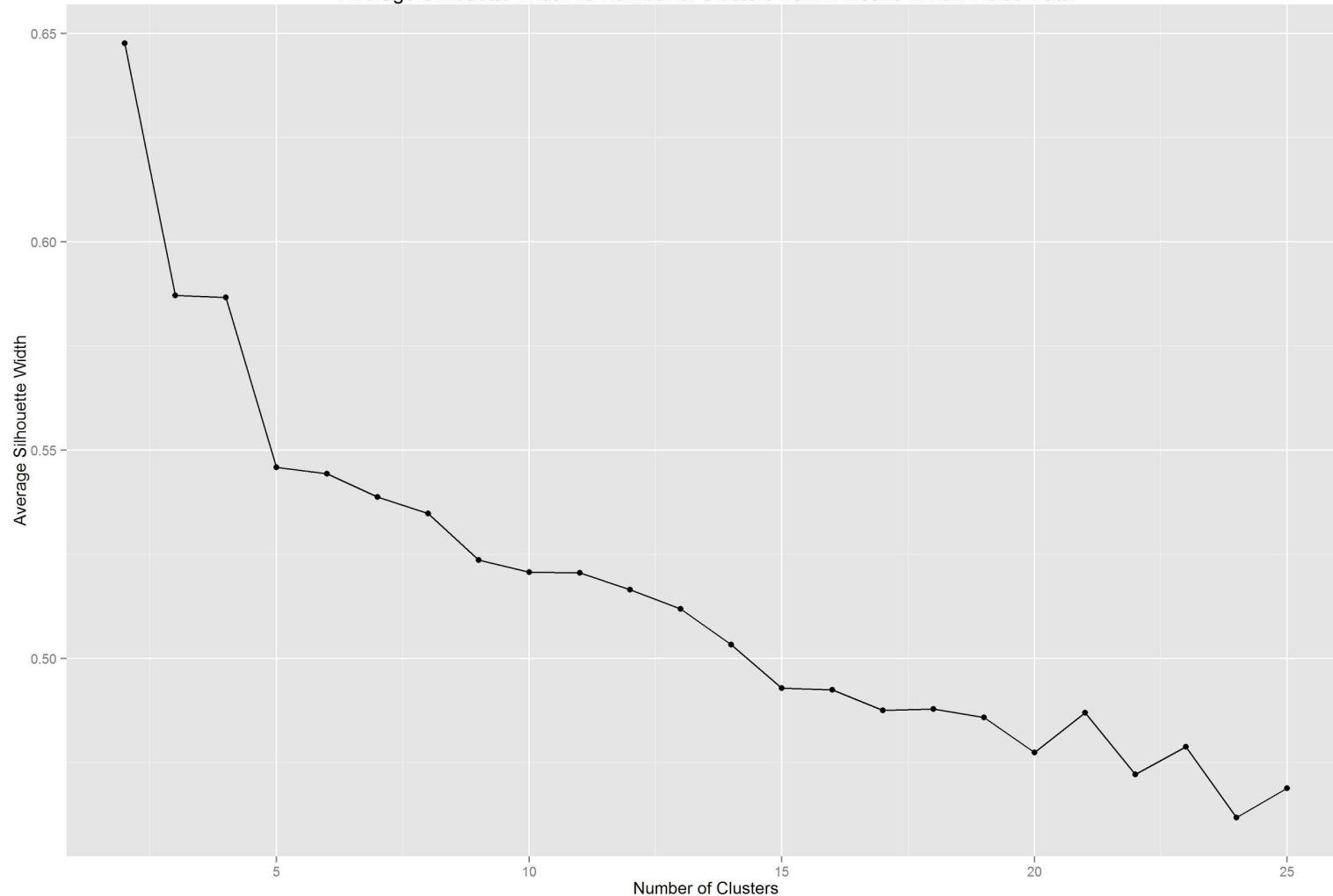
Principal Component Eigenvector (Width Proportional to Eigenvalue)

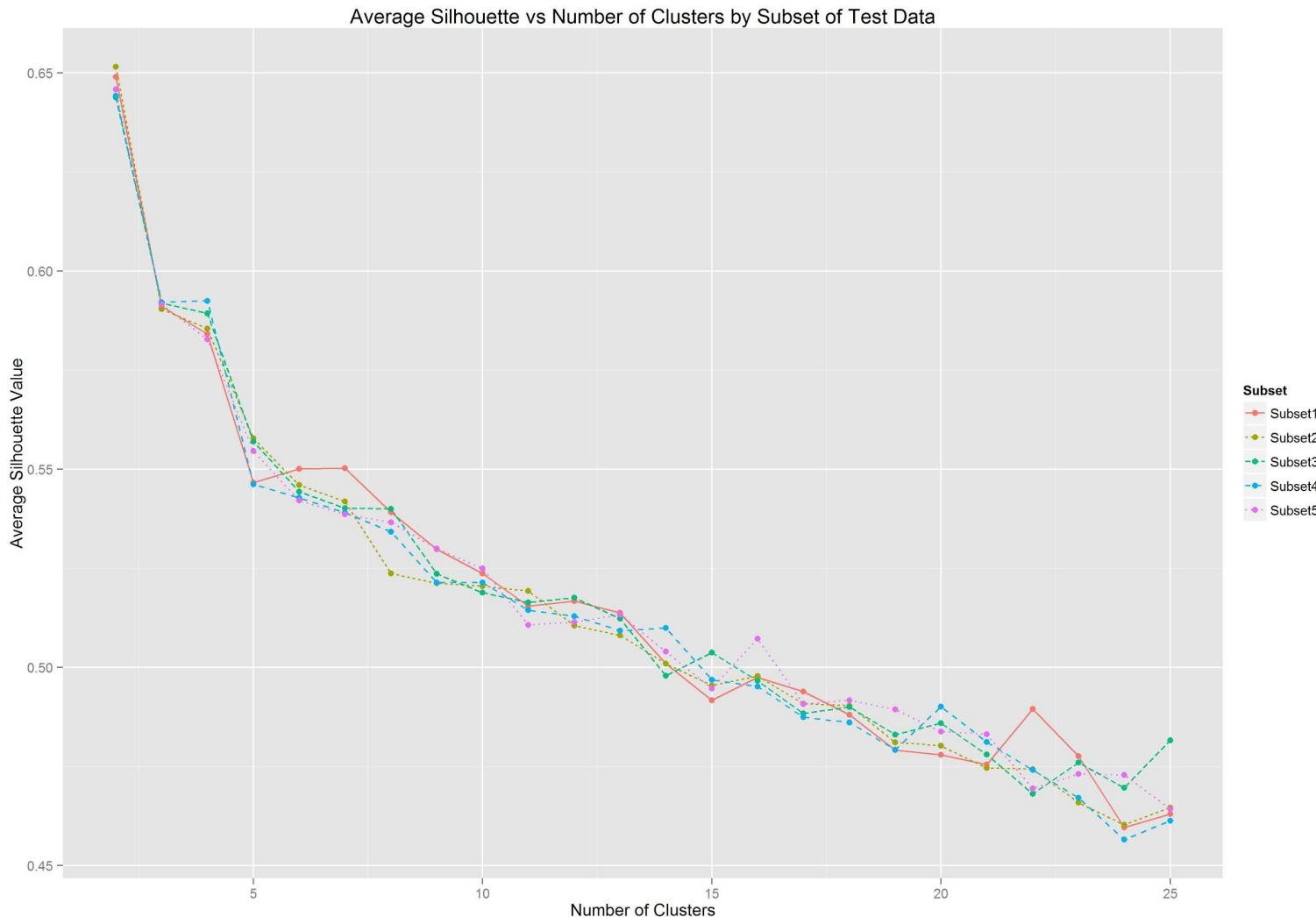
# DATA FOR CLUSTERING

Arbitrarily chose 28 features (both low noise and test data) to reduce dimensionality.

