

Solution to Problem Set 5

Fall 2012

Issued: Tues. Oct. 16, 2012 **Due:** Tues. Oct. 30, 2012

Reading: Chapters 9, 10, 11, 12

Problem 5.1

EM algorithm for hidden Markov models

- (a) Implement the EM algorithm for HMMs with Gaussian emission probabilities $p(y_t | x_t)$, where $y_t \in \mathbb{R}^2$ and $x_t \in \{0, 1, \dots, m-1\}$. Restrict the covariance matrices to be isotropic (i.e., $\Sigma = \sigma^2 I$).

Solution: Following equations (12.48)–(12.50) in the book with emission probability

$$p(y_t | x_t = i) = \frac{1}{2\pi\sigma_i^2} \exp\left(-\frac{(y_t - \mu_i)^\top (y_t - \mu_i)}{2\sigma_i^2}\right),$$

we can write the complete log likelihood of the data as

$$\begin{aligned} \ell(\theta; x, y) &= \log p(x_0) + \sum_{t=0}^{T-1} \log p(x_{t+1} | x_t) + \sum_{t=0}^{T-1} \log p(y_t | x_t) \\ &= \sum_{i=1}^M q_0^i \log \pi_i + \sum_{t=0}^{T-1} \sum_{i,j=1}^M q_t^i q_{t+1}^j \log a_{ij} \\ &\quad + \sum_{t=0}^T \sum_{i=1}^M q_t^i \left(-\log 2\pi - \log \sigma_i^2 - \frac{1}{2\sigma_i^2} (y_t - \mu_i)^\top (y_t - \mu_i) \right). \end{aligned}$$

Note that there is no $1/2$ in front of the $\log 2\pi$ and $\log \sigma_i^2$ since the Gaussians are 2-dimensional. This affected some people's M-step for updating the σ^2 's, and caused their σ^2 updates to go to infinity.

We see that for the E-step we need to calculate $\gamma_t^i = p(x_t = i | y, \theta^{(p)})$ and $\xi_{t,t+1}^{i,j} = p(x_t = i, x_{t+1} = j | y, \theta^{(p)})$. We compute these using the forward-backward algorithm, as described in Chapter 12.

The M-step for updating A is the same as in the book. The M-step for updating the mean and covariance parameters is identical to the

M-step for a Gaussian (non HMM) mixture model. The M-steps are thus:

$$\begin{aligned}\widehat{\pi}_i^{(p+1)} &= \gamma_0^i \\ \widehat{a}_{ij}^{(p+1)} &= \frac{\sum_{t=0}^{T-1} \xi_{t,t+1}^{i,j}}{\sum_{t=0}^{T-1} \gamma_t^i} \\ \widehat{\mu}_i^{(p+1)} &= \frac{\sum_{t=0}^T \gamma_t^i y_t}{\sum_{t=0}^T \gamma_t^i} \\ \widehat{\sigma}_i^2{}^{(p+1)} &= \frac{\sum_{t=0}^T \gamma_t^i (y_t - \widehat{\mu}_i^{(p+1)})^\top (y_t - \widehat{\mu}_i^{(p+1)})}{2 \sum_{t=0}^T \gamma_t^i}.\end{aligned}$$

All code for this problem set is attached at the end of this document.

- (b) Fit an HMM with $m = 4$ states to the two-dimensional data in *hmm-gauss.dat* and evaluate the log likelihood on the training and test data in *hmm-test.dat*. Plot the data together with the means of the component densities.

Solution: We run the EM algorithm for 5 trials (with random initializations) and choose the highest log likelihood of the training data. We then compute the log likelihood of the test data using the resulting parameters. The plots of the training and test data are shown in Figure 1, along with the component means and the circles indicating one standard deviation away from the mean. We also do the same procedure for the Gaussian mixture model in part (c) below; the plots are also shown in Figure 1.

The resulting estimated parameters for HMM are:

$$\begin{aligned}\pi &= \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, & \sigma^2 &= \begin{pmatrix} 0.9902 \\ 0.9966 \\ 0.8624 \\ 1.0268 \end{pmatrix}, & \mu &= \begin{pmatrix} 1.9896 & -1.9655 \\ -2.0398 & 1.9596 \\ 1.9667 & 2.0889 \\ -1.9979 & -2.1501 \end{pmatrix}, \\ A &= \begin{pmatrix} 0.7266 & 0.0742 & 0.0543 & 0.1448 \\ 0.1332 & 0.6289 & 0.1161 & 0.1217 \\ 0.1211 & 0.0833 & 0.6524 & 0.1432 \\ 0.0870 & 0.0956 & 0.0994 & 0.7180 \end{pmatrix}.\end{aligned}$$

