

FSMEM for MoG

Daniel Wagenaar, June 2000

FSMEM, or *free split/merge expectation maximization*, is a modification of the SMEM algorithm presented by Ueda et al in [1]. Unlike SMEM, FSMEM can modify the number of clusters dynamically. To compare likelihoods across solutions with different number of clusters, a minimal description length term is used.

I will be using the FSMEM algorithm exclusively for optimizing mixture of Gaussians, although the algorithm can be applied equally to other clustering problems. The following symbols will be used:

K refers to the number of clusters. k is a dummy index enumerating clusters.

N refers to the number of data point. n enumerates data points.

D refers to the dimensionality of the space.

Vectors in data space will be denoted $\vec{x} = (x^1, \dots, x^D)$. Inner product notation will be used whenever convenient:

$$\vec{x} \cdot \Sigma \vec{x} \equiv \sum_{d_1, d_2} x^{d_1} \Sigma_{d_1 d_2} x^{d_2}.$$

1 Mixture of Gaussians and EM

Mixture of Gaussians models a probability density function of the form

$$p(\vec{x}) = \sum_k \pi_k p(\vec{x}; \Theta_k),$$

where π_k are non-negative and sum to one, and $\Theta_k = \{\vec{\mu}_k, \Sigma_k\}$ are the parameters for a Gaussian:

$$p(\vec{x}; \Theta_k) = \frac{1}{\sqrt{(2\pi)^D \det \Sigma}} e^{-\frac{1}{2} \vec{x} \cdot \Sigma^{-1} \vec{x}}$$

The standard EM learning algorithm for MoG consists of iterating the following steps:

1. Calculate responsibilities:

$$R_{kn} = \frac{\pi_k p(\vec{x}_n; \Theta_k)}{\sum_{k'} \pi_{k'} p(\vec{x}_n; \Theta_{k'})}$$

2. Update parameters:

$$\begin{aligned} \vec{\mu}_k' &= \frac{\sum_n R_{kn} \vec{x}_n}{\sum_n R_{kn}} \\ \Sigma_k' &= \frac{\sum_n R_{kn} R_{kn} (\vec{x}_n - \mu_k') (\vec{x}_n - \mu_k')^T}{\sum_n R_{kn}} \\ \pi_k &= \sum_n R_{kn}. \end{aligned}$$

These steps are to be iterated until the log likelihood $L = \sum_n p(\vec{x}_n)$ converges.

2 SMEM by Ueda et al

This is the original SMEM algorithm by Ueda et al:

Initialize parameters as for EM.

Run EM until convergence.

repeat

Collect a list of merge candidates: pairs of clusters with their merge merits, and a list of split candidates: clusters with their split merits.

Sort these lists, then

for all promising triplets $\{k_{m_1}, k_{m_2}; k_s\}$ **do**

Perform the split and merge operations.

Reinitialize the affected clusters.

Run partial EM on the affected clusters until convergence.

Run full EM on all clusters until convergence.

if the likelihood is better than before the merge/split **then**

Ignore the other candidates, and go back to the candidate collection step.

else

Restore parameters.

until no candidate produced better result than old one

The split candidates are ordered by their *split merit*:

$$J_k^{(\text{split})} = \int d\vec{x} f_k(x) \log \frac{f_k(x)}{p(\vec{x}_n; \Theta_k)},$$

where

$$f_k(x) = \frac{\sum_n \delta(\vec{x} - \vec{x}_n) R_{kn}}{\sum_n R_{kn}}.$$

$J_k^{(\text{split})}$ is a measure of the KL distance between the local observed distribution $f_k(\vec{x})$ and the model distribution $p(\vec{x}; \Theta_k)$.

Similarly, merge candidates are ordered by their merge merit:

$$J_{kk'}^{(\text{merge})} = \sum_n R_{kn} R_{k'n}.$$

Note that split and merge merits are measured on different scales, so it is not possible to sort split and merge candidates together. Ueda et al state that it is only necessary to test about 5 promising candidates at each iteration, but they do not specify how they weigh split and merge merits together.

Initialization of parameters after a merge is as follows:

$$\vec{\mu}' = \frac{1}{2}(\vec{\mu}_1 + \vec{\mu}_2)$$

$$\Sigma' = \frac{1}{2}(\Sigma_1 + \Sigma_2)$$

$$\pi' = \pi_1 + \pi_2.$$

Note that no attempt is made to weight the parameters.

Initialization of parameters after a split is as follows:

$$\vec{\mu}'_{1,2} = \vec{\mu} + \vec{v}_{1,2},$$

where \vec{v}_k is a small noise term. I choose to pick \vec{v} from $\mathcal{G}(0, \Sigma)$.

$$\Sigma'_{1,2} = \sigma^2 \mathbf{I},$$

where σ^2 is the largest eigenvalue of Σ . Ueda et al use $\sigma^2 = (\det \Sigma)^{1/D}$, but I had problems with that approach in cases where Σ possessed some very small eigenvalues.

$$\pi_{1,2} = \frac{1}{2} \pi.$$

Partial EM is defined as EM acting on a subset K' of clusters. The equations are the same as for EM, except

$$R_{kn} = \frac{\pi_k p(\vec{x}_n; \Theta_k)}{\sum_{k' \in K'} \pi_{k'} p(\vec{x}_n; \Theta_{k'})} R_{kn}^{(0)}$$

where $R_{kn}^{(0)}$ are the responsibilities just after the split or merge. This modification ensures that $\sum_{k \in K'} R_{kn}$ remains constant during the partial EM, thus preserving normalization.

3 FSMEM

The proposed new FSMEM algorithm is subtly different:

Initialize parameters as for EM.

Run EM until convergence.

Let *state* be MERGE.

Let *fail* be zero.

repeat

if *state* = MERGE **then**

 Collect a list of merge candidates.

 Sort this list.

for all promising pairs $\{k_{m_1}, k_{m_2}\}$ **do**

 Perform the merge operation. {This reduces K by 1.}

 Reinitialize the affected cluster.

 Run partial EM on the affected cluster until convergence.

 Run full EM on all clusters until convergence.

if the MDL likelihood is better than before the merge **then**

 Ignore the other candidates, let *fail* be zero again, and go back to the outer loop.

else

 Restore parameters.

 Let *state* be SPLIT.

 Increment *fail* by one.

else

 {So *state* = SPLIT}

Collect a list of split candidates.

