

# Statistical inference with probabilistic graphical models

Angélique Drémeau\*, Christophe Schülke†, Yingying Xu‡, Devavrat Shah§

September 18, 2014

*These are notes from the lecture of Devavrat Shah given at the autumn school “Statistical Physics, Optimization, Inference, and Message-Passing Algorithms”, that took place in Les Houches, France from Monday September 30th, 2013, till Friday October 11th, 2013. The school was organized by Florent Krzakala from UPMC & ENS Paris, Federico Ricci-Tersenghi from La Sapienza Roma, Lenka Zdeborová from CEA Saclay & CNRS, and Riccardo Zecchina from Politecnico Torino.*

---

\*École Normale Supérieure, France

†Université Paris Diderot, France

‡Tokyo Institute of Technology, Japan

§Massachusetts Institute of Technology, USA

# Contents

<b>1</b>	<b>Introduction to Graphical Models</b>	<b>3</b>
1.1	Inference . . . . .	3
1.2	Graphical models . . . . .	3
1.2.1	Directed GMs . . . . .	3
1.2.2	Undirected GMs . . . . .	4
1.2.3	Cliques . . . . .	5
1.3	Factor graphs . . . . .	7
1.3.1	Image processing . . . . .	8
1.3.2	Crowd-sourcing . . . . .	8
1.4	MAP and MARG . . . . .	9
<b>2</b>	<b>Inference Algorithms: Elimination, Junction Tree and Belief Propagation</b>	<b>10</b>
2.1	The elimination algorithm . . . . .	10
2.2	Junction Tree property and chordal graphs . . . . .	12
2.2.1	Junction Tree (JCT) property . . . . .	12
2.2.2	Chordal graph . . . . .	13
2.2.3	Procedure to find a JCT . . . . .	14
2.2.4	Tree width . . . . .	14
2.3	Belief propagation (BP) algorithms . . . . .	15
2.3.1	Factor graphs . . . . .	16
<b>3</b>	<b>Understanding Belief Propagation</b>	<b>17</b>
3.1	Existence of a fixed point . . . . .	17
3.2	Nature of the fixed points . . . . .	18
3.2.1	Background on Nonlinear Optimization . . . . .	19
3.2.2	Belief Propagation as a variational problem . . . . .	19
3.3	Can the fixed points be reached? . . . . .	20
3.3.1	The hardcore model . . . . .	22
<b>4</b>	<b>Learning Graphical Models</b>	<b>22</b>
4.1	Parameter learning . . . . .	22
4.1.1	Single parameter learning . . . . .	22
4.1.2	Directed graphs . . . . .	23
4.1.3	Undirected graphs . . . . .	24
4.2	Graphical model learning . . . . .	24
4.2.1	Directed graphs . . . . .	25
4.2.2	Undirected graphs . . . . .	25
4.3	Latent Graphical Model learning: the Expectation-maximization algorithm . . . . .	26
	<b>References</b>	<b>27</b>

# 1 Introduction to Graphical Models

## 1.1 Inference

Consider two random variables  $A$  and  $B$  with a joint probability distribution  $P_{A,B}$ . From the observation of the realization of one of those variables, say  $B = b$ , we want to infer the one that we did not observe. To that end, we compute the conditional probability distribution  $P_{A|B}$ , and use it to obtain an estimate  $\hat{a}(b)$  of  $a$ .

To quantify how good this estimate is, we introduce the **error probability**:

$$\begin{aligned} P_{error} &\triangleq P(A \neq \hat{a}(b)|B = b) \\ &= 1 - P(A = \hat{a}(b)|B = b), \end{aligned} \tag{1}$$

and we can see from the second equality that minimizing this error probability is equivalent to the following maximization problem, called maximum a posteriori (**MAP**) problem:

$$\hat{a}(b) = \arg \max_a P_{A|B}(a|b). \tag{2}$$

The problem of computing  $P_{A|B}(a|b)$  for all  $a$  given  $b$  is called the marginal (**MARG**) problem. When the number of random variables increases, the MARG problem becomes difficult, because an exponential number of combinations has to be calculated.

**Fano's inequality** provides us an information-theoretical way of gaining insight into how much information about  $a$  the knowledge of  $b$  can give us:

$$P_{error} \geq \frac{H(A|B) - 1}{\log|A|}, \tag{3}$$

with

$$\begin{aligned} H(A|B) &= \sum_b P_B(b)H(A|B = b), \\ H(A|B = b) &= \sum_a P_{A|B}(a|b) \log \left( \frac{1}{P_{A|B}(a|b)} \right). \end{aligned}$$

Fano's inequality formalises only a theoretical bound that does not tell us how to actually make an estimation. From a practical point of view, graphical models (GM) constitute here a powerful tool allowing us to write algorithms that solve inference problems.

## 1.2 Graphical models

### 1.2.1 Directed GMs

Consider  $N$  random variables  $X_1 \cdots X_N$  on a discrete alphabet  $\mathcal{X}$ , and their joint probability distribution  $P_{X_1 \cdots X_N}$ . We can always factorize this joint distribution in the following way:

$$P_{X_1 \cdots X_N} = P_{X_1} P_{X_2|X_1} \cdots P_{X_N|X_1 \cdots X_{N-1}} \tag{4}$$

and represent this factorized form by the following directed graphical model:

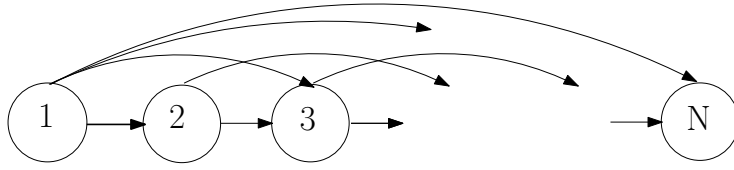


Figure 1: A directed graphical model representing the factorized form (4).

In this graphical model, each node is affected to a random variable, and each directed edge represents a conditioning. The way that we factorized the distribution, we obtain a complicated graphical model, in the sense that it has many edges. A much simpler graphical model would be:

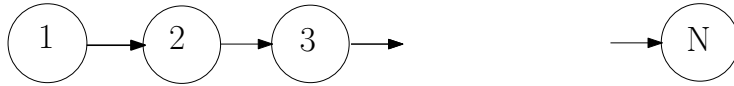


Figure 2: A simpler graphical model representing the factorized form (5).