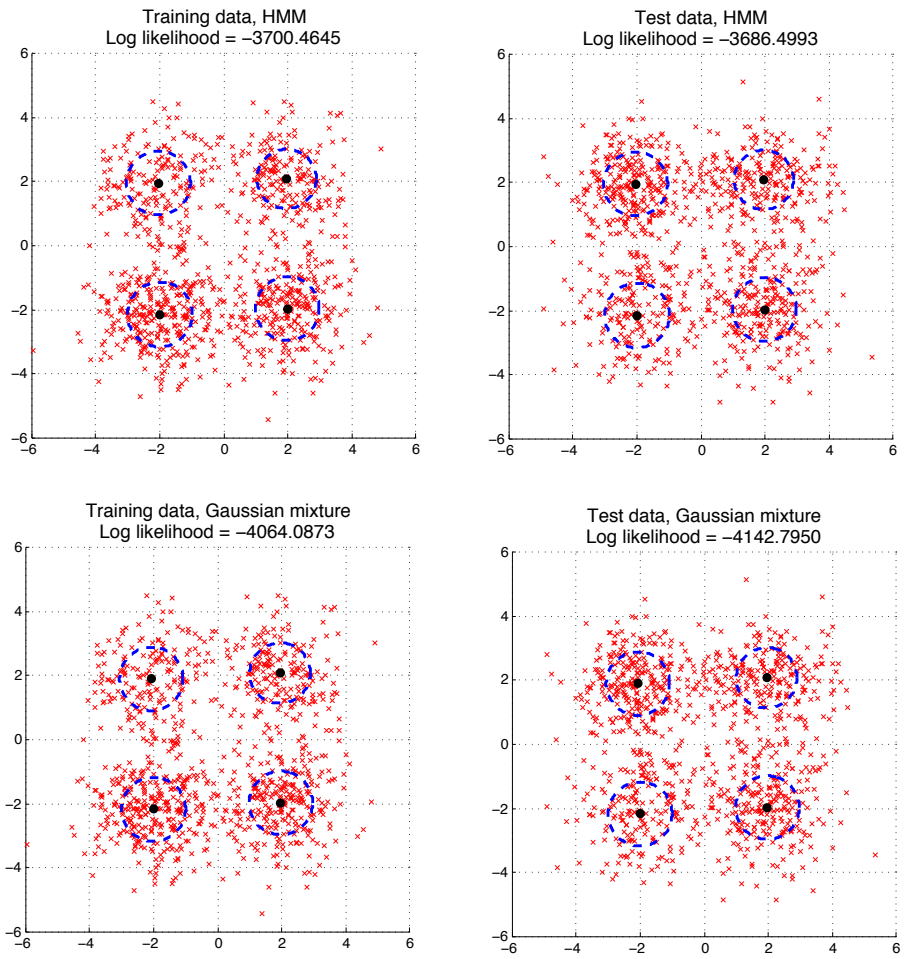


Figure 1: Results of the HMM and Gaussian mixture models.

- (c) Fit a Gaussian mixture model with 4 states to the same data (again with isotropic covariance matrices $\sigma^2 I$). Compare the performance with that of the HMM.

Solution: The complete log likelihood for the Gaussian mixture model is

$$\begin{aligned} \ell(\theta; x, y) &= \sum_{t=0}^T \sum_{i=1}^M x_t^i \log \pi_i - \sum_{t=0}^T \sum_{i=1}^M x_t^i \frac{(y_t - \mu_i)^\top (y_t - \mu_i)}{2\sigma_i^2} \\ &\quad - \sum_{t=0}^T \sum_{i=1}^M x_t^i \log \sigma_i^2 - T \log(2\pi), \end{aligned}$$

so the expected complete log likelihood is

$$\begin{aligned} \mathbb{E}[\ell(\theta; x, y)] &= \sum_{t=0}^T \sum_{i=1}^M \gamma_t^i \log \pi_i - \sum_{t=0}^T \sum_{i=1}^M \gamma_t^i \frac{(y_t - \mu_i)^\top (y_t - \mu_i)}{2\sigma_i^2} \\ &\quad - \sum_{t=0}^T \sum_{i=1}^M \gamma_t^i \log \sigma_i^2 - T \log(2\pi). \end{aligned}$$

E step. The posterior probabilities are

$$\gamma_t^i = p(x_t = i \mid y_t) = \frac{\pi_i \frac{1}{2\pi\sigma_i^2} \exp\left(-\frac{(y_t - \mu_i)^\top (y_t - \mu_i)}{2\sigma_i^2}\right)}{\sum_{j=1}^M \pi_j \frac{1}{2\pi\sigma_j^2} \exp\left(-\frac{(y_t - \mu_j)^\top (y_t - \mu_j)}{2\sigma_j^2}\right)}.$$

M step. By the same argument as in part (a), we have the update rules

$$\begin{aligned} \widehat{\pi}^{(p+1)} &= \frac{1}{T} \sum_{t=0}^T \gamma_t \\ \widehat{\mu}_i^{(p+1)} &= \frac{\sum_{t=0}^T \gamma_t^i y_t}{\sum_{t=0}^T \gamma_t^i} \\ \widehat{\sigma}_i^{2(p+1)} &= \frac{\sum_{t=0}^T \gamma_t^i (y_t - \widehat{\mu}_i^{(p+1)})^\top (y_t - \widehat{\mu}_i^{(p+1)})}{2 \sum_{t=0}^T \gamma_t^i}. \end{aligned}$$

The resulting log likelihoods are:

	training	test
HMM	-3700.46	-3686.50
Gaussian mixture	-4064.09	-4142.80

We see that the HMM yields higher log likelihood, which is to be expected since Gaussian mixture model is a special case of the HMM (when the transition probability $p(x_{t+1} | x_t)$ is independent of x_t). That is, HMM is a more expressive model, so it should be able to fit the data better.

Problem 5.2

EM and missing values: Suppose you have a random sample of twins and are interested in studying *identical* twins. However, you observe only:

- $m \equiv$ the total number of male twins (both identical and fraternal)
- $f \equiv$ the total number of female twins, and
- $b \equiv$ the number of twins of opposite gender.

Let θ be the probability that a pair of twins are identical. Assume that, given identical twins, the probability the twins are male is p . Given fraternal twins, assume the number of males is Binomial(2, q).

Give an algorithm for calculating the MLEs for θ , p , and q . (*Hint:* If you knew exactly how many identical male and female twins there are, then the MLEs would be easy to calculate.)

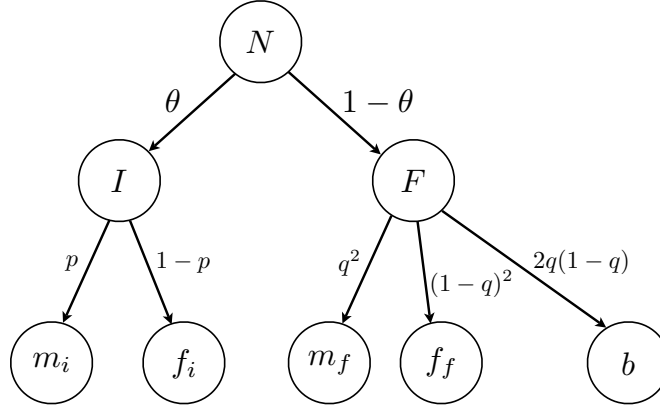
Solution: We introduce the following latent variables:

$$\begin{aligned}
 m_i &= \# \text{ identical male twins} \\
 f_i &= \# \text{ identical female twins} \\
 m_f &= \# \text{ fraternal male twins} \\
 f_f &= \# \text{ fraternal female twins.}
 \end{aligned}$$

For convenience, we write

$$\begin{aligned}
 I &= m_i + f_i = \# \text{ identical twins} \\
 F &= m_f + f_f = \# \text{ fraternal twins} \\
 N &= m + f + b = \# \text{ twins in total.}
 \end{aligned}$$

The generative process is described by the following diagram:



Note that this is not a graphical model describing the distribution. The graphical model has “cliques” $\{N, I, F\}$, $\{I, m_i, f_i\}$, and $\{F, m_f, f_f, b\}$.

