

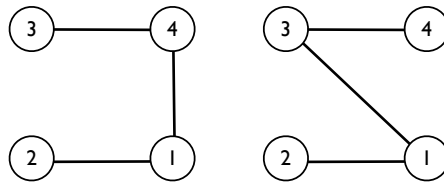
Homework 2

Covers lecture slides

3. Markov property
4. Elimination algorithm
5. Max-product algorithm

Problem 2.1

In this problem, we will show that when the distribution $\mu(x)$ is not strictly positive (i.e. $\mu(x) = 0$ for some x), then the I-map for this distribution is not unique. Consider a distribution of 4 binary random variables x_1, x_2, x_3 , and x_4 such that $\mu(x_1 = x_2 = x_3 = x_4 = 1) = 0.5$ and $\mu(x_1 = x_2 = x_3 = x_4 = 0) = 0.5$. The following two undirected graphical models are both minimal I-maps for this distribution, hence it is not unique.



- (a) Prove that the two undirected graphical models above are minimal I-maps for the distribution $\mu(x)$. You need to show that both graphs are I-maps for the given distribution $\mu(x)$ and that removing any edge results in introducing independencies that are not implied by the distribution $\mu(x)$.
- (b) Now, we show that starting with a complete graph and eliminating edges that are pairwise conditionally independent does not always give you an I-map (minimal or not). Start with a complete graph K_4 . For each pair of nodes, eliminate the edge between this pair if they are conditionally independent given the rest of the nodes in the graph. Continue this procedure for all pairs of nodes and examine the resulting graph. Is this an I-map of the distribution $\mu(x_1, x_2, x_3, x_4)$?

Recall from class, that a distribution over x is (globally) Markov with respect to $G = (V, E)$ if, for any disjoint subsets of nodes A, B, C such that B separates A from C , $x_A - x_B - x_C$ is satisfied. Recall two other notions of Markovity. A distribution is pairwise Markov with respect to G if, for any two nodes i and j not directly linked by an edge in G , the corresponding variables x_i and x_j are independent conditioned on all of the remaining variables, i.e. for all $(i, j) \notin E$,

$$x_i - x_{V \setminus \{i, j\}} - x_j$$

A distribution is locally Markov with respect to G if any node i , when conditioned on the variables on the neighbors of i , is independent of the remaining variables, i.e. for all $i \in V$,

$$x_i - x_{\partial i} - x_{V \setminus \{i, \partial i\}}$$

- (c) Using the example of distribution on 4 random variables as a counter example, prove that a distribution is pairwise Markov w.r.t. G does not always imply that it is locally Markov w.r.t. the same graph G . (However, if the distribution is positive, pairwise Markovity implies local and global Markovity.)
- (d) Using the definitions of Markov properties, prove that if a distribution is globally Markov with respect to G , then it is locally Markov with respect to G .
- (e) (Optional) Using the definitions of Markov properties, prove that if a distribution is locally Markov with respect to G , then it is pairwise Markov with respect to G .

Problem 2.2

Consider a stochastic process that transitions among a finite set of states s_1, \dots, s_k over time steps $i = 1, \dots, N$. The random variables X_1, \dots, X_N representing the state of the system at each time step are generated as follows:

- Sample the initial state $X_1 = s$ from an initial distribution p_1 , and set $i := 1$.
- Repeat the following:
 - Sample a duration d from a duration distribution p_D over the integers $\{1, \dots, M\}$, where M is the maximum duration.
 - Remain in the current state s for the next d time steps, i.e., set

$$X_i := X_{i+1} := \dots := X_{i+d-1} := s$$

- Sample a successor state s' from a transition distribution $p_T(\cdot|s)$ over the other states $s' \neq s$ (so there are no self-transitions).
- Assign $i := i + d$ and $s := s'$.

This process continues indefinitely, but we only observe the first N time steps. You need not worry about the end of the sequence to do any of the problems. As an example calculation with this model, the probability of the sample state sequence $s_1, s_1, s_1, s_2, s_3, s_3$ is

$$p_1(s_1)p_D(3)p_T(s_2|s_1)p_D(1)p_T(s_3|s_2) \sum_{2 \geq d \leq M} p_D(d) .$$

Finally, we do not directly observe the X_i 's, but instead observe emissions y_i at each step sampled from a distribution $p_{Y_i|X_i}(y_i|x_i)$.

