

# CS228 Homework 4

Instructor: Stefano Ermon – [ermon@stanford.edu](mailto:ermon@stanford.edu)

Available: 02/17/2017; Due: 03/3/2017

---

1. **[20 points]** We have a data association problem where there are  $K$  objects and we are given  $K$  observations. Each observation corresponds to a single object, and we are given one observation for each object. However, we don't know which observation corresponds to which object, and we would like to infer that using a probabilistic model relating observations to objects. Specifically we have

- $K$  objects  $u_1, \dots, u_K$
- observations  $v_1, \dots, v_K$ , where  $\text{Val}(v_i) = \{a_1, \dots, a_L\}$  (so  $v_i$  is a discrete random variable), where *each observation corresponds to the appearance of one object and there is exactly one observation of each object*
- correspondence variables  $C_1, \dots, C_K$ , where  $\text{Val}(C_i) = \{1, \dots, K\}$ ;  $C_i = k$  denotes that measurement  $v_i$  is derived from object  $u_k$
- a *known* appearance model for each object  $u_k, P_k(v_i = a_l | C_i = k)$ .

Note that because of the mutex constraints, the correspondence variables  $C_1, \dots, C_K$  will be a permutation over  $1, \dots, K$ . We also assume for simplicity that all permutations are equally likely *a priori*.

We wish to compute the marginals  $P(C_i | v_1, \dots, v_K)$ , for  $i = 1, \dots, K$ , using Metropolis-Hastings (MH) to sample the correspondence variables. We will start with an arbitrary assignment to  $C_1, \dots, C_K$ , and take MH-steps. The proposal distribution that we will use randomly picks two correspondence variables  $C_i, C_j$  from a uniform distribution over all pairs of correspondence variables, and swaps their assignments.

- (a) **[10 points]** Compute the acceptance probability for each MH step.
  - (b) **[5 points]** Suppose we have run the MH sampler for a long time and collected  $M$  samples  $(C_1[m], \dots, C_K[m])$  for  $m = 1, \dots, M$  after the chain has mixed. Give an explicit expression for estimating the marginal  $P(C_i | v_1, \dots, v_K)$ .
  - (c) **[5 points]** Your friend Geoff Gibbs hears about your MH algorithm and suggests that you can also consider using Gibbs sampling to compute your marginals. Briefly explain why this will or will not work.
2. **[20 points] Multi-conditional Parameter Learning, Markov Networks**

In this problem, we will consider the problem of learning parameters for a Markov network using a specific objective function. In particular assume that we have two sets of variables  $\mathbf{Y}$  and  $\mathbf{X}$ , and a dataset  $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^M, \mathbf{y}^M)\}$ . We will estimate the model parameters  $\boldsymbol{\theta} = [\theta_1 \dots \theta_n]$  by maximizing the following objective function:

$$g(\boldsymbol{\theta}; \mathcal{D}) = (1 - \alpha)\ell_{\mathbf{Y}|\mathbf{X}}(\boldsymbol{\theta}; \mathcal{D}) + \alpha\ell_{\mathbf{X}|\mathbf{Y}}(\boldsymbol{\theta}; \mathcal{D})$$

where  $\ell_{\mathbf{X}|\mathbf{Y}}(\boldsymbol{\theta}; \mathcal{D})$  means the conditional log-likelihood of the dataset  $\mathcal{D}$  using the distribution  $P_{\boldsymbol{\theta}}(\mathbf{X} | \mathbf{Y})$  defined by the Markov network with parameters  $\boldsymbol{\theta}$  (similarly for  $\ell_{\mathbf{Y}|\mathbf{X}}$ ). Thus, our objective is a mixture of two conditional log-likelihoods ( $0 < \alpha < 1$ ). As usual, we consider a log-linear parameterization of a Markov network, using a set of  $n$  features  $f_i(\mathbf{X}_i, \mathbf{Y}_i)$  where  $\mathbf{X}_i$  and  $\mathbf{Y}_i$  are some (possibly empty) subsets of the variables  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively.