The complete data likelihood is

$$\begin{aligned}
 L(\theta, p, q \mid \mathcal{D}) &= \mathbb{P}(I, F \mid N) \mathbb{P}(m_i, f_i \mid I, N) \mathbb{P}(m_f, f_f, b \mid F, N) \\
 &= \left\{ \binom{N}{m_i + f_i} \theta^{m_i + f_i} (1 - \theta)^{m_f + f_f + b} \right\} \left\{ \binom{I}{m_i} p^{m_i} (1 - p)^{f_i} \right\} \\
 &\quad \times \left\{ \binom{F}{m_f, f_f, b} q^{2m_f} (1 - q)^{2f_f} (2q(1 - q))^b \right\} \\
 &= C \theta^{m_i + f_i} (1 - \theta)^{m_f + f_f + b} p^{m_i} (1 - p)^{f_i} q^{2m_f + b} (1 - q)^{2f_f + b},
 \end{aligned}$$

where $C \equiv C(m_i, f_i, m_f, f_f, b)$ does not depend on the parameters θ, p, q . Taking the log, we see that the log likelihood is linear in m_i, f_i, m_f, f_f, b , so they are an appropriate set of sufficient statistics. Note that we have already written the likelihood in a form where the parameters are “decoupled” to make maximization easy.

E-step: For the E-step, we need to calculate the expected value of each of these sufficient statistics. For notational convenience, assume all expectations are taken with respect to the current estimate of the parameters.

Reading from the diagram or by algebraic manipulation, we see that

$$\begin{aligned}
m_i &\sim \text{Multinomial}(m + f + b, \theta p) \\
f_i &\sim \text{Multinomial}(m + f + b, \theta(1 - p)) \\
m_f &\sim \text{Multinomial}(m + f + b, (1 - \theta)q^2) \\
f_f &\sim \text{Multinomial}(m + f + b, (1 - \theta)(1 - q)^2) \\
b &\sim \text{Multinomial}(m + f + b, 2(1 - \theta)q(1 - q)).
\end{aligned}$$

This implies for any particular set of twins,

$$\mathbb{P}(\text{identical} \mid \text{male}) = \frac{\mathbb{P}(\text{male AND identical})}{\mathbb{P}(\text{male})} = \frac{p\theta}{p\theta + q^2(1 - \theta)}.$$

We can calculate similar results given the twins are female or mixed male/female, so we obtain

$$\begin{aligned}
\mathbb{E}[m_i \mid m] &= \frac{mp\theta}{p\theta + q^2(1 - \theta)} \\
\mathbb{E}[f_i \mid f] &= \frac{m(1 - p)\theta}{(1 - p)\theta + (1 - q)^2(1 - \theta)} \\
\mathbb{E}[m_f \mid m] &= m - \mathbb{E}[m_i \mid m] \\
\mathbb{E}[f_f \mid f] &= f - \mathbb{E}[f_i \mid f].
\end{aligned}$$

It is important to actually calculate these expectations and, in this case, to give them in closed form. Some of the main reasons one uses EM are that (1) it is often easy to implement, and (2) the iterations are fast so EM runs in a reasonable amount of time. Leaving the E-step in terms of the conditional probabilities (the maximizing auxiliary distribution) means one has to sum over a set of values, which makes EM too slow.

Of those that calculated expectations, there were two common mistakes. The first was saying $\mathbb{E}[m_i \mid m, f, b] = (m + f + b)p\theta$. This is incorrect since it only uses information about the total number of twins and does not make use of the additional information provided by the split into m, f, b . The second mistake was saying $\mathbb{E}[m_i \mid m, f, b] = m\theta$. There is some subtlety in the generative process since fraternal male twins require two ‘‘coin flips’’ versus one for the identical male twins, and the coin flips may have different probabilities.

M-step: Returning to the complete data likelihood, we see that it has the same form as the product of three binomials with parameters θ, p, q .

Therefore, maximizing the expected complete data (log-)likelihood is trivial:

$$\begin{aligned}\theta &\leftarrow \frac{\mathbb{E}[m_i | m] + \mathbb{E}[f_i | f]}{m + f + b} \\ p &\leftarrow \frac{\mathbb{E}[m_i | m]}{\mathbb{E}[m_i | m] + \mathbb{E}[f_i | f]} \\ q &\leftarrow \frac{2\mathbb{E}[m_f | m] + b}{2\mathbb{E}[m_f | m] + 2\mathbb{E}[f_f | f] + 2b} = \frac{\text{expected \# fraternal males}}{\text{expected \# fraternal siblings}}.\end{aligned}$$

A common mistake was to forget about or incorrectly including b in the update for q .

Problem 5.3

(*Social network analysis and IPF:*) Frank is studying dependencies in voting patterns of a collection of d US senators. For any given bill, he collects a vector $x \in \{0, 1\}^d$, where $x_i = 1$ means that senator i voted yes on that bill. He models the random vector (X_1, X_2, \dots, X_d) as a pairwise Markov random field

$$\mathbb{P}_\theta(x_1, x_2, \dots, x_d) \propto \prod_{(s,t) \in E} \exp(\theta_{st}(x_s, x_t)).$$

- (a) For a set $d = 4$ senators, the data file `Pairwise.dat` contains a 4×30 matrix, summarizing the data from $n = 30$ bills that were voted on in the senate. Implement and apply the IPF updates to estimate the model parameters for each of the following graphs: (i) the graph with edge set $E = \{(12), (23), (34), (14)\}$; and (ii) the graph with edge set $E = \{(12), (23), (13), (14)\}$; and (iii) the fully connected graph with all $\binom{4}{2}$ edges.

Solution: The code for the implementation of the IPF algorithm is provided at the end of this document. We choose to initialize the potential functions as $\theta_{st}(x_s, x_t) = 0$ for all $(s, t) \in E$ so that the normalization constant can be computed explicitly, $Z = 2^4 = 16$. Note that to compute the marginals $p^{(t)}(x_s, x_t)$ we probably should use algorithms such as junction tree, but considering that we are only working with graphs with 4 nodes, we can simply perform the marginalization directly. We also note that the values of θ_{st} are unnormalized, i.e. we can add any constant to any θ_{st} without changing the probability distribution.

(i) $E = \{(12), (23), (34), (14)\}$.

Algorithm converges after 6 batch iterations.¹ Log likelihood = -79.6370 .

