

## 4F13 Machine Learning: Course work #2: Variational and Sampling Methods

Zoubin Ghahramani & Carl Edward Rasmussen

Due: 4pm Fri Mar 14, 2008 to Rachel Fogg, room 37 (Baker)

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

$$p(\mathbf{s}|\boldsymbol{\pi}) = p(s_1, \dots, s_K|\boldsymbol{\pi}) = \prod_{i=1}^K p(s_i) = \prod_{i=1}^K \pi_i^{s_i} (1 - \pi_i)^{(1-s_i)}$$
$$p(\mathbf{y}|\mathbf{s}_1, \dots, \mathbf{s}_K, \boldsymbol{\mu}, \sigma^2) = \mathcal{N}\left(\sum_i \mathbf{s}_i \boldsymbol{\mu}_i, \sigma^2 \mathbf{I}\right)$$

where  $\mathbf{y}$  is a  $D$ -dimensional vector and  $\mathbf{I}$  is the  $D \times D$  identity matrix. Assume you have a data set of  $N$  i.i.d. observations of  $\mathbf{y}$ , i.e.  $\mathbf{Y} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(N)}\}$ . More details are provided in Appendix A.

**General Matlab hint:** wherever possible, avoid looping over the data points. Many (but not all) of these functions can be written using matrix operations. In Matlab it's much faster.

**Warning:** Each question depends on earlier questions. Start as soon as possible.

**Hand in:** Derivations, code and plots.

5% Gibbs sampling relies on computing conditional probabilities of certain variables given all other variables (read section 29.5 in MacKay's textbook if you haven't done so already). Derive the conditional probability:

$$p(s_i^{(n)} | s_1^{(n)}, \dots, s_{i-1}^{(n)}, s_{i+1}^{(n)}, \dots, s_K^{(n)}, \mathbf{y}^{(n)}, \boldsymbol{\theta})$$

for the binary latent variable model, which you will need for Gibbs sampling. Describe an algorithm for drawing samples from this distribution using the `rand` function in Matlab which gives you uniformly distributed random numbers between 0 and 1.

15% In this exercise you will implement the fully factored (a.k.a. mean-field) **variational approximation** described in the course notes. That is, for each data point  $\mathbf{y}^{(n)}$ , the code 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. Specifically, you should write a Matlab function:

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

where `lambda` is  $N \times K$ , `F` is the lower bound on the likelihood, `Y` is the  $N \times D$  data matrix, `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  $\epsilon$  the iterations halt.

15% Using the conditional probability derived in the first question you will now implement **Gibbs sampling** to approximate the posterior distribution over the hidden states given the data. Specifically, write a function:

```
[S] = Gibbs(Y, mu, sigma, pie, S0, nsamples)
```

where `S` is a  $N \times K \times \text{nsamples}$  array of samples over the hidden variables, `S0` is an  $N \times K$  matrix of initial settings for the hidden variables.

5% We have derived the M step for this model in terms of the quantities: `Y`, `ES` =  $E_q[\mathbf{s}]$ , which is an  $N \times K$  matrix of expected values, and `ESS`, which is an  $N \times K \times K$  array of expected values  $E_q[\mathbf{s}\mathbf{s}^T]$  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.

5% Using the above, we have implemented a function:

```
[mu, sigma, pie] = MStep(Y, 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.

7% Examine the data `images.jpg` shown on the web site (Do **not** look at `genimages.m` yet!). This shows 100 grayscale  $4 \times 4$  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 mixture of Gaussians? Explain your reasoning.

12% Put the E step and M step code together into a function:

```
[mu, sigma, pie] = LearnBinFactors(Y,K,iterations,gibbsflag)
```

where  $K$  is the number of binary factors, `iterations` is the maximum number of iterations of EM, and `gibbsflag = 1` means use Gibbs sampling, otherwise, use mean field. For the mean field algorithm, make sure  $F$  always increases (this is a good debugging tool).

7% 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 do the Gibbs and variational algorithms differ? Which do you prefer? Note that `nsamples` might have to be large for Gibbs sampling to work well. 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`).

7% Make sure you understand the idea of *convergence* of Markov chains (29.8 and 29.9 in MacKay's book). Given the parameters you have learned in the previous problem, and given just the first data point in the data set  $\mathbf{y}^{(1)}$ , i.e.  $N = 1$ , run your Gibbs sampling code with various initial conditions of the hidden state  $\mathbf{S}_0$ . In particular, you might want to try starting with  $\mathbf{S}_0$  as all 1's, as all 0's, or somewhere in between. Plot various statistics of  $\mathbf{S}$  (such as the sum over all  $K$  states, which should range between 0 and  $K$ , or the cumulative average over all  $K$  states) as a function of sampling iteration for Gibbs sampling for a large number of samples. Can you assess (visually) how long it takes for the Gibbs sampler to converge? How is this affected by increases and decreases in `sigma`? Why? Support your arguments.

7% For the same setting of the parameters as in the previous problem and again given just the first data point in the data set  $\mathbf{y}^{(1)}$ , i.e.  $N = 1$ , run the variational approximation. 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.

10% Given known values of  $\sigma^2$  and  $\boldsymbol{\pi}$ , and a sample of  $\mathbf{s}$  and  $\mathbf{y}$ , what is the conjugate prior for the  $\boldsymbol{\mu}$ s? Implement a sampling procedure for  $\boldsymbol{\mu}$  given  $\sigma^2$ ,  $\boldsymbol{\pi}$ ,  $\mathbf{s}$  and  $\mathbf{y}$ . You might want to use Gibbs sampling or Metropolis. Show some samples from the posterior distribution of  $\boldsymbol{\mu}$ .

5% Describe a Bayesian method for selecting  $K$ , the number of hidden binary variables. Does your method pose any computational difficulties and if so how would you tackle them?

## Appendix: M-step for Course work #2

Iain Murray, Dec 2003

### A Background

The generative model under consideration has a vector of  $K$  binary latent variables  $\mathbf{s}$ . Each  $D$ -dimensional data point  $\mathbf{y}^{(n)}$  is generated using a new hidden vector,  $\mathbf{s}^{(n)}$ . Each  $\mathbf{s}^{(n)}$  is identically and independently distributed according to:

$$P(\mathbf{s}^{(n)}|\boldsymbol{\pi}) = \prod_{i=1}^K \pi_i^{s_i^{(n)}} (1 - \pi_i)^{(1-s_i^{(n)})}. \quad (1)$$

Once  $\mathbf{s}^{(n)}$  has been generated, the data point is created according to the Gaussian distribution:

$$p(\mathbf{y}^{(n)}|\mathbf{s}^{(n)}, \boldsymbol{\mu}, \sigma^2) = (2\pi\sigma^2)^{-D/2} \exp \left[ -\frac{1}{2\sigma^2} \left( \mathbf{y}^{(n)} - \sum_{i=1}^K s_i^{(n)} \boldsymbol{\mu}_i \right)^\top \left( \mathbf{y}^{(n)} - \sum_{i=1}^K s_i^{(n)} \boldsymbol{\mu}_i \right) \right]. \quad (2)$$

When this process is repeated we end up obtaining a set of visible data  $Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(N)}\}$  generated by a set of  $N$  binary vectors  $S = \{\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)}\}$  and some model parameters  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \sigma^2, \boldsymbol{\pi}\}$ , which are constant across all the data. Given just  $Y$ , both  $S$  and  $\boldsymbol{\theta}$  are unknown. We might want to find the set of parameters that maximise the likelihood function  $P(Y|\boldsymbol{\theta})$ ; “the parameters that make the data probable”. EM is an approach towards this goal which takes our knowledge about the uncertain  $S$  into account.

In the EM algorithm we optimise the objective function

$$\begin{aligned} \mathcal{F}(\mathbf{q}, \boldsymbol{\theta}) &= \langle \log p(S, Y|\boldsymbol{\theta}) \rangle_{\mathbf{q}(S)} - \langle \log q(S) \rangle_{\mathbf{q}(S)} \\ &= \sum_{\mathbf{n}} \langle \log p(\mathbf{s}^{(n)}, \mathbf{y}^{(n)}|\boldsymbol{\theta}) \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} - \sum_{\mathbf{n}} \langle \log q(\mathbf{s}^{(n)}) \rangle_{\mathbf{q}(\mathbf{s}^{(n)})}, \end{aligned} \quad (3)$$

