

- (a) First, write down the log-joint probability for a single observation-source pair $\log(p(\mathbf{s}, \mathbf{x}))$. Rearrange the terms to form a sum of log-factors on \mathbf{s} (assuming \mathbf{x} is observed), each defined either on a single source variable, or on a pair:

$$\log(p(\mathbf{s}, \mathbf{x})) = \sum_i \log f_i(s_i) + \sum_{ij} \log g_{ij}(s_i, s_j).$$

Relate your result to the Boltzmann Machine. [Remember that, since the sources s are binary, $s_i^2 = s_i$.]

- (b) Next, derive a message passing scheme to find iterative approximations \tilde{f}_i and \tilde{g}_{ij} to each factor. Start your derivation from the KL divergence $\mathbf{KL}[p\|q]$ and identify clearly each time you make an approximate step. You don't need to make all of the EP approximations: which one(s) is(are) missing?

Give the final message-passing scheme in terms of updates to the natural parameters of the site approximations. There will be two different types of update: for the \tilde{f}_i and the \tilde{g}_{ij} respectively.

- (c) Rewrite your message passing approximation to use factored approximate messages. Explain how this leads to a loopy BP algorithm.
 (d) Describe a Bayesian method for selecting K , the number of hidden binary variables using EP. Does your method pose any computational difficulties and if so how would you tackle them?

Implement the EP/loopy-BP algorithm that you derived in the previous question, and compare your results to those of the variational mean-field algorithm.

$$* \log P(x|s) = \log (2\pi 6^2)^{-D/2} \exp \left\{ \left(\underline{x} - \sum_{k=1}^K s_k \mu_k \right)^T \left(\underline{x} - \sum_{k=1}^K s_k \mu_k \right) / 26^2 \right\}$$

$$= -\frac{D}{2} \log 2\pi 6^2 - \frac{1}{26^2} \left\{ \underline{x}^T \underline{x} + \sum_{m,k=1}^K \mu_k^T \mu_m \cdot s_k \cdot s_m - 2 \sum_{k=1}^K \mu_k^T \underline{x} s_k \right\}.$$

$$= -\frac{D}{2} \log 2\pi 6^2 - \frac{1}{26^2} \left\{ \underline{x}^T \underline{x} + \sum_{m \neq k}^K \sum_{k=1}^K s_k \cdot s_m \mu_k^T \mu_m + \sum_{k=1}^K s_k \mu_k^T \mu_k - 2 \sum_{k=1}^K \mu_k^T \underline{x} s_k \right\}$$

$$* \log P(s|\pi) = \log \prod_{k=1}^K \pi_k^{s_k} (1-\pi_k)^{1-s_k}$$

$$= \sum_{k=1}^K s_k \log \pi_k + (1-s_k) \log (1-\pi_k)$$

$$\log P(s, x) = \log P(x|s) + \log P(s|\pi)$$

$$= -\frac{1}{26^2} \left\{ \sum_{m \neq k}^K \sum_{k=1}^K \mu_k^T \mu_m s_k s_m \right\} + \sum_{k=1}^K \left\{ -\frac{1}{26^2} (\mu_k^T \mu_k - 2 \mu_k^T \underline{x}) + \log \pi_k - \log (1-\pi_k) \right\} s_k.$$

+ constant.

$$= \underbrace{\left\{ \sum_{m \neq k}^K \sum_{k=1}^K \left[-\frac{1}{26^2} \mu_k^T \mu_m \right] s_k s_m \right\}}_{①} + \underbrace{\left\{ \sum_{k=1}^K \left[\frac{1}{6^2} \mu_k^T \underline{x} - \frac{1}{26^2} \mu_k^T \mu_k + \log \pi_k - \log (1-\pi_k) \right] s_k \right\}}_{②}$$

+ constant.