Estimated parameters:

		(x_s, x_t)			
		(0,0)	(0,1)	(1,0)	(1,1)
$\theta_{st}(x_s, x_t)$	(12)	0.2880	-0.2237	-0.2237	0.0650
	(23)	0.1180	-0.1338	-0.8478	0.4521
	(34)	-0.0114	0.0116	0.0077	-0.0076
	(14)	0.1183	-0.1342	-0.1549	0.1341

(ii) $E = \{(12), (23), (13), (14)\}$.

Algorithm converges after 21 batch iterations. Log likelihood = -78.3564 .

Estimated parameters:

		(x_s, x_t)			
		(0,0)	(0,1)	(1,0)	(1,1)
$\theta_{st}(x_s, x_t)$	(12)	0.3923	-0.3874	-0.4051	0.2113
	(23)	0.2049	-0.2405	-1.0487	0.5105
	(13)	-0.4326	0.2454	0.3998	-0.3247
	(14)	0.1178	-0.1335	-0.1542	0.1335

(iii) $E = \{(12), (13), (14), (23), (24), (34)\}$.

Algorithm converges after 32 batch iterations. Log likelihood = -76.5029 .

Estimated parameters:

		(x_s, x_t)			
		(0,0)	(0,1)	(1,0)	(1,1)
$\theta_{st}(x_s, x_t)$	(12)	0.4706	-0.5140	-0.5232	0.3103
	(13)	-0.7243	0.4423	0.2229	-0.2071
	(14)	0.2392	-0.2826	-0.3229	0.2668
	(23)	0.4821	-0.4835	-0.9560	0.3961
	(24)	-0.4571	0.3348	0.3720	-0.5313
	(34)	0.2654	-0.2393	-0.1632	0.1728

¹One batch iteration is a cycle through all the edges and all the possible configurations of the edges.

- (b) Of models (i) and (ii), which model has a higher likelihood?

Solution: Model (ii) has a higher likelihood.

- (c) Of all three models (i), (ii) and (iii), which has the highest likelihood? Do you think that it is the “best” model?

Solution: Model (iii) has the highest likelihood. However, this is somewhat vacuous; since it contains all of the edges in both models (i) and (ii), plus additional edges, it is a strictly more expressive model than either of the other two and can capture more complex relationships in the training data. This says nothing about how well it generalizes, as we may have overfit by using such a powerful model without any prior or form of regularization. In practice, we would usually use criteria like the BIC or the AIC to balance the conflicting considerations of maximizing the data likelihood and choosing a model that will generalize well, or employ some expert information to start with a graph that makes practical sense.

Problem 5.4

Model selection for trees: Recall that for a given tree T with edge set $E(T)$, the MLE for the exponential parameters takes the form

$$\begin{aligned}\hat{\theta}_s(x_s) &= \log \hat{\mu}_s(x_s) \quad \text{for all } s \in V, \text{ and} \\ \hat{\theta}_{st}(x_s, x_t) &= \log \frac{\hat{\mu}_{st}(x_s, x_t)}{\hat{\mu}_s(x_s)\hat{\mu}_t(x_t)} \quad \text{for all } (s, t) \in E(T).\end{aligned}$$

Here $\hat{\mu}$ are the empirical marginals computed from the data (e.g., $\hat{\mu}_{st}(j, k) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{st;jk}(x_{is}, x_{it})$), and we assume that they take strictly positive values.

- (a) Define the rescaled log likelihood $\ell(\theta) = \frac{1}{n} \sum_{i=1}^n \log \mathbb{P}_\theta(x_i)$. Letting $\hat{\theta}(T)$ denote the MLE for trees, show that $\ell(\hat{\theta}(T))$ depends on the empirical marginals only via

$$\text{Singleton entropy: } H(\hat{\mu}_s) = - \sum_{x_s} \hat{\mu}_s(x_s) \log \hat{\mu}_s(x_s) \quad \forall s \in V, \text{ and}$$

$$\text{Joint edge entropy: } H(\hat{\mu}_{st}) = - \sum_{x_s, x_t} \hat{\mu}_{st}(x_s, x_t) \log \hat{\mu}_{st}(x_s, x_t) \quad \forall (s, t) \in E(T).$$

Solution: Suppose all the random variables X_s take values in a finite set \mathcal{X} . Note that for each $1 \leq i \leq n$ and $s \in V$, we can write

$$\hat{\theta}_s(x_{is}) = \log \hat{\mu}_s(x_{is}) = \sum_{j \in \mathcal{X}} \mathbb{I}_{s;j}(x_{is}) \log \hat{\mu}_s(j),$$

so

$$\frac{1}{n} \sum_{i=1}^n \hat{\theta}_s(x_{is}) = \frac{1}{n} \sum_{j \in \mathcal{X}} \sum_{i=1}^n \mathbb{I}_{s;j}(x_{is}) \log \hat{\mu}_s(j) = \sum_{j \in \mathcal{X}} \hat{\mu}_s(j) \log \hat{\mu}_s(j) = -H(\hat{\mu}_s).$$

Similarly, since for $1 \leq i \leq n$ and $(s, t) \in E$ we have

$$\begin{aligned} \hat{\theta}_{st}(x_{is}, x_{it}) &= \log \frac{\hat{\mu}_{st}(x_{is}, x_{it})}{\hat{\mu}_s(x_{is}) \hat{\mu}_t(x_{it})} \\ &= \sum_{j, k \in \mathcal{X}} \mathbb{I}_{st;jk}(x_{is}, x_{it}) \log \hat{\mu}_{st}(j, k) - \sum_{j \in \mathcal{X}} \mathbb{I}_{s;j}(x_{is}) \log \hat{\mu}_s(j) \\ &\quad - \sum_{k \in \mathcal{X}} \mathbb{I}_{t;k}(x_{it}) \log \hat{\mu}_t(k), \end{aligned}$$

we can also write

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{st}(x_{is}, x_{it}) &= \frac{1}{n} \sum_{j, k \in \mathcal{X}} \sum_{i=1}^n \mathbb{I}_{st;jk}(x_{is}, x_{it}) \log \hat{\mu}_{st}(j, k) \\ &\quad - \frac{1}{n} \sum_{j \in \mathcal{X}} \sum_{i=1}^n \mathbb{I}_{s;j}(x_{is}) \log \hat{\mu}_s(j) - \frac{1}{n} \sum_{k \in \mathcal{X}} \sum_{i=1}^n \mathbb{I}_{t;k}(x_{it}) \log \hat{\mu}_t(k) \\ &= \sum_{j, k \in \mathcal{X}} \hat{\mu}_{st}(j, k) \log \hat{\mu}_{st}(j, k) - \sum_{j \in \mathcal{X}} \hat{\mu}_s(j) \log \hat{\mu}_s(j) - \sum_{k \in \mathcal{X}} \hat{\mu}_t(k) \log \hat{\mu}_t(k) \\ &= -H(\hat{\mu}_{st}) + H(\hat{\mu}_s) + H(\hat{\mu}_t). \end{aligned}$$