- (a) **[10 points]** Write down the full objective function  $g(\boldsymbol{\theta}; \mathcal{D})$  in terms of the features  $f_i$  and weights  $\theta_i$
- (b) **[10 points]** Derive  $\frac{\partial}{\partial \theta_i} g(\boldsymbol{\theta}; \mathcal{D})$ : the derivative of the objective with respect to a weight  $\theta_i$ . Write your final answer in terms of feature expectations  $\mathbf{E}_Q[f_i]$ , where  $Q$  is either: the empirical distribution of our dataset  $\hat{P}$ ; or a conditional distribution of the form  $P_{\boldsymbol{\theta}}(\mathbf{W} \mid \mathbf{Z} = \mathbf{z})$  (for some sets of variables  $\mathbf{W}, \mathbf{Z}$ , and assignment  $\mathbf{z}$ .)

3. **[10 points] Expectation Maximization in a Naive Bayes Model**

Consider the Naive Bayes model with class variable  $C$  and discrete evidence variables  $X_1, \dots, X_n$ . The CPDs for the model are parameterized by  $P(C = c) = \theta_c$  and  $P(X_i = x \mid C = c) = \theta_{x_i|c}$  for  $i = 1, \dots, n$ , and for all assignments  $x_i \in \text{Val}(X_i)$  and classes  $c \in \text{Val}(C)$ .

Now given a data set  $\mathcal{D} = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$ , where each  $\mathbf{x}[m]$  is a complete assignment to the evidence variables,  $X_1, \dots, X_n$ , we can use EM to learn the parameters of our model. Note that the class variable,  $C$ , is never observed.

Show that if we initialize the parameters uniformly,

$$\theta_c^0 = \frac{1}{|\text{Val}(C)|} \quad \text{and} \quad \theta_{x_i|c}^0 = \frac{1}{|\text{Val}(X_i)|},$$

for all  $x_i, c$ , then the EM algorithm converges in one iteration, and give a closed form expression for the parameter values at this convergence point.

[60 points] Programming Assignment <sup>1</sup>

In this homework, you will apply Gibbs sampling to a simple Markov random field model for image restoration. In an image restoration problem, you are given an image corrupted by noise  $X$  and you want to recover the original image  $Y$  (see Figure 1 and Lecture 4).

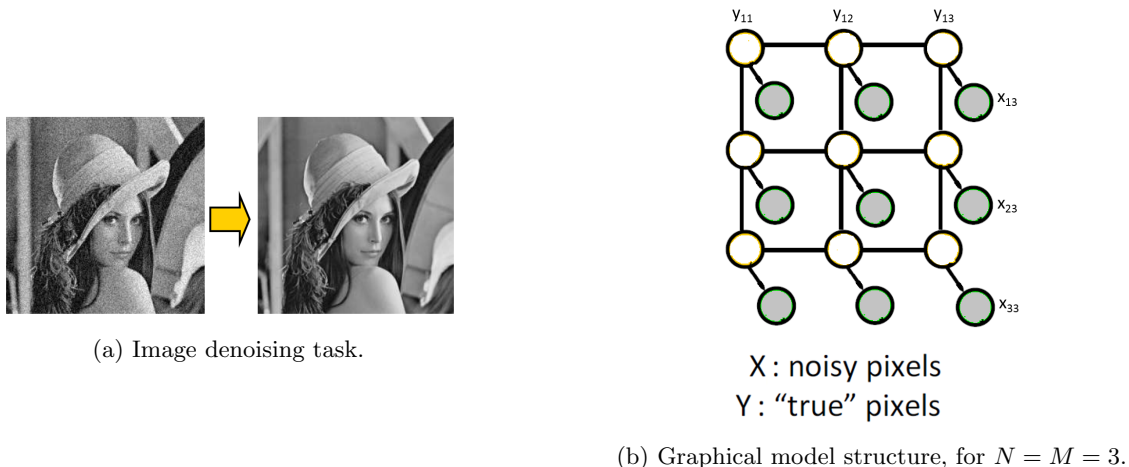


Figure 1: Image denoising with graphical models.

Let  $\mathbf{x} = \{x_{ij}\}$  denote the observed image, with  $x_{ij} \in \{-1, +1\}$  representing the pixel at row  $i$  and column  $j$ . Assume a black-and-white image, with -1 corresponding to white and +1 to black. The image has dimensions  $N \times M$ , so that  $1 \leq i \leq N$  and  $1 \leq j \leq M$ . Assume a set of (unobserved) variables  $\mathbf{y} = \{y_{ij}\}$  representing the true (unknown) image, with  $y_{ij} \in \{-1, +1\}$  indicating the value of  $x_{ij}$  before noise was added. Each (internal)  $y_{ij}$  is linked with four immediate neighbors,  $y_{i-1,j}$ ,  $y_{i+1,j}$ ,  $y_{i,j-1}$ , and  $y_{i,j+1}$ , which together are denoted  $y_{N(i,j)}$ . Pixels at the borders of the image (with  $i \in \{1, N\}$  or  $j \in \{1, M\}$ ) also have neighbors denoted  $y_{N(i,j)}$ , but these sets are reduced in the obvious way. We denote  $E$  the corresponding set of edges. For example, the pair  $((1, 1), (1, 2)) \in E$ , but the pair  $((1, 1), (2, 2)) \notin E$ . The joint probability of  $\mathbf{y}$  and  $\mathbf{x}$  can be written (with no prior preference for black or white):

