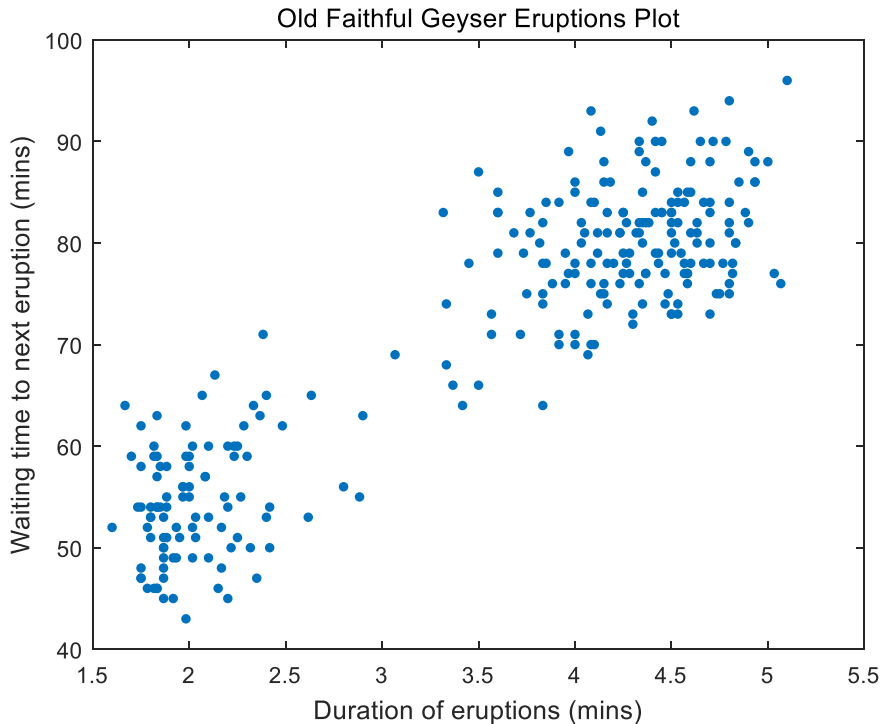


## EE511 Project 3

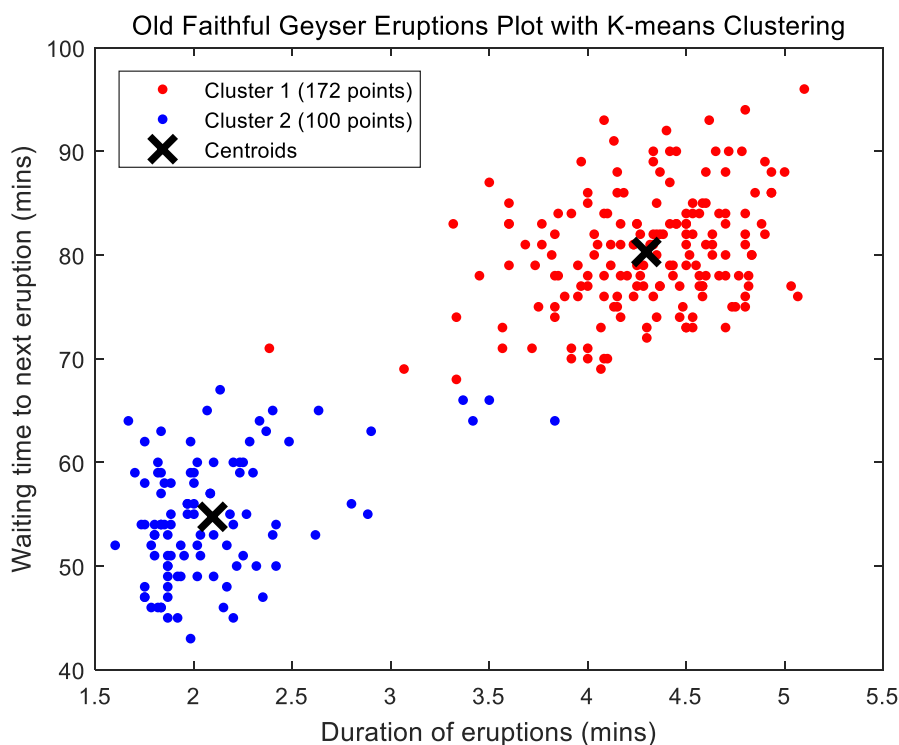
This project is implemented using MATLAB. Codes are attached in the end of this document.

### [Testing Faith]

In this problem, I use textscan function to read data from text files and convert to a 272 by 2 matrix. The first column is duration of eruption and the second column is the waiting time to next eruption. Then generate the scatter plot as following.