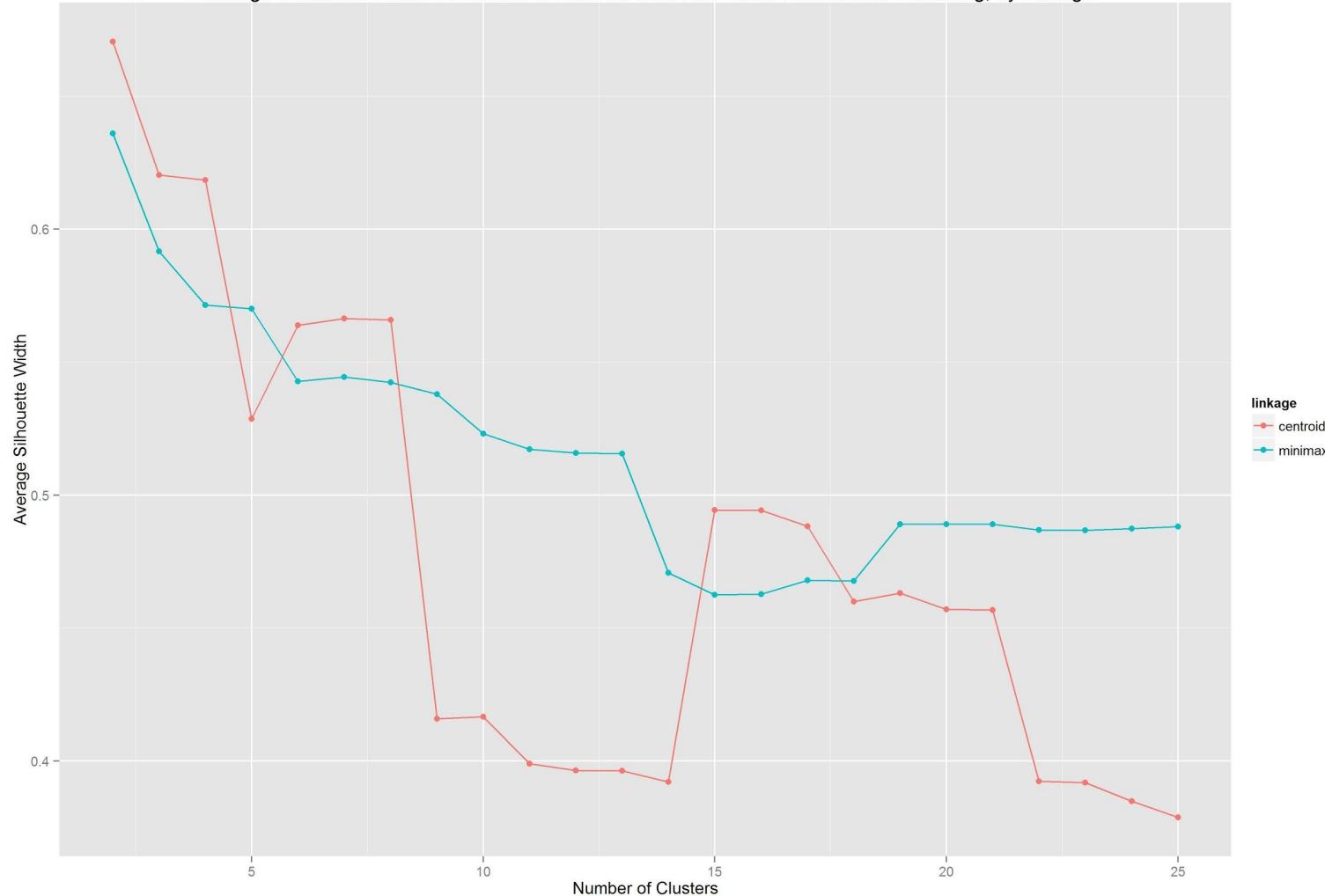
Test data set divided into 5 random subsets (approximately equal size) for ease of computation.

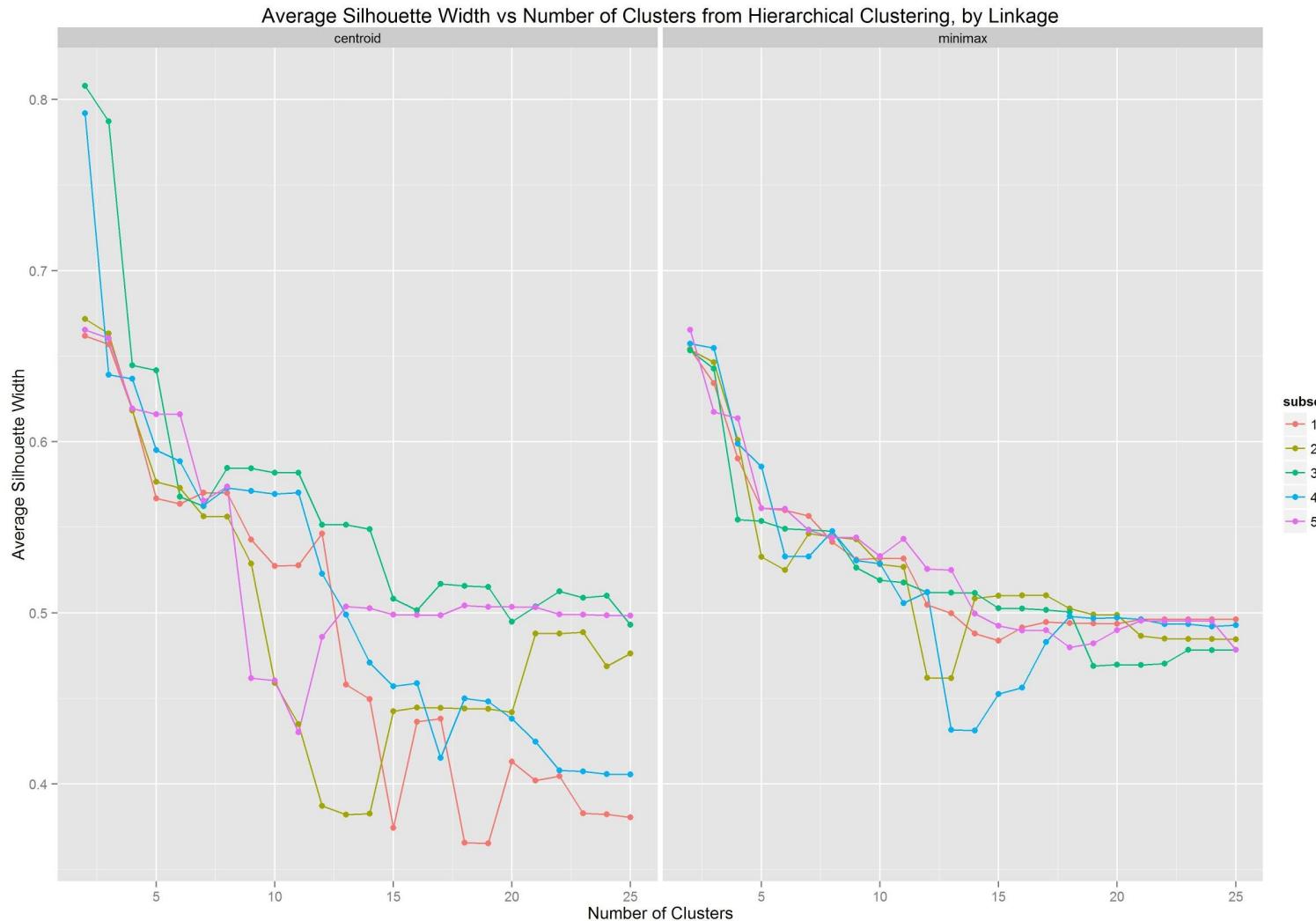
Average Silhouette Width vs Number of Clusters from k-means in Low Noise Data





Average Silhouette Width vs Number of Clusters from Low Noise Hierarchical Clustering, by Linkage