The latter graphical model corresponds to a factorization in which each of the probability distributions in the product is conditioned on only one variable:

$$P_{X_1 \dots X_N} = P_{X_1} P_{X_2 | X_1} \dots P_{X_N | X_{N-1}} \quad (5)$$

In the most general case, we can write a distribution represented by a directed graphical model in the factorized form:

$$P_{X_1 \dots X_N} = \prod_i P_{X_i | X_{\Pi_i}}, \quad (6)$$

where  $X_{\Pi_i}$  is the set containing the parents of  $X_i$  (the vertices from which an edge points to  $i$ ).

The following **notations** will hold for the rest of this chapter:

- random variables are capitalized:  $X_i$ ,
- realizations of random variables are lower case:  $x_i$ ,
- a set of random variables  $\{X_1 \dots X_N\}$  is noted  $\underline{X}$ ,
- a set of realizations of  $\underline{X}$  is noted  $\underline{x}$ ,
- the subset of random variables of indices in  $S$  is noted  $X_S$ .

### 1.2.2 Undirected GMs

Another type of graphical model is the undirected graphical model. In that case, we define the graphical model not through the **factorization**, but through **independence**.

Let:

$\mathcal{G}(\mathcal{V}, \mathcal{E})$  be an undirected graphical model, where  
 $\mathcal{V} = \{1, \dots, N\}$  is the set of vertices, and  
 $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  is the set of edges.

Each vertex  $i \in \mathcal{V}$  of this GM represents one random variables  $X_i$ , and each edge  $(i, j) \in \mathcal{E}$  represents a conditional dependence. As the GM is undirected, we have  $(i, j) \equiv (j, i)$ .

We define:

$$N(i) \triangleq \{j \in \mathcal{V} | (i, j) \in \mathcal{E}\} \quad \text{the set containing the neighbours of } i. \quad (7)$$

Undirected graphical model captures following dependence:

$$P_{X_i | X_{\mathcal{V} \setminus \{i\}}} \equiv P_{X_i | X_{N(i)}}, \quad (8)$$

meaning that only variables connected by edges have a conditional dependence.

Let  $A \subset \mathcal{V}$ ,  $B \subset \mathcal{V}$ ,  $C \subset \mathcal{V}$ . We write that  $X_A \perp X_B | X_C$  if  $A$  and  $B$  are disjoint and if all paths leading from one element of  $A$  to one element of  $B$  lead over an element of  $C$ , as is illustrated in Fig. 3. In other words, if we remove  $C$ , then  $A$  and  $B$  are unconnected (Fig. 4).

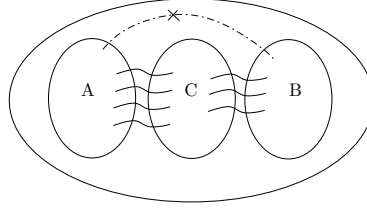


Figure 3: Schematic view of a graphical model in which  $X_A \perp X_B | X_C$ . All paths leading from  $A$  to  $B$  go through  $C$ .

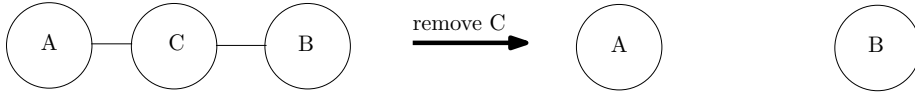


Figure 4: Simple view showing the independence of  $A$  and  $B$  conditioned on  $C$ .

Undirected GMs are also called **Markov random fields** (MRF).

### 1.2.3 Cliques

**(Definition)** A clique is a subgraph of a graph in which all possible pairs of vertices are linked by an edge. A maximal clique is a clique that is contained by no other clique.

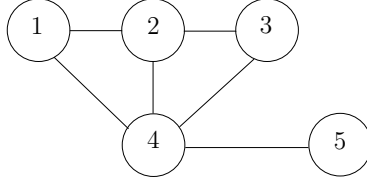


Figure 5: In this graphical model, the maximal cliques are  $\{1, 2, 4\}$ ,  $\{2, 3, 4\}$  and  $\{4, 5\}$ .

**Theorem 1** ([4]) Given a MRF  $\mathcal{G}$  and a probability distribution  $P_{\underline{X}}(\underline{x}) > 0$ . Then:

$$P_{\underline{X}}(\underline{x}) \propto \prod_{C \in \mathcal{C}} \phi_C(x_C) \quad (9)$$

where  $\mathcal{C}$  is the set of cliques of  $\mathcal{G}$ .

**Proof 1** ([3]) for  $\mathcal{X} = \{0, 1\}$ .

We will show the following, equivalent formulation:

$$P_{\underline{X}}(\underline{x}) \propto e^{\sum_{C \in \mathcal{C}} V_C(x_C)} \quad (10)$$

by exhibiting the solution:

$$V_C(x_C) = \begin{cases} Q(C) & \text{if } x_C = \mathbb{1}_C, \\ 0 & \text{otherwise,} \end{cases} \quad (11)$$

with

$$Q(C) = \sum_{A \subseteq C} (-1)^{|C-A|} \underbrace{\ln P_{\underline{X}}(x_A = \mathbb{1}_A, x_{V \setminus A} = \mathbb{0})}_{\triangleq G(A)}. \quad (12)$$

Suppose we have an assignment  $\underline{X} \mapsto N(\underline{X}) = \{i | x_i = 1\}$ . We want to prove that:

$$\begin{aligned} G(N(\underline{X})) &\triangleq \ln P_{\underline{X}}(\underline{x}), \\ &= \sum_{C \in \mathcal{C}} V_C(x_C), \\ &= \sum_{C \subseteq N(\underline{x})} Q(C). \end{aligned} \quad (13)$$

This is equivalent to proving the two claims:

$$C1: \quad \forall S \subset \mathcal{C}, \quad G(S) = \sum_{A \subseteq S} Q(A)$$

$$C2: \quad \text{if } A \text{ is not a clique,} \quad Q(A) = 0$$

Let us begin by proving C1:

$$\begin{aligned} \sum_{A \subseteq S} Q(A) &= \sum_{A \subseteq S} \sum_{B \subseteq A} (-1)^{|A-B|} G(B) \\ &= \sum_{B \subseteq S} G(B) \left( \sum_{B \subseteq A \subseteq S} (-1)^{|A-B|} \right) \end{aligned} \quad (14)$$

where we note that the term in brackets is zero except when  $B = S$ , because we can rewrite it as

$$\sum_{0 \leq l \leq k} (-1)^l \binom{l}{k} = (-1 + 1)^k = 0. \quad (15)$$

Therefore  $G(S) = \sum_{A \subseteq S} Q(A)$ .