Then use k-mean clustering function to partition the data into two clusters ( $k = 2$ ). This k-means function uses the squared Euclidean distance measure and the k-means++ algorithm for cluster center initialization. The number of points in cluster and centroids keep the same in each run.

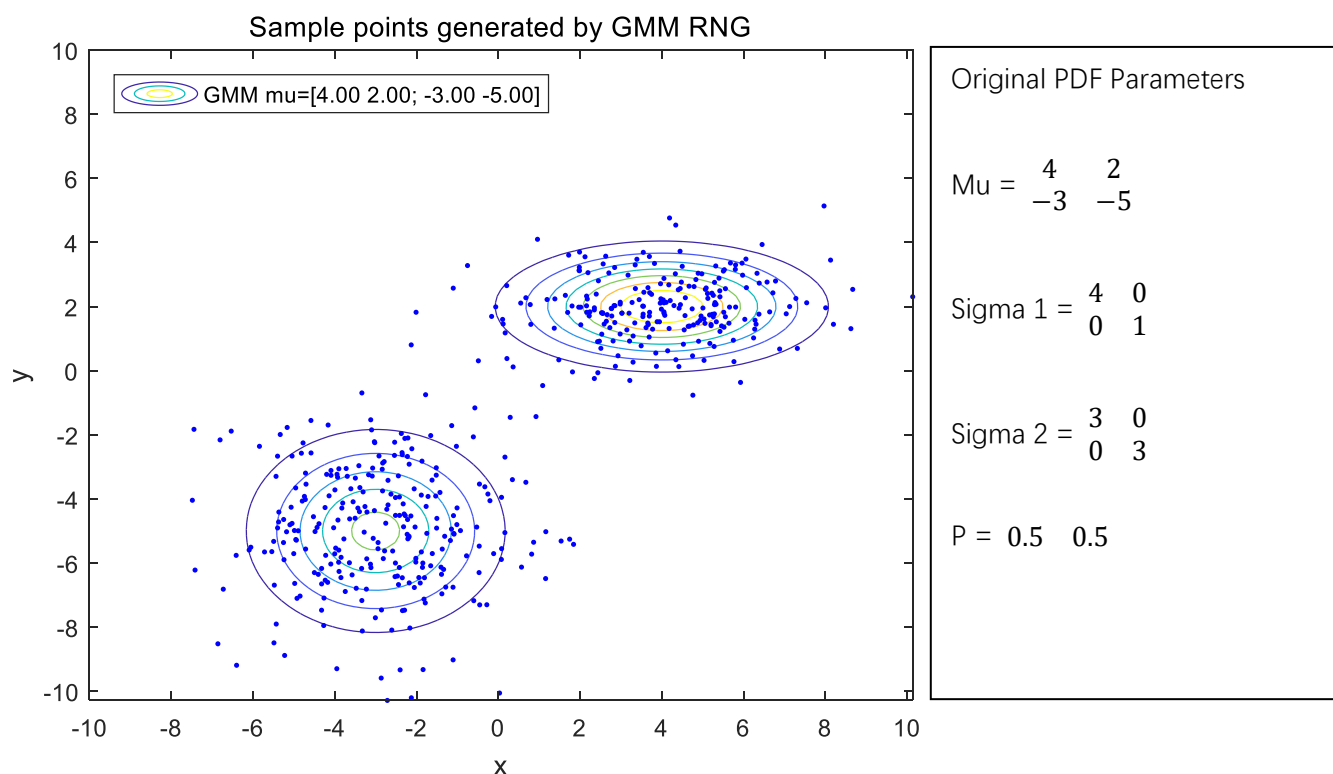


### [EM]

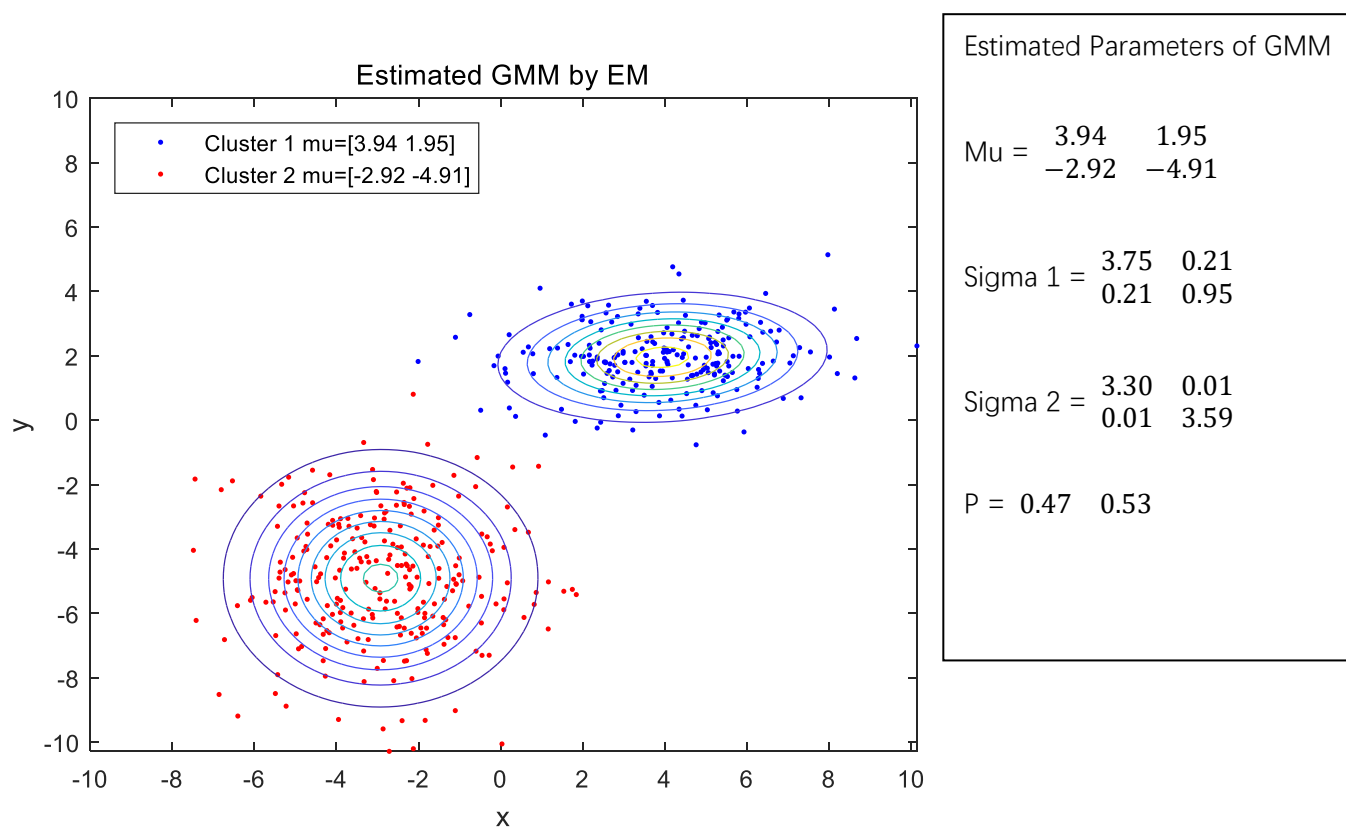
GMM RNG with 2 subpopulations:

The below plot is the contour of pdf and sampling points generated by this RNG.