alternately increasing  $\mathcal{F}$  by changing the distribution  $\mathbf{q}(S)$  in the “E-step”, and the parameters in the “M-step”. This document gives a derivation and Matlab implementation of the M-step. In this assignment you will implement two approximate E-steps and apply this EM algorithm to a data set.

### B M-step derivation

Here we maximise  $\mathcal{F}$  with respect to each of the parameters using differentiation. This only requires the term with  $\boldsymbol{\theta}$  dependence:

$$\sum_{\mathbf{n}} \langle \log p(\mathbf{s}^{(n)}, \mathbf{y}^{(n)}|\boldsymbol{\theta}) \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} = \sum_{\mathbf{n}} \langle \log p(\mathbf{y}^{(n)}|\mathbf{s}^{(n)}, \boldsymbol{\theta}) + \log P(\mathbf{s}^{(n)}|\boldsymbol{\theta}) \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \quad (4)$$

Substituting the given distributions from equations ?? and ?? gives:

$$\begin{aligned} &= -\frac{ND}{2} \log 2\pi - ND \log \sigma \\ &\quad - \frac{1}{2\sigma^2} \left[ \sum_{\mathbf{n}=1}^N \mathbf{y}^{(n)\top} \mathbf{y}^{(n)} + \sum_{i,j} \boldsymbol{\mu}_i^\top \boldsymbol{\mu}_j \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} - 2 \sum_i \boldsymbol{\mu}_i^\top \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \mathbf{y}^{(n)} \right] \\ &\quad + \sum_{i=1}^K \left[ \log \pi_i \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} + \log(1 - \pi_i) \left( N - \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \right) \right]. \end{aligned} \quad (5)$$

From which we can obtain all the required parameter settings:

$$\frac{\partial \mathcal{F}}{\partial \pi_i} = \frac{1}{\pi_i} \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} + \frac{1}{1 - \pi_i} \left[ \sum_{\mathbf{n}=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} - N \right] = 0 \quad (6)$$

$$\Rightarrow \boxed{\boldsymbol{\pi} = \frac{1}{N} \sum_{n=1}^N \langle \mathbf{s}^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})}} , \quad (7)$$

$$\frac{\partial \mathcal{F}}{\partial \boldsymbol{\mu}_i} = -\frac{1}{\sigma^2} \sum_{n=1}^N \left[ \sum_j \langle s_i^{(n)} s_j^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} - \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \mathbf{y}^{(n)} \right] \quad (8)$$

$$\sum_j \sum_{n=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \boldsymbol{\mu}_j = \sum_{n=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \mathbf{y}^{(n)}$$

$$\Rightarrow \boxed{\boldsymbol{\mu}_j = \sum_i \left[ \sum_{n=1}^N \langle \mathbf{s}^{(n)} \mathbf{s}^{(n)\top} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \right]_{ji}^{-1} \sum_{n=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \mathbf{y}^{(n)}} \quad (9)$$

and

$$\frac{\partial \mathcal{F}}{\partial \sigma} = 0 \Rightarrow \boxed{\sigma^2 = \frac{1}{ND} \left[ \sum_{n=1}^N \mathbf{y}^{(n)\top} \mathbf{y}^{(n)} + \sum_{i,j} \boldsymbol{\mu}_i^\top \boldsymbol{\mu}_j \sum_{n=1}^N \langle s_i^{(n)} s_j^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} - 2 \sum_i \boldsymbol{\mu}_i^\top \sum_{n=1}^N \langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})} \mathbf{y}^{(n)} \right]} . \quad (10)$$

Note that the required sufficient statistics of  $\mathbf{q}(\mathbf{S})$  are  $\langle s_i^{(n)} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})}$  and  $\sum_{n=1}^N \langle \mathbf{s}^{(n)} \mathbf{s}^{(n)\top} \rangle_{\mathbf{q}(\mathbf{s}^{(n)})}$ . In the code these are known as ES and ESS.

All of the sums above can be interpreted as matrix multiplication or trace operations. This means that each of the boxed parameters above can neatly be computed in one line of Matlab.

## C M-step code

MStep.m

---

```
% [mu, sigma, pie] = MStep(Y,ES,ESS)
%
% Inputs:
% _____
%       Y NxD data matrix
%       ES NxK E_q[s]
%       ESS KxK sum over data points of E_q[ss'] (NxKxK)
%           if E_q[ss'] is provided, the sum over N is done for you.
%
% Outputs:
% _____
%       mu DxK matrix of means in p(y|{s_i},mu,sigma)
%       sigma 1x1 standard deviation in same
%       pie 1xK vector of parameters specifying generative distribution for s
%
function [mu, sigma, pie] = MStep(Y,ES,ESS)

[N,D] = size(Y);
if (size(ES,1)~=N), error('ES must have the same number of rows as Y'); end;
K = size(ES,2);
if (isequal(size(ESS),[N,K,K])), ESS = shiftdim(sum(ESS,1),1); end;
if (~isequal(size(ESS),[K,K]))
    error('ESS must be square and have the same number of columns as ES');
end;

mu = (inv(ESS)*ES'*Y)';
sigma = sqrt((trace(Y'*Y)+trace(mu'*mu*ESS)-2*trace(ES'*Y*mu))/(N*D));
pie = mean(ES,1);
```

---