Using the relations above, we can write down the rescaled log-likelihood for the MLE as follows:

$$\begin{aligned} \ell(\hat{\theta}(T)) &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{Z} \prod_{s \in V} \exp(\hat{\theta}_s(x_{is})) \prod_{(s,t) \in E} \exp(\hat{\theta}_{st}(x_{is}, x_{it})) \right) \\ &= -\log Z + \frac{1}{n} \sum_{i=1}^n \left(\sum_{s \in V} \hat{\theta}_s(x_{is}) + \sum_{(s,t) \in E} \hat{\theta}_{st}(x_{is}, x_{it}) \right) \\ &= -\log Z - \sum_{s \in V} H(\hat{\mu}_s) + \sum_{(s,t) \in E} (-H(\hat{\mu}_{st}) + H(\hat{\mu}_s) + H(\hat{\mu}_t)) \\ &= -\log Z + \sum_{s \in V} (\deg(s) - 1) H(\hat{\mu}_s) - \sum_{(s,t) \in E} H(\hat{\mu}_{st}). \end{aligned}$$

In the equations above, Z is a normalizing constant (which is actually equal to 1 in this case) and $\deg(s)$ is the degree of vertex s in the tree.

This shows that $\ell(\hat{\theta}(T))$ only depends on the empirical marginals via the singleton and joint edge entropy. Alternatively, we can also write $\ell(\hat{\theta}(T))$ as

$$\ell(\hat{\theta}(T)) = -\log Z - \sum_{s \in V} H(\hat{\mu}_s) + \sum_{(s,t) \in E} I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t), \quad (1)$$

where $I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t) = H(\hat{\mu}_s) + H(\hat{\mu}_t) - H(\hat{\mu}_{st})$ is the mutual information.

- (b) In the model selection problem for trees, the goal is to choose, from *all trees on d nodes*, the highest likelihood tree i.e., $\hat{T} \in \arg \max_T \ell(\hat{\theta}(T))$. Show how this problem can be cast as a maximum weight spanning tree calculation. This is important, because it allows us to select the best tree by simple algorithms. (*Hint:* The mutual information $I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t) = H(\hat{\mu}_s) + H(\hat{\mu}_t) - H(\hat{\mu}_{st})$ should be relevant in your calculations.)

Solution: Using the representation (1), we see that the term $-\log Z - \sum_{s \in V} H(\hat{\mu}_s)$ is independent of the choice of the tree, so

$$\hat{T} \in \arg \max_T \ell(\hat{\theta}(T)) \iff \hat{T} \in \arg \max_T \sum_{(s,t) \in E} I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t).$$

Therefore, finding \hat{T} is equivalent to finding a maximum weight spanning tree in the complete graph with vertices V , in which each edge (s, t) is assigned weight $I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t)$.

- (c) Use the technique from (b) to select the best fitting tree for the data in `Pairwise.dat`.

Solution: The script for computing the mutual information $I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t)$ is attached at the end of this document. Figure 2 shows the weighted complete graph with the mutual information as the edge weights. The maximum weight spanning tree, which can be found by inspection (or by running Kruskal's algorithm), is shown with solid edges.

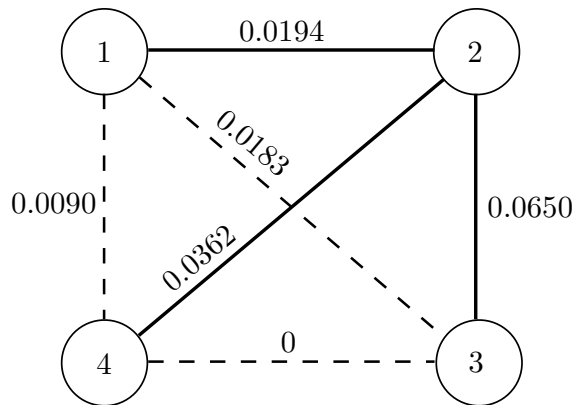


Figure 2: Weighted graph for Problem 5.4(c), with the mutual information $I(\hat{\mu}_{st}; \hat{\mu}_s, \hat{\mu}_t)$ as the edge weights. The maximum weight spanning tree is shown with solid edges.

Code for Problem 5.1

```

1 %% Script for running the computations for problem 5.1
2
3 load hmm-gauss.dat
4 load hmm-test.dat
5
6 numTrials = 5; % number of times to run EM
7
8
9 %% Fit a HMM with 4 states
10
11 % Initialize log likelihood to be a very small number
12 logLhoodHMM = -1e15;
13
14 % Run the EM algorithm numTrials times, choose the highest log likelihood
15 for i = 1:numTrials
16     [pi0_, A_, mu_, sigma2_, logLhood_] = EM_HMM(hmm-gauss, 4);
17     if (logLhood_ > logLhoodHMM)
18         logLhoodHMM = logLhood_;
19         pi0HMM = pi0_;
20         A = A_;
21         muHMM = mu_;
22         sigma2HMM = sigma2_;
23     end
24 end
25