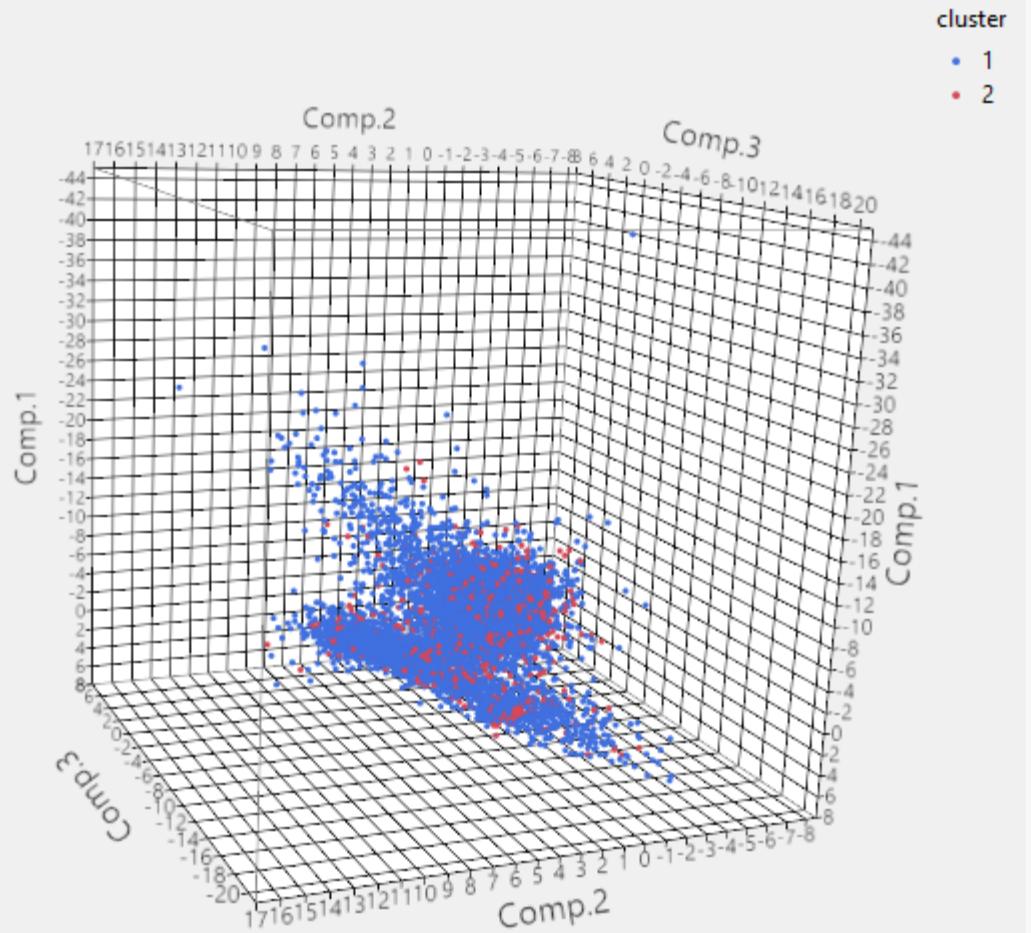
# OBSTACLES

- 1) Many observations had at least one missing feature
- 2) Some variables were constant across all observations
- 3) Clustering on all observations and features was computationally expensive
- 4) Visualizing clusters difficult (PCA, dendograms)

# FUTURE QUESTIONS/WORK

- 1) Other clustering methods (e.g., model-based)
- 2) Impute missing values
- 3) Find better ways to visualize data (e.g., use 3-d plots to display clustering forms)
- 4) Use COS Computing Cluster

*End*



Low Noise data set

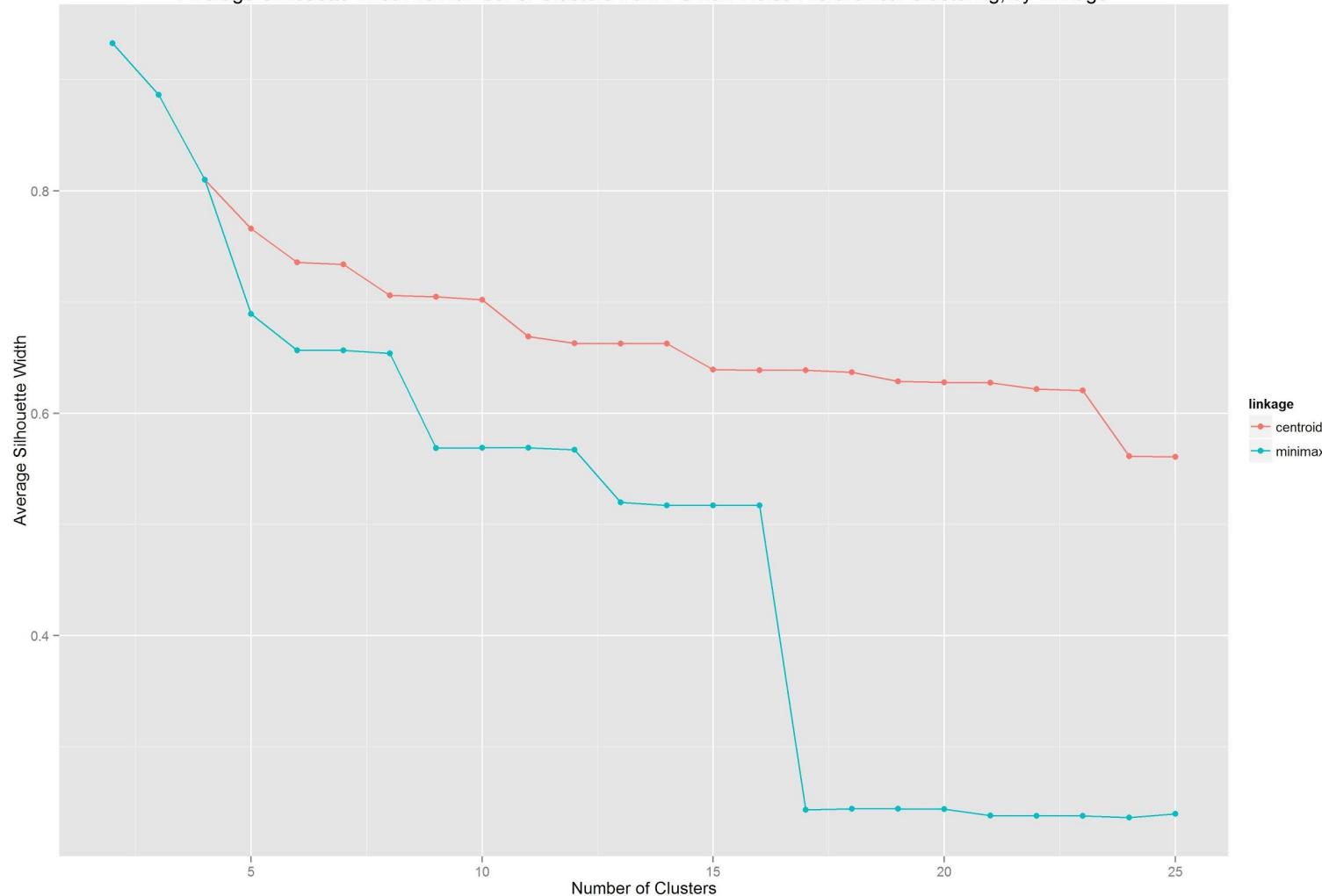
First three principal components

$k = 2$

cluster assignment from hierarchical clustering (28 features), minimax linkage

Cluster	1	2
Count	8759	1108

Average Silhouette Width vs Number of Clusters from PC Low Noise Hierarchical Clustering, by Linkage