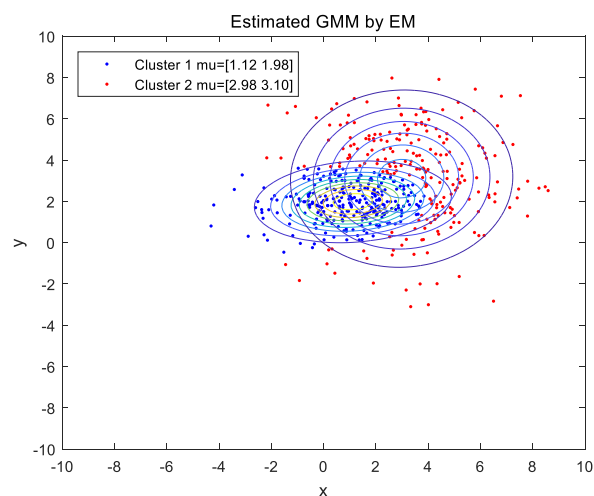
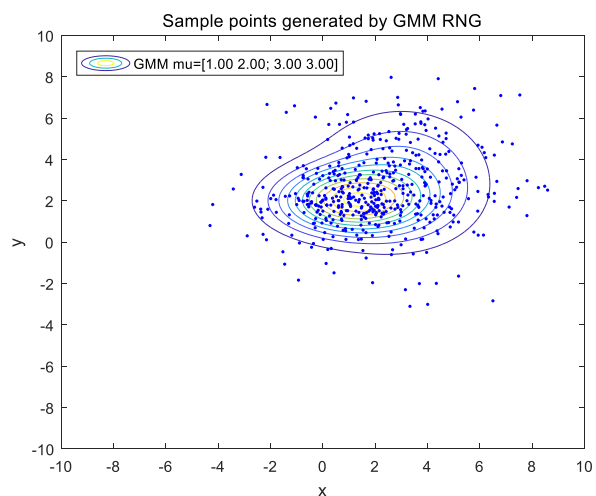
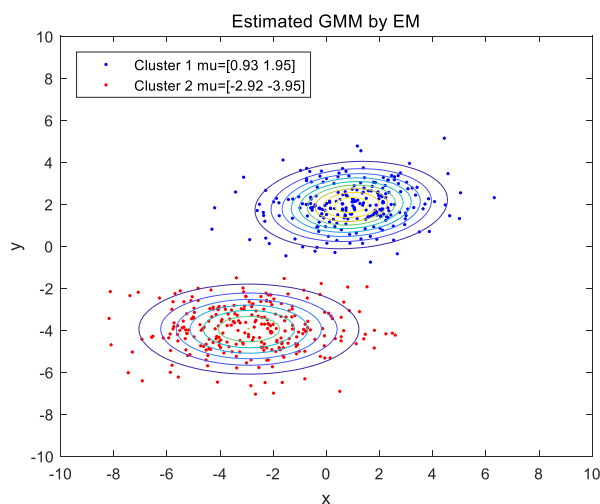
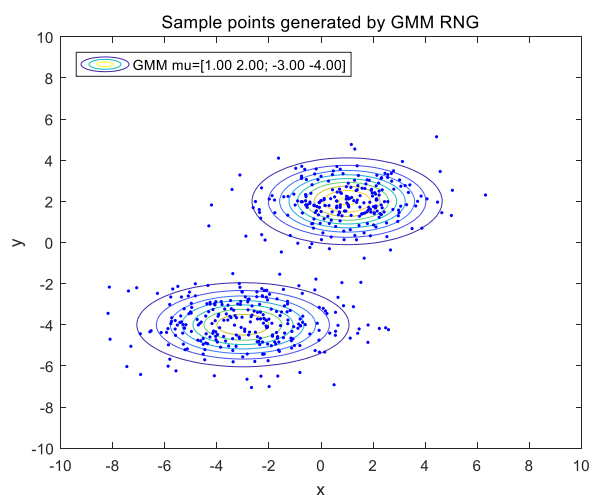
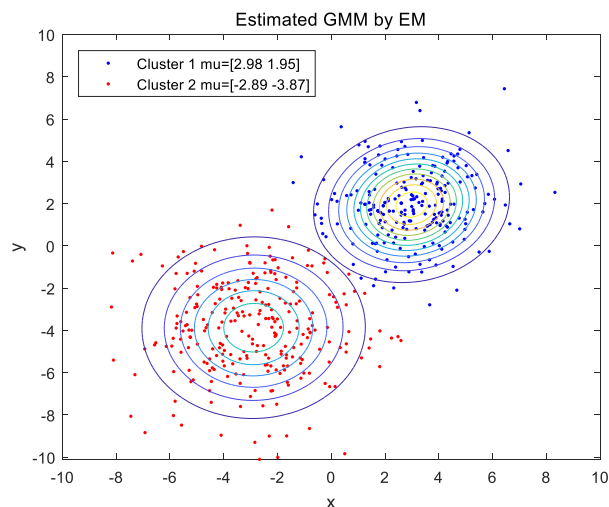
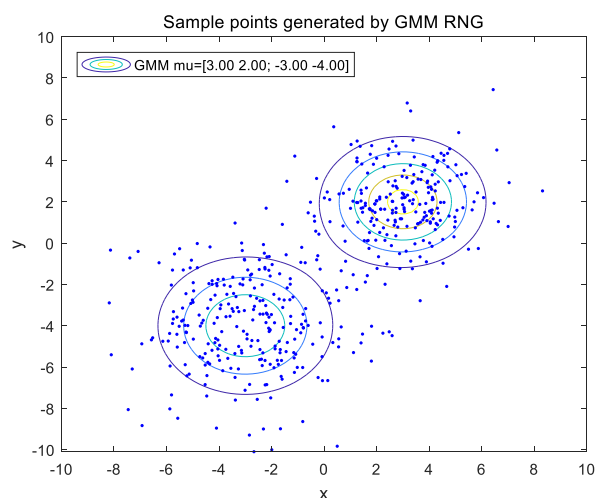
The original parameters of pdf used by this RNG are listed on right side.