$$= \left\{ \sum_{i \neq j} W_{ij} s_i s_j \right\} + \left\{ \sum_i \theta_i s_i \right\} \text{ Boltzmann Machine.}$$

$$\text{where } W_{ij} = -\frac{1}{26^2} \mu_i^T \mu_j, \quad \theta_i = \frac{1}{6^2} \mu_i^T \underline{x} - \frac{1}{26^2} \mu_i^T \mu_i + \log \pi_i - \log (1-\pi_i)$$

$$\Rightarrow P(x|s) \propto \exp \left(\sum_{i \neq j} W_{ij} s_i s_j \right) \cdot \exp \left(\sum_i \theta_i s_i \right). = \prod_{j \neq i} f_{ij}(s_i, s_j) \prod_i f_i(s_i).$$

$$= \prod_{i \neq j} \exp(W_{ij} s_i s_j) \cdot \prod_i \exp(\theta_i s_i).$$

$$f_{ij}(s_i, s_j) = \frac{1}{Z} \left(\frac{s_i}{s_j} \right)^{\frac{1}{26^2}}$$

$$Z = \prod_{i=1}^n \left(\frac{1}{s_i} \right)^{\frac{1}{26^2}}$$

(b) Note that the ② part in (a) can be regarded as a linear combination of s_k (sufficient statistics) and natural parameters. Therefore it is in exponential family. So we don't need to approximate f_i .

But the ① part in (a) is not in exponential family.

We use EP/loopy BP to approximate ① part.

Let $\tilde{g}_{ij}(s_i, s_j) = \exp(\beta_{ji}s_i + \beta_{ij}s_j)$ be Bernoulli distribution to approximate

$g_{ij}(s_i, s_j) = \exp(w_{ij}s_i s_j)$ to ensure tractability.

$$\text{Note } P(s) = \prod_i f_i(s_i) \cdot \prod_{i,j} g_{ij}(s_i, s_j)$$

$$\text{Rewrite } q(s) = \prod_i f_i(s_i) \prod_{i,j} \tilde{g}_{ij}(s_i, s_j).$$

The cavity distribution:

$$q_{\neg g_{ij}}(s) = \prod_k f_k(s_k) \cdot \prod_{(m,n) \neq (i,j)} \tilde{g}_{mn}(s_m, s_n)$$

$$\Rightarrow \tilde{g}_{ij}(s_i, s_j) = \underset{\tilde{g}}{\operatorname{argmin}} KL \left\{ g_{ij}(s_i, s_j) \| \tilde{g}_{ij}(s_i, s_j) \cdot q_{\neg g_{ij}} \right\}$$

Left part of KL:

$$g_{ij}(s_i, s_j) q_{\neg g_{ij}} = \exp \left\{ w_{ij}s_i s_j + \underbrace{\sum_k \theta_k s_k}_{\text{constant}} + \sum_{k \neq i,j} (\beta_{ki}s_i + \beta_{kj}s_j) \right\}$$

constant
cancel from s_i, s_j .

$$\propto \exp \left\{ w_{ij}s_i s_j + \theta_i s_i + \theta_j s_j + \sum_{k \neq i,j} (\beta_{ki}s_i + \beta_{kj}s_j) \right\}$$

Right part of KL:

$$\begin{aligned}
 \tilde{g}_{ij}(s_i, s_j) f_{\neg ij} &= \exp \left\{ \beta_{ji} s_i + \beta_{ij} s_j + \sum_k \theta_k s_k + \sum_{k \neq i, j} (\beta_{ki} s_i + \beta_{kj} s_j) \right\} \\
 &\propto \exp \left\{ \beta_{ji} s_i + \beta_{ij} s_j + \theta_i s_i + \theta_j s_j + \sum_{k \neq i, j} (\beta_{ki} s_i + \beta_{kj} s_j) \right\} \\
 &= \exp \left\{ \left(\beta_{ji} s_i + \theta_i s_i + \sum_{k \neq i, j} \beta_{ki} s_i \right) + \left(\beta_{ij} s_j + \theta_j s_j + \sum_{k \neq i, j} \beta_{kj} s_j \right) \right\} \\
 &= \exp \left\{ \left(\beta_{ji} + \theta_i + \sum_{k \neq i, j} \beta_{ki} \right) s_i + \left(\beta_{ij} + \theta_j + \sum_{k \neq i, j} \beta_{kj} \right) s_j \right\}
 \end{aligned}$$

Then do moment matching:

* In order to match s_i , we integrate out (sum over) s_j

in the left equation: note $s_j = 0$ or $s_j = 1$ (\sum_{s_j} means sum over $s_j=1$ and $s_j=0$)

$$\begin{aligned}
 \sum_{s_j} g_{ij}(s_i, s_j) f_{\neg ij} &= \exp \left\{ w_{ij} + \theta_i s_i + \theta_j + \sum_{k \neq i, j} (\beta_{ki} s_i + \beta_{kj}) \right\} + \exp \left\{ \theta_i s_i + \sum_{k \neq i, j} \beta_{ki} s_i \right\} \\
 &= \exp \left\{ \left(w_{ij} + \theta_i + \sum_{k \neq i, j} \beta_{ki} \right) s_i + \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + \exp \left\{ \left(\theta_i + \sum_{k \neq i, j} \beta_{ki} \right) s_i \right\}
 \end{aligned}$$

$$\begin{aligned}
 \sum_{s_i} \sum_{s_j} g_{ij}(s_i, s_j) f_{\neg ij} &= \exp \left\{ \left(w_{ij} + \theta_i + \sum_{k \neq i, j} \beta_{ki} \right) + \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + \exp \left\{ \left(\theta_i + \sum_{k \neq i, j} \beta_{ki} \right) \right\} \\
 &\quad + \exp \left\{ \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + 1 .
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow E_p(s_i) &= \frac{\sum_{s_i=1} \sum_{s_j} g_{ij}(s_i, s_j) f_{\neg ij}}{\sum_{s_i} \sum_{s_j} g_{ij}(s_i, s_j) f_{\neg ij}} \\
 &= \frac{\exp \left\{ \left(w_{ij} + \theta_i + \sum_{k \neq i, j} \beta_{ki} \right) \cancel{s_i} + \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + \exp \left\{ \left(\theta_i + \sum_{k \neq i, j} \beta_{ki} \right) \cancel{s_i} \right\}}{\exp \left\{ \left(w_{ij} + \theta_i + \sum_{k \neq i, j} \beta_{ki} \right) + \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + \exp \left\{ \left(\theta_i + \sum_{k \neq i, j} \beta_{ki} \right) \right\} + \exp \left\{ \left(\theta_j + \sum_{k \neq i, j} \beta_{kj} \right) \right\} + 1}
 \end{aligned}$$

Note $\underline{\mathbb{E}_f(S_i)} = \underline{\mathbb{E}_p(S_i)}$, where under the distribution f ,

$\tilde{g}_{ij}(S_i, S_j) f_{\pi_{ij}}$ is Bernoulli distribution by our assumption.

$$\begin{aligned} \text{and } \tilde{g}_{ij}(S_i, S_j) f_{\pi_{ij}} &= \exp \left\{ (\beta_{ji} + \theta_i + \sum_{k \neq i,j} \beta_{ki}) S_i + (\beta_{ij} + \theta_j + \sum_{k \neq i,j} \beta_{kj}) S_j \right\} \\ &= \exp \left\{ (\beta_{ji} + \theta_i + \sum_{k \neq i,j} \beta_{ki}) S_i \right\} \cdot \exp \left\{ (\beta_{ij} + \theta_j + \sum_{k \neq i,j} \beta_{kj}) S_j \right\}. \end{aligned}$$

Match the part with S_i Natural Parameter

Recall for Bernoulli Distribution: $p(S|\pi) = \pi^S (1-\pi)^{1-S}$

$$\begin{aligned} &= \exp \left\{ \log \left(\frac{\pi}{1-\pi} \right) S + \log(1-\pi) \right\} \\ &\quad \text{Natural Parameter.} \end{aligned}$$

Also note $\mathbb{E}_p(S) = \pi$

Therefore the update for natural parameters under

$$\begin{aligned} \hat{\beta}_{ji}^* + \theta_i + \sum_{k \neq i,j} \beta_{ki} &= \log \frac{\mathbb{E}_f(S_i)}{1 - \mathbb{E}_f(S_i)}. \end{aligned}$$