# **Clustering Astronomy Data**

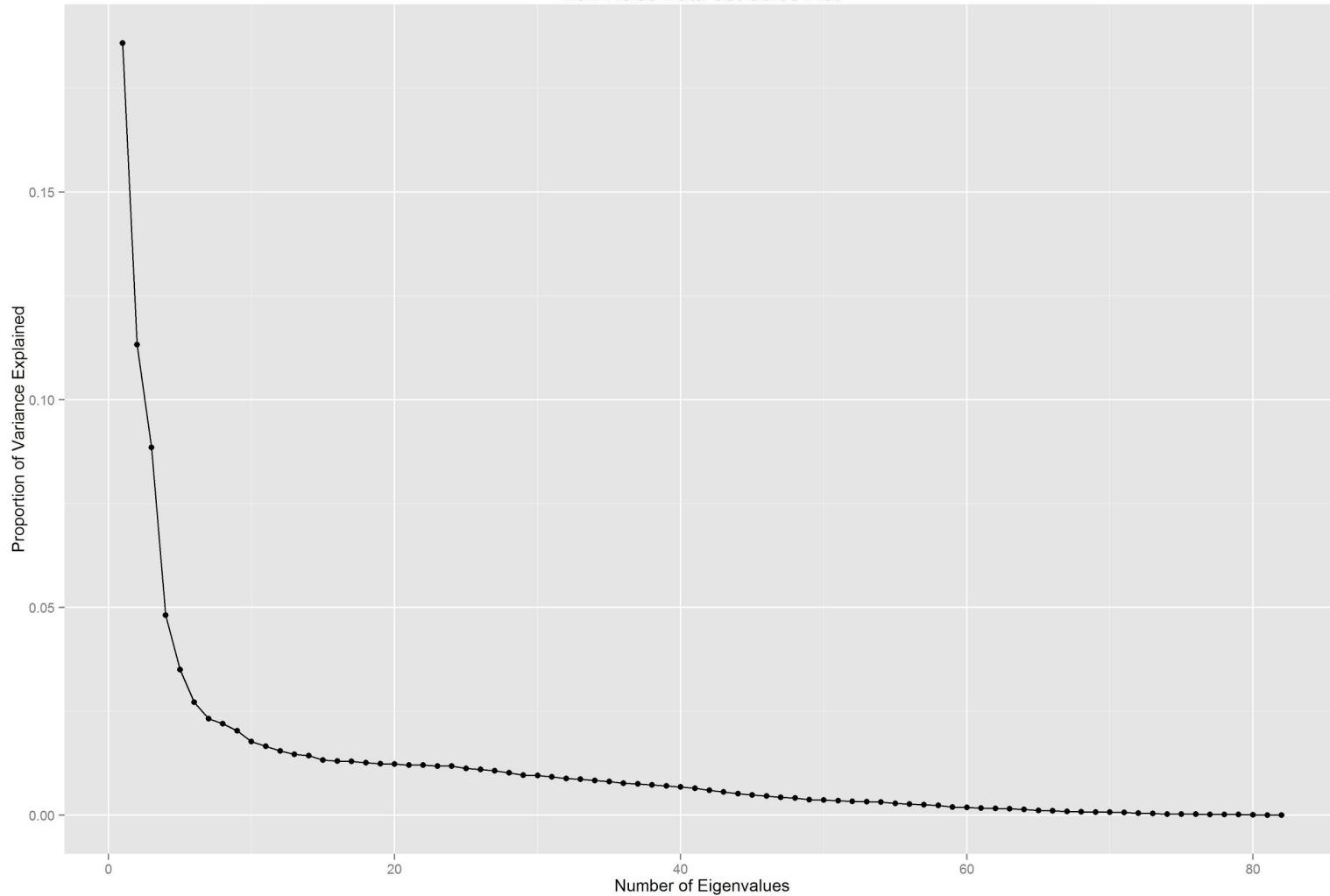
Katherine Eng, Jie Hu, Vanessa Lepe, Kalbi  
Zongo

# QUESTIONS OF INTEREST

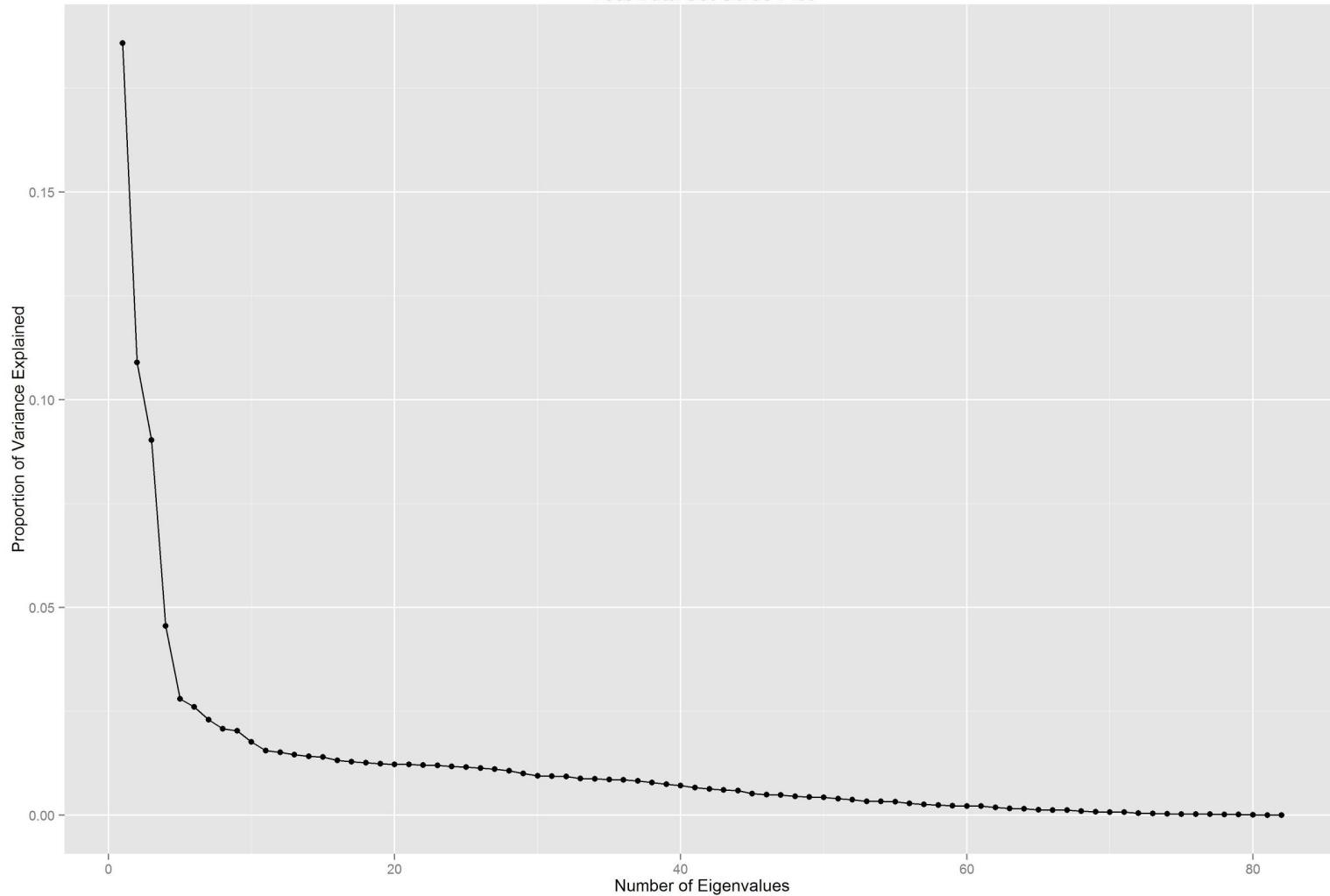
1. How does the principal component analysis (PCA) for low noise data compare to the PCA for the full test data set?
2. How do clustering methods (k-means vs hierarchical) differ in the low noise data set?
3. How do clustering methods (k-means vs hierarchical) differ in the test data set?