- (a) For this part only, suppose $M = 2$, and $p_D(d) = \begin{cases} 0.6 & \text{for } d = 1 \\ 0.4 & \text{for } d = 2 \end{cases}$, and each X_i takes on a value from an alphabet $\{a, b\}$. Draw a minimal directed I-map for the first five time steps using the variables $(X_1, \dots, X_5, Y_1, \dots, Y_5)$. Explain why none of the edges can be removed. [Note: you do not need to solve part (a) in order to solve part (b) and (c).]
- (b) This process can be converted to an HMM using an *augmented state representation*. In particular, the states of this HMM will correspond to pairs (x, t) , where x is a state in the original system, and t represents the time elapsed in that state. For instance, the state sequence $s_1, s_1, s_1, s_2, s_3, s_3$ would be represented as $(s_1, 1), (s_1, 2), (s_1, 3), (s_2, 1), (s_3, 1), (s_3, 2)$. the transition and emission distribution for the HMM take the forms

$$\tilde{p}_{X_{i+1}, T_{i+1}|X_i, T_i}(x_{i+1}, t_{i+1}|x_i, t_i) = \begin{cases} \phi(x_i, x_{i+1}, t_i) & \text{if } t_{i+1} = 1 \text{ and } x_{i+1} \neq x_i \\ \xi(x_i, t_i) & \text{if } t_{i+1} = t_i + 1 \text{ and } x_{i+1} = x_i \\ 0 & \text{otherwise} \end{cases}$$

and $\tilde{p}_{Y_i|X_i,T_i}(y_i|x_i,t_i)$, respectively. Express $\phi(x_i, x_{i+1}, t_i)$, $\xi(x_i, t_i)$, and $\tilde{p}_{Y_i|X_i,T_i}(y_i|x_i, t_i)$ in terms of parameters $p_1, p_D, p_T, p_{Y_i|X_i}, k, N$, and M of the original model.

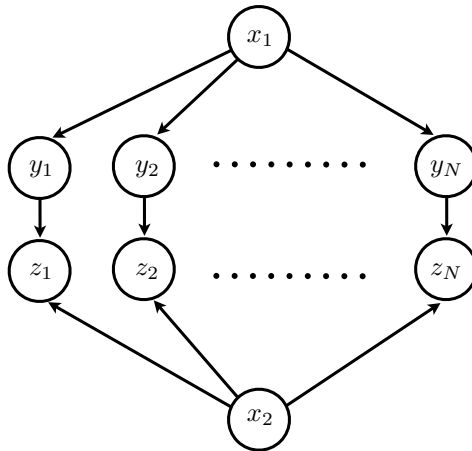
- (c) We wish to compute the marginal probability for the final state X_N given the observations Y_1, \dots, Y_N . If we naively apply the sum-product algorithm to the construction in part (b), the computational complexity is $O(Nk^2M^2)$. Show that by exploiting additional structure in the model, it is possible to reduce the complexity to $O(N(k^2 + kM))$. In particular, give the corresponding rules for computing the forward messages $\nu_{i+1 \rightarrow i+2}(x_{i+1}, t_{i+1})$ from the previous message $\nu_{i \rightarrow i+1}(x_i, t_i)$. Do not worry about the beginning or the end of the sequence and restrict your attention to $2 \leq i \leq N - 1$. [Hint: substitute your solution from part (b) into the standard update rule for HMM messages and simplify as much as possible.]
 [Note: If you cannot fully solve this part of the problem, you can receive substantial partial credit by constructing an algorithm with complexity $O(Nk^2M)$.]

Problem 2.3

Consider random variables $X_1, X_2, Y_1, \dots, Y_N, Z_1, \dots, Z_N$ distributed according to

$$p_{X_1, X_2, Y, Z}(x_1, x_2, y, z) = p_{X_1}(x_1)p_{X_2}(x_2) \prod_{i=1}^N \left[p_{Y|X_1}(y_i|x_1)p_{Z|Y, X_2}(z_i|y_i, x_2) \right],$$

where $X_1, Y_1, \dots, Y_N, Z_1, \dots, Z_N$ take on values in $\{1, 2, \dots, K\}$ and X_2 instead takes on a value in $\{1, 2, \dots, N\}$. A minimal directed I-map for the distribution is as follows:



Assume throughout this problem that the complexity of table lookups for $p_{X_1}, p_{X_2}, p_{Y|X_1}$, and $p_{Z|Y, X_2}$ are $O(1)$.

