



Mean-field learning

Consider a binary latent factor model. This is a model with a vector \mathbf{s} of K binary latent variables, $\mathbf{s} = (s_1, \dots, s_K)$, a real-valued observed vector \mathbf{x} and parameters $\boldsymbol{\theta} = \{\{\boldsymbol{\mu}_i, \pi_i\}_{i=1}^K, \sigma^2\}$. The model is described by:

$$p(\mathbf{s}|\boldsymbol{\pi}) = p(s_1, \dots, s_K|\boldsymbol{\pi}) = \prod_{i=1}^K p(s_i|\pi_i) = \prod_{i=1}^K \pi_i^{s_i} (1 - \pi_i)^{(1-s_i)}$$

$$p(\mathbf{x}|s_1, \dots, s_K, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}\left(\sum_i s_i \boldsymbol{\mu}_i, \sigma^2 I\right)$$

where \mathbf{x} is a D -dimensional vector and I is the $D \times D$ identity matrix. Assume you have a data set of N i.i.d. observations of \mathbf{x} , i.e. $\mathcal{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$.

We will implement generalized EM learning using the fully factored (a.k.a. mean-field) **variational approximation** for the model above. That is, for each data point $\mathbf{x}^{(n)}$, we will approximate the posterior distribution over the hidden variables by a distribution:

$$q_n(\mathbf{s}^{(n)}) = \prod_{i=1}^K \lambda_{in}^{s_i^{(n)}} (1 - \lambda_{in})^{(1-s_i^{(n)})}$$

and find the $\boldsymbol{\lambda}^{(n)}$'s that maximize \mathcal{F}_n holding $\boldsymbol{\theta}$ fixed.

- (a) Write a Matlab function:

```
[lambda,F] = MeanField(X,mu,sigma,pie,lambda0,maxsteps)
```

where `lambda` is $N \times K$, `F` is the lower bound on the likelihood, `X` is the $N \times D$ data matrix (\mathcal{X}), `mu` is the $D \times K$ matrix of means, `pie` is the $1 \times K$ vector of priors on \mathbf{s} , `lambda0` are initial values for `lambda` and `maxsteps` are maximum number of steps of the fixed point equations. You might also want to set a convergence criterion so that if `F` changes by less than some very small number ϵ the iterations halt.

- (b) We have derived the M step for this model in terms of the quantities: \mathbf{X} , $\underline{\mathbf{E}}\mathbf{s} = E_q[\mathbf{s}]$, which is an $N \times K$ matrix of expected values, and $\underline{\mathbf{E}}\mathbf{s}\mathbf{s}^\top$, which is an $N \times K \times K$ array of expected values $E_q[\mathbf{s}\mathbf{s}^\top]$ for each n . The full derivation is provided in Appendix B. Write two or three sentences discussing how the solution relates to linear regression and why.

- (c) Using the above, we have implemented a function:

```
[mu, sigma, pie] = MStep(X,ES,ESS)
```

This can be implemented either taking in `ESS` = a $K \times K$ matrix summing over N the `ESS` array as defined above, or taking in the full $N \times K \times K$ array. This code can be found in Appendix C and can also be found on the web site. Study this code and figure out what the computational complexity of the code is in terms of N , K and D for the case where `ESS` is $K \times K$. Write out and justify the computational complexity; don't assume that any of N , K , or D is large compared to the others.

- (d) Examine the data `images.jpg` shown on the web site (Do **not** look at `genimages.m` yet!). This shows 100 greyscale 4×4 images generated by randomly combining several features and adding a little noise. Try to guess what these features are by **staring at the images**. How many are there? Would you expect factor analysis to do a good job modelling this data? How about ICA? mixture of Gaussians? Explain your reasoning.
- (e) Put the E step and M step code together into a function:
- ```
[mu, sigma, pie] = LearnBinFactors(X,K,iterations)
```
- where  $K$  is the number of binary factors, and `iterations` is the maximum number of iterations of EM. Include a check that `F` increases at every iteration (this is a good debugging tool).
- (f) Run your algorithm for learning the binary latent factor model on the data set generated by `genimages.m`. What features `mu` does the algorithm learn (rearrange them into  $4 \times 4$  images)? How could you improve the algorithm and the features it finds? Explain any choices you make along the way and the rationale behind them (e.g. what to set  $K$ , how to initialize parameters, hidden states, and `lambdas`).
- (g) For the setting of the parameters learned in the previous step, run the variational approximation for just the first data point (i.e. to find  $q_1(\mathbf{s}^{(1)})$ ) (i.e. set  $N = 1$ ). Convergence of a variational approximation results when the value of  $\lambda$ 's as well as  $F$  stops changing. Plot `F` and `log(F(t)-F(t-1))` as a function of iteration number  $t$  for `MeanField`. How rapidly does it converge? Plot `F` for three widely varying `sigmas`. How is this affected by increases and decreases of `sigma`? Why? Support your arguments.
- (h) Describe a (variational) Bayesian method for selecting  $K$ , the number of hidden binary variables in this model. Does your method pose any computational difficulties and if so how would you tackle them?

