

Clustering Pseudogenes



USING K-MEANS WITH PCA

by Hugo Y. K. Lam

Introduction



SECTION ONE

Pseudogene



- “false” genes , which look like real genes but have no apparent function
- First recognized and dubbed pseudogenes during the late 1970s, when early gene hunters began trying to pinpoint the chromosomal locations associated with production of important molecules

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogene



- For example, while seeking the gene responsible for making beta globin (a key component of the hemoglobin protein that transports oxygen through the bloodstream) scientists identified a DNA sequence that looked like a globin gene but could not possibly give rise to a protein.
- Essential functional parts of the gene’s anatomy were disabled by mutations, making it impossible for cellular machinery to translate the gene into a useful molecule.

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogenome

- Pseudogenes are the molecular remains of broken genes, which are unable to function because of lethal injury to their structures.
- The great majority of pseudogenes are damaged copies of working genes and serve as genetic fossils that offer insight into gene evolution and genome dynamics.
- Identifying pseudogenes involves intensive data mining to locate gene-like sequences and analysis to establish whether they function.
- Recent evidence of activity among pseudogenes, and their potential resurrection, suggests some are not entirely dead after all

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Birth of Pseudogene

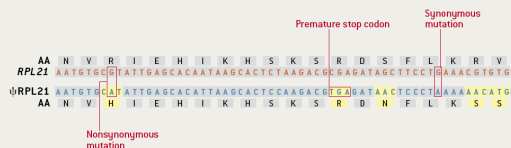
PSEUDOGENE BIRTH AND GENE DEATH

Two distinct processes can duplicate genes, and together they allow genomes to grow and diversify over evolutionary time. If errors in a copy destroy its ability to function as a gene, however, it becomes a pseudogene instead (*right*). The mutations that can kill a gene (*below*) range from gross deletions (such as the loss of the promoter region that signals the start of a gene sequence) to minute changes in the DNA sequence that skew the meaning of the gene's protein-encoding segments, called exons.

GENE DEATH

Genes die and become pseudogenes when mutations generated during the gene-copying process or accumulated over time render them incapable of giving rise to a protein. Cellular machinery reads the DNA alphabet of nucleotide bases (abbreviated A, C, G, T) in three-base increments called codons, which name an amino acid building block in a protein sequence or encode "stop" signals indicating the end of a gene. Even single-base mutations in codons

can alter their amino acid meaning, and base deletions or insertions can affect neighboring codons by shifting the cellular machinery's reading frame. The alignment shown here of a partial sequence for a human gene (*RPL21*) against one of its pseudogene copies (*ψRPL21*), along with each codon's corresponding amino acid (AA), illustrates some of the disabling mutations typically found in pseudogenes.



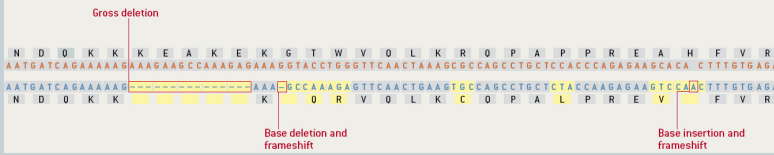
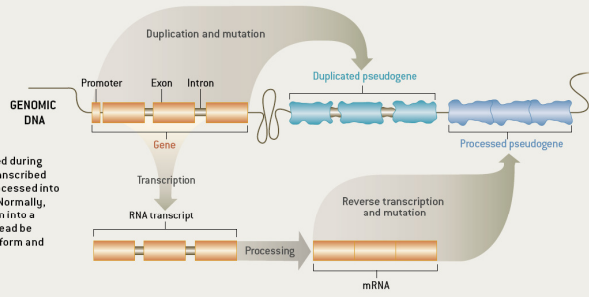
Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Birth of Pseudogene