- (a) A Bayesian network represented by a directed acyclic graph can be turned into a Markov random field by *moralization*. The moralized counterpart of a directed acyclic graph is formed by connecting all pairs of nodes that have a common child, and then making all edges in the graph undirected. Draw the moral graph over random variables $X_1, X_2, Y_1, \dots, Y_N$ conditioned on Z_1, \dots, Z_N . In other words, find an undirected I-map for the distribution of random variables $X_1, X_2, Y_1, \dots, Y_N$ conditioned on Z_1, \dots, Z_N .

Provide a good elimination ordering for computing marginals of $X_1, X_2, Y_1, \dots, Y_N$ conditioned on Z_1, \dots, Z_N . For your elimination ordering, determine α and β such that complexity of computing $p_{X_1|Z_1, \dots, Z_N}$ using the associated elimination algorithm is $O(N^\alpha K^\beta)$.

- (b) For the remainder of this problem, suppose that we also have the following *context-dependent* conditional independencies: Y_i is conditionally independent of Z_i given $X_2 = c$ for all $i \neq c$. For fixed z_1, \dots, z_N, x_1 , and c , show that

$$p_{Z_1, \dots, Z_N | X_1, X_2}(z_1, \dots, z_N | x_1, c) = \eta(x_1, c, z_c) \lambda(c, z_1, \dots, z_{c-1}, z_{c+1}, \dots, z_N)$$

for some function $\eta(x_1, c, z_c)$ that can be evaluated in $O(K)$ operations for fixed (x_1, c, z_c) , and some function $\lambda(c, z_1, \dots, z_N)$ that can be evaluated in $O(N)$ operations for fixed (c, z_1, \dots, z_N) . Express $\eta(x_1, c, z_c)$ in terms of $p_{Y|X_1}$ and $p_{Z|Y, X_2}$, and $\lambda(c, z_1, \dots, z_{c-1}, z_{c+1}, \dots, z_N)$ in terms of $p_{Z|X_2}$.

Problem 2.4 [Optional]

The graph G is a perfect undirected map for some strictly positive distribution $\mu(x)$ over a set of random variables $x = (x_1, \dots, x_n)$, each of which takes values in a discrete set \mathcal{X} . Choose some variable x_i and let x_A denote the rest of the variables in the model, i.e., $\{x_i, x_A\} = \{x_1, \dots, x_N\}$. Construct the graph G' from G by removing the node x_i and all its edges. Let some value $c \in \mathcal{X}$ be given. Show that G' is not necessarily a perfect map for the conditional distribution $\mathbb{P}_{x_A | x_i}(\cdot | c)$ by giving a counterexample.

Problem 2.5

Consider the (parallel) sum-product algorithm on an undirected tree $T = (V, E)$ with compatibility functions ψ_{ij} such that $\mu(x) = \prod_{(i,j) \in E} \psi_{ij}(x_i, x_j)$. Consider any initialization of messages, which is denoted by $\nu_{i \rightarrow j}^{(0)}(x_i)$ for all directions $i \rightarrow j$ and all states x_i . Messages at step $t \geq 1$ are denoted by $\nu_{i \rightarrow j}^{(t)}(x_i)$. In this problem, we will prove by induction that the sum-product algorithm, with the parallel schedule, converges in at most diameter of the graph iterations. (Diameter of the graph is the length of the longest path.)

- (a) For $D = 1$, the result is immediate. Consider a graph of diameter D . At each time step the message that each of the leaf nodes sends out to its neighbors is constant because it does not depend on messages from any other nodes. Construct a new undirected graphical model $T' = (V', E')$ by stripping each of the leaf nodes from the original graph T . Let $\psi'_{ij}(x_i, x_j)$ be the compatibility functions for the new graphical model, and $\nu'_{i \rightarrow j}(x_i)$ be the messages of (parallel) sum-product algorithm on the new graphical model. Let L be the set of leaves in T and L' be the set of nodes that is adjacent to a node in L . For the new graphical model, we add, for all $i \in L'$,

$$\psi'_i(x_i) = \psi_i(x_i) \prod_{k \in \partial i \cap L} \sum_{x_k} \nu_{k \rightarrow i}^{(0)}(x_k) \psi_{ki}(x_k, x_i)$$

where $\psi_i(x_i) = 1$ if $\psi_i(x_i)$ is not defined for the original graph G and for all other edges we keep the original compatibility functions

$$\psi'_{ij}(x_i, x_j) = \psi_{ij}(x_i, x_j).$$

Also we initialize the messages as

$$\nu'_{i \rightarrow j}(x_i) = \nu_{i \rightarrow j}^{(1)}(x_i).$$

Show that $\nu'_{i \rightarrow j}(x_i) = \nu_{i \rightarrow j}^{(t+1)}(x_i)$ for all $(i, j) \in E'$ and all $t \geq 0$.