# **RESULTS**

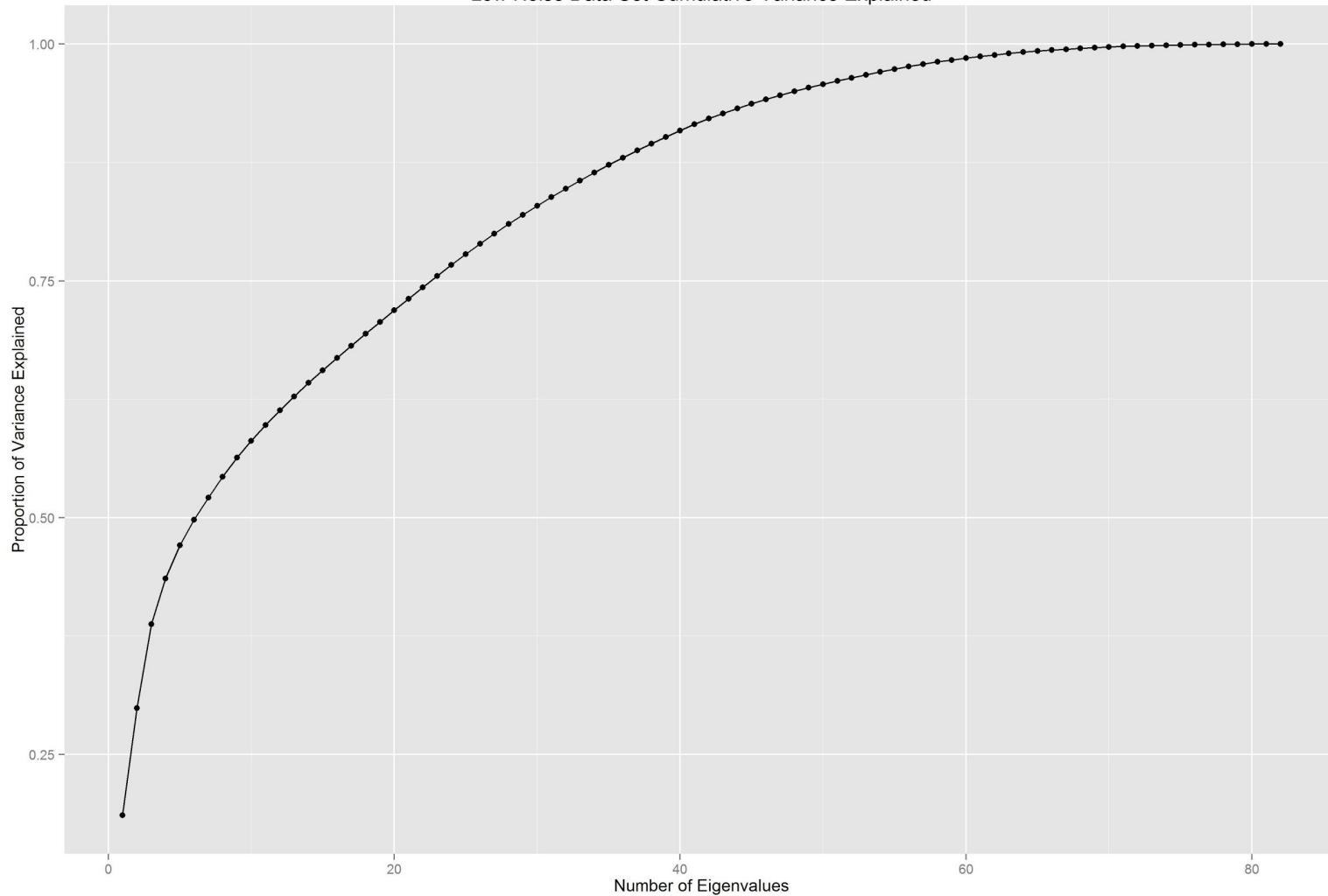
Low Noise Data Set Scree Plot



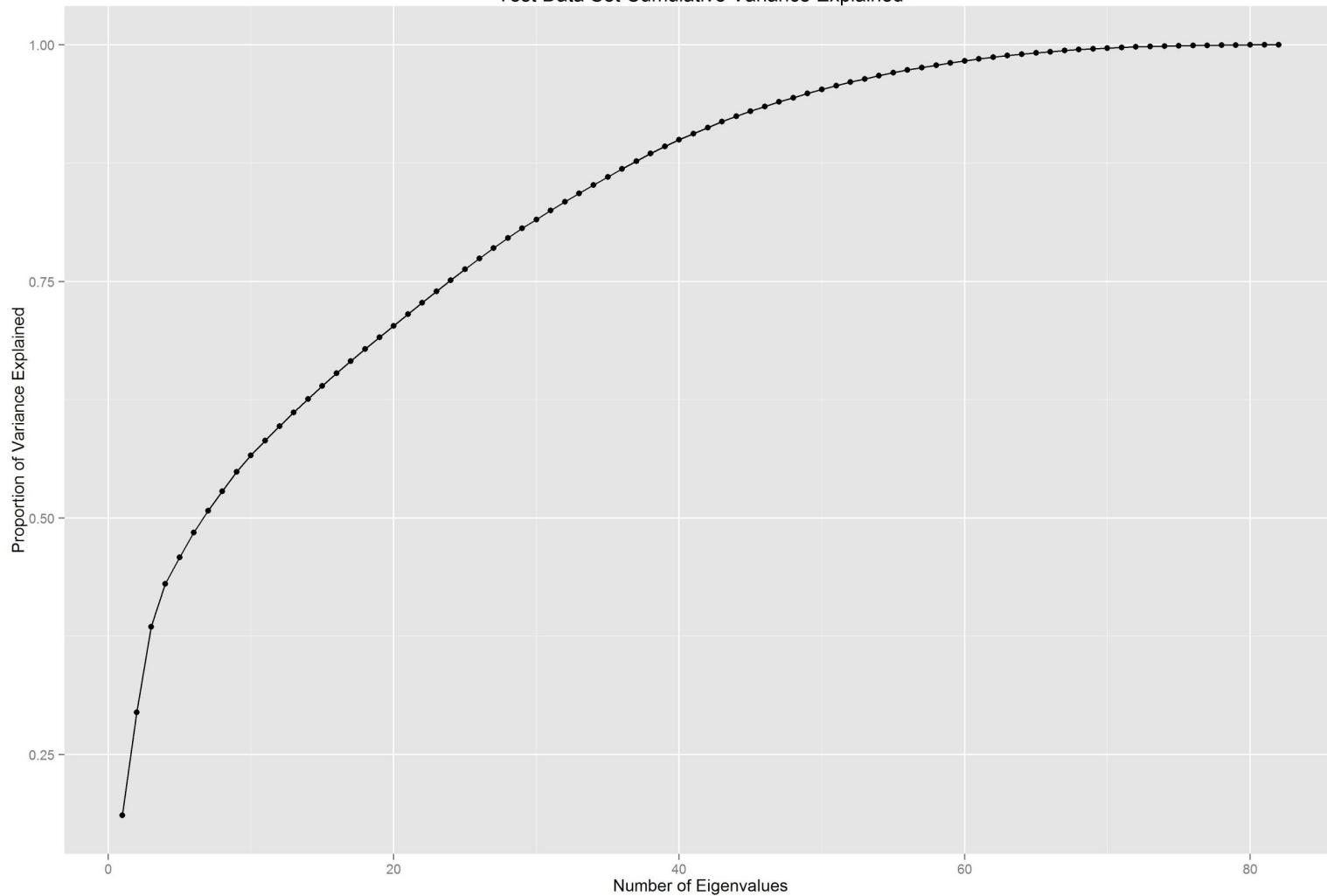
Test Data Set Scree Plot



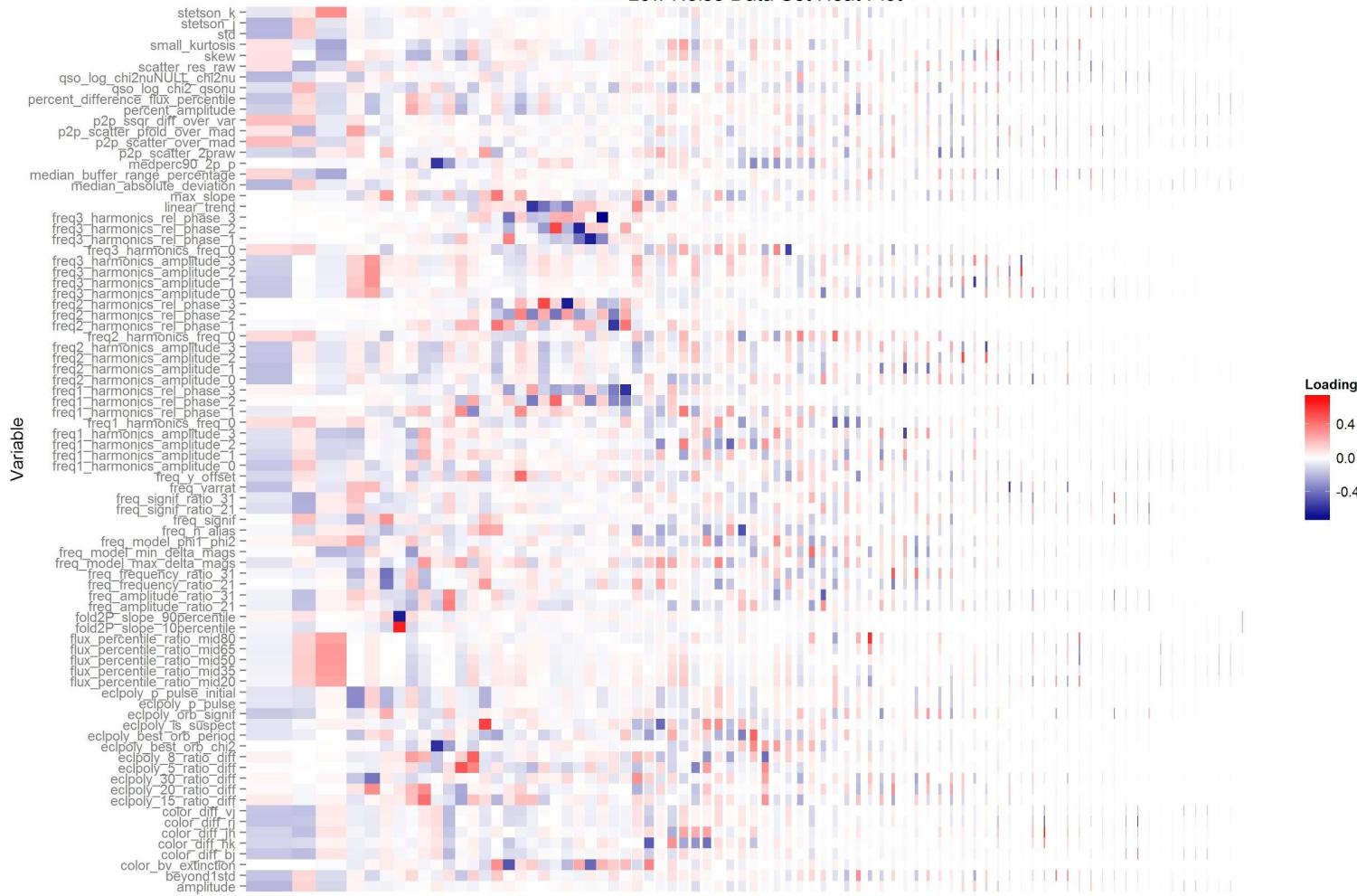
Low Noise Data Set Cumulative Variance Explained