Sort this list.

for all promising candidates k_s **do**

 Perform the split operation. {This increments K by 1.}

 Run partial EM on the affected clusters until convergence.

 Run full EM on all clusters until convergence.

if the MDL likelihood is better than before the merge **then**

 Ignore the other candidates, let *fail* be zero again, and go back to the outer loop.

else

 Restore parameters.

 Let *state* be MERGE.

 Increment *fail* by one.

until *fail* equals two.

Thus split and merge operations can now be executed independently. The algorithm starts by looking for merge candidates, and continues to do that until it doesn't find any. Then it looks for split candidates, and continues to do that until there aren't any successful ones. Then it returns to scanning merge candidates, and so on until neither split nor merge finds an improvement in the MDL likelihood

$$L_{\text{MDL}} = L - \frac{1}{2} \log(N)K(1 + D + \frac{1}{2}D(D + 1)).$$

The second term here expresses the entropy in the parameters, and corrects for the scaling of L with number of parameters.

4 Stability issues

4.1 Numerical

In several equations I found that numerical problems occurred when $p(\vec{x})$ was so small it was rounded to zero in the computer algorithms. This caused trouble when none of the clusters was near a particular point. I solved these problems by adding a vanishingly small constant $\epsilon = 10^{-300}$ to some equations.

4.2 Fundamental

There is another, more fundamental stability issue identified by Ueda et al.: there are global optima in parameter space corresponding to the appearance of clusters with vanishing variance. To avoid being sucked into those, they suggest the following alternative for the update step for the variance matrix:

$$\Sigma'_k = \frac{\sum_n R_{kn}(\vec{x}_n - \mu'_k)(\vec{x}_n - \mu'_k)^T + \lambda I}{\sum_n R_{kn} + \lambda}.$$

(In their paper the λ in the denominator is replaced by 1, but I believe that to be a misprint.) I found good results for $\lambda \sim .01$.

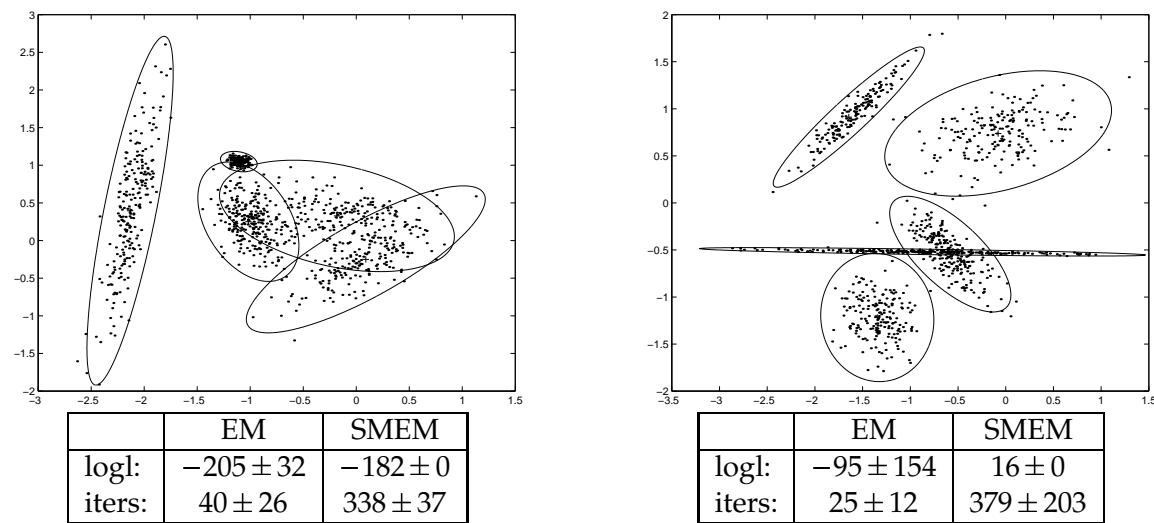
5 Results on artificial data

The SMEM algorithm was able to correctly find clusters in many cases where plain EM was trapped in a local optimum. It could not always correctly reconstruct clusters that were so much on top of each other that taking them apart didn't yield much of an improvement in log likelihood.

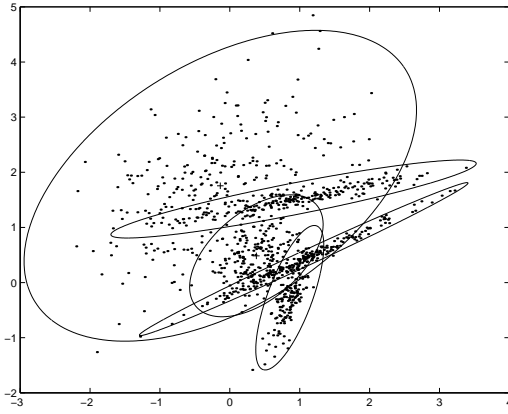
The FSMEM algorithm did just as well, with the added bonus that it could determine the correct number of clusters automatically. Predictably, this worked best if the target clusters were well-separated.

5.1 Two dimensional data

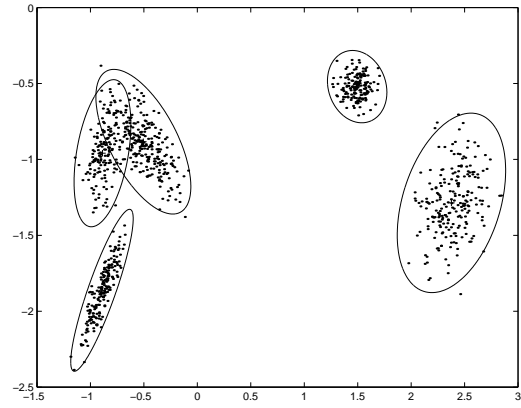
To gain some preliminary insight in the performance of the SMEM, I performed some runs on two dimensional data. The original data are shown below together with the log likelihoods¹ found by the various algorithms and the number of iterations they took. All results based on five runs per source. Reported are mean and spread, even though the distributions are far from Gaussian.