```

```

26 % Compute log likelihood of the test data
27 [¬, ¬, logLhoodHMM.test] = FwBw(hmm_test, pi0HMM, A, muHMM, sigma2HMM);
28
29
30 %% Fit a Gaussian mixture model with 4 states
31
32 % Initialize log likelihood to be a very small number
33 logLhoodMix = -1e15;
34
35 % Run the EM algorithm numTrials times, choose the highest log likelihood
36 for i = 1:numTrials
37     [pi0_, mu_, sigma2_, logLhood_] = EM.GaussianMix(hmm_gauss, 4);
38     if (logLhood_ > logLhoodMix)
39         logLhoodMix = logLhood_;
40         pi0Mix = pi0_;
41         muMix = mu_;
42         sigma2Mix = sigma2_;
43     end
44 end
45
46 % Compute log likelihood of the test data
47 logLhoodMix.test = computeLogLhoodGaussianMix(hmm_test, pi0Mix, muMix, sigma2Mix);
48
49
50 %% Plot training and test data
51
52 % Plot training data for HMM
53 figure;
54 scatter(hmm_gauss(:,1), hmm_gauss(:,2), 'xr');
55 title(sprintf('Training data, HMM\nLog likelihood = %.4f', ...
56     logLhoodHMM), 'FontSize', 13);
57 hold on; grid;
58 axis([-3.5 3.5 -3.5 3.5]); axis square;
59 scatter(muHMM(:,1), muHMM(:,2), 50, 'ok', 'filled');
60 for i = 1:4
61     z = (0:255)*2*pi/255;
62     x = sqrt(sigma2HMM(i))*cos(z) + muHMM(i,1);
63     y = sqrt(sigma2HMM(i))*sin(z) + muHMM(i,2);
64     plot(x, y, '--b', 'LineWidth', 2);
65 end
66
67
68 % Plot test data for HMM
69 figure;
70 scatter(hmm_test(:,1), hmm_test(:,2), 'xr');
71 title(sprintf('Test data, HMM\nLog likelihood = %.4f', ...
72     logLhoodHMM.test), 'FontSize', 13);
73 hold on; grid;
74 axis([-3.5 3.5 -3.5 3.5]); axis square;

```

```

75 scatter(muHMM(:,1), muHMM(:,2), 50, 'ok', 'filled');
76 for i = 1:4
77     z = (0:255)*2*pi/255;
78     x = sqrt(sigma2HMM(i))*cos(z) + muHMM(i,1);
79     y = sqrt(sigma2HMM(i))*sin(z) + muHMM(i,2);
80     plot(x, y, '--b', 'LineWidth', 2);
81 end
82
83
84 % Plot training data for Gaussian mixture model
85 figure;
86 scatter(hmm_gauss(:,1), hmm_gauss(:,2), 'xr');
87 title(sprintf('Training data, Gaussian mixture\nLog likelihood = %.4f', ...
88     logLhoodMix), 'FontSize', 13);
89 hold on; grid;
90 axis([-3.5 3.5 -3.5 3.5]); axis square;
91 scatter(muMix(:,1), muMix(:,2), 50, 'ok', 'filled');
92 for i = 1:4
93     z = (0:255)*2*pi/255;
94     x = sqrt(sigma2Mix(i))*cos(z) + muMix(i,1);
95     y = sqrt(sigma2Mix(i))*sin(z) + muMix(i,2);
96     plot(x, y, '--b', 'LineWidth', 2);
97 end
98
99
100 % Plot test data for Gaussian mixture model
101 figure;
102 scatter(hmm_test(:,1), hmm_test(:,2), 'xr');
103 title(sprintf('Test data, Gaussian mixture\nLog likelihood = %.4f', ...
104     logLhoodMix_test), 'FontSize', 13);
105 hold on; grid;
106 axis([-3.5 3.5 -3.5 3.5]); axis square;
107 scatter(muMix(:,1), muMix(:,2), 50, 'ok', 'filled');
108 for i = 1:4
109     z = (0:255)*2*pi/255;
110     x = sqrt(sigma2Mix(i))*cos(z) + muMix(i,1);
111     y = sqrt(sigma2Mix(i))*sin(z) + muMix(i,2);
112     plot(x, y, '--b', 'LineWidth', 2);
113 end

```

The code for implementing EM for HMM:

```

1 function [pi0, A, mu, sigma2, logLhood] = EM_HMM(data, M)
2 % M is the number of states in the latent variable
3
4 % data is T-by-K
5 T = size(data,1); % number of iid observations
6 K = size(data,2); % dimensionality of the Gaussian

```

```

7
8 % Initialize cluster probabilities to be uniform 1/M
9 pi0 = ones(M,1)/M;
10
11 % Initialize transition probabilities to be uniform
12 A = ones(M,M)/M;
13
14 % Initialize mean values to be random data points -- mu is M-by-K
15 mu = data(randi(T,M,1), :);
16
17 % Initialize sigma^2 to be the average distance between the data points and
18 % one random data point
19 centered = data - data(randi(T,1,1)*ones(T,1), :);
20 sigma2_init = (sum(sum(centered.^2))/(2*T));
21 sigma2 = sigma2_init*ones(M,1);
22
23 % Compute gamma and xi via forward-backward
24 [gamma, xi, logLhood] = FwBw(data, pi0, A, mu, sigma2);
25
26 % Maximum number of iterations
27 maxIter = 1000;
28
29 % Iterate over the E and M steps
30 for iter = 1:maxIter
31     % Update parameters using gamma and xi
32     pi0 = gamma(:,1);
33
34     % Update transition probabilities matrix
35     for i = 1:M
36         for j = 1:M
37             A(i,j) = sum(xi(i,j,:))/sum(gamma(i,1:end-1));
38         end
39     end
40
41     % Update mean parameters
42     for i = 1:M
43         mu(i,:) = sum(data.*gamma(i*ones(1,K),:))'/sum(gamma(i,:));
44     end
45
46     % Update variance parameters - use the new mu values
47     for i = 1:M
48         centered = data - mu(i*ones(T,1), :); % centered data
49         sigma2(i) = sum(sum(centered.^2,2).*gamma(i,:))'/(2*sum(gamma(i,:)));
50     end
51
52     % Then recompute gamma, xi, and log likelihood
53     logOld = logLhood;
54     [gamma, xi, logLhood] = FwBw(data, pi0, A, mu, sigma2);
55

```



```

56     % Check for convergence
57     if (logLhood - logOld < 1e-14)
58         fprintf('Converges after %d batch iterations!\n', iter);
59         break;
60     elseif (iter == maxIter)
61         fprintf('EM did not converge: final  $\Delta = %.8f$ \n', logLhood - logOld);
62     end
63 end