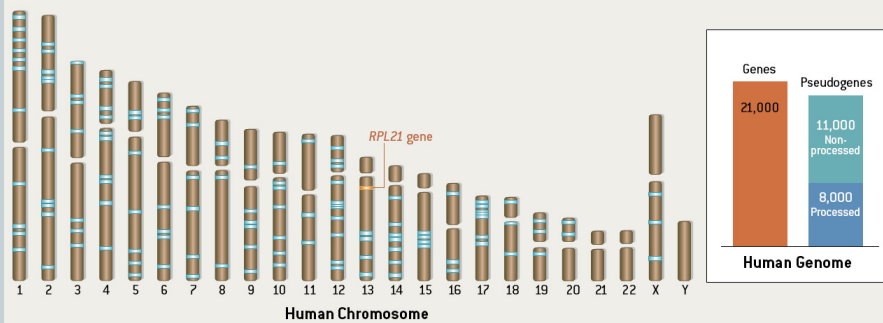
FLAWED COPIES

A "duplicated" pseudogene arises when a cell is replicating its own DNA and inserts an extra copy of a gene into the genome in a new location.

A "processed" pseudogene is formed during gene expression, when a gene is transcribed into RNA, then that transcript is processed into a shorter messenger RNA (mRNA). Normally, the mRNA is destined for translation into a protein—but sometimes it can instead be reverse-transcribed back into DNA form and inserted in the genome.



Human Pseudogenome



PSEUDOGENE DESCENDANTS (blue) of the ribosomal protein gene *RPL21* (orange) are scattered across the human chromosomal landscape. Overall distribution of pseudogenes in the human genome appears to be completely random, although some local genome regions tend to contain more pseudogenes. Those DNA regions may be analogous to certain geochemical environments that better

preserve mineral fossils. Identification of genes and pseudogenes is an ongoing process, but to date more than 19,000 pseudogenes have been identified in the human genome—only slightly less than the current tally of around 21,000 human genes (inset). About 8,000 of our pseudogenes are processed; the rest include duplicated pseudogenes and other nonprocessed subcategories.

Gerstein, M & Zheng, D. The real life of pseudogenes. *Sci Am* 295: 48-55 (2006).

Pseudogenomics



- **Hints about Life Histories**

- Often, genes involved in an organism's response to its environment are subject to extensive duplication and diversification over time, leading to large gene families, such as the olfactory receptor repertoire.
- Many dead-on-arrival pseudogene copies are an immediate byproduct of this process.

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogenomics



- But the subsequent death of additional duplicates, which gives rise to new pseudogenes, is also frequently connected to changes in an organism's environment or its circumstances.
- Consequently, differences in the pseudogenes of animals offer hints about their diverse life histories that are not as easily detected in comparisons of working genes, which are strongly constrained by their function.

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogenomics



- **May not be totally dead genes**
 - A few pseudogenes appear to be better preserved than one would expect if their sequences were drifting neutrally.
 - Recent experiments by Thomas Gingeras of Affymetrix and by Michael Snyder of Yale University have found that a significant fraction of the intergenic regions in the human genome are actively transcribed.

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogenomics



- **Pseudogene conversion**
 - Some evidence also exists for the possibility of pseudogene resurrection — a dead gene turning back into a living one that makes a functional protein product.
 - Careful sequence comparisons have shown that one cow gene for a ribonuclease enzyme was a pseudogene for much of its history but appears to have been reactivated during recent evolutionary time

Gerstein, M & Zheng, D. The real life of pseudogenes. Sci Am 295: 48-55 (2006).

Pseudogenomics



- **Improves gene annotation**
 - Early efforts to catalogue pseudogenes were largely driven by the need to distinguish them from true genes when annotating genome sequences.
 - Identifying pseudogenes is not as straightforward as recognizing genes
 - Establishing a suspected pseudogene's inability to function is more challenging.

Gerstein, M & Zheng, D. The real life of pseudogenes. *Sci Am* 295: 48-55 (2006).