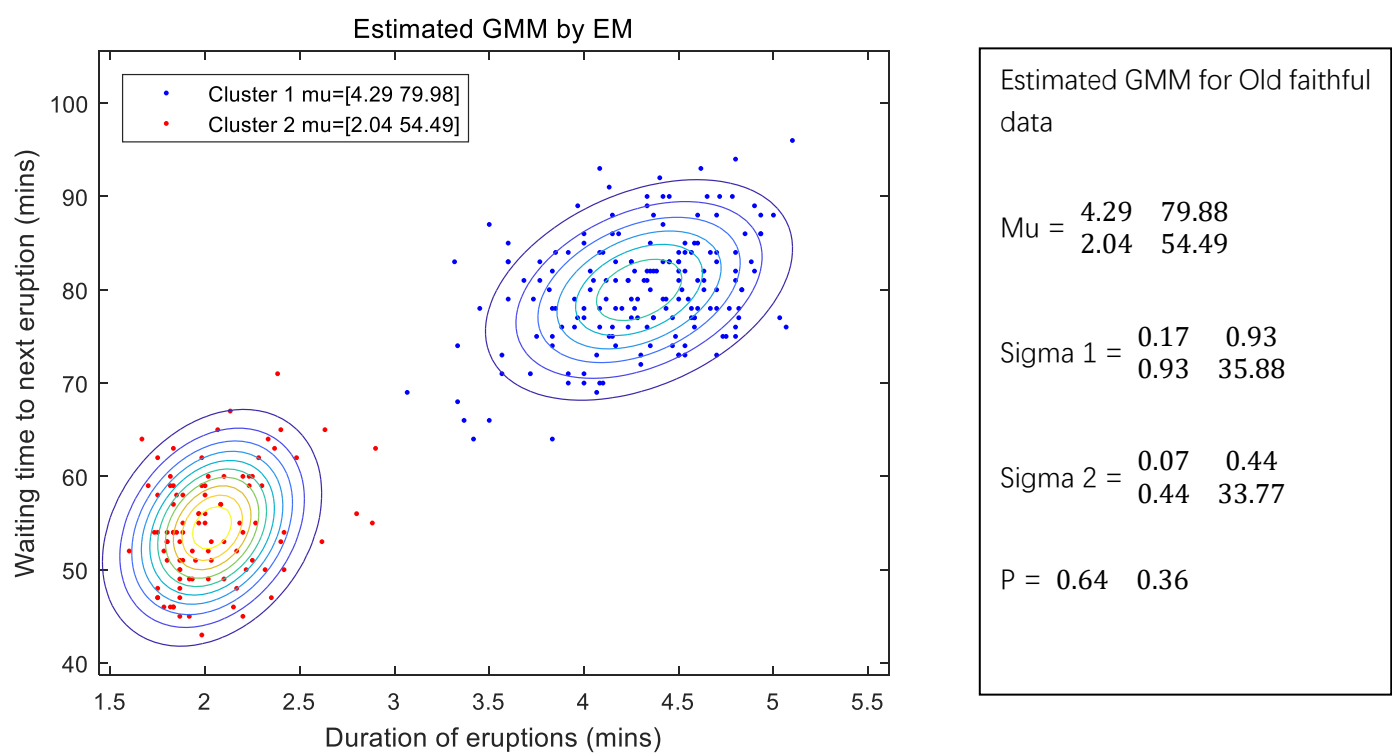
Then use expectation maximization (EM) algorithm to estimate the pdf of sampling points. The below plot is the contour of estimated pdf. The estimated parameters are listed on right side. Comparing with original parameters and the shape of pdf, the parameters and shape conform to the original approximately.



The following plots are estimated pdf of spherical covariance matrices, ellipsoidal covariance matrices and poorly-separated subpopulations. Left side are original pdf and right side are estimated pdf. For the quality, spherical and ellipsoidal covariance are better than poorly-separated. For speed of estimation, spherical and ellipsoidal covariance are much faster than poorly-separated. Spherical and ellipsoidal cases generally take dozens of iterations, but poorly-separated case may take hundreds of iterations.