[Continued . . .]

(a) The Free Energy

$$\begin{aligned} \mathcal{F} &= \sum_{n=1}^N \left\langle \log P(x^{(n)}, s^{(n)} | \theta, \pi) \right\rangle_{q(s^{(n)})} - \left\langle \log q(s^{(n)}) \right\rangle_{q(s^{(n)})} \\ &= \sum_{n=1}^N \underbrace{\left\langle \log \hat{P}(x^{(n)} | s^{(n)}, \theta) \right\rangle_{q(s^{(n)})}}_{①} + \underbrace{\left\langle \log P(s^{(n)} | \pi) \right\rangle_{q(s^{(n)})}}_{②} - \underbrace{\left\langle \log q(s^{(n)}) \right\rangle_{q(s^{(n)})}}_{③}. \end{aligned}$$

Maximizing  $\mathcal{F}$  to get expression of  $\lambda_i$  in

by taking derivative w.r.t.  $\lambda_i$  and set it to 0 to get fixed point.

$$\begin{aligned} ② &\equiv \sum_{n=1}^N \left\langle \log \prod_{k=1}^K \pi_k^{s_k^{(n)}} (1 - \pi_k)^{(1-s_k^{(n)})} \right\rangle_{q(s^{(n)})} = \left\langle \sum_{k=1}^K s_k^{(n)} \log \pi_k + (1 - s_k^{(n)}) \log (1 - \pi_k) \right\rangle_{q(s^{(n)})} \\ &= \sum_{n=1}^N \sum_{k=1}^K \left\langle s_k^{(n)} \right\rangle_{q(s^{(n)})} \log \pi_k + \left\langle 1 - s_k^{(n)} \right\rangle_{q(s^{(n)})} \log (1 - \pi_k) \\ &= \sum_{n=1}^N \sum_{k=1}^K \lambda_k^{(n)} \log \pi_k + (1 - \lambda_k^{(n)}) \log (1 - \pi_k). \end{aligned}$$

$$\begin{aligned} ③ &\equiv \sum_{n=1}^N \left\langle \log \prod_{k=1}^K \lambda_k^{(n)s_k^{(n)}} (1 - \lambda_k)^{(1-s_k^{(n)})} \right\rangle_{q(s^{(n)})} = \sum_{n=1}^N \sum_{k=1}^K \left\langle s_k^{(n)} \right\rangle_{q(s^{(n)})} \log \lambda_k^{(n)} + \left\langle 1 - s_k^{(n)} \right\rangle_{q(s^{(n)})} \log (1 - \lambda_k^{(n)}) \\ &= \sum_{n=1}^N \sum_{k=1}^K \lambda_k^{(n)} \log \lambda_k^{(n)} + (1 - \lambda_k^{(n)}) \log (1 - \lambda_k^{(n)}) \end{aligned}$$

$$\begin{aligned} ① &\equiv \sum_{n=1}^N \left\langle \log \hat{P}(x^{(n)} | s^{(n)}, \theta) \right\rangle_{q(s^{(n)})} = \sum_{n=1}^N \left\langle \log (2\pi\sigma^2)^{-\frac{D}{2}} \exp \left\{ \left( \underline{x} - \sum_{k=1}^K s_k^{(n)} \underline{\mu}_k \right)^T \left( \underline{x} - \sum_{k=1}^K s_k^{(n)} \underline{\mu}_k \right) \right\} \right\rangle_{q(s^{(n)})} \\ &= -\frac{ND}{2} \log 2\pi - ND \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N \left\{ \underline{x}^T \underline{x} + \sum_{m \neq k}^K \underline{\mu}_k^T \underline{\mu}_m \left\langle s_k^{(n)} s_m^{(n)} \right\rangle_{q(s^{(n)})} - 2 \sum_{k=1}^K \underline{\mu}_k^T \underline{x} \left\langle s_k^{(n)} \right\rangle_{q(s^{(n)})} \right\} \end{aligned}$$

where  $\underline{x} \in \mathbb{R}^{D \times 1}$ ,  $\underline{\mu}_k \in \mathbb{R}^{D \times 1}$

$$\text{Note } \left\langle s_k s_m \right\rangle_{q(s)} = \begin{cases} \lambda_k \lambda_m & \text{if } m \neq k \\ \lambda_k & \text{if } m = k. \end{cases}$$

$$= -\frac{ND}{2} \log 2\pi - ND \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N \left\{ \underline{x}^T \underline{x} + \sum_{m \neq k}^K \sum_{k=1}^K \lambda_k \lambda_m \underline{\mu}_k^T \underline{\mu}_m + \sum_{k=1}^K \lambda_k \underline{\mu}_k^T \underline{\mu}_k - 2 \sum_{k=1}^K \underline{\mu}_k^T \underline{x} \lambda_k \right\}$$

$$\frac{\partial \mathcal{F}}{\partial \lambda_i^{(n)}} = 0, \text{ where } i \in \{1, \dots, K\}, n \in \{1, \dots, N\}.$$