PseudoPipe



- **An automated pseudogene identification pipeline**
 - identify all the regions in the genome that share sequence similarity with any known protein, using BLAST
 - resolve the paternity ambiguity of the pseudogenes, i.e. determine among the paralogous query proteins which one most likely gave rise to the pseudogene.

Zhang, Z, Carriero, N, Zheng, D, Karro, J, Harrison, PM & Gerstein, M. PseudoPipe: an automated pseudogene identification pipeline. *Bioinformatics* 22: 1437-9 (2006).

Pseudogene Classification

- **Processed Pseudogenes**
 - Lack introns, have small flanking direct repeats and a 3' polyadenine tail
- **Duplicated Pseudogenes**
 - Distinguished from processed pseudogenes by a combination of these features, with the emphasis on the evidence of ancient introns
- **Pseudogene Fragments**
 - are protein/chromosome homologies that have high sequence similarity, but are too decayed to be reliably assessed as processed or duplicated (i.e. <70% coverage of the parent gene)

Zhang, Z, Carriero, N, Zheng, D, Karro, J, Harrison, PM & Gerstein, M. PseudoPipe: an automated pseudogene identification pipeline. *Bioinformatics* 22: 1437-9 (2006).

Pseudogene Data Set

SECTION TWO

Pseudogene Dataset

Pseudogene Classic Pipeline HS 42_36d Data

Query:

Available variables:

Chromosome	Parent chromosome	NumDels	Number of deletions
ChrStart	Start of gene	NumShifts	Number of shifts
ChrEnd	End of gene	NumStops	Number of stops
Strand	(Type?)	Expect	Expectation score
Parent	Gene parent	PctIdent	Percent identical
QueryStart	Start of matched parent region	PolyA	Poly A signal
QueryEnd	End of matched parent region	Disable	Disablement class
QueryLength	Length of parent matched	NumExons	Number of exons
Frac	Fraction of parent spanned	BasicClass	Basic gene classification
NumIns	Number of insertions	RefinedClass	Refined gene classification

Examples:

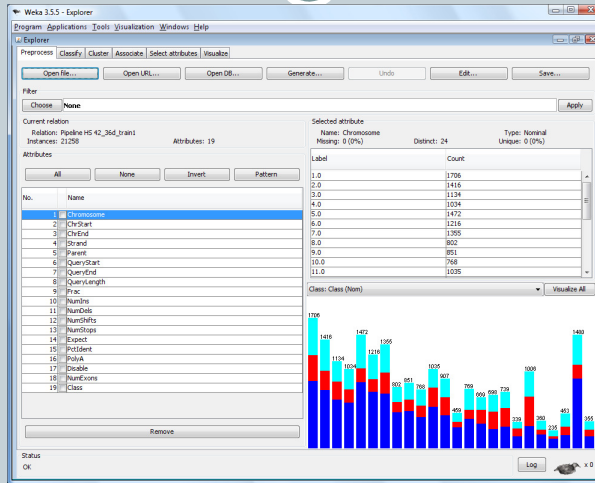
Query	Description
<code>*</code>	Returns all data, displays the first few entries
<code>BasicClass startsWith('DUP') or BasicClass startsWith('PSS')</code>	Returns just the pseudogenes (no fragments, false positives, or 'real' genes)
<code>BasicClass startsWith('DUP') or BasicClass startsWith('PSS') and Chromosome == '3' and 10000000 < ChrStart <= 20000000</code>	Returns the pseudogenes on chromosome 3 starting somewhere in the range from 10,000,000 to 20,000,000 bp.