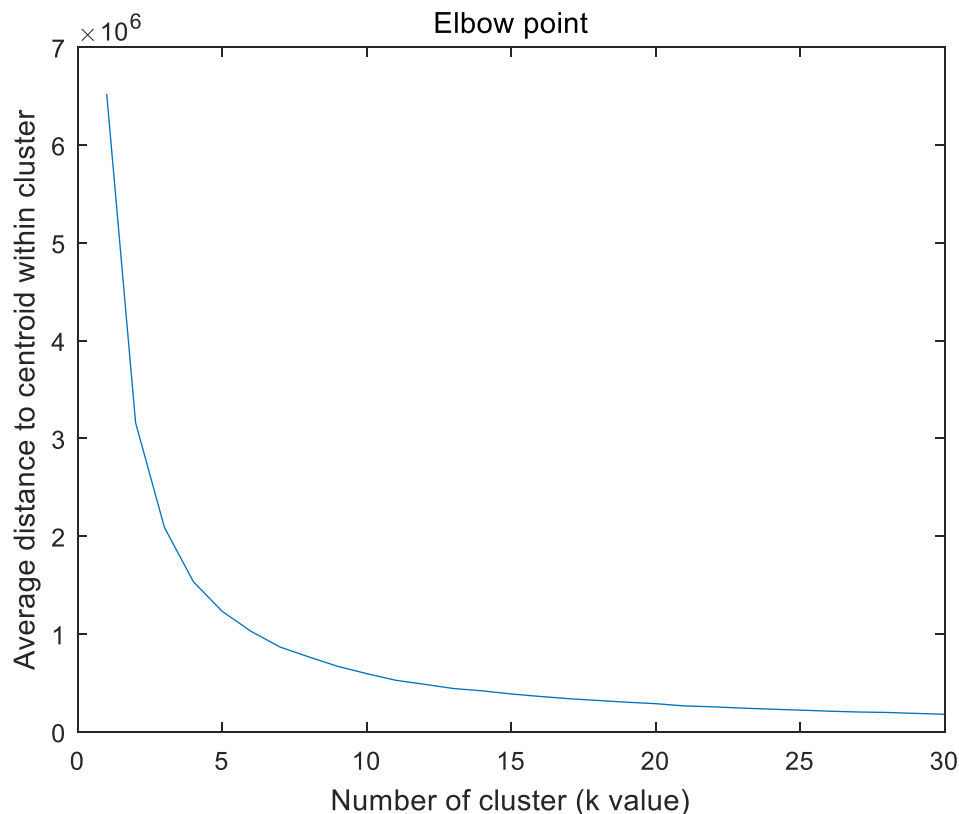


The below plot is the GMM-EM estimation for old faithful data set.



## [Clusters of Text]

I tried different k values and record the average distance to centroid within cluster. The below plot it the relationship between k value and average distance to centroid with cluster. I choose k = 12 because the average distance changes little when k goes greater than 12.



Cluster 1 (238 docs)

-----

1987\_10 1987\_14 1987\_17 1987\_19 1987\_21 1987\_22 1987\_23 1987\_24 1987\_25 1987\_28  
1987\_30 1987\_31 1987\_32 1987\_34 1987\_40 1987\_42 1987\_43 1987\_45 1987\_48 1987\_49  
1987\_52 1987\_54 1987\_58 1987\_59 1987\_61 1987\_64 1987\_68 1987\_70 1987\_79 1987\_80  
1987\_81 1987\_82 1987\_85 1987\_86 1987\_87 1987\_90 1988\_11 1988\_31 1988\_32 1988\_33  
1988\_34 1988\_35 1988\_37 1988\_38 1988\_41 1988\_48 1988\_50 1988\_51 1988\_52 1988\_56  
1988\_59 1988\_61 1988\_62 1988\_63 1988\_64 1988\_65 1988\_67 1988\_69 1988\_71 1988\_72  
1988\_73 1988\_74 1988\_76 1988\_77 1988\_78 1988\_79 1988\_80 1988\_81 1988\_82 1988\_83  
1988\_84 1988\_85 1988\_86 1988\_88 1988\_90 1988\_91 1988\_92 1988\_93 1988\_94 1988\_95  
1989\_1 1989\_5 1989\_6 1989\_8 1989\_12 1989\_14 1989\_20 1989\_28 1989\_32 1989\_34  
1989\_35 1989\_37 1989\_38 1989\_41 1989\_44 1989\_46 1989\_48 1989\_52 1989\_71 1989\_72  
1989\_78 1989\_81 1989\_85 1989\_89 1989\_91 1989\_92 1989\_93 1989\_94 1989\_95 1989\_96  
1989\_97 1989\_100 1990\_1 1990\_2 1990\_4 1990\_7 1990\_8 1990\_9 1990\_10 1990\_11  
1990\_12 1990\_15 1990\_17 1990\_18 1990\_19 1990\_22 1990\_24 1990\_33 1990\_35 1990\_39  
1990\_41 1990\_42 1990\_44 1990\_45 1990\_48 1990\_49 1990\_50 1990\_51 1990\_53 1990\_54  
1990\_55 1990\_57 1990\_61 1990\_62 1990\_66 1990\_69 1990\_71 1990\_75 1990\_80 1990\_84  
1990\_87 1990\_88 1990\_101 1990\_117 1990\_118 1990\_120 1990\_129 1990\_135 1990\_136 1990\_139  
1990\_140 1990\_141 1990\_142 1990\_143 1991\_1 1991\_3 1991\_8 1991\_9 1991\_11 1991\_13  
1991\_24 1991\_25 1991\_27 1991\_28 1991\_33 1991\_35 1991\_43 1991\_44 1991\_46 1991\_49  
1991\_50 1991\_52 1991\_53 1991\_56 1991\_57 1991\_59 1991\_66 1991\_73 1991\_75 1991\_76  
1991\_77 1991\_78 1991\_83 1991\_92 1991\_94 1991\_95 1991\_97 1991\_98 1991\_99 1991\_100