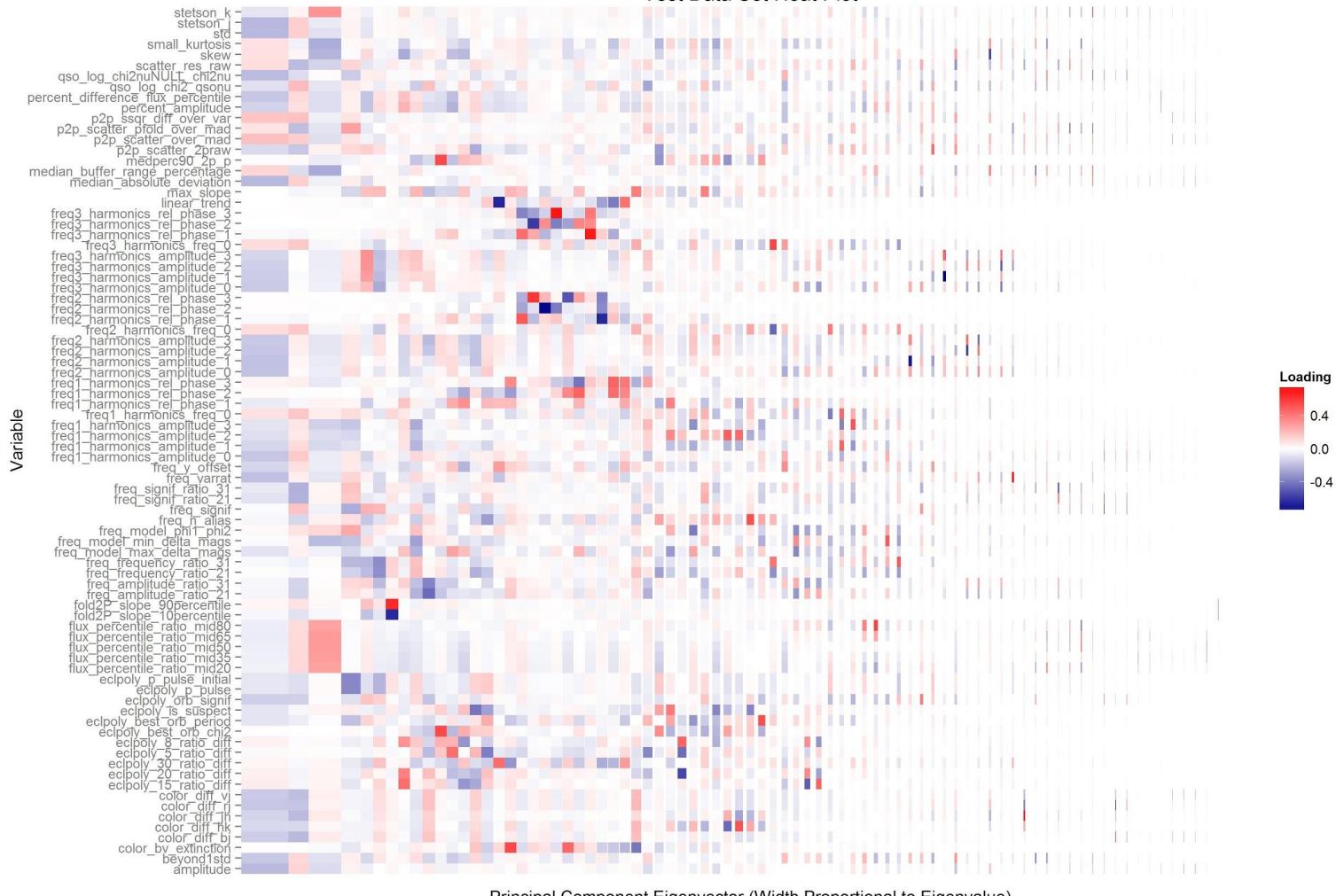
Test Data Set Cumulative Variance Explained



Low Noise Data Set Heat Plot



Test Data Set Heat Plot



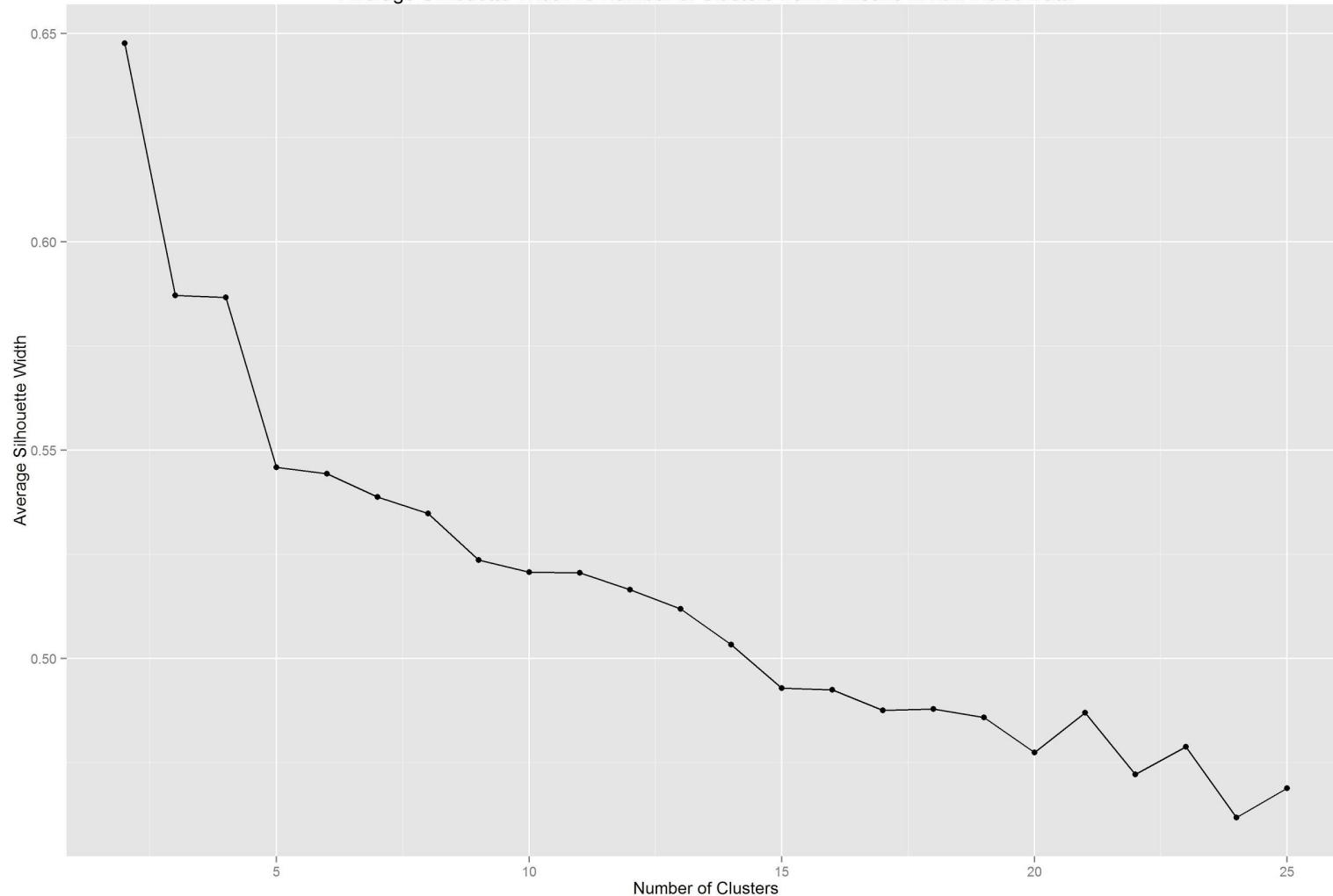
Principal Component Eigenvector (Width Proportional to Eigenvalue)