$$p(\mathbf{y}, \mathbf{x}) = \frac{1}{Z} \left\{ \prod_{i=1}^N \prod_{j=1}^M \exp^{\eta y_{ij} x_{ij}} \right\} \times \left\{ \prod_{((i,j),(i',j')) \in E} \exp^{\beta y_{ij} y_{i'j'}} \right\} \quad (1)$$

$$= \frac{1}{Z} \exp \left\{ \eta \sum_{i=1}^N \sum_{j=1}^M y_{ij} x_{ij} + \beta \sum_{((i,j),(i',j')) \in E} y_{ij} y_{i'j'} \right\} \quad (2)$$

where

$$Z = \sum_{\mathbf{y}, \mathbf{x}} \exp \left\{ \eta \sum_{i,j} y_{ij} x_{ij} + \beta \sum_{((i,j),(i',j')) \in E} y_{ij} y_{i'j'} \right\} \quad (3)$$

(Notice in particular that each pair of neighbors,  $y_{ij}$  and  $y_{i'j'}$ , factors into the formula only once, despite that each variable is a neighbor of the other. Failing to account for this will lead to double counting of  $\beta$  values.) This is equivalent to a Boltzmann (sometimes called Gibbs) distribution with "energy":

$$E(\mathbf{y}, \mathbf{x}) = -\eta \sum_{i,j} y_{ij} x_{ij} - \beta \sum_{((i,j),(i',j')) \in E} y_{ij} y_{i'j'} \quad (4)$$

<sup>1</sup>Assignment adapted from Cornell's BTRY 6790, instructed by Adam Siepel

The system will have lower energy, and hence higher probability, in states in which neighboring  $y_{ij}$  variables, and neighboring  $y_{ij}$  and  $x_{ij}$  variables, tend to have the same value (assuming  $\eta$  and  $\beta$  are positive). This captures the fact that each noisy pixel  $x_{ij}$  is likely to be similar to the corresponding “true” pixel  $y_{ij}$ , and that images tend to be “smooth”.

There are algorithms for deterministically estimating  $\mathbf{y}$  given an image  $\mathbf{x}$  but we will here use the alternative approach: we will devise a Markov Chain Monte Carlo (MCMC) algorithm to sample values of  $\mathbf{y}$  conditional on  $\mathbf{x}$ . Here are some advantages over the deterministic algorithms:

- i. It is very general, and can easily be extended to more complex graphs.
- ii. It provides great flexibility for quantifying the uncertainty of  $\mathbf{y}$  (and, potentially, for the parameters  $\eta$  and  $\beta$ ).
- iii. It is relatively straightforward in this setting to derive the exact conditional distributions for nodes given the Markov blanket, so Gibbs sampling is possible, and one need not worry about the acceptance rate for proposed samples.

You will apply your methods to two small, black-and-white images that have been made available with the problem set. These two noisy images, and the original, undistorted image from which they derive, are available both in PNG format and in a simple text format that lists each coordinate pair  $(i, j)$  and the corresponding value of  $x_{ij}$ . You may find it useful to convert between this text representation and a viewable image format.