1991\_101 1991\_106 1991\_126 1991\_142 1992\_4 1992\_13 1992\_16 1992\_27 1992\_43 1992\_48  
1992\_50 1992\_55 1992\_64 1992\_68 1992\_69 1992\_70 1992\_71 1992\_77 1992\_81 1992\_82  
1992\_89 1992\_92 1992\_95 1992\_97 1992\_98 1992\_99 1992\_101 1992\_102 1992\_105 1992\_108  
1992\_110 1992\_112 1992\_116 1992\_120 1992\_121 1992\_122 1992\_123 1992\_125

Cluster 2 (1 docs)

-----

1991\_84

Cluster 3 (1 docs)

-----

1989\_83

Cluster 4 (5 docs)

-----

1987\_41 1989\_21 1989\_24 1990\_132 1992\_18

Cluster 5 (24 docs)

-----

1987\_8 1987\_13 1987\_38 1988\_5 1988\_21 1989\_42 1989\_77 1990\_119 1991\_58 1991\_103  
1991\_104 1991\_112 1991\_117 1991\_122 1991\_129 1991\_132 1991\_140 1992\_20 1992\_21 1992\_29  
1992\_59 1992\_61 1992\_65 1992\_103

Cluster 6 (117 docs)

-----

1987\_6 1987\_18 1987\_37 1987\_50 1987\_65 1987\_89 1988\_4 1988\_7 1988\_8 1988\_16  
1988\_20 1988\_30 1988\_58 1988\_68 1989\_56 1989\_59 1989\_60 1989\_62 1989\_65 1989\_66  
1989\_67 1989\_68 1989\_69 1989\_74 1989\_75 1989\_76 1989\_82 1989\_88 1989\_99 1990\_25  
1990\_72 1990\_90 1990\_91 1990\_92 1990\_94 1990\_95 1990\_96 1990\_98 1990\_99 1990\_100  
1990\_102 1990\_103 1990\_106 1990\_108 1990\_110 1990\_111 1990\_113 1990\_114 1990\_115 1990\_121  
1990\_122 1990\_124 1990\_125 1990\_127 1990\_128 1990\_131 1990\_133 1990\_134 1991\_45 1991\_48  
1991\_55 1991\_63 1991\_72 1991\_79 1991\_81 1991\_82 1991\_85 1991\_86 1991\_102 1991\_105  
1991\_109 1991\_110 1991\_111 1991\_113 1991\_115 1991\_118 1991\_123 1991\_124 1991\_125 1991\_127  
1991\_128 1991\_130 1991\_131 1991\_134 1991\_135 1991\_136 1991\_137 1991\_139 1991\_141 1991\_143  
1992\_1 1992\_2 1992\_7 1992\_10 1992\_19 1992\_25 1992\_28 1992\_32 1992\_44 1992\_49  
1992\_53 1992\_56 1992\_57 1992\_58 1992\_60 1992\_62 1992\_66 1992\_67 1992\_72 1992\_73  
1992\_74 1992\_75 1992\_83 1992\_90 1992\_91 1992\_94 1992\_109

Cluster 7 (165 docs)