Denote $\theta_i + \sum_{k \neq i,j} \beta_{ki} = N_i$

$$\Rightarrow \hat{\beta}_{ji}^* = \log \frac{\mathbb{E}_f(S_i)}{1 - \mathbb{E}_f(S_i)} - \underbrace{\theta_i + \sum_{k \neq i,j} \beta_{ki}}_{= \log(\exp(-N_i))} \quad \theta_j + \sum_{k \neq i,j} \beta_{kj} = N_j$$

$$= \log \left(\frac{\frac{\exp(w_{ij} + N_i + N_j) + \exp(N_i)}{1 + \exp(w_{ij} + N_i + N_j) + \exp(N_i) + \exp(N_j)}}{1 - \frac{\exp(w_{ij} + N_i + N_j) + \exp(N_i)}{1 + \exp(w_{ij} + N_i + N_j) + \exp(N_i) + \exp(N_j)}} \right) + \log(\exp(-N_i))$$

$$= \log \left(\frac{\exp(w_{ij} + N_i + N_j) + \exp(N_i)}{1 + \exp(N_j)} \right) + \log \exp(-N_i)$$

$$= \log \left(\frac{\exp(w_{ij} + N_j) + 1}{1 + \exp(N_j)} \right)$$

Conclude :

$$\beta_{ij}^* = \log \left(\frac{\exp(w_{ij} + \nu_j) + 1}{\exp(\nu_j) + 1} \right), \quad \beta_{ji}^* = \log \left(\frac{\exp(w_{ij} + \nu_i) + 1}{\exp(\nu_i) + 1} \right)$$

$$\text{where: } \nu_i = \theta_i + \sum_{k \neq i, j} \beta_{ki}$$

$$\nu_j = \theta_j + \sum_{k \neq i, j} \beta_{kj}$$

(c) Write pairwise factor $g_{ij}(s_i, s_j)$ into message passing:

$$g_{ij}(s_i, s_j) \approx \tilde{g}_{ij}(s_i, s_j) = M_{i \rightarrow j}(s_j) M_{j \rightarrow i}(s_i)$$

Rewrite $g(s)$ into beliefs

$$g(s) = \prod_i f_i(s_i) \prod_{i,j} \tilde{g}_{ij}(s_i, s_j) = \prod_i \left(f_i(s_i) \prod_{i,j} M_{i \rightarrow j}(s_j) M_{j \rightarrow i}(s_i) \right) = \prod_i b_i(s_i)$$

Cavity distribution can be rewrite into:

$$\begin{aligned} \tilde{g}_{ij}(s_i, s_j) &= \int \prod_k f_k(s_k) \prod_{(m,n) \neq (i,j)} \tilde{g}_{mn}(s_m, s_n) \\ &= f_i(s_i) f_j(s_j) \prod_{\alpha \in \text{ne}(i) \setminus j} M_{\alpha \rightarrow i}(s_i) \prod_{\beta \in \text{ne}(j) \setminus i} M_{\beta \rightarrow j}(s_j) \end{aligned}$$

$$\Rightarrow \left\{ M_{i \rightarrow j}^{\text{new}}, M_{j \rightarrow i}^{\text{new}} \right\} = \underset{M_{i \rightarrow j}, M_{j \rightarrow i}}{\text{argmin}} \text{KL} \left\{ g_{ij}(s_i, s_j) \middle\| \tilde{g}_{ij}(s_i, s_j) \right\} \left\| M_{i \rightarrow j}(s_j) M_{j \rightarrow i}(s_i) \middle\| \tilde{g}_{ij}(s_i, s_j) \right\}$$

$$\text{where } \tilde{g}_{ij}(s_i, s_j) = \exp(\beta_{ji}s_i + \beta_{ij}s_j) = \exp(\beta_{ij}s_j) \cdot \exp(\beta_{ji}s_i)$$

$$= M_{i \rightarrow j}(s_j) \cdot M_{j \rightarrow i}(s_i)$$

This is the same as:

$$\tilde{g}_{ij}(s_i, s_j) = \underset{g}{\text{argmin}} \text{KL} \left\{ g_{ij}(s_i, s_j) \middle\| \tilde{g}_{ij}(s_i, s_j) \right\}$$

We then get updated messages by:

$$\sum_{s_i} f_i(s_i) g_{ij}(s_i, s_j) \prod_{\alpha \in \text{ne}(i) \setminus j} M_{\alpha \rightarrow i}(s_i), \text{ project into exponential family, get } M_{i \rightarrow j}^{\text{new}}$$

Similar for $M_{j \rightarrow i}^{\text{new}}$

(d) By the Automatic Relevance determination (ARD) described in lecture.