- (b) Argue that T' has diameter strictly less than $D - 1$.

- (c) By the induction assumption that the parallel sum-product algorithm converges to a fixed point after at most d time steps when the diameter is $d \leq D - 1$, the sum-product algorithm on T' converges after at most $D - 2$ time steps. Show that if we add back the leaf nodes into T' and run (parallel) sum-product algorithm for one more time step, then all messages will have converged to a fixed point.

Problem 2.6

For $\ell \in \mathbb{N}$, let $G_\ell = (V_\ell, E_\ell)$ be an $\ell \times \ell$ two-dimensional grid¹. We consider an Ising model on G_ℓ with parameters $\theta = \{\theta_{ij}, \theta_i : (i, j) \in E_\ell, i \in V_\ell\}$. This is the probability distribution over $x \in \{+1, -1\}^{V_\ell}$

$$\mu(x) = \frac{1}{Z_G} \exp \left\{ \sum_{(i,j) \in E_\ell} \theta_{ij} x_i x_j + \sum_{i \in V_\ell} \theta_i x_i \right\} \quad (1)$$

- (a) Write the belief propagation (BP) update equations for this model. Also write the update equation for the log-likelihood ratio

$$L_{i \rightarrow j}^{(t)} = \frac{1}{2} \log \left(\frac{\nu_{i \rightarrow j}^{(t)}(+1)}{\nu_{i \rightarrow j}^{(t)}(-1)} \right)$$

- (b) We give a MATLAB implementation of the code `bpsol.m` from the course website. Make yourself familiar with it to answer the following questions.
- (c) Consider the case $\ell = 10$ (and hence $n = 100$ nodes). For each $\beta \in \{0.2, 0.4, \dots, 2.8, 3.0\}$, generate an instance by drawing θ_i, θ_{ij} uniformly random in $[0, \beta]$. Run the BP iteration and monitor convergence by computing the quantity

$$\Delta(t) \equiv \frac{1}{|\vec{E}_\ell|} \sum_{(i,j) \in \vec{E}_\ell} |\nu_{i \rightarrow j}^{(t+1)}(+1) - \nu_{i \rightarrow j}^{(t)}(+1)|. \quad (2)$$

Here \vec{E}_ℓ denotes the set of directed edges in G_ℓ , in particular $|\vec{E}_\ell| = 2|E_\ell|$.

Plot $\Delta(t = 15)$ and $\Delta(t = 25)$ versus β , for the random instances generated with $\beta \in \{0.2, 0.4, \dots, 2.8, 3.0\}$. Comment on the results.

- (d) Repeat the calculation at the previous point, with now θ_i, θ_{ij} uniformly random in $[-\beta, +\beta]$, with $\beta \in \{0.2, 0.4, \dots, 2.8, 3.0\}$. Comment on the results.

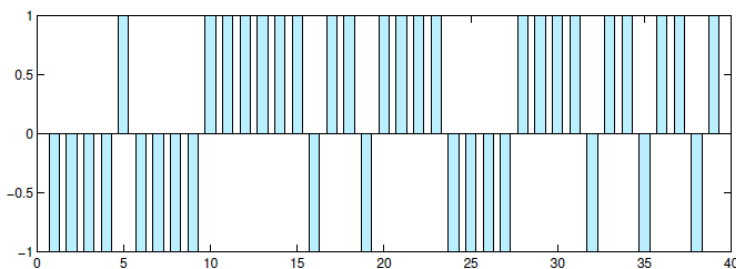
Problem 2.7

In this problem, you will implement the sum-product algorithm on a line graph and analyze the behavior of S&P 500 index over a period of time. The following figure shows the price of S&P 500 index from January 2, 2009 to September 30, 2009 (<http://finance.yahoo.com>).