$$\frac{\partial \mathcal{F}}{\partial \lambda_i^{(n)}} = \log \frac{\pi_i}{1-\pi_i} - \log \frac{\lambda_i^{(n)}}{1-\lambda_i^{(n)}} - \frac{1}{6^2} \left\{ \left( 2 \sum_{\substack{j=1 \\ j \neq i}}^K \lambda_j^{(n)} \mu_i^\top \mu_j \right) + \mu_i^\top \mu_i - 2 \mu_i^\top \tilde{x}^{(n)} \right\} = 0$$

$$\begin{aligned} \Leftrightarrow \log \frac{\lambda_i^{(n)}}{1-\lambda_i^{(n)}} &= \log \frac{\pi_i}{1-\pi_i} - \frac{1}{6^2} \left\{ \left( \sum_{\substack{j=1 \\ j \neq i}}^K \lambda_j^{(n)} \mu_i^\top \mu_j \right) + \frac{1}{2} \mu_i^\top \mu_i - \mu_i^\top \tilde{x}^{(n)} \right\} \\ &= \log \frac{\pi_i}{1-\pi_i} - \frac{1}{6^2} \left\{ \left( \sum_{j=1}^K \lambda_j^{(n)} \mu_i^\top \mu_j \right) - \lambda_i^{(n)} \mu_i^\top \mu_i + \frac{1}{2} \mu_i^\top \mu_i - \mu_i^\top \tilde{x}^{(n)} \right\} \\ &= \log \frac{\pi_i}{1-\pi_i} - \frac{\mu_i^\top}{6^2} \left\{ \left( \sum_{j=1}^K \lambda_j^{(n)} \mu_j \right) + \left( \frac{1}{2} - \lambda_i^{(n)} \right) \mu_i - \tilde{x}^{(n)} \right\}. \end{aligned}$$

Denote the right hand side by  $\alpha$

$$\begin{aligned} \text{then } \lambda_i^{(n)} &= \text{sigmoid}(\alpha) = \frac{e^\alpha}{1+e^\alpha} = \frac{1}{e^{-\alpha}+1} \\ &= \text{sigmoid} \left[ \log \frac{\pi_i}{1-\pi_i} - \frac{\mu_i^\top}{6^2} \left\{ \left( \sum_{j=1}^K \lambda_j^{(n)} \mu_j \right) + \left( \frac{1}{2} - \lambda_i^{(n)} \right) \mu_i - \tilde{x}^{(n)} \right\} \right] \\ &\quad \text{where } \lambda_j \in \mathbb{R}^{D \times 1} \end{aligned}$$

Rewrite  $\lambda$  and  $\mathcal{F}$  into matrix form:

$$\begin{aligned} \lambda_i^{(n)} &= \text{sigmoid} \left[ \log \frac{\pi_i}{1-\pi_i} - \frac{\mu_i^\top}{6^2} \left\{ \left( \sum_{j=1}^K \lambda_j^{(n)} \mu_j \right) + \left( \frac{1}{2} - \lambda_i^{(n)} \right) \mu_i - \tilde{x}^{(n)} \right\} \right] \\ \lambda_i &= \text{sigmoid} \left[ \log \frac{\pi_i^\top}{1-\pi_i^\top} - \frac{1}{6^2} \left\{ \left( \sum_{j=1}^N \lambda_j^\top \mu_j \right) + \left( \frac{1}{2} - \lambda_i^\top \right) \mu_i^\top - \tilde{x}^\top \right\} \mu_i \right] \\ \mathcal{F} &= \langle \log p(S, X | \theta) \rangle_{q(S)} - \langle \log q(S) \rangle_{q(S)} = \sum_{n=1}^N \langle \log p(S^{(n)}, X^{(n)} | \theta) \rangle_{q(S^{(n)})} - \sum_{n=1}^N \langle \log q(S^{(n)}) \rangle_{q(S^{(n)})} \\ &= -\frac{ND}{2} \log 2\pi - ND \log 6 - \frac{1}{26^2} \left[ \sum_{n=1}^N \tilde{x}^{(n)\top} \tilde{x}^{(n)} + \sum_{i,j=1}^K \mu_i^\top \mu_j \sum_{n=1}^N \langle S_i^{(n)} S_j^{(n)} \rangle_{q(S^{(n)})} - 2 \sum_{i=1}^K \mu_i^\top \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} \tilde{x}^{(n)} \right] \\ &\quad + \sum_{i=1}^K \left[ \log \pi_i \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} + \log (1-\pi_i) \left( N - \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} \right) \right] - \sum_{i=1}^K \sum_{n=1}^N \left[ \lambda_i^{(n)} \log \lambda_i^{(n)} + (1-\lambda_i^{(n)}) \log (1-\lambda_i^{(n)}) \right] \\ &= -\frac{ND}{2} \log 2\pi - ND \log 6 - \frac{1}{26^2} \left[ \sum_{n=1}^N \tilde{x}^{(n)\top} \tilde{x}^{(n)} + \sum_{i,j=1}^K \mu_i^\top \mu_j \sum_{n=1}^N \langle S_i^{(n)} S_j^{(n)} \rangle_{q(S^{(n)})} - 2 \sum_{i=1}^K \mu_i^\top \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} \tilde{x}^{(n)} \right] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N \left[ \lambda_i^{(n)} \log \pi_i + (1-\lambda_i^{(n)}) \log (1-\pi_i) \right] - \sum_{i=1}^K \sum_{n=1}^N \left[ \lambda_i^{(n)} \log \lambda_i^{(n)} + (1-\lambda_i^{(n)}) \log (1-\lambda_i^{(n)}) \right] \end{aligned}$$