[N.J. Carrara](#) [Help](#)

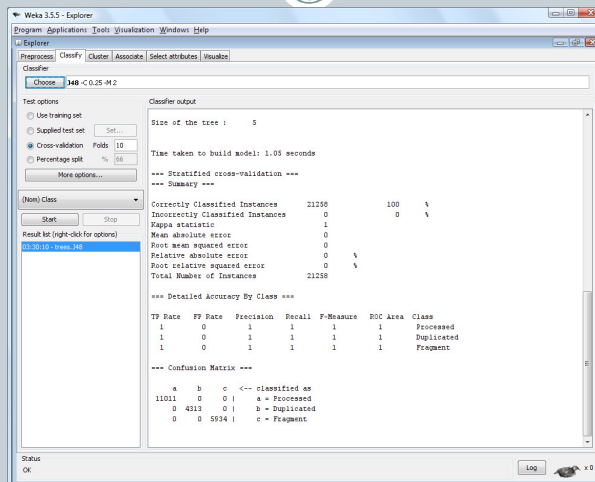
Pseudogene Dataset

Chromosome	ChrStart	ChrEnd	Strand	Parent	QueryStart	QueryEnd	QueryLength	Frac	NumIns	NumDels	NumShifts	NumStops	Expect	PctIdent	PolyA	Disable	NumExons	BasicClass
1	484	639	-	ENSP00000323932	5	56	194	2.700000e-01	0	0	0	0	1.000000e-08	7.120000e-01	0	0	1	FP
1	817	1367	+	ENSP00000371947	1	181	181	1.000000e+00	4	0	1	1	2.000000e-24	9.290000e-01	0	0	1	PSSDDO
1	2845	3533	+	ENSP00000228264	821	970	150	1.500000e-01	0	4	0	0	1.000000e-28	7.800000e-01	0	0	0	3 DUP
1	3887	4009	+	ENSP00000295199	1	41	41	1.000000e+00	0	0	0	0	4.000000e-17	9.760000e-01	0	0	0	1 GENE_SINGLE
1	4215	43196	+	ENSP00000373977	7	307	315	9.600000e-01	8	9	7	7	6.000000e-55	6.260000e-01	0	0	1	PSSDDO
1	118927	123443	-	ENSP00000306241	1	134	149	9.000000e-01	0	0	0	1	2.000000e-24	9.100000e-01	3	0	2	DUP
1	125776	127971	-	ENSP00000372165	1	262	262	1.000000e+00	2	0	0	1	2.000000e-70	9.700000e-01	0	0	0	2 DUP
1	125730	126242	+	ENSP00000343312	8	185	215	8.300000e-01	10	13	4	4	6.000000e-11	5.190000e-01	0	0	1	PSSDDO
1	128560	128999	-	ENSP00000302684	14	128	129	8.900000e-01	1	0	1	1	4.000000e-20	6.900000e-01	0	0	0	2 DUP
1	218158	218517	-	ENSP000003027694	1	121	121	1.000000e+00	0	1	0	0	1.000000e-61	9.920000e-01	0	0	0	1 GENE_SINGLE
1	218384	218624	+	ENSP0000034263	254	333	350	2.300000e-01	3	0	2	1	9.000000e-11	8.020000e-01	0	0	0	1 FRAG
1	224218	224322	-	ENSP00000307202	1	35	114	3.100000e-01	0	0	0	0	3.000000e-10	8.860000e-01	0	0	0	1 FP
1	314823	315148	+	ENSP00000302684	1	110	129	8.500000e-01	1	1	1	1	1.000000e-14	5.550000e-01	0	0	1	PSSDDO
1	316287	317502	+	ENSP00000372165	1	262	262	1.000000e+00	2	0	0	1	8.000000e-71	9.700000e-01	0	0	0	2 DUP
1	316996	317421	-	ENSP00000343312	1	135	215	6.300000e-01	11	2	2	2	5.000000e-11	5.210000e-01	0	0	0	1 FRAG
1	319633	324186	+	ENSP00000306241	1	134	149	9.000000e-01	0	0	0	1	6.000000e-11	9.300000e-01	0	0	0	2 DUP