```

A helper function that computes the γ_t and $\xi_{t,t+1}$ variables, as well as the log likelihood of the data, using the alpha-gamma recursion.

```

1 function [gamma, xi, logLhood] = FwBw(data, pi0, A, mu, sigma2)
2 % Computes posterior probabilities using alpha-gamma recursion
3
4 M = length(pi0); % number of states in the latent variables
5 T = size(data,1); % number of observations
6
7 % Compute the alpha variables
8 % Note: we always normalize so sum of alpha(q_t) over q_t is = 1
9 % But keep track of the normalizers to compute p(y)
10 alpha = zeros(M,T);
11 normalizer = zeros(T,1);
12
13 % First initialize alpha(q_0)
14 centered = data(ones(M,1),:) - mu;
15 alpha(:,1) = (pi0.*exp(-sum(centered.^2,2)/(2*sigma2)))/(2*pi*sigma2);
16 normalizer(1) = sum(alpha(:,1));
17 alpha(:,1) = alpha(:,1)/normalizer(1);
18
19 % Then move forward in time
20 for t = 1:T-1
21     % Compute p(y_{t+1}|q_{t+1}=i) for i=1..M
22     centered = data((t+1)*ones(M,1),:) - mu;
23     post = exp(-sum(centered.^2,2)/(2*sigma2))/(2*pi*sigma2);
24
25     % Compute \sum_{q_t} alpha(q_t)*a_{q_t,q_{t+1}=i} for i=1..M
26     c = sum(alpha(:,t*ones(1,M)).*A,1)';
27
28     % Store value of alpha
29     alpha(:,t+1) = post.*c;
30     normalizer(t+1) = sum(alpha(:,t+1));
31     alpha(:,t+1) = alpha(:,t+1)/normalizer(t+1);
32 end
33
34 % Compute the gamma variables
35 gamma = zeros(M,T);
36 gamma(:,T) = alpha(:,T);

```

```

37
38 % Then move backward in time
39 for t = (T-1):(-1):1
40     % Compute  $\alpha(q_t) * a_{\{q_t, q_{t+1}\}}$  and normalize the sum
41     c = alpha(:, t * ones(1, M)) .* A;
42     sum_c = sum(c, 1);
43     c = c ./ sum_c(ones(M, 1), :);
44
45     % Store value of gamma
46     gamma(:, t) = sum(c .* gamma(:, (t+1) * ones(1, M))', 2);
47 end
48
49 % Compute the marginal log likelihood
50 logLhood = sum(log(normalizer));
51
52 % Compute the xi variables
53 xi = zeros(M, M, T-1);
54 for t = 1:T-1
55     % Compute  $p(y_{t+1} | q_{t+1}=j)$  for  $j=1..M$ 
56     centered = data((t+1) * ones(M, 1), :) - mu;
57     post = exp(-sum(centered.^2, 2) ./ (2 * sigma2)) ./ (2 * pi * sigma2);
58
59     for i = 1:M
60         xi(i, :, t) = (alpha(i, t) * (post .* gamma(:, t+1) .* A(i, :)') ./ alpha(:, t+1))';
61     end
62
63     % Need to normalize the xi since we rescaled alpha
64     xi(:, :, t) = xi(:, :, t) / sum(sum(xi(:, :, t)));
65 end

```

The code for implementing EM for Gaussian mixture model:

```

1 function [pi0, mu, sigma2, logLhood] = EM_GaussianMixture(data, M)
2 % M is the number of states in the latent variable
3
4 % data is T-by-K
5 T = size(data, 1); % number of iid observations
6 K = size(data, 2); % dimensionality of the Gaussian
7
8 % Initialize cluster probabilities to be uniform 1/M
9 pi0 = ones(M, 1) / M;
10
11 % Initialize mean values to be random data points -- mu is M-by-K
12 mu = data(randi(T, M, 1), :);
13
14 % Initialize sigma^2 to be the average distance between the data points and
15 % one random data point
16 centered = data - data(randi(T, 1, 1) * ones(T, 1), :);

```

```

17 sigma2_init = (sum(sum(centered.^2))/(2*T));
18 sigma2 = sigma2_init*ones(M,1);
19
20 % Compute log likelihood of the data
21 logLhood = computeLogLhoodGaussianMix(data, pi0, mu, sigma2);
22
23 % Maximum number of iterations
24 maxIter = 1000;
25
26 % Iterate over the E and M steps
27 for iter = 1:maxIter
28     % Compute the posterior probabilities, each row for one q_t
29     post = zeros(T,M);
30     for t = 1:T
31         % Compute y_t-mu_i for i = 1..M
32         centered = data(t*ones(M,1),:) - mu;
33
34         % Compute joint probability for y_t:
35         % p(y_t, q_t=i) = p(q_t=i)*p(y_t|q_t=1) for i=1..M
36         joint_yt = (pi0.*exp(-sum(centered.^2,2)/(2*sigma2)))/(2*pi*sigma2);
37
38         % Normalize and store to post
39         post(t,:) = joint_yt'/sum(joint_yt);
40     end
41
42     % Then update parameters
43     pi0 = sum(post,1)'/T;
44
45     % Update mean parameters
46     for i = 1:M
47         mu(i,:) = sum(data.*post(:,i*ones(1,K)))/sum(post(:,i));
48     end
49
50     % Update variance parameters - use the new mu values
51     for i = 1:M
52         centered = data - mu(i*ones(T,1), :); % centered data
53         sigma2(i) = sum(sum(centered.^2,2).*post(:,i))/(2*sum(post(:,i)));
54     end
55
56     % Compute new log likelihood of the data
57     logOld = logLhood;
58     logLhood = computeLogLhoodGaussianMix(data, pi0, mu, sigma2);
59
60     % Check for convergence
61     if (logLhood - logOld < 1e-14)
62         fprintf('Converges after %d batch iterations!\n', iter);
63         break;
64     elseif (iter == maxIter)
65         fprintf('EM did not converge: final Δ = %.8f\n', logLhood - logOld);

```

```

66     end
67 end

```

A helper function that computes the log likelihood of the data given the parameters. Note that in this case the log likelihood is easy to compute since the observations are i.i.d., so

$$p(y) = \prod_{t=0}^T p(y_t) = \prod_{t=0}^T \sum_{i=1}^M p(y_t, x_t = i).$$

```

1 function logLhood = computeLogLhoodGaussianMix(data, pi0, mu, sigma2)
2
3 % Candidate for the log likelihood log(p(data))
4 logLhood = 0;
5
6 % Compute marginal p(y-t) for each t separately
7 for t = 1:size(data,1)
8     % Compute y-t - mu_i for i = 1..M
9     centered = data(t*ones(size(mu,1),1), :) - mu;
10
11     % Compute joint probability for y-t:
12     % p(y-t, q-t=i) = p(q-t=i)*p(y-t|q-t=1) for i=1..M
13     joint_yt = (pi0.*exp(-sum(centered.^2,2)/(2*sigma2)))./(2*pi*sigma2);
14
15     % Then compute marginal probability p(y-t) and add to total
16     logLhood = logLhood + log(sum(joint_yt));
17 end