For C2, suppose that  $A$  is not a clique, which allows us to choose  $(i, j) \in A$  with  $(i, j) \notin \mathcal{E}$ . Then

$$Q(A) = \sum_{B \subseteq A \setminus \{i, j\}} (-1)^{|A-B|} [G(B) - G(B+i) + G(B+i+j) - G(B+j)].$$

Let us show that the term in brackets is zero by showing

$$G(B+i+j) - G(B+j) = G(B+i) - G(B)$$

or equivalently

$$\ln \frac{P_X(x_B = \mathbb{1}_B, x_i = 1, x_j = 1, x_{\mathcal{V} \setminus \{i, j, B\}} = 0)}{P_X(x_B = \mathbb{1}_B, x_i = 0, x_j = 1, x_{\mathcal{V} \setminus \{i, j, B\}} = 0)} = \ln \frac{P_X(x_B = \mathbb{1}_B, x_i = 1, x_j = 0, x_{\mathcal{V} \setminus \{i, j, B\}} = 0)}{P_X(x_B = \mathbb{1}_B, x_i = 0, x_j = 0, x_{\mathcal{V} \setminus \{i, j, B\}} = 0)},$$

where  $\mathcal{V} \setminus \{i, j, B\}$  stands for the set of all vertices except  $i, j$  and those in  $B$ . We see that the only difference between the left-hand side and the right-hand side is the value taken by  $x_j$ . Using Bayes' rule, we can rewrite both the right-hand side and the left-hand side under the form

$$\ln \frac{P_X(X_i = 1 | X_j = \pm 1, X_B = \mathbb{1}_B, X_{\mathcal{V} \setminus \{i, j, B\}} = 0)}{P_X(X_i = 0 | X_j = \pm 1, X_B = \mathbb{1}_B, X_{\mathcal{V} \setminus \{i, j, B\}} = 0)}.$$

As  $(i, j) \notin \mathcal{E}$ , the conditional probabilities on  $X_i$  do not depend on the value taken by  $X_j$ , and therefore the right-hand side equals the left-hand side,  $Q(A) = 0$  and C2 is proved.

### 1.3 Factor graphs

Thanks to the Hammersley-Clifford theorem, we know that we can write a probability distribution corresponding to a MRF  $\mathcal{G}$  in the following way

$$P_{\underline{X}}(\underline{x}) \propto \prod_{C \in \mathcal{C}^*} \phi_C(x_C) \quad (16)$$

where  $\mathcal{C}^*$  is the set of maximal cliques of  $G$ . In a general definition, we can also write

$$P_{\underline{X}}(\underline{x}) \propto \prod_{F \in \mathcal{F}} \phi_F(x_F) \quad (17)$$

where the partition  $\mathcal{F} \subseteq 2^{\mathcal{V}}$  has nothing to do with any underlying graph.

In what follows, we give two examples in which introducing factor graphs is a natural approach to an inference problem.

### 1.3.1 Image processing

We consider an image with binary pixels ( $\mathcal{X} = \{-1, 1\}$ ), and a probability distribution:

$$p(\underline{x}) \propto e^{\sum_{i \in V} \theta_i x_i + \sum_{(i,j) \in E} \theta_{ij} x_i x_j} \quad (18)$$

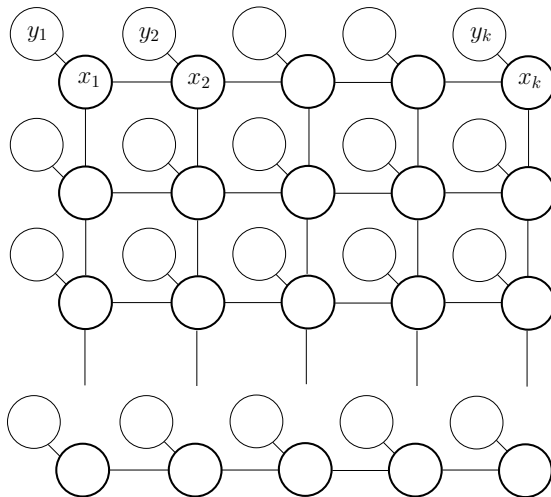


Figure 6: Graphical model representing a 2D image. The fat circles correspond to the pixels of the image  $x_k$ , and each one is linked to a noisy measurement  $y_k$ . Adjacent pixels are linked by edges that allow modelling the assumed smoothness of the image.

For each pixel  $x_k$ , we record a noisy version  $y_k$ . We consider natural images, in which big jumps in intensity between two neighbouring pixels are unlikely. This can be modelled with:

$$a \sum_i x_i y_i + b \sum_{(i,j) \in \mathcal{E}} x_i x_j \quad (19)$$

This way, the first term pushes  $x_k$  to match the measured value  $y_k$ , while the second term favours piecewise constant images. We can identify  $\theta_i \equiv a y_i$  and  $\theta_{ij} \equiv b$ .

### 1.3.2 Crowd-sourcing

Crowd-sourcing is used for tasks that are easy for humans but difficult for machines, and that are as hard to verify as to evaluate. Crowd-sourcing then consists in assigning to each of  $M$  human “workers” a subset of  $N$  tasks to evaluate, and to collect their answers  $A$ . Each worker has a different error probability  $p_i \in \{\frac{1}{2}, 1\}$ : either he gives random answers, or he is fully reliable. The goal is to infer both the correct values of each task,  $t_j$ , and the  $p_i$  of each worker. The factor graph corresponding to that problem is represented in Fig 7.



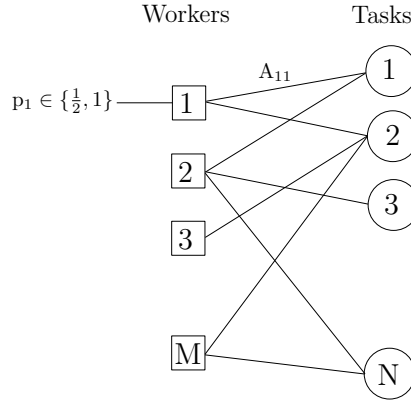


Figure 7: Graphical model illustrating crowd-sourcing. Each worker  $i$  is assigned a subset of the tasks for evaluation, and for each of those tasks  $a$ , his answer  $A_{ia}$  is collected.