-----

1987\_3 1987\_11 1987\_15 1987\_33 1987\_36 1987\_46 1987\_47 1987\_53 1987\_56 1987\_62  
1987\_66 1987\_72 1987\_77 1987\_78 1987\_83 1987\_88 1988\_1 1988\_2 1988\_3 1988\_6  
1988\_9 1988\_12 1988\_13 1988\_14 1988\_15 1988\_17 1988\_18 1988\_19 1988\_24 1988\_25  
1988\_26 1988\_27 1988\_28 1988\_29 1988\_36 1988\_39 1988\_40 1988\_42 1988\_57 1988\_60  
1988\_70 1988\_75 1988\_89 1989\_13 1989\_22 1989\_23 1989\_25 1989\_26 1989\_27 1989\_29  
1989\_30 1989\_31 1989\_33 1989\_39 1989\_43 1989\_45 1989\_49 1989\_50 1989\_51 1989\_53  
1989\_54 1989\_55 1989\_57 1989\_63 1989\_64 1989\_70 1989\_73 1989\_79 1989\_80 1989\_98  
1989\_101 1990\_5 1990\_20 1990\_21 1990\_23 1990\_26 1990\_27 1990\_28 1990\_29 1990\_30  
1990\_31 1990\_32 1990\_34 1990\_36 1990\_37 1990\_40 1990\_58 1990\_60 1990\_73 1990\_74

1990\_76 1990\_77 1990\_78 1990\_79 1990\_82 1990\_83 1990\_85 1990\_86 1990\_93 1990\_97  
1990\_104 1990\_105 1990\_107 1990\_116 1990\_123 1990\_126 1990\_137 1991\_17 1991\_18 1991\_19  
1991\_20 1991\_21 1991\_22 1991\_23 1991\_26 1991\_29 1991\_30 1991\_32 1991\_34 1991\_37  
1991\_38 1991\_40 1991\_42 1991\_54 1991\_60 1991\_61 1991\_62 1991\_67 1991\_68 1991\_71  
1991\_74 1991\_80 1991\_87 1991\_89 1991\_90 1991\_91 1991\_96 1991\_114 1991\_119 1991\_120  
1991\_121 1991\_133 1991\_138 1992\_5 1992\_6 1992\_12 1992\_14 1992\_15 1992\_17 1992\_22  
1992\_24 1992\_26 1992\_35 1992\_39 1992\_45 1992\_80 1992\_84 1992\_85 1992\_86 1992\_87  
1992\_88 1992\_93 1992\_106 1992\_107 1992\_117

#### Cluster 8 (70 docs)

-----

1987\_7 1987\_9 1987\_12 1987\_16 1987\_20 1987\_26 1987\_29 1987\_35 1987\_39 1987\_51  
1987\_69 1987\_71 1987\_73 1987\_74 1987\_75 1988\_22 1988\_43 1988\_44 1988\_45 1988\_46  
1988\_47 1988\_49 1988\_54 1988\_55 1989\_2 1989\_3 1989\_4 1989\_7 1989\_10 1989\_11  
1989\_15 1989\_16 1989\_17 1989\_18 1989\_19 1989\_36 1990\_3 1990\_6 1990\_38 1990\_43  
1990\_46 1990\_47 1990\_52 1990\_70 1991\_2 1991\_4 1991\_5 1991\_6 1991\_7 1991\_10  
1991\_12 1991\_15 1991\_16 1991\_47 1991\_51 1991\_88 1991\_93 1992\_46 1992\_47 1992\_51  
1992\_52 1992\_54 1992\_96 1992\_100 1992\_111 1992\_113 1992\_114 1992\_124 1992\_126 1992\_127

#### Cluster 9 (3 docs)

-----

1987\_76 1990\_81 1992\_104

#### Cluster 10 (11 docs)

-----

1987\_5 1988\_10 1989\_86 1990\_112 1990\_130 1991\_108 1991\_116 1992\_3 1992\_8 1992\_76  
1992\_79

#### Cluster 11 (55 docs)

-----

1987\_2 1987\_4 1987\_27 1987\_57 1987\_60 1987\_67 1987\_84 1988\_23 1988\_87 1989\_40  
1989\_47 1989\_58 1989\_61 1989\_84 1989\_87 1989\_90 1990\_16 1990\_56 1990\_59 1990\_63  
1990\_64 1990\_65 1990\_67 1990\_68 1990\_89 1990\_109 1990\_138 1991\_14 1991\_31 1991\_36  
1991\_39 1991\_41 1991\_64 1991\_65 1991\_69 1991\_70 1991\_107 1991\_144 1992\_9 1992\_11  
1992\_23 1992\_31 1992\_33 1992\_34 1992\_36 1992\_37 1992\_38 1992\_40 1992\_41 1992\_42  
1992\_63 1992\_78 1992\_115 1992\_118 1992\_119

#### Cluster 12 (9 docs)

-----

1987\_44 1987\_55 1987\_63 1988\_53 1988\_66 1989\_9 1990\_13 1990\_14 1992\_30