¹Modified by the MDL term discussed above for easy comparison with FSMEM results.



| | EM | SMEM |
|--------|-----------------|-----------------|
| logl: | -1450 ± 154 | -1278 ± 132 |
| iters: | 24 ± 17 | 269 ± 149 |



| | EM | SMEM |
|--------|---------------|--------------|
| logl: | 526 ± 165 | 613 ± 58 |
| iters: | 21 ± 14 | 114 ± 31 |

It is seen that SMEM does better than simple EM in all cases, and moreover produces much more predictable results. However, there is a price to pay: the number of iterations is much higher. (Iteration counts include iterations made in discarded steps.) These data sets were generated using `mog_gendata(2, 1000, 5, .3)`, i.e. $D = 2$, $N = 1000$, $K = 5$ and typical cluster size .3 spread over a typical area of size 1.

5.2 Four dimensional data

After these tests, a more extended test series was done on four dimensional data, for EM, SMEM and FSMEM together. In all cases FSMEM was able to find the correct number of clusters, whether started from $K = 5$ (the correct number of clusters), $K = 1$ or $K = 10$. These were the log likelihoods found in 5 runs for each of 10 sources:

| Run | EM | SMEM | FSMEM $K_0 = 1$ | FSMEM $K_0 = 5$ | FSMEM $K_0 = 10$ |
|-----|---------------------|---------------------|--------------------|--------------------|--------------------|
| 01 | 376.48 ± 55.70 | 417.24 ± 0.02 | 417.21 ± 0.09 | 417.21 ± 0.09 | 417.21 ± 0.09 |
| 02 | 782.42 ± 294.98 | 969.08 ± 280.43 | 1094.50 ± 0.00 | 1094.50 ± 0.00 | 1094.50 ± 0.00 |
| 03 | 29.59 ± 238.59 | 173.50 ± 0.00 | 173.50 ± 0.01 | 173.50 ± 0.01 | 173.50 ± 0.01 |
| 04 | -197.61 ± 48.78 | -161.98 ± 0.00 | -162.00 ± 0.01 | -162.00 ± 0.01 | -162.00 ± 0.01 |
| 05 | 406.12 ± 150.98 | 570.17 ± 0.00 | 570.12 ± 0.04 | 570.12 ± 0.04 | 570.12 ± 0.04 |
| 06 | 136.20 ± 297.57 | 471.60 ± 0.00 | 471.58 ± 0.04 | 471.58 ± 0.04 | 471.58 ± 0.04 |
| 07 | 471.46 ± 252.24 | 708.11 ± 0.00 | 708.11 ± 0.00 | 708.11 ± 0.00 | 708.11 ± 0.00 |
| 08 | 90.66 ± 131.87 | 323.93 ± 0.00 | 323.93 ± 0.00 | 323.93 ± 0.00 | 323.93 ± 0.00 |
| 09 | 104.33 ± 98.83 | 175.02 ± 0.00 | 175.02 ± 0.00 | 175.02 ± 0.00 | 175.02 ± 0.00 |
| 10 | 225.83 ± 81.66 | 262.40 ± 0.03 | 262.20 ± 0.16 | 262.20 ± 0.16 | 262.20 ± 0.16 |

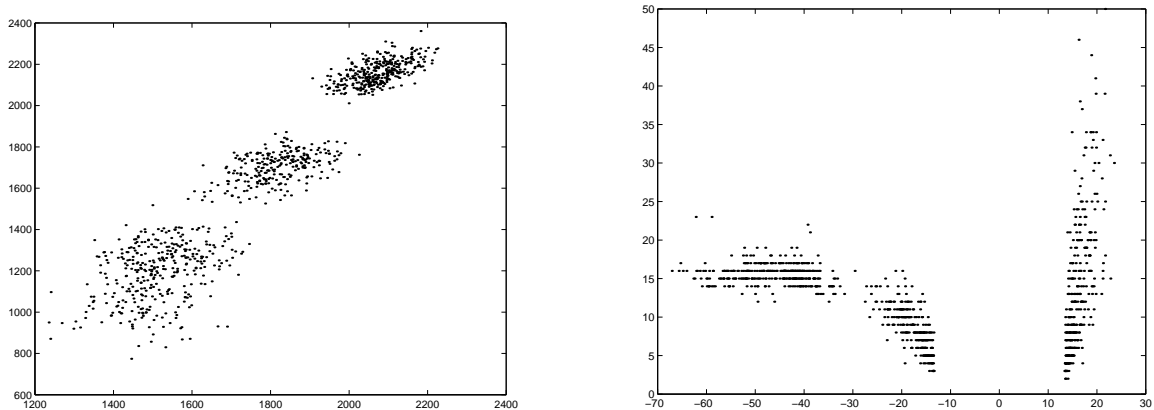
It is seen that FSMEM is at least as proficient than SMEM in avoiding local optima. Moreover, its run-time is much more predictable: these are the iteration count statistics:

| Run | EM | SMEM | FSMEM $K_0 = 1$ | FSMEM $K_0 = 5$ | FSMEM $K_0 = 10$ |
|-----|---------------------|-------------------|------------------|------------------|------------------|
| 01 | 376.48 ± 55.70 | 140.8 ± 29.7 | 289.0 ± 27.7 | 289.0 ± 27.7 | 289.0 ± 27.7 |
| 02 | 782.42 ± 294.98 | 179.4 ± 44.2 | 206.6 ± 24.6 | 206.6 ± 24.6 | 206.6 ± 24.6 |
| 03 | 29.59 ± 238.59 | 351.0 ± 149.5 | 263.4 ± 19.4 | 263.4 ± 19.4 | 263.4 ± 19.4 |
| 04 | -197.61 ± 48.78 | 509.0 ± 104.5 | 275.4 ± 36.9 | 275.4 ± 36.9 | 275.4 ± 36.9 |
| 05 | 406.12 ± 150.98 | 381.0 ± 200.4 | 228.2 ± 44.8 | 228.2 ± 44.8 | 228.2 ± 44.8 |
| 06 | 136.20 ± 297.57 | 259.8 ± 57.5 | 225.6 ± 22.1 | 225.6 ± 22.1 | 225.6 ± 22.1 |
| 07 | 471.46 ± 252.24 | 141.8 ± 21.0 | 194.6 ± 21.5 | 194.6 ± 21.5 | 194.6 ± 21.5 |
| 08 | 90.66 ± 131.87 | 321.6 ± 105.4 | 229.6 ± 21.5 | 229.6 ± 21.5 | 229.6 ± 21.5 |
| 09 | 104.33 ± 98.83 | 184.8 ± 89.3 | 227.2 ± 17.3 | 227.2 ± 17.3 | 227.2 ± 17.3 |
| 10 | 225.83 ± 81.66 | 239.0 ± 148.2 | 233.0 ± 37.0 | 233.0 ± 37.0 | 233.0 ± 37.0 |