The conditional probability distribution of  $\underline{t}$  and  $\underline{p}$  knowing the answers  $A$  reads

$$\begin{aligned} P_{\underline{t}, \underline{p} | A} &\propto P_{A | \underline{t}, \underline{p}} P_{\underline{t}, \underline{p}} \\ &\propto P_{A | \underline{t}, \underline{p}} \end{aligned} \quad (20)$$

where we assumed a uniform distribution on the joint probability  $P_{\underline{t}, \underline{p}}$ . Then

$$P_{A | \underline{t}, \underline{p}} = \prod_e P_{A_e | t_e, p_e} \quad (21)$$

with

$$P_{A_e | t_e, p_e} = \left( \left( \frac{p_e}{1 - p_e} \right)^{A_e t_e} (1 - p_e) p_e \right)^{\frac{1}{2}}. \quad (22)$$

#### 1.4 MAP and MARG

*MAP.* The MAP problem consists in solving:

$$\max_{\underline{x} \in \{0,1\}^N} \sum_i \theta_i x_i + \sum_{(i,j) \in \mathcal{E}} \theta_{ij} x_i x_j. \quad (23)$$

When  $\theta_{ij} \rightarrow -\infty$ , neighbouring nodes can not be in the same state anymore. This is the hard-core model, which is very hard to solve.

*MARG.* The MARG focuses on the evaluation of marginal probabilities, depending on only one random variable, for instance:

$$P_{X_1}(0) = \frac{Z(X_1 = 0)}{Z} \quad (24)$$

as well as conditional marginal probabilities:

$$P_{X_2|X_1}(X_2 = 0|X_1 = 0) = \frac{Z(X_1 = 0, X_2 = 0)}{Z(X_1 = 0)} \quad (25)$$

$$P_{X_N|X_1 \dots X_{N(1)}}(X_N = 0|X_1 \dots X_{N-1} = 0) = \frac{Z(\text{all } 0)}{Z(\text{all but } X_N \text{ are } 0)} \quad (26)$$

$$P_{X_1}(0) \times \dots \times P_{X_N|X_1 \dots X_{N-1}}(0) = \frac{1}{Z} \quad (27)$$

Both of these problems are computationally hard. Can we design efficient algorithms to solve them?

## 2 Inference Algorithms: Elimination, Junction Tree and Belief Propagation

In the MAP and MARG problems described previously, the hardness comes from the fact that with growing instance size, the number of combinations of variables over which to maximize or marginalize becomes quickly intractable. But when dealing with GMs, one can exploit the structure of the GM in order to reduce the number of combinations that have to be taken into account. Intuitively, the smaller the connectivity of the variables in the GM is, the smaller this number of combination becomes. We will formalize this by introducing the elimination algorithm, that gives us a systematic way of making fewer maximizations/marginalizations on a given graph. We will see how substantially the number of operations is reduced on a graph that is not completely connected.

### 2.1 The elimination algorithm

We consider the GM in Fig. 8 which is not fully connected. The colored subgraphs represent the maximal cliques.

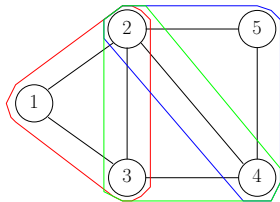


Figure 8: A GM and its maximal cliques.

Using decomposition (16), we can write

$$P_{\underline{X}}(\underline{x}) \propto \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot \phi_{245}(x_2, x_4, x_5). \quad (28)$$

We want to solve the MARG problem on this GM, for example for calculating the marginal probability of  $x_1$ :

$$P_{X_1}(x_1) = \sum_{x_2, x_3, x_4, x_5} P_{\underline{X}}(\underline{x}). \quad (29)$$

A priori, this requires to evaluate  $|\mathcal{X}|^4$  terms, each of them taking  $|\mathcal{X}|$  different values. In the end,  $3|\mathcal{X}||\mathcal{X}|^4$  operations are needed for calculating this marginal naively. But if we take advantage of the factorized form (28), we can eliminate some of the variables. The elimination process goes along these lines:

$$P_{X_1}(x_1) \propto \sum_{x_2, x_3, x_4, x_5} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot \phi_{245}(x_2, x_4, x_5), \quad (30)$$

$$\propto \sum_{x_2, x_3, x_4} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot \sum_{x_5} \phi_{245}(x_2, x_4, x_5), \quad (31)$$

$$\propto \sum_{x_2, x_3, x_4} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot m_5(x_2, x_4), \quad (32)$$

$$\propto \sum_{x_2, x_3} \phi_{123}(x_1, x_2, x_3) \cdot \sum_{x_4} \phi_{234}(x_2, x_3, x_4) \cdot m_5(x_2, x_4), \quad (33)$$

$$\propto \sum_{x_2, x_3} \phi_{123}(x_1, x_2, x_3) \cdot m_4(x_2, x_3), \quad (34)$$

$$\propto \sum_{x_2} \left( \sum_{x_3} \phi_{123}(x_1, x_2, x_3) m_4(x_2, x_3) \right), \quad (35)$$

$$\propto \sum_{x_2} m_3(x_1, x_2), \quad (36)$$

$$\propto m_2(x_1). \quad (37)$$

With this elimination process made, the number of operations necessary to compute the marginal scales as  $|\mathcal{X}|^3$  instead of  $|\mathcal{X}|^5$ , thereby greatly reducing the complexity of the problem by using the structure of the GM. Similarly, we can rewrite the MAP problem as follows

$$\max_{x_1, x_2, x_3, x_4, x_5} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot \phi_{245}(x_2, x_4, x_5), \quad (38)$$

$$= \max_{x_1, x_2, x_3, x_4} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot \max_{x_5} \phi_{245}(x_2, x_4, x_5), \quad (39)$$

$$= \max_{x_1, x_2, x_3, x_4} \phi_{123}(x_1, x_2, x_3) \cdot \phi_{234}(x_2, x_3, x_4) \cdot m_5^*(x_2, x_4), \quad (40)$$

$$= \max_{x_1, x_2, x_3} \phi_{123}(x_1, x_2, x_3) \cdot \max_{x_4} \phi_{234}(x_2, x_3, x_4) \cdot m_5^*(x_2, x_4), \quad (41)$$

$$= \max_{x_1, x_2, x_3} \phi_{123}(x_1, x_2, x_3) \cdot m_4^*(x_2, x_3), \quad (42)$$

$$= \max_{x_1, x_2} \left( \max_{x_3} \phi_{123}(x_1, x_2, x_3) m_4^*(x_2, x_3) \right), \quad (43)$$

$$= \max_{x_1, x_2} m_3^*(x_1, x_2), \quad (44)$$

$$= \max_{x_1} \left( \max_{x_2} m_3^*(x_1, x_2) \right), \quad (45)$$

leading to

$$x_1^* = \arg \max_{x_1} m_2^*(x_1). \quad (46)$$

Just like for the MARG problem, the complexity is reduced from  $|\mathcal{X}|^5$  (a priori) to  $|\mathcal{X}|^3$ . We would like to further reduce the complexity of the marginalizations (in  $|\mathcal{X}|^3$ ). One simple idea would be to reduce the GM into a linear



Figure 9: A linear graph. Each marginalization is computed in  $|\mathcal{X}|^2$  operations.

graph as in Fig. 9.