```

Code for Problem 5.3

```

1 function [logLhood Theta] = IPF(data, edges)
2 % Implement the IPF algorithm in the pairwise Markov random field model
3 % where the cliques are edges.
4
5 % data is M-by-N, assumed binary
6 % M is the number of nodes in the graph
7 % N is the number of iid observations
8 [M N] = size(data);
9
10 % edges is a cell, each entry is of the form (s,t) for some 1 ≤ s < t ≤ N
11 K = length(edges); % number of edges in the graph
12 for k = 1:K % make sure s < t
13     edges{k} = [min(edges{k}) max(edges{k})];
14 end
15

```

```

16 % Compute empirical distributions: Emp(i,j,k) is the fraction of the data
17 % that has (X_s, X_t) = (i-1, j-1), where (s,t) = edges{k}.
18 Emp = zeros(2,2,K);
19 for k = 1:K
20     % Find the nodes of the edge
21     s = edges{k}(1);
22     t = edges{k}(2);
23
24     % Iterate over the observations and update the appropriate entry of Emp
25     for n = 1:N
26         i = data(s,n) + 1;
27         j = data(t,n) + 1;
28         Emp(i,j,k) = Emp(i,j,k) + 1/N;
29     end
30 end
31
32 % Initialize the parameters: Theta(i,j,k) encodes theta_st(X_s=i-1,
33 % X_t=j-1), where (s,t) = edges{k}.
34 Theta = zeros(2,2,K);
35
36 % In this case the normalizer Z can be computed explicitly
37 Z = 2^M;
38
39 % Compute initial log likelihood
40 logLhood = N*sum(sum(sum(Emp.*Theta))) - N*log(Z);
41
42 % Iterate over all edges and configurations
43 maxIter = 100; % maximum number of batch iterations
44 for iter = 1:maxIter
45     % Iterate over the edges
46     for k = 1:K
47         % Find the nodes of the edge
48         s = edges{k}(1);
49         t = edges{k}(2);
50
51         % Iterate over the configurations of (s,t)
52         for i = 1:2
53             for j = 1:2
54                 % Generate possible configurations for the remaining M-2
55                 % nodes.
56                 % poss is 2^(M-2)-by-(M-2)
57                 poss = generatePoss(M-2);
58
59                 % Interleave with the configuration of (s,t)
60                 poss = [poss(:,1:s-1), ...
61                        (i-1) * ones(2^(M-2),1), ...
62                        poss(:,s:t-2), ...
63                        (j-1) * ones(2^(M-2),1), ...
64                        poss(:,t-1:end)];

```

```

65
66         % Compute marginal probability by summing over all possibilities
67         margPr = 0;
68         for w = 1:size(poss,1)
69             logPr = -log(Z);
70             for e = 1:length(edges)
71                 logPr = logPr + ...
72                     Theta(poss(w,edges{e}(1))+1, poss(w,edges{e}(2))+1, e);
73             end
74             margPr = margPr + exp(logPr);
75         end
76
77         % Update the corresponding entry of Theta
78         Theta(i,j,k) = Theta(i,j,k) + log(Emp(i,j,k)) - log(margPr);
79     end
80 end
81 end
82
83 % Compute log likelihood
84 logOld = logLhood;
85 logLhood = N*sum(sum(sum(Emp.*Theta))) - N*log(Z);
86
87 % Check for convergence
88 if (logLhood - logOld < 1e-14)
89     fprintf('Converges after %d batch iterations!\n', iter);
90     break;
91 elseif (iter == maxIter)
92     fprintf('IPF did not converge: final Δ = %.8f\n', logLhood - logOld);
93 end
94 end
95
96
97 function poss = generatePoss(m)
98 % Generates a (2^m)-by-m array consisting of all possible binary
99 % configurations of m nodes.
100
101 if (m == 1)
102     poss = [0; 1];
103 else
104     possRec = generatePoss(m-1);
105     poss = [zeros(size(possRec,1),1), possRec; ...
106             ones(size(possRec,1),1), possRec];
107 end

```

Code for Problem 5.4

```

1 % This script finds a maximum weight spanning tree from the Pairwise.dat

```

```

2 % dataset with edge weights given by the mutual information I(\mu_st;
3 % \mu_s, \mu_t).
4
5 load Pairwise.dat
6 d = size(Pairwise,1); % number of nodes in the graph
7 n = size(Pairwise,2); % number of observations
8
9 % Compute empirical vertex marginals \mu_s
10 mu_s = [(1/n) * sum(Pairwise==0, 2)'; ... % X_s = 0
11         (1/n) * sum(Pairwise==1, 2)']; % X_s = 1
12
13 % Compute vertex entropy H(\mu_s)
14 H_s = -sum(mu_s .* log(mu_s))';
15
16 % Compute empirical edge marginals \mu_st
17 mu_st = zeros(2, 2, d*(d-1)/2);
18 count = 1; % for ordering the edges from 1 to d*(d-1)/2
19 for i = 1:d
20     for j = i+1:d
21         ones_i = real(Pairwise(i,:) == 1); % when X_i = 1
22         ones_j = real(Pairwise(j,:) == 1); % when X_j = 1
23         mu_st(:,:,count) = (1/n) * ...
24             [(1-ones_i)*(1-ones_j)', (1-ones_i)*ones_j'; ...
25              ones_i*(1-ones_j)', ones_i*ones_j'];
26         count = count + 1;
27     end
28 end
29
30 % Compute edge entropy H(\mu_st)
31 H_st = zeros(d*(d-1)/2, 1);
32 for count = 1:d*(d-1)/2
33     mu_edge = mu_st(:,:,count);
34     H_st(count) = -sum(sum(mu_edge .* log(mu_edge)));
35 end
36
37 % Compute mutual information I(\mu_st; \mu_s, \mu_t)
38 I_st = zeros(d*(d-1)/2, 1);
39 count = 1; % for ordering the edges from 1 to d*(d-1)/2
40 for i = 1:d
41     for j = i+1:d
42         I_st(count) = H_s(i) + H_s(j) - H_st(count);
43         fprintf('I(%d,%d) = %.4f\n', i, j, I_st(count));
44         count = count + 1;
45     end
46 end

```