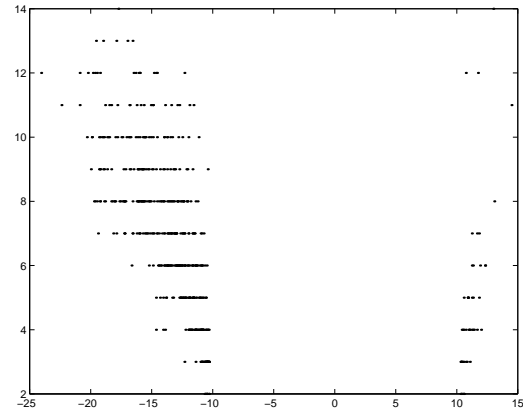
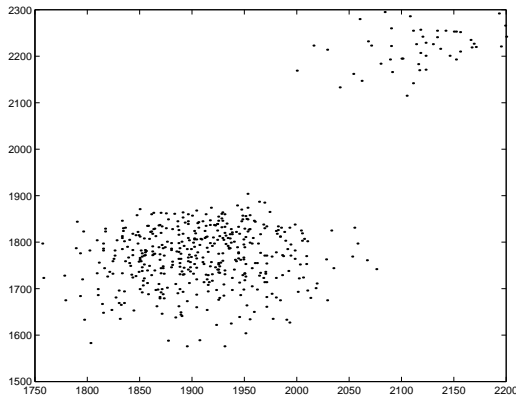
6 Results on real data

For a real test of the SMEM and FSMEM algorithms, I took a dataset of 15000 recorded action potentials from a multi-electrode array with 60 electrodes and several thousand neurons growing on them. These recordings consist of the electric potential measured over a period from 1 ms before to 2 ms after the peak of the action potential. Each electrode typically records from one to three neurons. MoG clustering can potentially be used to classify the recorded spikes according to responsible neuron.

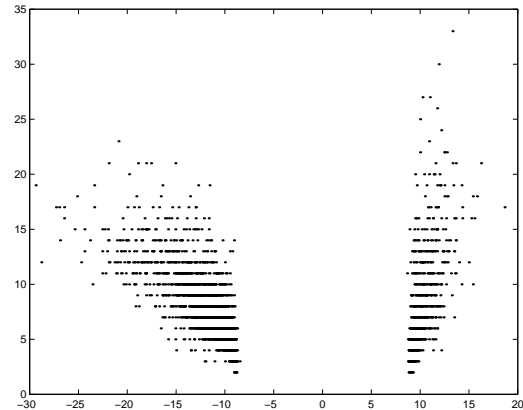
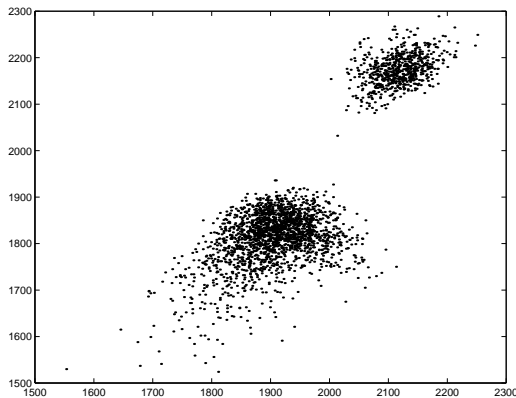
Here are a few examples of clustering results. For each, the clusters found are displayed in two forms: a plot of y_0 vs y_{-5} , i.e. the voltage right at the peak vs the voltage .2 ms earlier, and a plot of spike width versus spikes height. Results for EM and SMEM are not shown, because those algorithms are not able to judge the extract the number of clusters.



Log likelihood: -40204 ± 66 . Iterations: 73 ± 22 . Number of clusters: 3.



Log likelihood: -20170 ± 34 . Iterations: 44 ± 11 . Number of clusters: 2.



Log likelihood: -97930 ± 15 . Iterations: 68 ± 12 . Number of clusters: 3. (The blob on the left consists of two, partially overlapping, clusters.)

In a few cases the algorithm found different number of clusters on different runs (not shown).

Whether this way of clustering spike data is really useful will require some further study. For example, I suspect that in many cases the spikes with positive height (rightmost cluster in the righthand column) may not be physical, while some spikes detected with very low height are a noise artifact.

7 The code

This section presents all the matlab code written for this project. Electronic versions are available upon request.

The outer loop for SMEM:

```

1 function  $y = mog\_smem(X,K,epsi,lambda)$ 
2 % MOG.SMEM performs Ueda et al split/merge EM algorithm for mixture
3 % of gaussians.
4 % Input: X: DxN data
5 %       K: nr of clusters
6 %       epsi: relative change of log likelihood used for termination

```



```

7 %      lambda: fudge parameter to prevent zero–variance attractor
8 % Output: structure with members
9 % p: 1xK vector of mixing coefficients
10 % mu: DxK matrix of means
11 % sig: K cell vector of DxD matrices of variances
12 % likelies : successive local optima of log likelihood.
13 % iters : successive iteration counts of em runs: even numbers
14 %      refer to partial em runs. Discarded runs are included.
15
16 % Some constants:
17 max_iter = 100;
18 max_cands = 5;
19 split_init_epsilon = .1;
20
21 par = mog_init(X,K);
22 likelies = zeros(0,1);
23 iters = zeros(0,1);
24
25 par = mog_fullem(X,par,epsi,lambda);
26 iters = cat(1, iters , par.iters );
27
28 FIG = 1;
29
30 fprintf(1,'First run converged after %i iterations\n',par.iters);
31 figure(FIG);
32 plot_mog(par,X);
33
34 for iter=1:max_iter;
35     FIG = 3 – FIG;
36     figure(FIG);
37     fprintf(1,'mog_smem: major iteration %i\n',iter);
38     merits = mog_merits(X,par);
39
40     likelies = cat(1, likelies , par.likely );
41
42 for c=1:max_cands
43     m1 = merits.Jmerge(c,2);
44     m2 = merits.Jmerge(c,3);
45     for d=1:K
46         s = merits.Jsplit(d,2);
47         if ((s ~=m1) & (s ~=m2))
48             break; % found a candidate
49         end
50     end
51     % So now m1,m2,s are a merge/split triplet
52     disp(m1); disp(m2); disp(s);
53     para = mog_mergeinit(par , m1, m2);
54     para = mog_splitinit(para , s , m2 , split_init_epsilon );
55
56     fprintf(1,'Trying candidate %i:\n',c);
57     plot_mog(para,X);