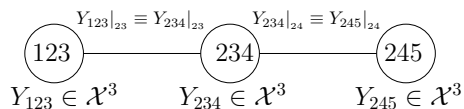


Figure 10: Linear GM obtained by grouping variables.

By grouping variables in the GM (Fig. 8), it is in fact possible to obtain a linear graph, as shown in Fig. 10, with the associated potentials  $\phi_{123}(Y_{123})$ ,  $\phi_{234}(Y_{234})$  and  $\phi_{245}(Y_{245})$  and the consistency constraints  $Y_{123}|_{23} \equiv Y_{234}|_{23}$  and  $Y_{234}|_{24} \equiv Y_{245}|_{24}$ . For other GMs, the simplest graph achievable by grouping variables might be a tree instead of a simple chain. But not all groupings of variables will lead to a tree graph that correctly represents the problem. In order for the grouping of variables to be correct, we need to build the tree attached to the maximal cliques, and we have to resort to the Junction Tree property.

## 2.2 Junction Tree property and chordal graphs

The Junction Tree property allows us to find groupings of variables under which the GM becomes a tree (if such groupings exist). On this tree, the elimination algorithm will need a lower number of maximizations/marginalizations than on the initial GM. However, there is a remaining problem: running the algorithm on the junction tree does not give a straightforward solution to the initial problem, as the variables on the junction tree are groupings of variables of the original problem. This means that further maximizations/marginalizations are then required to have a solution in terms of the variables of the initial problem.

### 2.2.1 Junction Tree (JCT) property

**(Definition)** A graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  is said to possess the JCT property if it has a Junction Tree  $\mathcal{T}$  which is defined as follows: it is a tree graph such that

- its nodes are maximal cliques of  $\mathcal{G}$
- an edge between nodes of  $\mathcal{T}$  is allowed only if the corresponding cliques share at least one vertex
- for any vertex  $v$  of  $\mathcal{G}$ , let  $\mathcal{C}_v$  denote set of all cliques containing  $v$ . Then  $\mathcal{C}_v$  forms a connected sub-tree of  $\mathcal{T}$ .

Two questions then arise

- Do all graphs have a JCT?

- If a graph has a JCT, how can we find it?

### 2.2.2 Chordal graph

**(Definition)** A graph is chordal if all of its loops have chords. Fig. 11 gives an illustration of the concept.

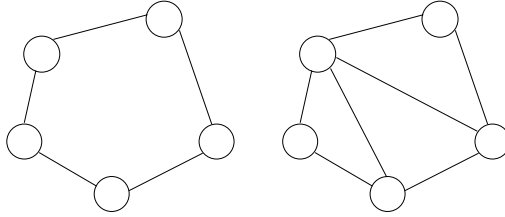


Figure 11: The graph on the left is not chordal, the one on the right is.

**Proposition 1**  $\mathcal{G}$  has a junction tree  $\Leftrightarrow \mathcal{G}$  is a chordal graph.

**Proof 2** of the implication  $\Leftarrow$ . Let us take a chordal graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  that is not complete, as represented in Fig. 12.

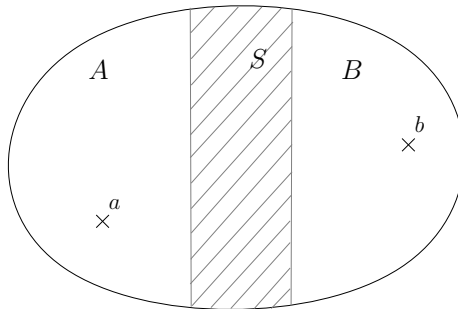


Figure 12: On a chordal graph that is not complete, two vertices  $a$  and  $b$  that are not connected, separated by a subgraph  $S$  that is fully connected.

We will use the two following lemmas that can be shown to be true:

1. If  $\mathcal{G}$  is chordal, has at least three nodes and is not fully connected, then  $\mathcal{V} = \mathcal{A} \cup \mathcal{B} \cup \mathcal{S}$ , where all three sets are disjoint and  $\mathcal{S}$  is a fully connected subgraph that separates  $\mathcal{A}$  from  $\mathcal{B}$ .
2. If  $\mathcal{G}$  is chordal and has at least two nodes, then  $\mathcal{G}$  has at least two nodes each with all neighbors connected. Furthermore, if  $\mathcal{G}$  is not fully connected, then there exist two nonadjacent nodes each with all its neighbors connected.

The property “If  $\mathcal{G}$  is a chordal graph with  $N$  vertices, then it has a junction tree.” can be shown by induction on  $N$ . For  $N = 2$ , the property is trivial. Now, suppose that the property is true for all integers up to  $N$ . Consider a chordal graph with  $N + 1$  nodes. By the second lemma,  $\mathcal{G}$  has a node  $a$  with all its neighbors connected. Removing it creates a graph  $\mathcal{G}'$  which is chordal, and therefore has a JCT,  $\mathcal{T}'$ . Let  $C$  be the maximal clique that  $a$  participates in. Either  $C \setminus a$  is a maximal-clique node in  $\mathcal{T}'$ , and in this case adding  $a$  to this clique node results in a junction tree  $\mathcal{T}$  for  $\mathcal{G}$ . Or  $C \setminus a$  is not a maximal-clique node in  $\mathcal{T}'$ . Then,  $C \setminus a$  must be a subset of a maximal-clique node  $D$  in  $\mathcal{T}'$ . Then, we add  $C$  as a new maximal-clique node in  $\mathcal{T}'$ , which we connect to  $D$  to obtain a junction tree  $\mathcal{T}$  for  $\mathcal{G}$ .