For each week, we measure the price movement relative to the previous week and denote it using a binary variable (+1 indicates up and -1 indicates down). The price movements from week 1 (the week of January 5) to week 39 (the week of September 28) are plotted below:

Consider a hidden Markov model in which x_t denotes the economic state (good or bad) of week t and y_t denotes the price movement (up or down) of the S&P 500 index. We assume that $x_{t+1} = x_t$ with

¹Namely $V_\ell = [\ell] \times [\ell]$ and, for any two vertices $i, j \in V_\ell$, $i = (i_1, i_2)$, $j = (j_1, j_2)$, $i_1, i_2, j_1, j_2 \in [\ell]$, $(i, j) \in E_\ell$ if and only if $i_1 = j_1$ and $|i_2 - j_2| = 1$, or $i_2 = j_2$ and $|i_1 - j_1| = 1$.



probability 0.8, and $\mathbb{P}_{Y_t|X_t}(y_t = +1|x_t = \text{'good'}) = \mathbb{P}_{Y_t|X_t}(y_t = -1|x_t = \text{'bad'}) = q$. In addition, assume that $\mathbb{P}_{X_1}(x_1 = \text{'bad'}) = 0.8$. Download the file `sp500.mat` from course website, and load it into MATLAB. The variable `price_move` contains the binary data above. Implement the (sequential) sum-product algorithm and submit the code with the homework solutions.

- Assume that $q = 0.7$. Plot $\mathbb{P}_{X_t|Y}(x_t = \text{'good'}|y)$ for $t = 1, 2, \dots, 39$. What is the probability that the economy is in a good state in the week of September 28, 2009 (week 39)?
- Repeat (a) for $q = 0.9$. Compare the results of (a) and (b).

Problem 2.8 [Optional]

Consider a hidden Markov model (HMM) with binary states $x_i \in \{0, 1\}$ for $i \in \{1, \dots, n\}$ and observations y_i 's. For simplicity, let us assume that the model is homogeneous, i.e., $\psi_{i,i+1}(x_i, x_{i+1}) = \psi(x_i, x_{i+1})$ and $\phi_i(x_i, y_i) = \phi(x_i, y_i)$. Given the observations y_i 's we are interested in state estimates $\hat{x}_i(y_1, \dots, y_n)$ that maximizes the probability that at least one of those state estimates \hat{x}_i is correct.

- The desired state estimates can be expressed in the form

$$(\hat{x}_1, \dots, \hat{x}_n) \in \arg \min \mathbb{P}(X_1 = f(\hat{x}_1) \wedge \dots \wedge X_n = f(\hat{x}_n) | y_1, \dots, y_n).$$

Determine the function $f(\cdot)$.

- (b) Show that if only the marginal distributions $\mu(x_i|y_1 \dots, y_n)$, $i \in \{1, \dots, n\}$ for the model are available, the desired state estimates cannot be determined. In particular, construct two HMMs whose marginals coincide, but whose state estimates differ.
 [Hint: it suffices to consider a model with $n = 2$, and in which the observations are independent of the states thus can be ignored. Accordingly, express your answer in the form of two compatibility functions $\psi(x_1, x_2)$ and $\psi'(x_1, x_2)$.]
- (c) Construct an example of an HMM in which our desired estimates are not the same as the MAP estimates obtained from running the max-product algorithm on our model. The same hint in part (b) applies, so again give your answer in the form of a compatibility function $\psi(x_1, x_2)$.
- (d) Let's assume that you are given two pieces of code (e.g., matlab scripts).

The first routine implements the sum-product algorithm, taking as input the potential functions that describe a homogeneous HMM, and an associated list of n observations. It produces as output the list of marginal distributions for each associated n states conditioned on the full set of n observations, for the specified HMM.

The second routine implements the max-product algorithm, taking the same inputs as sum-product algorithm, but producing as output the max-marginals for each associated n states conditioned on the full set of n observations, for the specified HMM.

Describe how to use one or both of these routines to compute the desired estimates $\hat{x}_i(y_1, \dots, y_n)$ for $i \in \{1, \dots, n\}$ for our model of interest, assuming that the potentials are strictly positive. You are free to use these routines with any input values you like (whether or not related to the model of interest), and you can further process the outputs of these routines to compute the desired state estimates. However, in such further processing, you are not allowed to (re)use the model's potential functions or observations.