- (a) **[5 points]** Derive an expression for the conditional probability that pixel  $(i, j)$  is black given its Markov blanket, i.e.  $p(y_{ij} = 1 | y_{M(i,j)})$ , where  $y_{M(i,j)}$  denotes the variables in the Markov blanket of  $y_{ij}$  (but you should be explicit about which variables are included). Your expression should take the form of a logistic function and should depend only on  $\eta, \beta$ , and  $y_{M(i,j)}$ .
- (b) **[10 points]** Outline a Gibbs sampling algorithm (in pseudocode) that iterates over the pixels in the image and samples each  $y_{ij}$  given its Markov blanket. Use the simple approach of sweeping across the image in row-major fashion on every iteration of the algorithm. Thus, an “iteration” will generate a complete new sample of  $\mathbf{y}$ . Allow for a burn-in of  $B$  iterations, followed by draws of  $S$  samples. You may assume  $\eta$  and  $\beta$  are fixed constants. How can we show in our case that the equilibrium distribution is in fact the posterior distribution  $p(\mathbf{y} | \mathbf{x})$ ?
- (c) **[15 points]** Implement your algorithm and apply it to the image with 20% noise (noisy 20.png,txt). Use values of  $\eta = 1$ ,  $\beta = 1$ ,  $B = 100$ , and  $S = 1000$ . On each iteration of your algorithm, compute the energy  $E(\mathbf{y}, \mathbf{x})$  for the current sample of  $\mathbf{y}$  and output it to a log file, keeping track of which values correspond to the burn-in. Run your algorithm with three different initializations - one in which each  $y_{ij}$  is initialized to  $x_{ij}$ , one in which each  $y_{ij}$  is initialized to  $-x_{ij}$ , and one in which the  $y_{ij}$  are set to  $-1$  or  $+1$  at random. Plot the energy of the model as a function of the iteration number for all three chains and visually inspect these traces for signs of convergence. Do all three seem to be converging to the same general region of the posterior, or are some obviously suboptimal? Does the burn-in seem to be adequate in length? Is there substantial fluctuation from iteration to iteration, indicating that the chain is mixing well, or does it become stuck at particular energies for several iterations at a time?
- (d) **[10 points]** Have your program output a restored image after completing its sampling iterations, by thresholding the estimated posterior probabilities for the  $y_{ij}$  variables at 0.5 - i.e., by estimating the “true” color of each pixel  $(i, j)$  as:

$$\hat{y}_{ij} = \begin{cases} +1 & \text{if } p(y_{ij} = 1 | \mathbf{x}) > 0.5 \\ -1 & \text{otherwise} \end{cases}$$

To estimate the required posterior probabilities, store a running count  $c_{ij}$  of the number of (retained) samples for which each  $y_{ij} = 1$ , and then use the Monte Carlo estimate:

$$p(y_{ij} = 1|\mathbf{x}) \approx \frac{1}{S} \sum_t 1(y_{ij}^{(t)} = 1) = \frac{1}{S} c_{ij} \quad (5)$$

where  $y_{ij}^{(t)}$  represents the  $t^{th}$  sample of  $y_{ij}$ . Restore both the 10% - and 20% -noise images in this way, using the same values of  $\eta$ ,  $\beta$ ,  $B$ , and  $S$  as above. Evaluate the quality of the restoration by computing the fraction of all pixels that differ between the restored images and the original image. Prepare a figure for each the the two images, showing the original, the noisy version, and the restoration side by side.

- (e) **[10 points]** If you have implemented your algorithm correctly, your restored images should be quite close to the original. But is this because you have a clever algorithm or just because the problem is easy? To examine this question, implement a trivial reconstruction algorithm that sets each  $y_{ij}$  equal to the consensus (majority) of its neighbors (including  $x_{ij}$ ), and iterates a few times until convergence (use sequential rather than batch updates, as in Gibbs sampling. This algorithm need not converge in theory, but quickly do quite often in practice. To be safe, you can force it to terminate after, say, 30 iterations.) You should be able to get this program working quickly and easily by reusing code from your Gibbs sampler. However, note that in this case, you should not average over samples (there are no samples here) but instead should use the final value of the  $y_{ij}$  variables for your restored image. Run this program on both images and compute its restoration error. Include figures for the images restored in this way. Does the Gibbs sampler do better than the trivial algorithm? Why or why not?
- (f) **[10 points]** While the Gibbs sampler is useful for obtaining marginal posterior probabilities of interest, much of its appeal derives from its flexibility in estimating posterior distributions for more complex features of the model. To get a sense for its flexibility, use your Gibbs sampler to estimate the posterior distribution over the number of pixels in the “Z” in the image, which approximately falls in the rectangle from  $(i = 125, j = 143)$  to  $(i = 162, j = 174)$ . Using the same parameters as above, simply count the number of cases of  $y_{ij} = +1$  within this rectangle for each retained sample, output one count per iteration as your sampler runs, then use their relative frequencies as an estimate of the posterior distribution of interest. Plot a histogram showing these relative frequencies for both images and comment on any differences between the two estimated posterior distributions.