### 2.2.3 Procedure to find a JCT

Let  $G$  be the initial GM, and  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  be the GM in which  $\mathcal{V}$  is the set of maximal cliques of  $G$  and  $(c_1, c_2) \in \mathcal{E}$  if the maximal cliques  $c_1$  and  $c_2$  share a vertex. Let us take  $e = (c_1, c_2)$  with  $c_1, c_2 \in \mathcal{V}$  and define the weight of  $e$  as  $w_e = |c_1 \cap c_2|$ . Then, finding a junction tree of  $G$  is equivalent to finding the max-cut spanning tree of  $\mathcal{G}$ . Denoting by  $T$  the set of edges in a tree, we define the weight of the tree as

$$\begin{aligned} W(T) &= \sum_{e \in T} w_e \\ &= \sum_{e \in T} |c_1 \cap c_2| \\ &= \sum_{v \in V} \sum_{e \in T} \mathbb{1}_{\{v \in e\}}. \end{aligned} \tag{47}$$

and we claim that  $W(T)$  is maximal when  $T$  is a JCT.

**Procedure** to get the maximum weighted spanning tree

- List all edges in a decreasing order,
- Include  $e_i$  in  $\mathcal{E}_{i-1}$  if you can.

what we are left with at the end of the algorithm is the maximal weight spanning tree.

### 2.2.4 Tree width

**(Definition)** The width of a tree decomposition is the size of its maximal clique minus one.

*Toy examples*

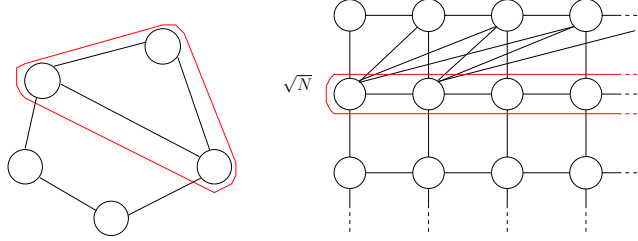


Figure 13: tree width = 2 (left), tree width =  $\sqrt{N}$  (right)

### 2.3 Belief propagation (BP) algorithms

Until now, everything we have done was exact. The elimination algorithm is an exact algorithm. But as we are interested in **efficient algorithms**, as opposed to exact algorithms with too high complexities to actually end in reasonable time, we will from now on **introduce approximations**.

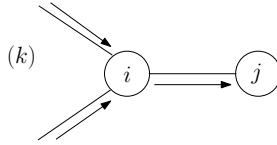


Figure 14: Message passing on a graph.

Coming back to the elimination algorithm (30)-(37), we can generalize the notations used as

$$m_i(x_j) \propto \sum_{x_i} \phi_i(x_i) \cdot \phi_{i,j}(x_i, x_j) \cdot \prod_k m_k(x_i). \quad (48)$$

Considering now the same but oriented GM (arrows on figure above), we get

$$m_{i \rightarrow j}(x_j) \propto \sum_{x_i} \phi_i(x_i) \cdot \phi_{i,j}(x_i, x_j) \cdot \prod_{k \in N(i) \setminus j} m_{k \rightarrow i}(x_i), \quad (49)$$

where  $N(i)$  is the neighbourhood of  $x_i$ .

The MARG problem can then be solved using the **sum-product** procedure.

---

### Sum-product BP

- $t = 0$ ,

$$\forall (i, j) \in E, \forall (x_i, x_j) \in \mathcal{X}^2 : m_{i \rightarrow j}^0(x_j) = m_{j \rightarrow i}^0(x_i) = 1. \quad (50)$$

- $t > 0$ ,

$$m_{i \rightarrow j}^{t+1}(x_j) \propto \sum_{x_i} \phi_i(x_i) \cdot \phi_{ij}(x_i, x_j) \cdot \prod_{k \in N(i) \setminus j} m_{k \rightarrow i}^t(x_i), \quad (51)$$

$$P_{X_i}^{t+1}(x_i) = \prod_{k \in N(i)} m_{k \rightarrow i}^{t+1}(x_i). \quad (52)$$


---

While, for the MAP problem, the **max-sum** procedure is considered.

---

### Max-sum BP

- $t = 0$ ,

$$m_{i \rightarrow j}^0(x_j) = m_{j \rightarrow i}^0(x_i) = 1. \quad (53)$$

- $t > 0$ ,

$$m_{i \rightarrow j}^{t+1}(x_j) \propto \max_{x_i} \phi_i(x_i) \cdot \phi_{ij}(x_i, x_j) \cdot \prod_{k \in N(i) \setminus j} m_{k \rightarrow i}^t(x_i), \quad (54)$$

$$x_i^{t+1} = \arg \max_{x_i} \phi_i(x_i) \cdot \prod_{k \in N(i)} m_{k \rightarrow i}^{t+1}(x_i). \quad (55)$$


---

Note: here, we use only the potentials of pairs. But in case of cliques, we have to consider the JCT and iterate on it. To understand this point, let us apply the sum-product algorithm on factor graphs.

#### 2.3.1 Factor graphs

Considering the general notations in Fig. 15, the sum-product BP algorithm is particularized such that

$$m_{i \rightarrow f}^{t+1}(x_i) = \prod_{f' \in N(i) \setminus f} m_{f' \rightarrow i}^t(x_i), \quad (56)$$

$$m_{f \rightarrow i}^{t+1}(x_i) = \sum_{x_j, j \in N(f) \setminus i} f(x_i, x_j) \prod_{j \in N(f) \setminus i} m_{j \rightarrow f}^t(x_j). \quad (57)$$

On a tree, the leaves are sending the right messages at time 1 already, and



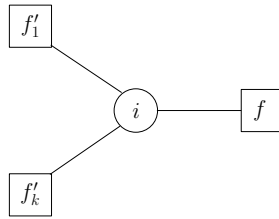


Figure 15: A simple factor graph.

after a number of time steps proportional to the tree diameter<sup>1</sup>, all messages are correct: the steady point is reached and the algorithm is exact. Therefore, **BP is exact on trees**. The JCT property discussed before is therefore useful, and can in certain cases allow us to construct graphs on which we know that BP is exact. However, the problem mentioned before remains: if BP is run on the JCT of a GM, subsequent maximizations/marginalizations will be necessary to recover the solution in terms of the initial problem's variables.

### 3 Understanding Belief Propagation

We have seen how to use the (exact) elimination algorithm in order to design the BP algorithms max-product and sum-product, that are exact only on trees. The JCT property has taught us how to group variables of an initial loopy GM such that the resulting GM is a tree (when it is possible), on which we can then run BP with a guarantee of an exact result. However, the subsequent operations that are necessary to obtain the solution in terms of the initial problem's variables can be a new source of intractability. Therefore, we would like to know what happens if we use BP on the initial (loopy) graph anyway. The advantage is that BP remains tractable because of the low number of operations per iteration. The danger is that BP is not exact anymore and therefore we need to ask ourselves the following 3 questions:

1. Does the algorithm have fixed points?
2. What are those fixed points?
3. Are they reached?

The analysis will be made with the sum-product BP algorithm, but could be carried out similarly for the max-product version.

