

Homework 5

Assigned: 24/10/2011.

Due: 4/11/2011.

Exercise 1. BLAHUT-ARIMOTO ALGORITHM

In this exercise, you will analyze a special case of the EM algorithm, which can be shown to converge to a unique optimum. This variant is older than EM itself, and is known as Blahut-Arimoto algorithm. You can find more information on the role of this algorithm for rate distortion theory in the textbook by Cover and Thomas.

Recall mixtures models from the lecture. For latent variable x , observed variable y (both discrete with finite range), the log marginal likelihood is

$$\log P(y = \tilde{y}) = \log \sum_{x=1}^K \pi_x P(y = \tilde{y}|x).$$

In the following, we write $P(\tilde{y})$ instead of $P(y = \tilde{y})$, but keep in mind the difference between y (a variable) and \tilde{y} (a fixed observed value for it). For the purpose of this exercise, we will work with this single datapoint \tilde{y} , but all you'll show here holds for i.i.d. datasets just as well. In this exercise, assume that $P(\tilde{y}|x)$ is fixed and known, and the goal is maximize $\log P(\tilde{y})$ w.r.t. the distribution $P(x) = (\pi_x)$.

If (p_i) , (q_i) are distributions over a finite set, the relative entropy is defined as

$$D[\mathbf{q} \parallel \mathbf{p}] := \sum_i q_i \log \frac{q_i}{p_i}.$$

The function $(\mathbf{q}, \mathbf{p}) \mapsto D[\mathbf{q} \parallel \mathbf{p}]$ is convex, although you don't need this fact here. Given two convex sets of distributions \mathcal{P} , \mathcal{Q} , Csiszár and Tusnády showed that the minimum

$$D[\mathbf{q}_* \parallel \mathbf{p}_*] = \min_{\mathbf{p} \in \mathcal{P}} \min_{\mathbf{q} \in \mathcal{Q}} D[\mathbf{q} \parallel \mathbf{p}]$$

is found by a simple *alternating minimization* algorithm: $\mathbf{q}_{j+1} = \operatorname{argmin}_{\mathbf{q}} D[\mathbf{q} \parallel \mathbf{p}_j]$, $\mathbf{p}_{j+1} = \operatorname{argmin}_{\mathbf{p}} D[\mathbf{q}_{j+1} \parallel \mathbf{p}]$, $j = 1, 2, \dots$

- [4 points]** Bring the EM problem $\min_{(\pi_x)} (-\log P(\tilde{y}))$ into the form of the alternating minimization algorithm. What is \mathcal{P} ? What is \mathcal{Q} ?
Hint: Use the log partition function bound from the lecture. Use families \mathcal{P} , \mathcal{Q} of *joint* distributions over (x, y) , try the form $q_{(x,y)} = Q(x) \mathbb{I}_{\{y=\tilde{y}\}}$.
- [2 points]** Show that the sets \mathcal{P} , \mathcal{Q} you determined above (for the EM problem) are both convex.
- [2 points]** Show that with your re-formulation of the EM problem, the E step update is equivalent to $\mathbf{q}_{j+1} = \operatorname{argmin}_{\mathbf{q}} D[\mathbf{q} \parallel \mathbf{p}_j]$, and the M step update is equivalent to $\mathbf{p}_{j+1} = \operatorname{argmin}_{\mathbf{p}} D[\mathbf{q}_{j+1} \parallel \mathbf{p}]$.
This means that in this restricted case, EM is just alternating minimization between two convex sets, therefore converges to a unique minimum.

4. **[2 points]** Prove the information inequality: for any two distributions \mathbf{p}, \mathbf{q} , the relative entropy $D[\mathbf{q} \parallel \mathbf{p}]$ is always nonnegative. Moreover, $D[\mathbf{q} \parallel \mathbf{p}] = 0$ if and only if $q_i = p_i$ for all i .

Hint: Use Jensen's inequality, in much the same way as in the lecture.

Exercise 2. RAUCH-TUNG-STRIEBEL SMOOTHER

In this exercise, you will work out the RTS smoother, based on an alternative view of HMM inference. Recall the definition of an HMM, based on observation probabilities $P(\mathbf{y}|\mathbf{x})$ and transition probabilities $P(\mathbf{x}_t|\mathbf{x}_{t-1})$. Given data $D = \{\mathbf{y}_1, \dots, \mathbf{y}_T\}$, recall the difference between filtering (computing $P(\mathbf{x}_t|\mathbf{y}_{\leq t})$, where $\mathbf{y}_{\leq t} = \{\mathbf{y}_1, \dots, \mathbf{y}_t\}$) and smoothing (computing the marginal posteriors $P(\mathbf{x}_t|D)$). Since filtering is important per se (for online prediction, given an incoming data stream), most smoothing algorithms make use of filtering code as far as possible. An example is the two-filter smoother discussed in the lecture. The RTS smoother is different, in that the backward pass is simpler, and the data is not required anymore. On the other hand, since the backward pass can start only once the forward (filtering) pass has been completed, the RTS smoother is less suitable for parallel hardware.

1. **[2 points]** Let us first consider a general HMM, not specifying the CPT types. You've learned about belief propagation in the lecture. For a chain, compute all messages, then combine them at each node. For LDS models, BP corresponds to the two-filter smoother. In this exercise, you will work out an alternative to BP (for Markov chains), that will lead to RTS smoothing.

Suppose we have determined $P(\mathbf{x}_t|\mathbf{y}_{\leq t})$, $t = 1, \dots, T$, in a forward pass. Instead of running an equivalent backward pass, we seek a recursion for the marginals $P(\mathbf{x}_t|D)$ directly. The start is easy: $P(\mathbf{x}_T|D) = P(\mathbf{x}_T|\mathbf{y}_{\leq T})$. Now assume that we have determined $P(\mathbf{x}_t|D)$, $t > 1$. Develop a recursion of the form

$$P(\mathbf{x}_{t-1}|D) = \int f_t(\mathbf{x}_{t-1}, \mathbf{x}_t) P(\mathbf{x}_t|D) d\mathbf{x}_t$$

(if \mathbf{x} is discrete, then \int becomes \sum), so that f_t can be computed easily (by local computations that do not need any further global propagation) from the filtering distributions $P(\mathbf{x}_{t'}|\mathbf{y}_{\leq t'})$ we have already determined. State what f_t is, and how to compute it.

Hint: Draw the local part of the model at $\mathbf{x}_{t-1}, \mathbf{x}_t$. Look for conditional independencies that help you. It is helpful to think of $P(\mathbf{x}_{t-1}|\mathbf{y}_{\leq(t-1)})$ as a "prior" in this context (state why).

2. **[4 points]** The RTS smoother is what you get if you apply this idea to the LDS model discussed in the lecture. Develop the RTS backward pass equations (for $P(\mathbf{x}_t|D)$, $t = 1, \dots, T$), given that you have the Kalman filter outcomes available. Use the normal form (moment parameters), it is simpler this way. Do state the final recursion of means and covariance matrices.

Hint: Yes, you need those Gaussian formulae. $f_t(\mathbf{x}_{t-1}, \mathbf{x}_t)$ above will be some simplification of $P(\mathbf{x}_{t-1}|\mathbf{x}_t, D)$. Work out its mean and covariance from the local filtering distributions' moments, and recall that the former is affine, the latter independent of \mathbf{x}_t . Then, average over $P(\mathbf{x}_t|D)$.