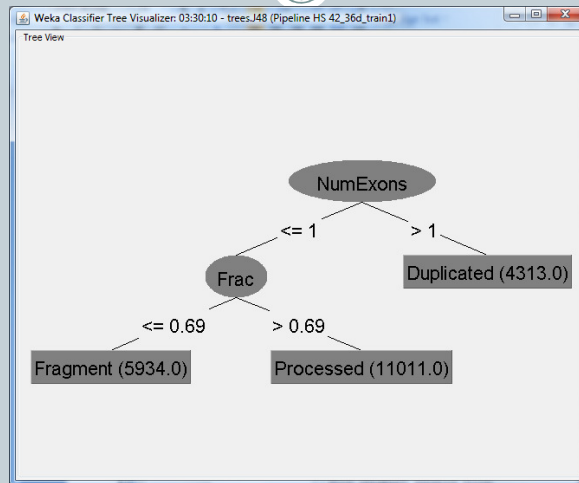
Weka: Datamining Tool



Reproduce the Model



The Pseudogene Classifier



Question

- Is there any way to classify the pseudogenes into processed or duplicated pseudogenes even they are pseudogene fragments that do not have enough information to make a significant conclusion just based on the classification model?

Solution

- Use data mining techniques
 - Cluster the data into two groups
 - Label the clusters as processed and duplicated pseudogenes
 - Remap the pseudogene fragments to the clusters
- Clustering Technique
 - K-means

The Clustering

SECTION THREE

K-means



- K-means is one of the simplest unsupervised learning algorithms that solve the well known clustering problem.
- The procedure follows a simple and easy way to classify a given data set through a certain number of clusters (assume k clusters)
- The main idea is to define k centroids, one for each cluster

http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/kmeans.html

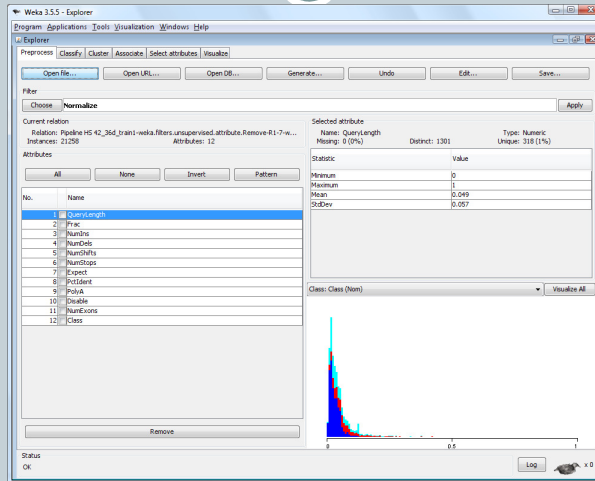
K-means Algorithm



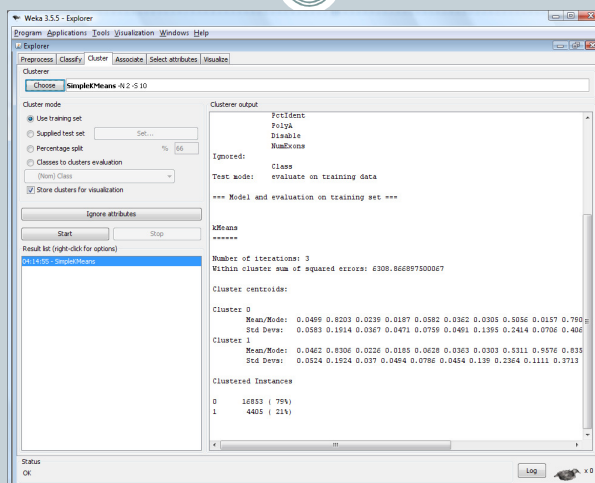
1. Place K points into the space represented by the objects that are being clustered. These points represent initial group centroids.
2. Assign each object to the group that has the closest centroid.
3. When all objects have been assigned, recalculate the positions of the K centroids.
4. Repeat Steps 2 and 3 until the centroids no longer move. This produces a separation of the objects into groups from which the metric to be minimized can be calculated.