Apply it to loopy BP

Then we do a "Hyper" M-step to optimise hyperparameter k

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

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

Solution for computational difficulty:

Introduce additionally latent variables

3. Implement EP.

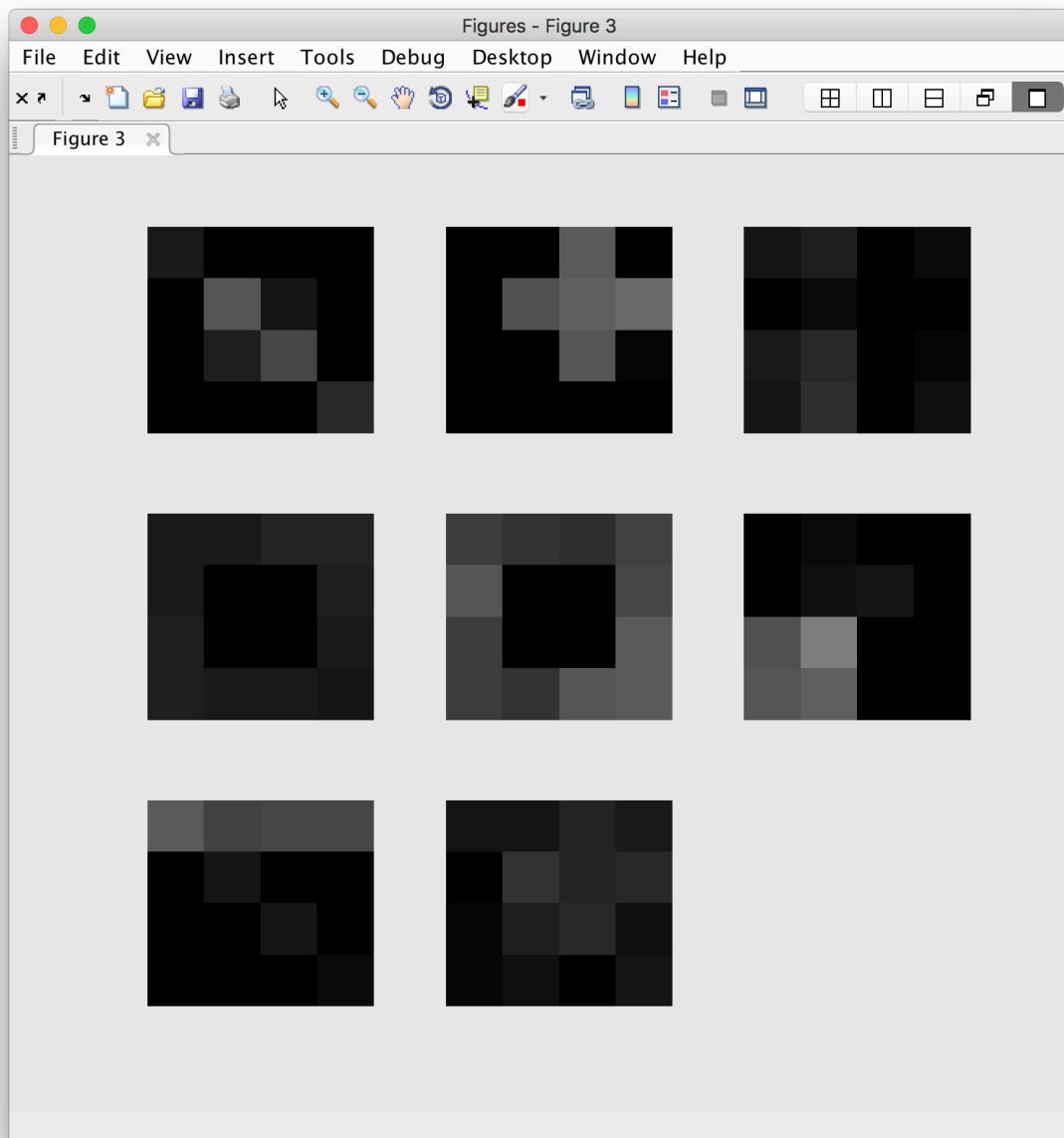
As we derived in 2(b), we have:

$$\beta_{ji}^* + \theta_i + \sum_{k \neq i,j} \beta_{ki} = \log \frac{E_g(s_i)}{1 - E_g(s_i)}.$$