#### 3.1 Existence of a fixed point

The algorithm is of the type

$$\underline{m}^{t+1} = F(\underline{m}^t) \quad \text{with} \quad \underline{m}^t \in [0, 1]^{2|\mathcal{E}||\mathcal{X}|} \quad (58)$$

and **the existence of a fixed point is guaranteed** by a theorem.

<sup>1</sup>The eccentricity of a vertex  $v$  in a graph is the maximum distance from  $v$  to any other vertex. The diameter of a graph is the maximum eccentricity over all vertices in a graph.

### 3.2 Nature of the fixed points

Let us remind that we had factorized  $P_{\underline{X}}(\underline{x})$  in this way:

$$\begin{aligned} P_{\underline{X}}(\underline{x}) &\propto \prod_{i \in \mathcal{V}} \phi_i(x_i) \prod_{(i,j) \in \mathcal{E}} \psi_{ij}(x_i, x_j) \\ &= \frac{1}{Z} e^{Q(\underline{x})}. \end{aligned} \quad (59)$$

The fixed points are a solution of the following problem

$$P_{\underline{X}} \in \arg \max_{\mu \in M(\mathcal{X}^{\mathcal{N}})} \mathbb{E}[Q(X)] + H(\mu) \quad (60)$$

with

$$\mathbb{E}_{\mu}[Q(X)] + H(\mu) = \sum_{\underline{x}} \mu(\underline{x}) Q(\underline{x}) - \sum_{\underline{x}} \mu(\underline{x}) \log \mu(\underline{x}) = F(\mu). \quad (61)$$

Let us find a bound for this quantity. From (59), we get  $Q(\underline{x}) = \log P_{\underline{X}}(\underline{x}) + \log Z$ . Then

$$\begin{aligned} F(\mu) &= \left( \sum_{\underline{x}} \mu(\underline{x}) \log Z \right) + \left( \sum_{\underline{x}} \mu(\underline{x}) \log \frac{P_{\underline{X}}(\underline{x})}{\mu(\underline{x})} \right) \\ &= \log Z + \mathbb{E}_{\mu} \left[ \log \frac{P_{\underline{X}}}{\mu(\underline{x})} \right] \\ &\leq \log Z + \log \mathbb{E}_{\mu} \left[ \frac{P_{\underline{X}}}{\mu} \right] \quad \text{using Jensen's inequality} \\ &\leq \log Z \end{aligned} \quad (62)$$

and the equality is reached when the distributions  $\mu$  and  $P$  are equal.

This maximization in equation (60) is made over the space of all possible distributions, which is a far too big search space. But if we restrict ourselves to trees, we know that  $\mu$  has the form:

$$\mu \propto \prod_i \mu_i \prod_{(i,j)} \frac{\mu_{ij}}{\mu_i \mu_j} \quad (63)$$

BP has taught us that:

$$\mu_i \propto \phi_i \prod_{k \in N(i)} m_{k \rightarrow i} \quad (64)$$

$$\mu_{ij} \propto \prod_{k \in N(i) \setminus j} m_{k \rightarrow i} \phi_i \psi_{ij} \phi_j \prod_{l \in N(j) \setminus i} m_{l \rightarrow j} \quad (65)$$

If we marginalize  $\mu_{ij}$  with respect to  $x_j$ , we should obtain  $\mu_i$ :  $\sum_{x_j} \mu_{ij}(x_i, x_j) = \mu_i(x_i)$ . Writing this out, we obtain:

$$\prod_{k \in N(i) \setminus j} m_{k \rightarrow i} \phi_i \left( \sum_{x_j} \psi_{ij} \phi_j \prod_{l \in N(i) \setminus j} m_{l \rightarrow j} \right) = \phi_i \prod_{k \in N(i)} m_{k \rightarrow i} \quad (66)$$

and this should lead us to what we believe from the fixed points of BP. Let us make a recharacterization in terms of the fixed points. In order to lighten notations, we will write  $\phi$  instead of  $\log \phi$  and  $\psi$  instead of  $\log \psi$ :

$$F_{\text{Bethe}}(\mu) = \mathbb{E}_{\mu} \left[ \sum_i \phi_i + \sum_{i,j} \psi_{ij} \right] - \mathbb{E}_{\mu} [\log \mu] \quad (67)$$

We now use following factorization

$$\mathbb{E}_{\mu} [\log \mu] = - \sum_i \mathbb{E}_{\mu_i} [\log \mu_i] - \sum_{ij} \left( \mathbb{E}_{\mu_{ij}} [\log \mu_{ij}] - \mathbb{E}_{\mu_i} [\log \mu_i] - \mathbb{E}_{\mu_j} [\log \mu_j] \right) \quad (68)$$

and obtain a new expression for the Bethe free energy

$$F_{\text{Bethe}} = \sum_i (1-d_i) \left( H_{\mu_i} + \mathbb{E}_{\mu_i} [\phi_i] \right) + \sum_{ij} \left( H(\mu_{ij}) + \mathbb{E}_{\mu_{ij}} [\psi_{ij} + \phi_i + \phi_j] \right), \quad (69)$$

where  $d_i$  is the degree of node  $i$ .

### 3.2.1 Background on Nonlinear Optimization

The problem

$$\max_q G(q) \quad \text{s.t.} \quad Aq = b \quad (70)$$

can be expressed in a different form by using Lagrange multipliers  $\lambda$

$$L(q, \lambda) = G(q) + \lambda^T (Aq - b) \quad (71)$$

and maximizing

$$\begin{aligned} \max_q L(q, \lambda) &= M(\lambda) \leq G(q^*) \\ \inf_{\lambda} M(\lambda) &\leq G(q^*). \end{aligned}$$

Let us look at all  $\lambda$  such that  $\nabla_q L(q) = 0$ . In a sense, BP is finding stationary points of this Lagrangian.