(b) Consider the Multivariate Linear Regression in Lecture 1:

$P(y|x, w, \Sigma_y) = \mathcal{N}(y|xw, \Sigma_y)$  where  $y$  is a function of  $x$ , with Gaussian noise, then the ML estimate for  $w$  is:

$$w^{ML} = \hat{w}^T = \left( \sum_i x_i x_i^T \right)^{-1} \sum_i y_i x_i$$

In our case,  $P(x|s_1, \dots, s_k, \mu, \sigma^2) = \mathcal{N}(x|s_i \mu_i, \sigma^2 I)$ , where  $x$  is a linear function of  $\mu_i$ , with Gaussian noise. And in the M-step, update of  $\mu_j$  has the same form of ML estimate.

i.e.  $\frac{\partial}{\partial \mu_j}$  Free Energy =  $\frac{\partial}{\partial \mu_j}$  Log likelihood

And we get the same form as linear regression:

$$\mu_j = \sum_i \left( \sum_n \langle s_i^{(n)} s_i^{(n)T} \rangle_{f(s^{(n)})} \right)^{-1} \sum_n \langle s_i^{(n)} \rangle_{f(s^{(n)})} x^{(n)}$$

$O(K^3) + O(K^2N) + O(KND)$ .

(c)  $\mu = (\text{inv}(E_{SS}) * E_{S'} * X)'$   
 $\sigma = \sqrt{(\text{trace}(X' * X) + \text{trace}(\mu' * \mu * E_{SS}) - 2 * \text{trace}(E_{S'} * X * \mu)) / (N * D)}$   
 $\pi = \text{mean}(E_S, 1)$

$K \times N \quad N \times D \quad D \times K$   
 $O(D^2N) \quad O(K^2D) + O(K^3) \quad O(KND) + O(K^2D)$

\* Inverse of Matrix:  $O(n^3)$ , where  $n$  is the dimension.

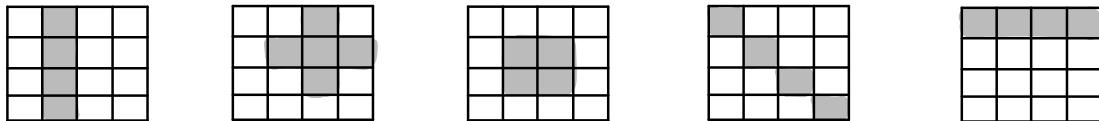
\* Matrix Multiplication:  $A \in \mathbb{R}^{x \times y}$ ,  $B \in \mathbb{R}^{y \times z}$ .

$AB : O(xyZ)$ .

Here we omit the small terms.

The total computational cost is about  $O(K^3) \approx O(KND) \approx O(K^2D)$   
 $\approx O(D^2N) \approx \dots$

(d) Features from staring at the graphs: 5 factors



\* Mixture of Gaussian: NO.

This requires that each cluster follows a circle or ellipsis

But the features doesn't follow this assumptions

Also the features are not continuous but discrete.

\* Factor Analysis: NO

FA requires  $p(x|s)$  and  $p(s)$  to be both Gaussian,

but in this case,  $p(s)$  is Bernoulli. So we can't use FA.

\* ICA: YES.

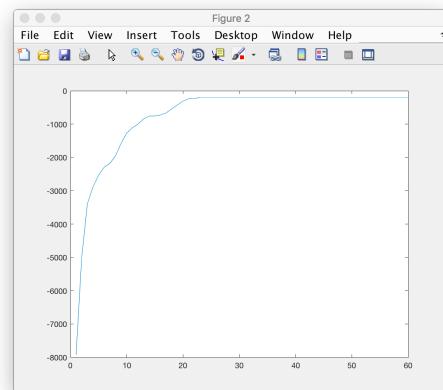
ICA requires that  $x|y$  is Gaussian and  $y$  is non Gaussian,  
which fits the model assumption here.