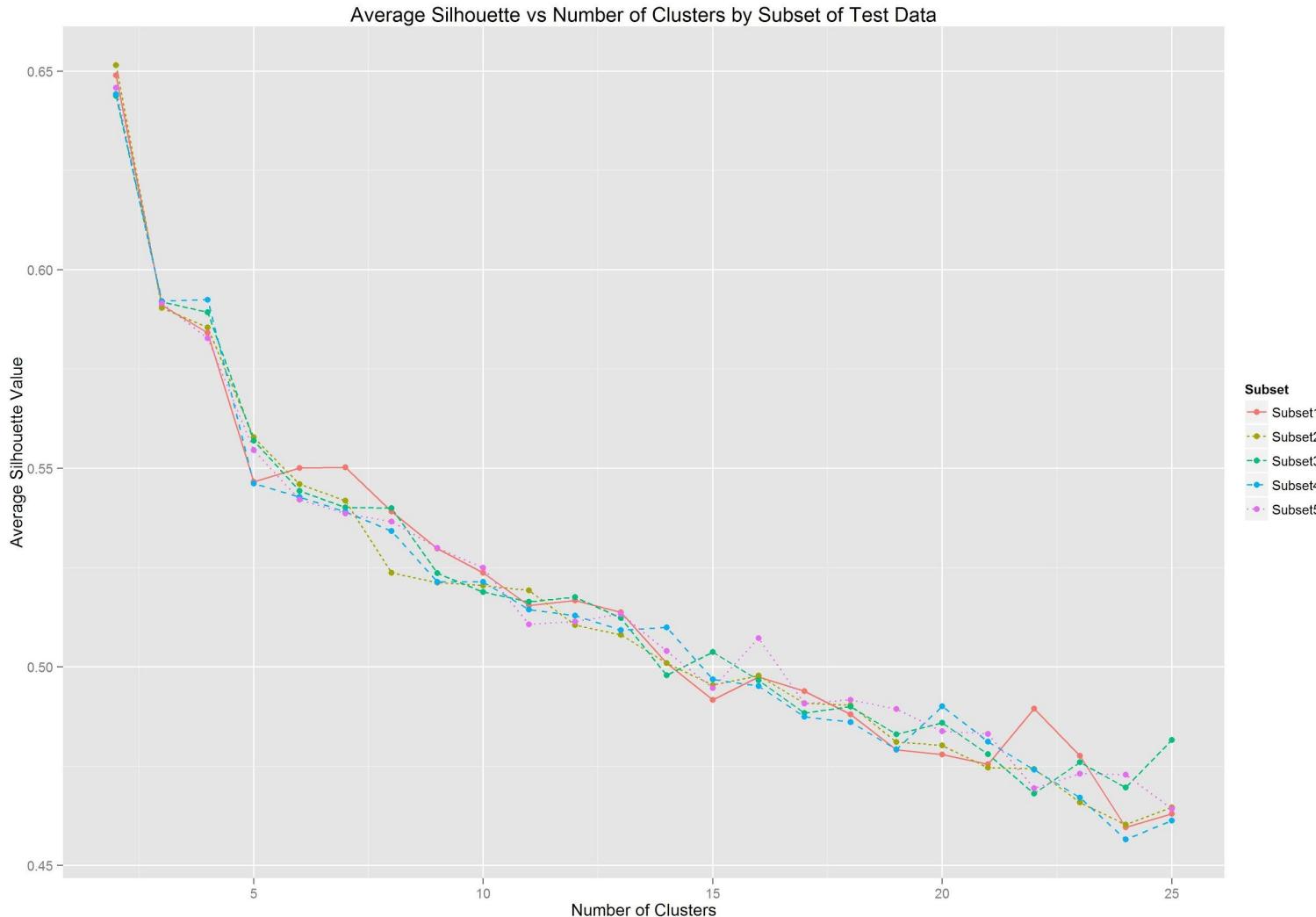
# DATA FOR CLUSTERING

Arbitrarily chose 28 features (both low noise and test data) to reduce dimensionality.

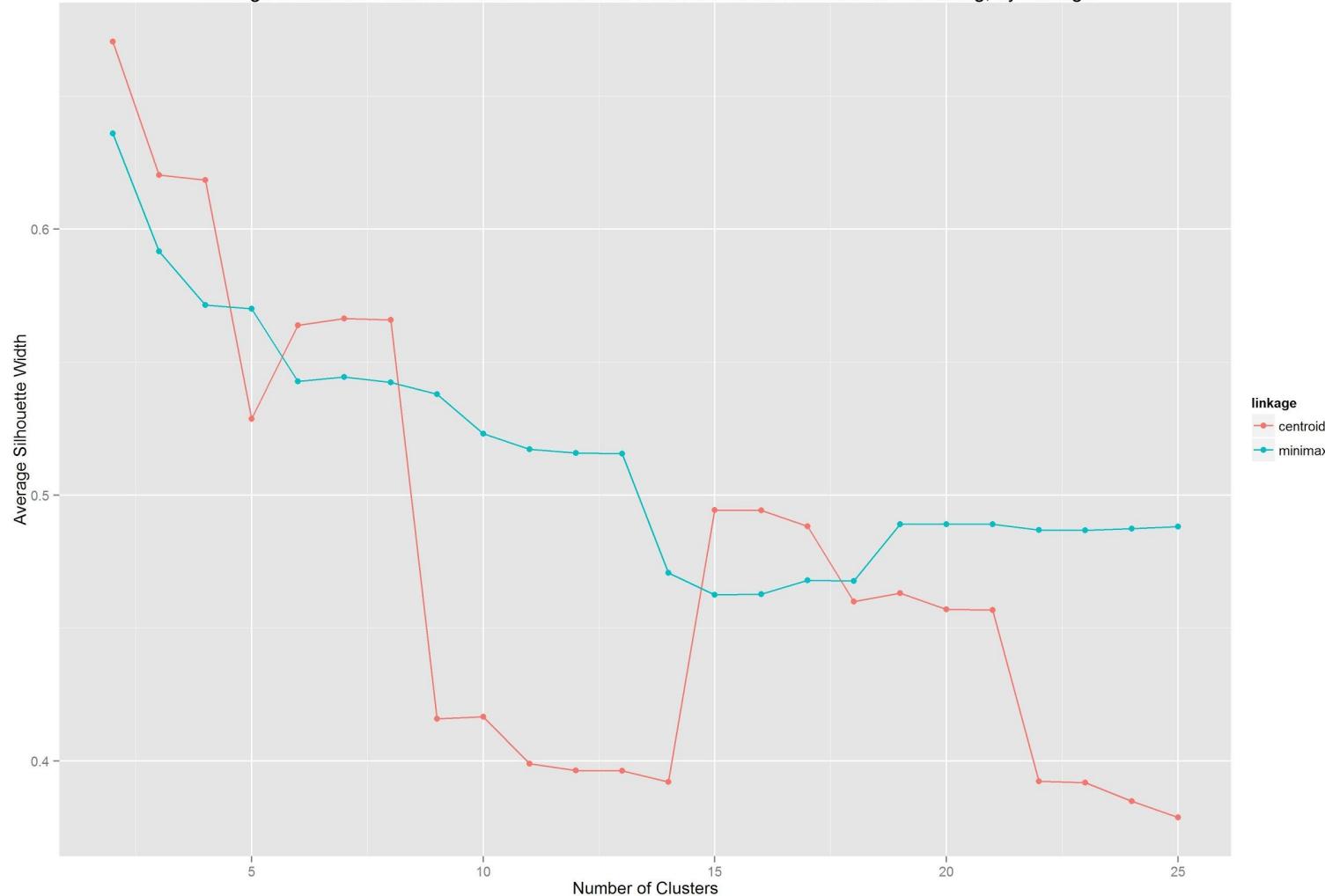
Test data set divided into 5 random subsets (approximately equal size) for ease of computation.