```

```

58
59  para = mog_partialem(X,para,epsi,lambda,[m1,m2,s]);
60  iters = cat(1, iters , para. iters );
61
62  fprintf(1,'Partial EM converged after %i\n',para.iters);
63  plot_mog(para,X);
64
65  para = mog_fullem(X,para,epsi,lambda);
66  iters = cat(1, iters , para. iters );
67
68  fprintf(1,'Full EM converged after %i\n',para.iters);
69  plot_mog(para,X);
70
71  if (para.likely > par.likely )
72      fprintf(1,'Improvement found from candidate %i\n',c);
73      break; % out of candidate test loop
74  end
75  end % for candidates
76  if (para.likely <= par.likely )
77      fprintf('No more improvements available after %i major iters\n',iter);
78      break; % out of main loop: split/merge didn't produce better results
79  end
80  par = para;
81  end
82
83  fprintf(1,'Final result:\n');
84  plot_mog(par,X);
85
86  y=par;
87  y.likelies = likelies ;
88  y.iters = iters ;

```

The outer loop for FSMEM:

```

1  function y = mog_fsmem(X,K,epsi,lambda)
2  % MOG.FSMEMDW performs free split/merge EM: DWs modification of
3  % Ueda et al split/merge EM algorithm for mixture of gaussians.
4  % Input: X: DxN data
5  %       K: nr of clusters
6  %       epsi: relative change of log likelihood used for termination
7  %       lambda: fudge parameter to prevent zero-variance attractor
8  % Output: structure with members
9  % p: 1xK vector of mixing coefficients
10 % mu: DxK matrix of means
11 % sig: K cell vector of DxK matrices of variances
12 % likelies : successive local optima of log likelihood.
13 % iters : successive iteration counts of em runs: even numbers
14 %       refer to partial em runs. Discarded runs are included.
15
16 % Some constants:
17 max_iter = 100;
18 max_cands = 5;

```

```

19  split_init_epsilon = 1;
20
21  par = mog_init(X,K);
22  likelies = zeros(0,1);
23  iters = zeros(0,1);
24
25  par = mog_fullem(X,par,epsilon,lambda);
26  iters = cat(1, iters , par. iters );
27  likelies = cat(1, likelies , par. likely );
28
29  FIG = 1;
30
31  fprintf(1,'First run converged after %i iterations\n',par.iters);
32  figure(FIG);
33
34  next = 1;
35  fail = .5;
36
37  for iter=1:max_iter;
38      plot_mog(par,X);
39      FIG = 3 - FIG;
40      figure(FIG);
41      fprintf(1,'mog_smem: major iteration %i\n',iter);
42
43      merits = mog_merits(X,par);
44
45      if (next > 0)
46          % try to merge
47          merge_cands = size(merits.Jmerge,1);
48          if (merge_cands > max_cands)
49              merge_cands = max_cands;
50          end
51          para. likely = -1e9;
52          for c=1:merge_cands
53              m1 = merits.Jmerge(c,2);
54              m2 = merits.Jmerge(c,3);
55
56              para = mog_mergeinit(par , m1, m2);
57              para = mog_truncate(para , m2);
58
59              fprintf(1,'Trying merge candidate %i [%i %i]:\n',c,m1,m2);
60              plot_mog(para,X);
61
62              para = mog_partialem(X,para,epsilon,lambda,[m1]);
63              iters = cat(1, iters , para. iters );
64
65              fprintf(1,'Partial EM converged after %i\n',para.iters);
66              plot_mog(para,X);
67
68              para = mog_fullem(X,para,epsilon,lambda);
69              iters = cat(1, iters , para. iters );

```

```

70
71     fprintf(1,'Full EM converged after %i\n',para.itors);
72     plot_mog(para,X);
73
74     if (para.likely > par.likely )
75         fprintf(1,'Improvement found from merge candidate %i. New K is %i\n',c,length(para.p));
76         break; % out of candidate test loop
77     end
78     end % for merge candidates
79     if (para.likely > par.likely )
80         par = para;
81         likelies = cat(1, likelies , par.likely );
82         fprintf(1,'Setting par to para from merge. New K is %i\n',length(par.p));
83         fail = .5;
84     else
85         fail = fail + 1;
86         next = -1;
87     end
88     else
89         % try to split
90         split_cands = size(merits. Jsplit ,1);
91         if ( split_cands > max_cands)
92             split_cands = max_cands;
93         end
94         for c=1:split_cands
95             s = merits. Jsplit (c ,2);
96             para = mog_extend(par);
97             s2 = length(para.p);
98             para = mog_splitinit (para , s , s2 , split_init_epsilon );
99             fprintf(1,'Trying split candidate %i [%i]:\n',c,s);
100            plot_mog(para,X);
101
102            para = mog_partialem(X,para,epsilon ,lambda ,[s,s2 ]);
103            iters = cat(1, iters , para. iters );
104
105            fprintf(1,'Partial EM converged after %i\n',para.itors);
106            plot_mog(para,X);
107
108            para = mog_fullem(X,para,epsilon ,lambda);
109            iters = cat(1, iters , para. iters );
110
111            fprintf(1,'Full EM converged after %i\n',para.itors);
112            plot_mog(para,X);
113
114            if (para.likely > par.likely )
115                fprintf(1,'Improvement found from split candidate %i. new K is %i\n',c,s2);
116                break; % out of candidate test loop
117            end
118        end % for split candidates
119        if (para.likely > par.likely )
120            par = para;

```

```

121     likelies = cat(1, likelies , par.likely );
122     fprintf(1,'Setting par to para from split. New K is %i\n',length(par.p));
123     fail = .5;
124     else
125         fail = fail +1;
126         next = 1;
127     end
128 end
129
130 if ( fail > 2)
131     fprintf('No more improvements available after %i major iters\n',iter);
132     break; % out of main loop: split nor merge produced better results
133 end
134 end
135
136 fprintf(1,'Final result:\n');
137 plot_mog(par,X);
138
139 y=par;
140 y.likelies = likelies ;
141 y.iters = iters ;