(e) LearnBinFactors :

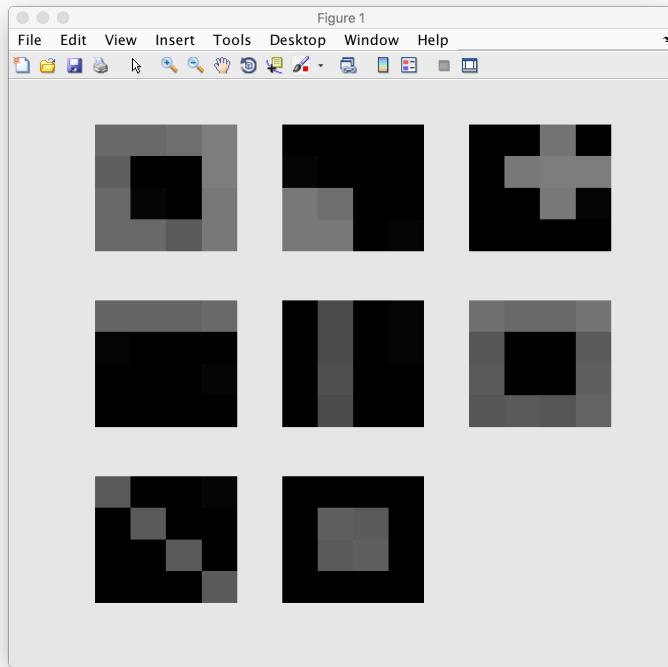
```
function [mu, sigma, pie, FF] = LearnBinFactors(X,K,iterations)
%% INPUT:
% K: number of latent factors
% X: observed data
%% Initialise mu, sigma, pie, ESS (lambda)
[N,D] = size(X);
mu = rand(D,K);
lambda = ones(N,K);
FF = zeros(iterations,1);
sigma = 1;
pie = ones(1,K) * (1/K);
maxsteps = 500; % max 500 steps in Esteps, but usually finishes much early than 100
for iter = 1:iterations
 [lambda,F] = MeanField(X,mu,sigma,pie,lambda,maxsteps);
 FF(iter) = F;
 % check whether Free Energy increases every time
 if (iter > 1) && (F - FF(iter-1)<0)
 disp('Free Energy decreases! Something is wrong!!!!')
 disp(F)
 else
 disp(F)
 end
 ES = lambda;
 [ESS] = CalculateESS(ES);
 [mu, sigma, pie] = MStep(X,ES,ESS);
end
=====

function [ESS] = CalculateESS(ES)
[N,K] = size(ES);
ESS = zeros(K,K,N);
for n = 1:N
 ESS(:,:,n) = ES(n,:)' * ES(n,:);
 for k = 1:K
 ESS(k,k,n) = ES(n,k);
 end
end
ESS = sum(ESS,3);
end
```

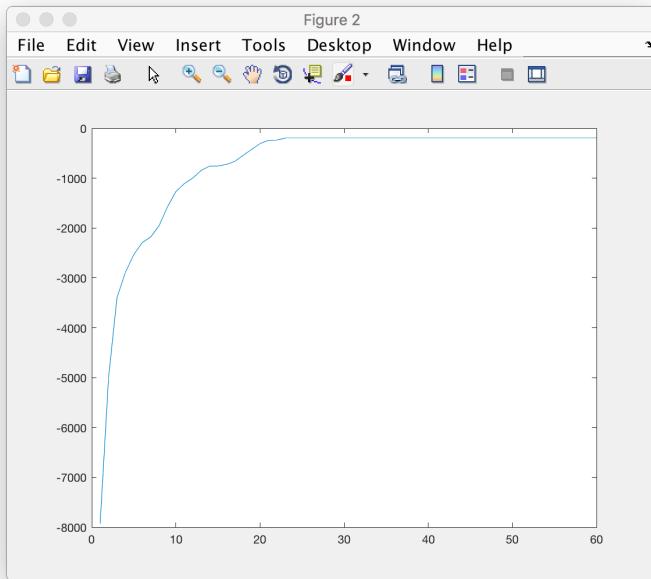
Also do a plot to check that  
↓ free energy always increase.



(f) learnt mu: Successfully learnt 7 out of 8 features.



A check that free energy always increase.



Parameters learnt for (f)

The screenshot shows the MATLAB desktop interface. The Command Window displays the following code and output:

```
>> mu_best
mu_best =
 0.8602 0.0139 -0.0171 0.8015 -0.0098 0.8930 0.7251 -0.0015
 0.7743 0.0633 0.0173 0.0709 -0.0786 0.6918 -0.0448 -0.0640
 0.8732 0.9389 -0.0062 0.0339 0.0015 0.7390 -0.0050 -0.0369
 0.8475 0.9570 -0.0092 0.0332 0.0020 0.7002 -0.0184 -0.0056
 0.8683 -0.0035 0.0055 0.8103 0.6132 0.8719 0.0059 -0.0018
 0.0333 -0.0289 0.9666 0.0137 0.6249 0.0169 0.7224 0.7528
 0.0678 0.8987 -0.0017 0.0075 0.6395 0.0309 0.0093 0.7254
 0.8571 0.9567 -0.0284 0.0491 0.6023 0.7377 -0.0111 -0.0100
 0.8871 -0.0209 0.9279 0.8077 0.0296 0.8174 0.0309 0.0486
 -0.0011 -0.0236 0.9772 0.0397 0.0159 -0.0803 0.0312 0.7231
 -0.0027 -0.0094 0.9681 0.0121 0.0158 -0.0078 0.7349 0.7608
 0.7389 0.0529 0.0082 0.0623 -0.0544 0.6950 -0.0047 -0.0498
 1.0219 0.0401 0.0295 0.8187 -0.0709 0.9217 0.0648 -0.0077
 0.9775 0.0558 0.9728 0.0571 0.0662 0.7438 0.0493 0.0236
 0.9576 0.0618 0.0644 0.0909 -0.0002 0.7621 0.0273 -0.0569
 0.9650 0.0868 0.0427 0.0548 0.0035 0.8122 0.7244 -0.0110
>> pie_best
pie_best =
 0.2350 0.2900 0.2400 0.2825 0.3050 0.0950 0.2925 0.2803
>> sigma_best
sigma_best =
 0.1893
>> |
```

The Current Folder browser on the left shows various MATLAB files in the directory.