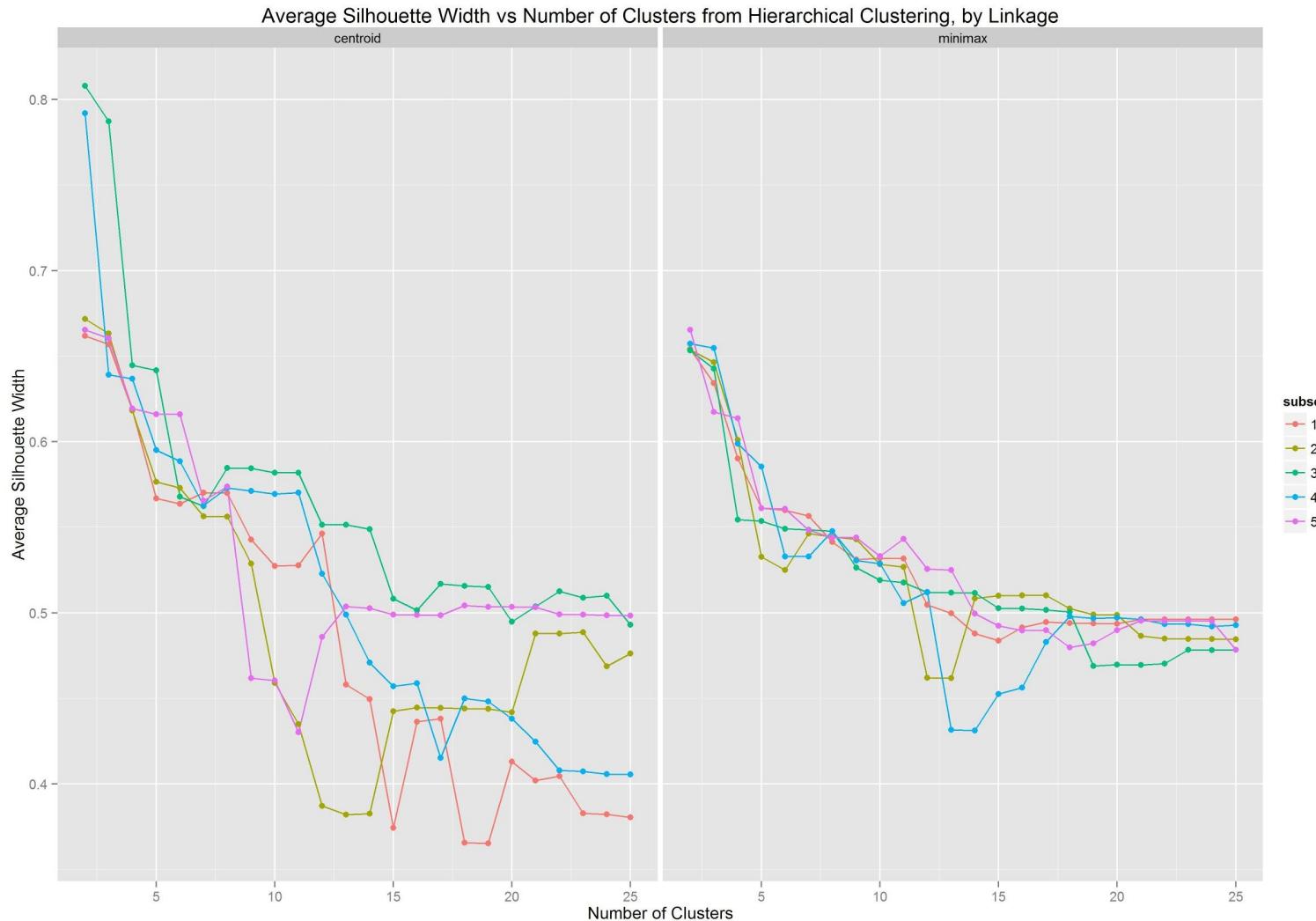
Average Silhouette Width vs Number of Clusters from k-means in Low Noise Data





Average Silhouette Width vs Number of Clusters from Low Noise Hierarchical Clustering, by Linkage





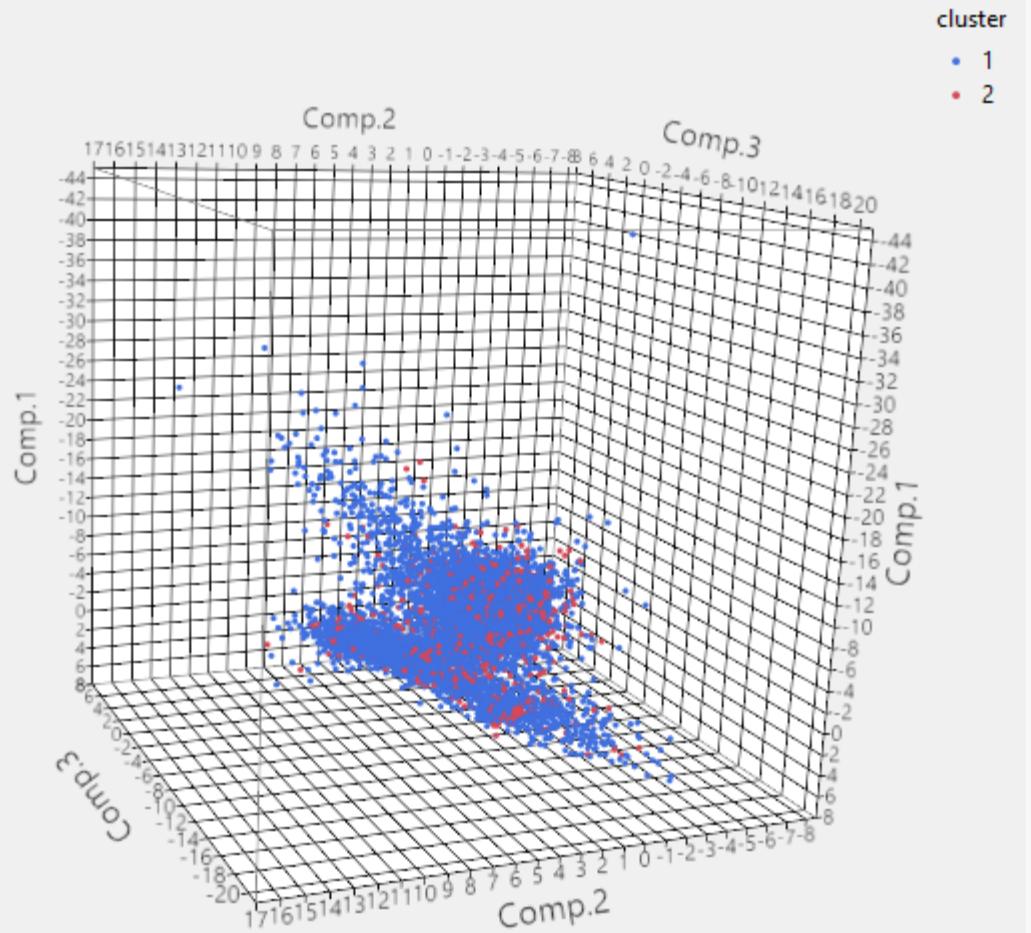
# OBSTACLES

- 1) Many observations had at least one missing feature
- 2) Some variables were constant across all observations
- 3) Clustering on all observations and features was computationally expensive
- 4) Visualizing clusters difficult (PCA, dendograms)

# FUTURE QUESTIONS/WORK

- 1) Other clustering methods (e.g., model-based)
- 2) Impute missing values
- 3) Find better ways to visualize data (e.g., use 3-d plots to display clustering forms)
- 4) Use COS Computing Cluster

*End*



Low Noise data set

First three principal components

$k = 2$

cluster assignment from hierarchical clustering (28 features), minimax linkage

Cluster	1	2
Count	8759	1108

Average Silhouette Width vs Number of Clusters from PC Low Noise Hierarchical Clustering, by Linkage