```

These use the following code for full EM:

```

1 function y = mog_fullem(X,par,epsi,lambda)
2 % MOG.FULLEM implements the EM algorithm for mixture of Gaussians.
3 % Input: X: D×N data
4 %     par: structure with members:
5 %         p: 1×K vector of mixing coefficients
6 %         mu: D×K matrix of means
7 %         sig: K cell vector of D×D matrices of variances
8 %         epsi: relative change of log likelihood used for termination
9 %         lambda: fudge parameter to prevent zero–variance attractor
10 % Output: structure with members as par, plus
11 %     likely: log likelihood at end of run
12 % Algorithm: Max Welling, in class notes for CS156b
13 % Coding: Daniel Wagenaar, April–May 2000
14
15 % some constants
16 max_iter = 100;
17
18 % initialize parameters
19 mu = par.mu;
20 p = par.p;
21 sig = par.sig;
22
23 N=size(X,2);
24 D=size(X,1);
25 K=size(p,2);
26
27 old_likely = -1e9;
28

```

```

29 for iter = 1:max_iter
30
31     lastiter = iter ;
32
33     % E step: compute responsibilities
34     norma = zeros(1,N);
35     R = zeros(N,K); % R(n,k) will be the responsibility of cluster k for point n
36     px = zeros(1,N);
37     for k=1:K
38         % This section is copied verbatim from mog_responsibility.m
39         s = sig{k};
40         siginv = inv(s);
41         detsig = det(s);
42         dx = X - repmat(mu(:,k),[1 N]);
43         sdx = siginv * dx;
44         expo = -.5 * sum(dx .* sdx ,1);
45         G = (1/sqrt(2*pi*detsig )) * exp(expo); % a 1xN vector of p(x|k)
46         pG = p(k)*G;
47         px = px + pG;
48         norma = norma + pG;
49         R(:,k) = pG';
50     % End of copied section
51     end
52     likely = sum(log(px+1e-300));
53
54     if (abs(( likely - old_likely )/ likely ) < epsi )
55         break;
56     end
57     old_likely = likely ;
58
59     for k=1:K
60         R(:,k) = R(:,k) ./ ( norma+1e-300)';
61     end
62
63     % M step: recompute mu, sig, p
64     for k=1:K
65         sumR = sum(R(:,k));
66         mu(:,k) = ( X*R(:,k))./ sumR; % DxN * Nx1 = Dx1
67         dx = X - repmat(mu(:,k),[1 N]);
68         Rdx = repmat(R(:,k)',[D 1]) .* dx; % DxN
69         sig{k} = (Rdx*dx' + lambda*eye(D))/(sumR+lambda);
70         % sig{k} = (Rdx * dx')/sumR; % DxN * NxD = DxD
71     end
72     p = mean(R,1);
73
74 end
75
76 % return parameters
77 y.p = p;
78 y.mu = mu;
79 y.sig = sig ;

```

```

80 y. likely = likely - .5*log(N)*K*(1+D+.5*D*(D+1)); % Subtract MDL term
81 y. iters = lastiter ;
82 return;

```

And this for partial EM:

```

1 function y = mog_partialem(X,par,epsi,lambda,idx)
2 % MOG_PARTIAL implements the partial update alg for MoG from Ueda et al
3 % Input: X: D×N data
4 %     par: structure with members:
5 %     p: 1×K vector of mixing coefficients
6 %     mu: D×K matrix of means
7 %     sig: K cell vector of D×D matrices of variances
8 %     epsi: relative change of log likelihood used for termination
9 %     lambda: fudge parameter to prevent zero–variance attractor
10 %     idx: index vector specifying which clusters should be updated
11 % Output: structure with members as par
12 % Algorithm: Max Welling, in class notes for CS156b and Ueda et al
13 % Coding: Daniel Wagenaar, April–May 2000
14
15 % some constants
16 max_iter = 100;
17
18 % initialize parameters
19 mu = par.mu;
20 p = par.p;
21 sig = par.sig;
22
23 N=size(X,2);
24 D=size(X,1);
25 K=size(p,2);
26
27 % Compute initial responsibilities
28 R = mog_responsibility (X,par);
29 primalR = sum(R(:,idx),2); % N×1 vector of primary total responsibilities
30
31 old_likely = -1e9;
32
33 for iter = 1: max_iter
34
35     lastiter = iter ;
36
37     % E step: compute responsibilities
38     norma = zeros(1,N);
39     R = zeros(N,K); % R(n,k) will be the responsibility of cluster k for point n
40     px = zeros(1,N);
41     for k = idx
42         % This section is copied verbatim from mog_responsibility.m
43         s = sig{k};
44         siginv = inv(s);
45         detsig = det(s);
46         dx = X - repmat(mu(:,k),[1 N]);

```

```

47   sdx = sginv * dx;
48   expo = -.5 * sum(dx .* sdx ,1);
49   G = (1/sqrt(2*pi*detsig )) * exp(expo); % a 1xN vector of p(x|k)
50   pG = p(k)*G;
51   px = px + pG;
52   norma = norma + pG;
53   R(:,k) = pG';
54   % End of copied section
55   end
56   likely = sum(log(px+1e-300)); % This is only the 'local' likelihood
57   % within the idx group
58
59   if (abs(( likely - old_likely )/ likely ) < epsi )
60       break;
61   end
62   old_likely = likely ;
63
64   norma = primalR ./ (norma+1e-300)';
65   for k = idx
66       R(:,k) = R(:,k) .* norma;
67   end
68
69   % M step: recompute mu, sig, p
70   for k = idx
71       sumR = sum(R(:,k));
72       mu(:,k) = ( X*R(:,k) ) ./ sumR; % DxN * Nx1 = Dx1
73       dx = X - repmat(mu(:,k),[1 N]);
74       Rdx = repmat(R(:,k)',[D 1]) .* dx; % DxN
75       sig{k} = ( Rdx*dx' +lambda*eye(D))/(sumR+lambda);
76       % sig{k} = (Rdx * dx')/sumR; % DxN * Nx1 = Dx1
77       p(k) = mean(R(:,k));
78   end
79
80 end
81
82 % return parameters
83 y.p = p;
84 y.mu = mu;
85 y.sig = sig;
86 y.iters = lastiter ;
87 return;

```

Initialization is done using:

```

1 function y = mog_init(X,K)
2 % MOG.INIT provides initialization for mixture of Gaussians EM.
3 % Input: X: DxN data
4 %       K: number of clusters
5 % Output: Structure that can be passed to mog_fullem
6 % Coding: Daniel Wagenaar, April–May 2000
7
8