\* How to improve the features it find?

Run multiple times with different random initialisation.

As it is likely to stuck at a local optimum.

Then choose the run with the largest free energy.

\* Above initialisations:

K : Test K with different numbers, from 5 to 10, then look at the features it find and whether the features are clear.

I start with K=5, which is the number of features I spotted

(or look at genimages.m, know that we have 8 features...)

$$\pi: \text{Set } \pi = \text{ones}(1, k) * (y_k) \\ = \left(\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k}\right)$$

to give each feature an "equal" chance to be discovered.

$\mu$ : Set  $\mu = \text{rand}(D, k)$ , take uniform values from  $[0, 1]$  at each element.

Because each "pixel" takes value between 0, 1, this is a good initialisation.

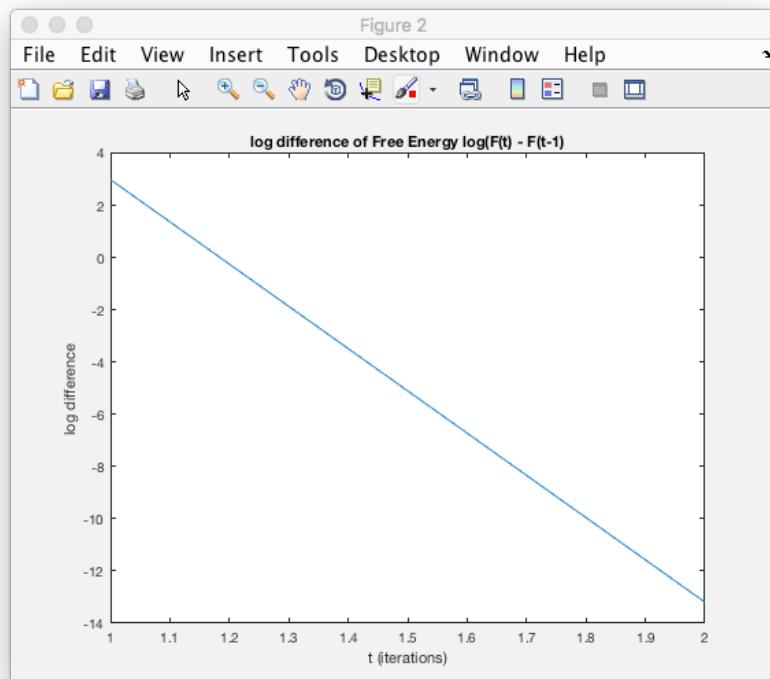
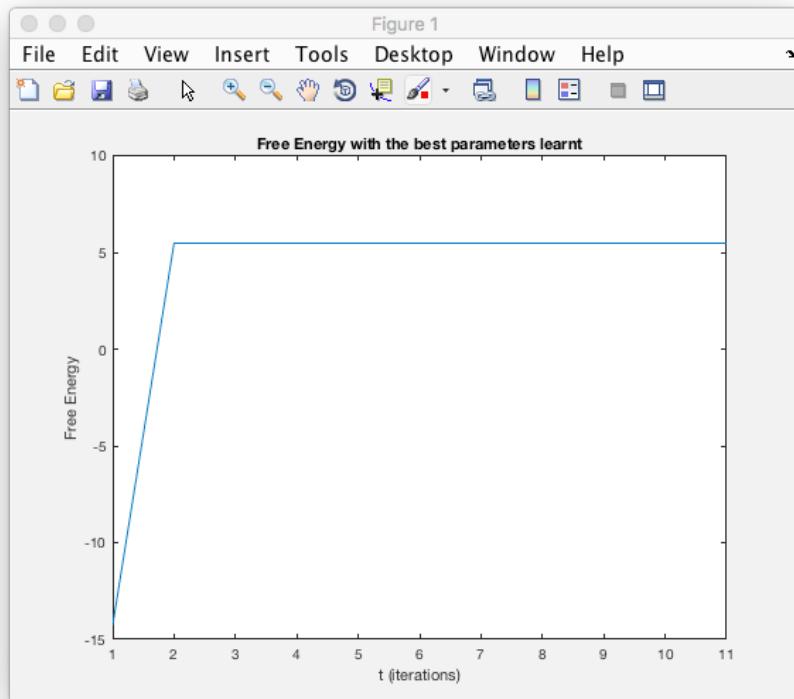
lambda: I tried  $\text{rand}(N, k)$ ,

$$\begin{matrix} \leftarrow & k & \rightarrow \\ \left( \begin{array}{c} \frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k} \\ \vdots \\ \frac{1}{k}, \dots, \frac{1}{k} \end{array} \right) & \xrightarrow{\quad N \quad} & \text{ones}(N, k) \\ \textcircled{1} & & \textcircled{2} \\ & & \textcircled{3} \end{matrix}$$

represent  $(\text{ones}(1, k) * (1/k), N, 1)$ .