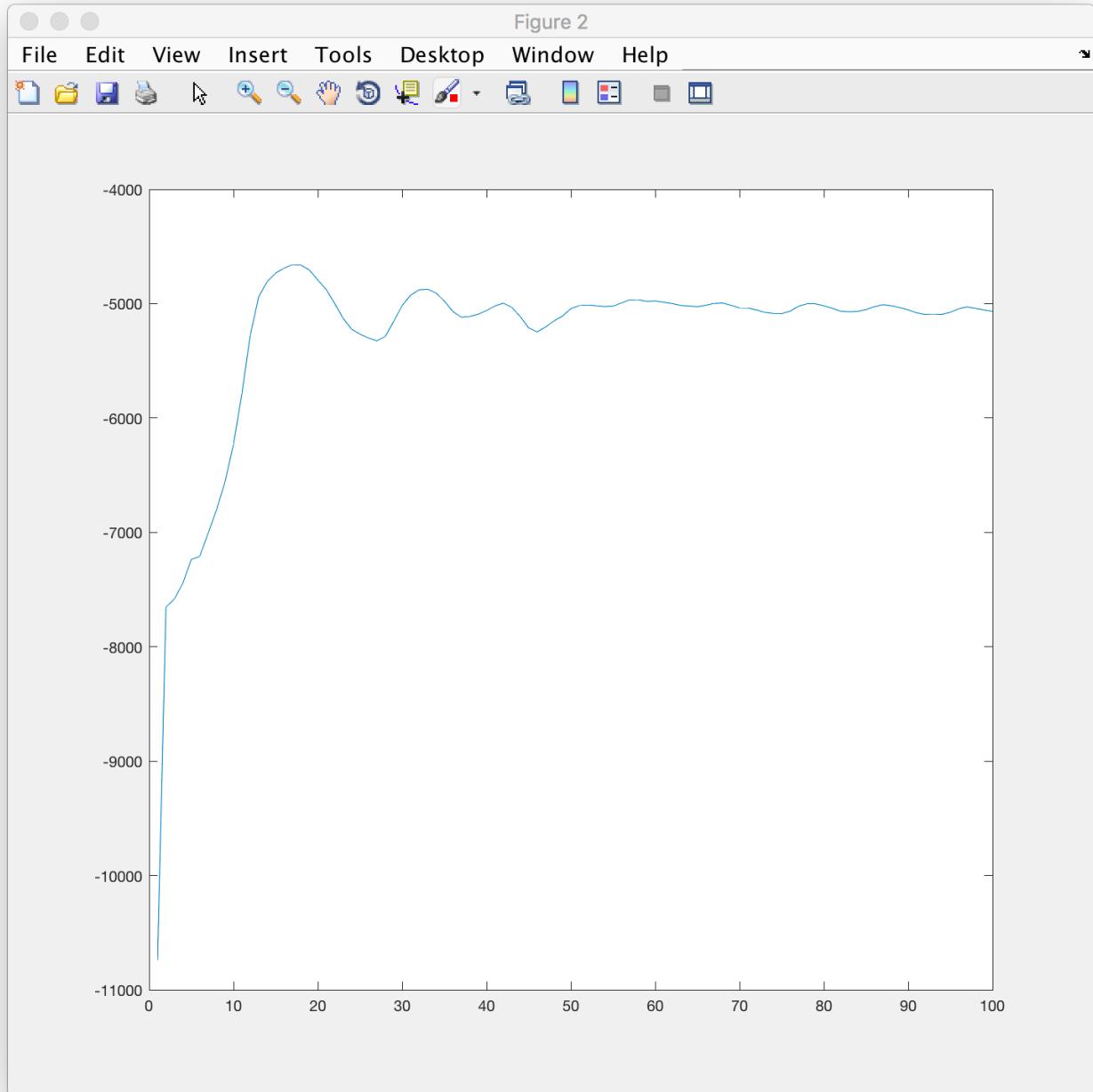
$$\begin{aligned} \text{Therefore } E_g(s_i) &= \text{sigmoid} \left(\beta_{ji}^* + \theta_i + \sum_{k \neq i,j} \beta_{ki} \right) \\ &= \text{sigmoid} \left(\underbrace{\beta_{ji}^* + \theta_i}_{XX(i)} \right). \end{aligned}$$

RESULTS:

1. Features Learnt from EP:

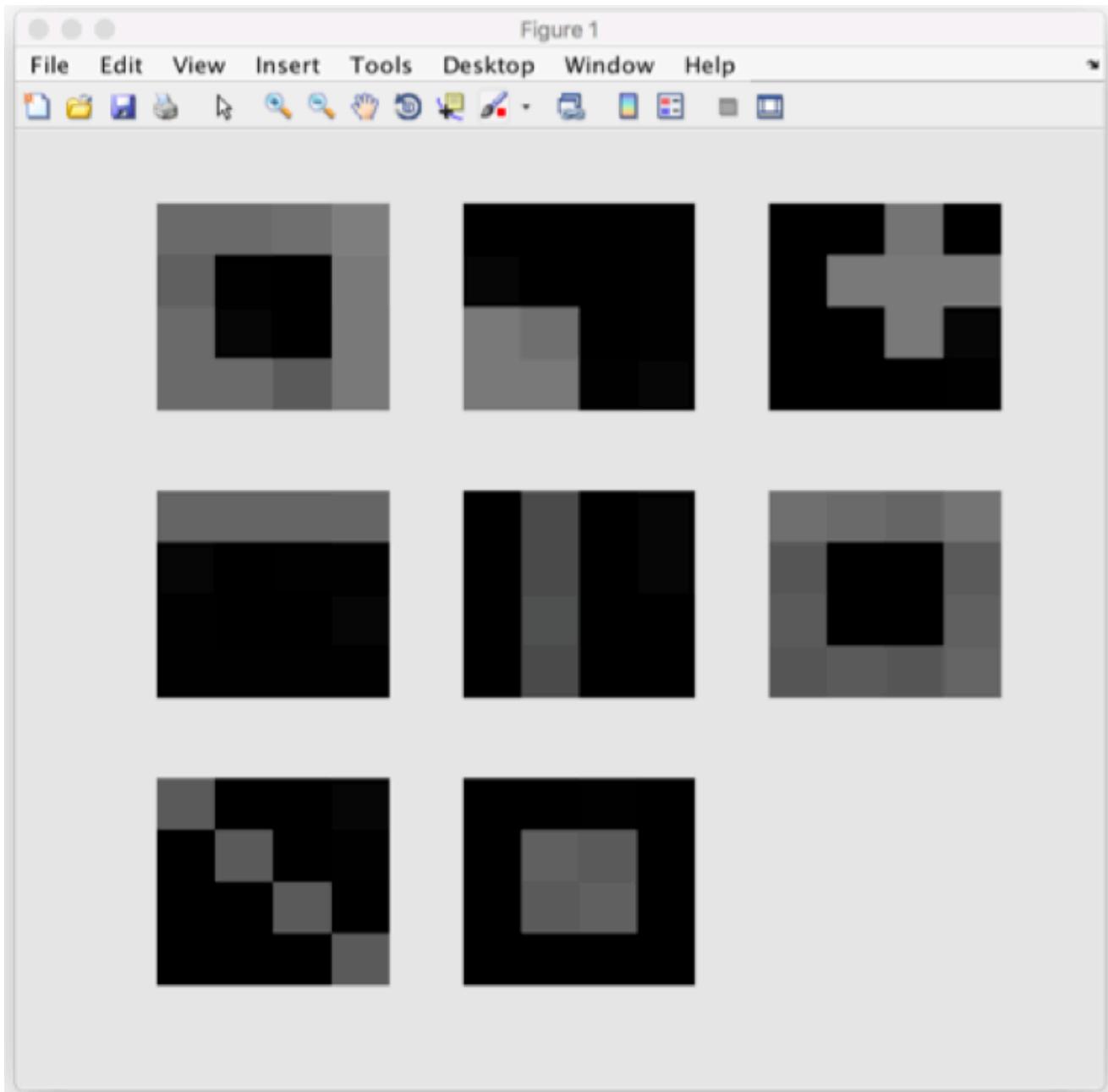


2. Corresponding Free energy with 100 iterations



The 'free energy' for EP does not always increase.

3. Compare the results with Mean Field:



Comparison with features learnt from EP:

It seems that my mean field feature is better than the EP feaures. This might be biased as I ran 100 times mean field in assignment 5 and chose the one with the best free energy. But on average, the features learnt by mean field is not very clear and very similar to the features we learnt by EP here.

In theory, EP should have better results. But here we might have some convergence issues. I have tried power EP update as described in the last slide in the EP lecture, which gives similar results.