http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/kmeans.html

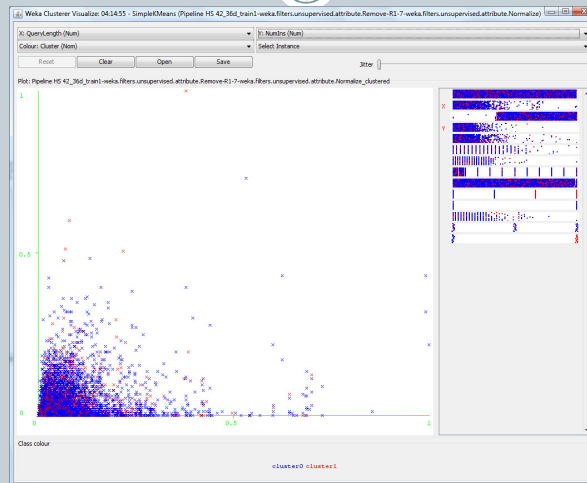
Data Pre-Processing



K-means Clustering



Clustering Result



K-means with PCA

- A weakness, which is common to clustering in general, concerns the visualization of the obtained clusters
- A possible solution is to preprocess the data using PCA
 - the PCA procedure is applied to the data. Using the principal components the data is mapped into the new feature space
 - the k-means algorithm is applied to the data in the feature space. The final objective is to be better able to distinguish the different clusters

<http://dataminingresearch.blogspot.com/search/label/PCA>

Principle Component Analysis

- PCA transforms the data to a new coordinate system such that the greatest variance by any projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on.
- PCA can be used for dimensionality reduction in a data set while retaining those characteristics of the data set that contribute most to its variance, by keeping lower-order principal components and ignoring higher-order ones. Such low-order components often contain the "most important" aspects of the data.

http://en.wikipedia.org/wiki/Karhunen-Lo%C3%A8ve_transform

PCA Attribute Selection

The screenshot shows the Weka 3.5.5 Explorer interface. The 'Attribute Selection' tab is active, displaying the results of a PCA analysis. The 'Attribute selection output' table lists 10 principal components (PC1 to PC10) with their corresponding loadings for 10 original attributes (V1 to V10). The 'Ranked attributes' list shows the top 10 attributes based on their variance explained, with PC1 being the most significant.

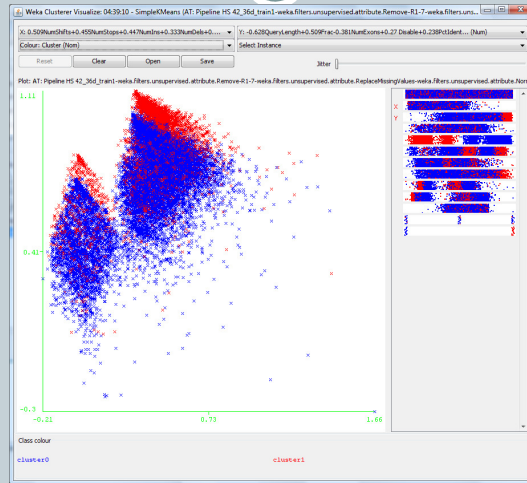
Attribute selection output										
EXPRESSIVE										
V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	
0.1041	-0.6283	0.2389	-0.1471	-0.0704	-0.2061	-0.2109	0.098	0.0305	-0.6149	Querylength
0.2012	0.5091	0.3245	0.1232	-0.0282	0.4346	0.1663	-0.0668	-0.14	-0.5861	Frac
0.4493	-0.1211	0.0136	0.0345	0.0281	0.0036	-0.0653	-0.6207	0.3026	-0.0454	NumIn
0.333	-0.1645	-0.0408	0.0346	0.4218	0.5031	-0.2409	0.5637	0.1082	0.1114	NumIn
0.5099	0.0507	-0.0093	-0.038	-0.0433	-0.0405	-0.2016	-0.2007	-0.0095	0.3347	NumIn
0.4548	0.0207	-0.0613	-0.0202	-0.1309	-0.2669	0.0004	0.1084	-0.7445	0.018	NumIn
-0.0487	-0.1256	-0.5299	-0.087	-0.6477	0.477	-0.1758	-0.0154	-0.0566	-0.0983	Expect
-0.2459	0.2384	0.4764	-0.1131	-0.2021	0.0145	-0.7296	-0.0019	-0.0412	0.1263	PostIdent
0.0195	0.1092	-0.0187	-0.063	0.1617	0.0009	0.1367	-0.0071	-0.0311	0.0068	Prn
0.324	0.2702	0.0472	-0.1009	-0.4577	-0.3118	0.1697	0.4622	0.5013	0.0004	Disable
0.0147	-0.3807	0.3631	-0.0101	-0.3172	0.3173	0.4084	-0.0082	-0.0971	0.3972	NumExms