I think the second one is the best, because it gives each feature the equal chance to be discovered.

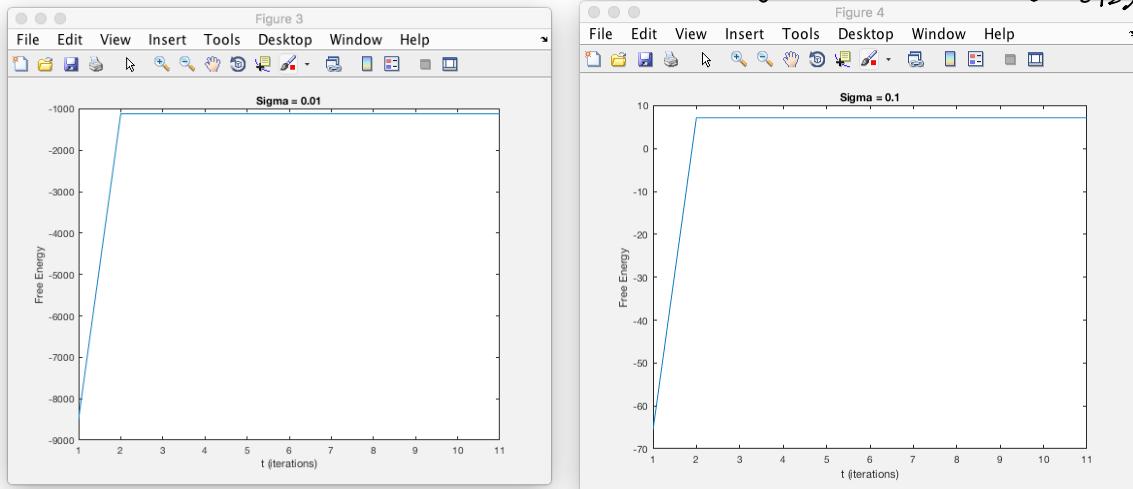
(g) (The parameters learnt are printed in the previous page ...)



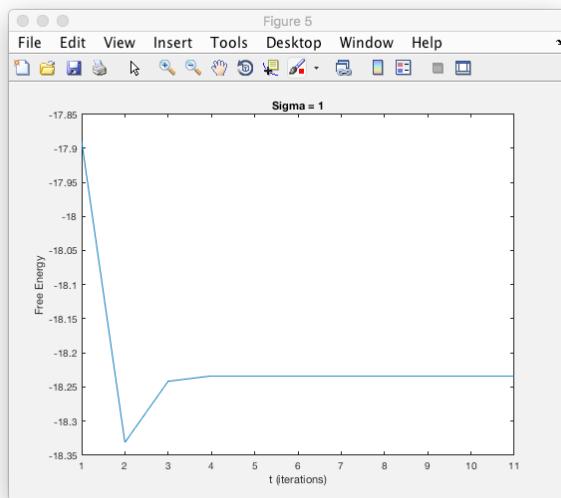
The convergence is RAPID, just at the second step  
 Because we use the parameters learnt from previous question, also we just  
 use a single delta, it is not surprising that the Variational E step  
 converges rapidly.

With different sigmas: 0.01, 0.1 and 1

close to the learnt sigma  
 $\downarrow$   
 $(0.1893)$



Free Energy is smaller with small sigma, but convergence is still rapid.



For a large sigma, it takes longer (4 steps here) to converge.

The free energy decreases to around -18 compared with the previous sigma

(b): By the Variational Bayesian Model Selection described in Lecture 10  
 We factor  $Q(S, \mu) = q(s) q(\mu)$  for variational Bayes.

The VB free Energy in this case:

$$F(q_s(s), q_\mu(\mu), \sigma, \pi) = \left\langle \log P(x, s | \mu, \sigma) + \log P(\mu | \pi) + \log(\sigma) \right\rangle_{q_s, q_\mu} + \dots$$

Then we optimise with respect to  $q(s)$  and  $q(\mu)$  to obtain a lower bound.

Then we do a "Hyper" M-step to optimise hyperparameter  $\lambda$

$$\pi \leftarrow \operatorname{argmax}_{\pi} \left\langle \log P(\mu | \pi) \right\rangle_{q_\mu}.$$

In this step, the dimension of  $s$  is selected, effectively learning the number of  $K$ .