```



```

9 % initialize parameters
10 N=size(X,2);
11 D=size(X,1);
12 p = (1/K) * ones(1,K);
13 idx = floor(rand(1,K)*N+1);
14 mu = X(:,idx);
15 datvar = diag(var(X')); % a DxD diagonal matrix containing the data variance
16 for k=1:K
17     sig{k} = datvar / 40; % This is arbitrary , and not very clever
18 end
19
20 % return parameters
21 y.p = p;
22 y.mu = mu;
23 y.sig = sig;
24 return;

```

This calculates the split criterion:

```

1 function J = mog_splitmerit(X,par,R)
2 % MOG_SPLITMERIT computes the split merit matrix for MoG SMEM
3 % Input: X, par: as for mog_fullem
4 %       R: NxK responsibility matrix as from mog_responsibility
5 % Output: Kx1 split merit vector
6 % Algorithm: Ueda et al
7 % Coding: DW
8
9 % Jspl_k := int dx f_k(x) log (f_k(x) / p(x|theta_k)),
10 % where f_k(x) := (sum_n delta(x-x_n) R_kn) / (sum_n R_kn)
11 % This can be re-expressed as
12 % Jspl_k := sum_n f_kn log (f_kn / p_nk), [1]
13 % where f_kn := R_kn / sum_m R_km and p_nk = p(x|theta_k).
14 % This is not quite true: int dx delta(x-y) log delta(x-y)*f(x) [2]
15 % isn't quite equal to log f(y), but I suspect Ueda et al haven't
16 % been careful about this either. Strictly , of course, the integral
17 % [2] diverges. I suspect that the divergent term originates from
18 % the fact that the KL distance between a continuous and a discreet
19 % pdf is necessarily divergent. I hope the above prescription [1] is
20 % the sensible one.
21
22 mu = par.mu;
23 p = par.p;
24 sig = par.sig;
25
26 N=size(X,2);
27 D=size(X,1);
28 K=size(p,2);
29
30 J=zeros(K,1);
31
32 for k=1:K
33     % This section is copied verbatim from mog_responsibility.m

```

```

34 s = sig{k};
35 siginv = inv(s);
36 detsig = det(s);
37 dx = X - repmat(mu(:,k),[1 N]);
38 sdx = siginv * dx;
39 expo = -.5 * sum(dx .* sdx ,1);
40 G = (1/sqrt(2*pi*detsig )) * exp(expo); % a 1xN vector of p(x|k)
41 % End of copied section
42
43 f = R(:,k) ./ sum(R(:,k));
44 G = G + .00001;
45 idx = find(f>.00001);
46 J(k) = sum(f(idx).*log(f(idx)./G(idx) '));
47 end

```

And this the merge criterion:

```

1 function J = mog_mergemerit(X,par,R)
2 % MOG_MERGEMERIT computes the merge merit matrix for MoG SMEM
3 % Input: X, par: as for mog_fullem
4 %       R: NxK responsibility matrix as from mog_responsibility
5 % Output: KxK symmetric matrix of merge merits
6 % Algorithm: Ueda et al
7 % Coding: DW
8
9 J = R'*R;

```

These are called through a small wrapper:

```

1 function y = mog_merits(X,par)
2 % MOG_MERITS returns sorted lists of split and merge candidates.
3 % Input: X, par as per mog_fullem
4 % Output: structure with members
5 %   Jsplitted : Kx2 matrix of [merit; clusterno] rows
6 %   Jmerge: (K*(K-1)/2)x3 matrix of [merit, clusterno_1, clusterno_2] rows
7 % By definition clusterno_1 < clusterno_2 for all reported
8 % candidates.
9 % Note that split and merge merits cannot be usefully added
10 % directly: they are unfortunately defined on different scales.
11
12 K=size(par.p,2);
13
14 R = mog_responsibility (X,par);
15
16 Jm = mog_mergemerit(X,par,R);
17 Jsplitted = mog_splitmerit (X,par,R);
18
19 Jmerge = zeros(K*(K-1)/2,3);
20 idx = 1;
21 for k=1:K
22   for l=(k+1):K
23     Jmerge(idx ,:) = [ Jm(k,l) , k , l ];
24     idx = idx + 1;
25   end

```

```

26 end
27
28 Jsplit = cat(2, Jsplit ,[1: K]');
29
30 [dummy, idx] = sort(-Jmerge (:,1));
31 y.Jmerge = Jmerge(idx ,:);
32 [dummy, idx] = sort(-Jsplit (:,1));
33 y.Jsplit = Jsplit (idx ,:);

```

Initialization after a merge is done by:

```

1 function y = mog_mergeinit(par , k , l , epsi)
2 % MOG_MERGEINIT(par,k,l,epsi) initializes the parameters for clusters k
3 % and l as the merges of cluster k and l into one. The result is
4 % stored in cluster k, and p(l) is set to zero.
5 % Input: par: parameters as for all mog fns
6 %       k: source cluster 1, also dest. cluster
7 %       l : source cluster 2
8 % Output: structure of updated parameters
9 % Alg: Ueda
10 % Coding: DW
11
12 y=par;
13
14 y.p(k) = par.p(k)+par.p(l);
15 y.mu(:,k) = ( par.mu(:,k) + par.mu(:,l ) ) ./ 2;
16 y.sig{k} = ( par.sig{k} + par.sig{l} ) ./ 2;
17 y.p(l) = 0;

```

And after a split by:

```

1 function y = mog_splitinit (par , k , l , epsi)
2 % MOG_SPLITINIT(par,k,l,epsi) initializes the parameters for clusters k
3 % and l as the split of cluster k into two. The original cluster l
4 % is destroyed. Normalization is lost unless par.p(l) was zero
5 % before the call.
6 % Input: par: parameters as for all mog fns
7 %       k: source cluster , also dest. cluster 1
8 %       l : destination cluster 2
9 %       epsi: scale of noise to be added to mean
10 % Output: structure of updated parameters
11 % Alg: Ueda / DW
12 % Coding: DW
13
14 D=size(par.mu,1);
15
16 y=par;
17 y.p(k) = par.p(k)/2;
18 y.p(l) = y.p(k);
19
20 [ U DD V ] = svd(par.sig{k});
21 stddev = sqrt(diag(DD));
22
23 % Move centres away according to principal components:

```

```

24 y.mu(:,k) = par.mu(:,k) + epsi.*U*(stddev.*randn(D,1));
25 y.mu(:,l) = par.mu(:,k) + epsi.*U*(stddev.*randn(D,1));
26
27 % Set sigma to maximal component (really primitive, I admit):
28 y.sig{k} = DD(1,1) * eye(D);
29 y.sig{l} = y.sig{k};

```

Helper functions to change the number of clusters are

```

1 function y = mog_extend(par)
2 % y=MOG_EXTEND(par) returns par unchanged, except that a new cluster
3 % with p=0, mu=0, sig=1 is added.
4 % Input: par: as for mog_fullem
5 % Output: as input
6 % Coding: DW
7
8 K = length(par.p);
9 D = size(par.mu,1);
10 y.p = cat(2,par.p,[0]);
11 y.mu = cat(2,par.mu,zeros(D,1));
12 for k=1:K
13     y.sig{k} = par.sig{k};
14 end
15 y.sig{K+1} = eye(D);

```

and

```

1 function y = mog_truncate(par,k0)
2 % y=MOG_TRUNCATE(par,k0) returns par unchanged, except that all
3 % information pertaining to cluster k0 is removed.
4 % Input: par: as for mog_fullem
5 %     k0: cluster to be removed
6 % Output: as input
7 % Coding: DW
8 % Note: the final cluster is renumbered as k0.
9
10 K = length(par.p);
11 Knew = K-1;
12 y.p = par.p(:,1:Knew);
13 y.mu = par.mu(:,1:Knew);
14 for k = 1:Knew
15     y.sig{k} = par.sig{k};
16 end
17 if (k0<K)
18     y.p(:,k0) = par.p(:,K);
19     y.mu(:,k0) = par.mu(:,K);
20     y.sig{k0} = par.sig{K};
21 end

```

This calculates the responsibilities for use outside the inner loops:

```

1 function R = mog_responsibility(X,par)
2 % MOG_RESPONSIBILITY performs the e step of mog. Using this
3 % function within a loop will be slow. However, it is nice and easy
4 % to use for non-looping calcs e.g. at the start of partialem and

```

```

5 % when calculating split/merge merits.
6 % Input: X, par as for mog_fullem
7 % Output: R: NxK responsibility matrix
8 % Algorithm: Max Welling
9 % Coding: DW
10
11 mu = par.mu;
12 p = par.p;
13 sig = par.sig;
14
15 N=size(X,2);
16 D=size(X,1);
17 K=size(p,2);
18
19 norma = zeros(1,N);
20 R = zeros(N,K); % R(n,k) will be the responsibility of cluster k for point n
21 px = zeros(1,N);
22 for k=1:K
23     s = sig{k};
24     siginv = inv(s);
25     detsig = det(s);
26     dx = X - repmat(mu(:,k),[1 N]);
27     sdx = siginv * dx;
28     expo = -.5 * sum(dx .* sdx ,1);
29     G = (1/sqrt(2*pi*detsig )) * exp(expo); % a 1xN vector of p(x|k)
30     pG = p(k)*G;
31     px = px + pG;
32     norma = norma + pG;
33     R(:,k) = pG';
34 end
35
36 for k=1:K
37     R(:,k) = R(:,k) ./ ( norma+1e-300)';
38 end

```

This makes a graphical representation of the current stage, and prints various parameters as well:

```

1 function plot_mog(y,X)
2 % PLOT_MOG(y,X) plots the MoG results in y. Use y=em_mog(X,K) to
3 % fill y. The plot is a projection onto the first two axes.
4
5 K = length(y.p);
6
7 displ(X,K,y.mu,y.sig); % see below
8 %return;
9 disp(y.p);
10 disp(y.mu);
11 sigs = zeros(size(y.mu,1),0);
12 for k=1:K
13     sigs = cat(2, sigs , y.sig{k});
14 end

```

```

15 disp(sigs);
16 fprintf(1,'Press enter to continue...\n');
17 pause;
18 return;
19
20 function displ(X,K,mu,sig)
21 % DISPL(X,K,mu,sig) is a helper function for plot_mog to
22 % plot the distribution.
23 % X: DxN
24 % K: 1x1
25 % mu: DxK
26 % sig: K cells of DxD
27
28
29 plot(X(1,:), X(2,:), 'b. ');
30 hold on;
31 for k=1:K
32     s = sig{k};
33     plotGauss(mu(1,k),mu(2,k),s(1,1), s(2,2), s(1,2));
34 end
35 hold off;
36 drawnow;
37 return;

```

And finally, this generates a test data set:

```

1 function y=mog_gendata(D,N,K,sig)
2 % y=MOG_GENDATA(D,N,K,sig) generates DxN data points from a mixture of
3 % K gaussians with typical variance set by sig.
4 % y.X is the DxN data
5 % y.par contains the parameters:
6 % p: 1xK
7 % mu: DxK
8 % sig: { K } DxD
9
10
11 X=zeros(D,0);
12 N0=1;
13 Nr=N;
14 Kr=K;
15
16 MEAN=zeros(D,0);
17 P=zeros(1,0);
18
19 for k=1:K
20     % Compute variance matrix
21     Rot = eye(D);
22     for d=1:D
23         for e=(d+1):D
24             theta = rand(1) * 2*pi;
25             rot = eye(D);
26             rot(d,d) = cos(theta);

```

```

27     rot(e,e) = cos(theta);
28     rot(d,e) = sin(theta);
29     rot(e,d) = -sin(theta);
30     Rot = rot*Rot;
31     end
32 end
33 Rot = sig*Rot*diag(randn(D,1) + 1);
34 SIG{k} = Rot*Rot';
35 % So now SIG is a randomly rotated covariance matrix
36 Mean = randn(D,1);
37 MEAN = cat(2,MEAN,Mean);
38 Nk = floor((rand(1)+1.5)*(.5*Nr/Kr));
39 if ((Nk>Nr) | (Kr==1))
40     Nk = Nr;
41 end
42 P = cat(2,P,Nk/N);
43 X = cat(2,X,Rot*randn(D,Nk)+repmat(Mean,1,Nk));
44 Nr = Nr - Nk;
45 Kr = Kr - 1;
46 end
47
48 par.mu=MEAN;
49 par.sig = SIG;
50 par.p = P;
51 figure(3);
52 plot_mog(par,X);
53 y.X = X
54 y.par = par;

```

8 References

[1]