Ranked attributes:

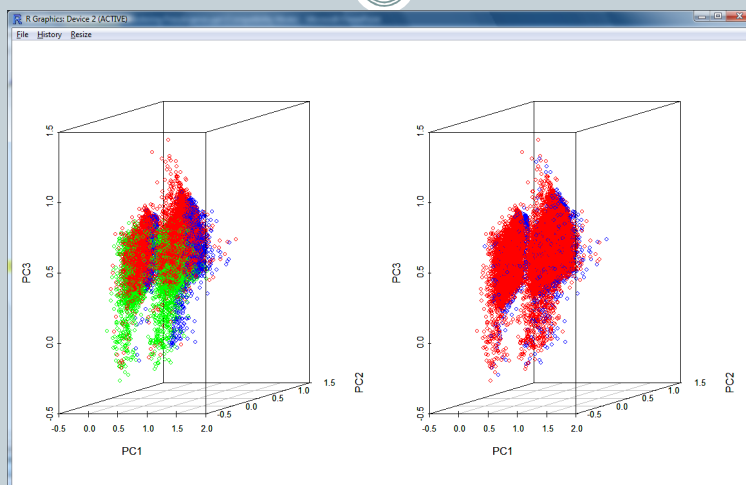
- 0.7851 1 0.5058NumIn+0.4518NumTop+0.4478NumIn+0.3338NumIn+0.3281Disable...
- 0.5677 2 -0.6283Querylength+0.5099Frac+0.3818NumExms+0.27 Disable+0.2398PostIdent...
- 0.482 3 0.5658NumExms+0.53Expect+0.4767PostIdent+0.325Frac+0.2397Querylength...
- 0.3958 4 -0.5639V9+0.1470Querylength+0.1525Frac+0.1159V10+0.1010Disable...
- 0.3079 5 -0.6488Expect+0.4858Disable+0.4228NumIn+0.3178NumExms+0.2022PostIdent...
- 0.2322 6 0.5028NumIn+0.477Expect+0.439Frac+0.3178NumExms+0.3125Disable...
- 0.1688 7 -0.7381V10+0.4018NumExms+0.3028NumIn+0.2418NumIn+0.2110Querylength...
- 0.107 8 -0.6118NumIn+0.5648NumIn+0.4622Disable+0.2888NumIn+0.1098NumTop...
- 0.058 9 -0.7548NumTop+0.3010Disable+0.3038NumIn+0.14Frac+0.1098NumIn...
- 0.0279 10 -0.6170Querylength+0.5367Frac+0.3978NumExms+0.3388NumIn+0.1263PostIdent...

Selected attributes: 1,2,3,4,5,6,7,8,9,10 : 10

Clustering PCA Result



A 3D View of the Clusters



Work to do



- Try more different parameters and clustering algorithms
- Relabel the pseudogene fragments according to the clusters
- Calculate the error rate
 - Use the training data set with just the processed and duplicated pseudogenes as a test set
 - Use the test set to evaluate the clustering
 - Calculate the error rate and plot the ROC curve