Computational Difficulties:

$$\text{Working with } P(x | z, \mu, \sigma^2) = \mathcal{N}\left(\sum_i s_i \mu_i, \sigma^2 I\right).$$

the computational complexity of inference scales cubically with the number of data points we have. (i.e  $O(n^3)$ ) because we have to do covariance matrix inversion

Solution:

Introduce additionally latent variables

M step derivation:

$$\text{Model: } \begin{cases} P(S^{(n)} | \pi) = \prod_{i=1}^K \pi_i^{S_i^{(n)}} \cdot (1 - \pi_i)^{1 - S_i^{(n)}} \\ P(X^{(n)} | S^{(n)}, \mu, \sigma^2) = (2\pi\sigma^2)^{-\frac{D}{2}} \exp \left[ -\frac{1}{2\sigma^2} \left( X^{(n)} - \sum_{i=1}^K S_i^{(n)} \mu_i \right)^T \left( X^{(n)} - \sum_{i=1}^K S_i^{(n)} \mu_i \right) \right] \end{cases}$$

Visible data  $X = \{X^{(1)}, \dots, X^{(N)}\}$  generated by  $N$  binary vectors:

$S = \{S^{(1)}, \dots, S^{(N)}\}$ , parameters  $\Theta = \{\mu, \sigma^2, \pi\}$  which are constant across all data. Target: We want to find the set of parameters that max likelihood  $P(X|\Theta)$

In EM, the free energy:

$$\begin{aligned} F(q, \Theta) &= \langle \log P(S, X | \Theta) \rangle_{q(S)} - \langle \log q(S) \rangle_{q(S)} \\ &= \sum_n \langle \log P(S^{(n)}, X^{(n)} | \Theta) \rangle_{q(S^{(n)})} - \sum_n \langle \log q(S^{(n)}) \rangle_{q(S^{(n)})} \end{aligned}$$

In M step: we max  $F$  w.r.t each of the params using differentiation only requires  $\Theta$  dependence

$$\sum_n \langle \log P(S^{(n)}, X^{(n)} | \Theta) \rangle_{q(S^{(n)})} = \sum_n \langle \log p(X^{(n)} | S^{(n)}, \Theta) + \log P(S^{(n)} | \Theta) \rangle_{q(S^{(n)})}$$

$$= -\frac{ND}{2} \log 2\pi - ND \log \sigma - (x^{-2})' = -2x^{-3} = \frac{2}{\sigma^3}$$

$$-\frac{1}{2\sigma^2} \left[ \sum_{n=1}^N X^{(n)T} X^{(n)} + \sum_{i,j} \mu_i^T \mu_j \sum_{n=1}^N \langle S_i^{(n)} S_j^{(n)} \rangle_{q(S^{(n)})} - 2 \sum_i \mu_i^T \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} X^{(n)} \right]$$

$$+ \sum_{i=1}^K \left[ \log \pi_i \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} + \log (1 - \pi_i) \left( N - \sum_{n=1}^N \langle S_i^{(n)} \rangle_{q(S^{(n)})} \right) \right]$$

$$* \frac{\partial \bar{F}}{\partial \pi_i} = \frac{1}{\pi_i} \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} + \frac{1}{1-\pi_i} \left[ \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} - N \right] = 0$$

$$\Rightarrow \bar{\pi}_i = \frac{1}{N} \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})}$$

$$* \frac{\partial \bar{F}}{\partial \mu_i} = -\frac{1}{6^2} \sum_{n=1}^N \left[ \sum_j \langle s_i^{(n)} s_j^{(n)} \rangle_{q(\xi^{(n)})} \mu_j - \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} \cdot \bar{x}^{(n)} \right] = 0$$

$$\Leftrightarrow \sum_j \sum_{n=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{q(\xi^{(n)})} \mu_j = \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} \bar{x}^{(n)}$$

$$\Rightarrow \mu_j = \sum_i \left[ \sum_{n=1}^N \langle \xi^{(n)} \xi^{(n)T} \rangle_{q(\xi^{(n)})} \right]^{-1}_{ji} \cdot \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} \bar{x}^{(n)}$$

$$* \frac{\partial \bar{F}}{\partial \sigma^2} = -\frac{ND}{6} + \frac{1}{6^3} \left[ \sum_{n=1}^N \bar{x}^{(n)T} \bar{x}^{(n)} + \sum_{i,j} \mu_i^T \mu_j \cdot \sum_{n=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{q(\xi^{(n)})} - 2 \sum_i \mu_i^T \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} \bar{x}^{(n)} \right] = 0$$

$$\Leftrightarrow \sigma^2 = \frac{1}{ND} \left[ \sum_{n=1}^N \bar{x}^{(n)T} \bar{x}^{(n)} + \sum_{i,j} \mu_i^T \mu_j \cdot \sum_{n=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{q(\xi^{(n)})} - 2 \sum_i \mu_i^T \sum_{n=1}^N \langle s_i^{(n)} \rangle_{q(\xi^{(n)})} \bar{x}^{(n)} \right]$$

Sufficient stats of  $q(\xi)$  are  $\underbrace{\langle \xi^{(n)} \rangle_{q(\xi^{(n)})}}_{ES}, \underbrace{\sum_{n=1}^N \langle \xi^{(n)} \xi^{(n)T} \rangle_{q(\xi^{(n)})}}_{ESS}$ .