### 3.2.2 Belief Propagation as a variational problem

In our case, here are the conditions we will enforce with Lagrange multipliers:

$$\mu_{ij}(x_i, x_j) \geq 0 \quad (72)$$

$$\sum_{x_i} \mu_i(x_i) = 1 \quad \rightarrow \lambda_i \quad (73)$$

$$\sum_{x_j} \mu_{ij}(x_i, x_j) = \mu_i(x_i) \quad \rightarrow \lambda_{j \rightarrow i}(x_i) \quad (74)$$

$$\sum_{x_i} \mu_{ij}(x_i, x_j) = \mu_j(x_j) \quad \rightarrow \lambda_{i \rightarrow j}(x_j) \quad (75)$$

The complete Lagrangian reads

$$\begin{aligned} \mathcal{L} = & F_{\text{Bethe}}(\mu) + \sum_i \lambda_i \left( \sum_{x_i} \mu_i(x_i) - 1 \right) \\ & + \sum_{ij} \left[ \left( \sum_{x_j} \mu_{ij}(x_i, x_j) - \mu_i(x_i) \right) \lambda_{j \rightarrow i}(x_i) \right. \\ & \left. + \left( \sum_{x_i} \mu_{ij}(x_i, x_j) - \mu_j(x_j) \right) \lambda_{i \rightarrow j}(x_j) \right]. \end{aligned} \quad (76)$$

We need to minimize this Lagrangian with respect to all possible variables, which we obtain by setting the partial derivatives to zero:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mu_i(x_i)} = 0 \quad (77) \\ = -(1 - d_i)(1 + \log \mu_i(x_i)) + (1 - d_i)\phi_i(x_i) + \lambda_i - \sum_{j \in N(i)} \lambda_{j \rightarrow i}(x_i) \end{aligned}$$

which imposes following equality for the distribution  $\mu_i$ :

$$\boxed{\mu_i(x_i) \propto e^{\phi_i(x_i) + \frac{1}{d_i - 1} \sum_{j \in N(i)} \lambda_{j \rightarrow i}(x_i)}} \quad (78)$$

Let us now use the transformation  $\lambda_{j \rightarrow i}(x_i) = \sum_{k \in N(i) \setminus j} \log m_{k \rightarrow i}(x_i)$ , and we obtain

$$\sum_{j \in N(i)} \lambda_{j \rightarrow i}(x_i) \equiv (d_i - 1) \sum_{j \in N(i)} \log m_{j \rightarrow i}(x_i). \quad (79)$$

In the same way, we can show that:

$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}(x_i, x_j)} = 0 \Rightarrow \boxed{\mu_{ij}(x_i, x_j) \propto e^{\phi_i(x_i) + \phi_j(x_j) + \psi_{ij}(x_i, x_j) + \lambda_{j \rightarrow i}(x_i) + \lambda_{i \rightarrow j}(x_j)}}$$

This way, we found the distributions  $\mu_i$  and  $\mu_{ij}$  that are the fixed points of BP.

### 3.3 Can the fixed points be reached?

We will now try to analyze if the algorithm can actually reach those fixed points that we have exhibited in the previous section. Let us look at the simple (but loopy) graph in Fig. 16. At time  $t = 1$ , we have

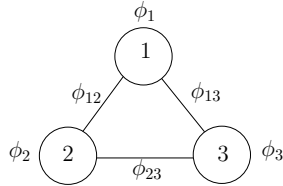


Figure 16: A simple loopy graph.

$$m_{2 \rightarrow 1}^1(x_1) \propto \sum_{x_2} \phi_2(x_2) \phi_{12}(x_1, x_2) \underbrace{m_{3 \rightarrow 2}^0(x_2)}_{=1} \quad (80)$$

and

$$m_{3 \rightarrow 1}^1 \propto \sum_{x_3} \phi_3 \phi_{13} \quad (81)$$

which also corresponds to the messages of the modified graph in Fig. 17.

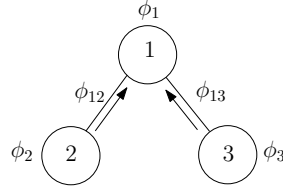


Figure 17: Graph seen by BP at time  $t = 1$ .

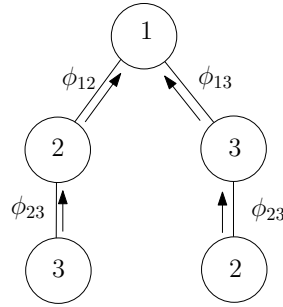


Figure 18: Graph seen by BP at time  $t = 2$ .

At time  $t = 2$ , the messages will be as

$$m_{2 \rightarrow 1}^2 \propto \sum_{x_2} \phi_2 \phi_{12} m_{3 \rightarrow 2}^1(x_2) \quad (82)$$

corresponding to the messages on the modified graph in Fig. 18. If we increase  $t$ , the corresponding non-loopy graph gets longer at each time step.

Another way of seeing this is by looking at the recursion equation:

$$F_{ij}(m^*) = m_{ij}^* \quad (83)$$

$$m_{ij}^{t+1} = F_{ij}(m^t)$$

$$\begin{aligned} |m_{ij}^{t+1} - m_{ij}^*| &= |F_{ij}(m^t) - F_{ij}(m^*)| \\ &= |\nabla F_{ij}(\theta)^T (m^t - m^*)| \quad (\text{mean value theorem}) \end{aligned}$$

$$|m^{t+1} - m^*|_\infty \leq |\nabla F_{ij}(\theta)|_1 |m^t - m^*|_\infty \quad (84)$$

From this last inequality, it is clear that if we can prove that  $|F_{ij}|_1$  is bounded by some constant  $\rho < 1$ , the convergence is proved. Unfortunately, it is not often easy to prove such a thing.

### 3.3.1 The hardcore model

In the hardcore model, we have

$$\phi_i(x_i) = 1 \quad \text{for all } x_i \in \{0, 1\} \quad (85)$$

$$\psi_{ij}(x_i, x_j) = 1 - x_i x_j. \quad (86)$$

Instead of using BP, let us do the following gradient-descent like algorithm:

$$y(t+1) = \left[ y(t) + \alpha(t) \frac{\partial F}{\partial y_i} \Big|_{y(t)} \right] \quad (87)$$

where the operator  $[\cdot]$  is a clipping function that ensures that the result stays in the interval  $(0, 1)$ . This is a projected version of a gradient algorithm with variable step size  $\alpha(t)$ . Choosing this step size with following rule:

$$\alpha(t) = \frac{1}{\sqrt{t}} \frac{1}{2^d} \quad (88)$$

then we can show that in a time  $T \sim n^2 2^d \frac{1}{\epsilon^4}$  we will find  $F_b$  up to  $\epsilon$ , and convergence is proved.

## 4 Learning Graphical Models

In this final section, we focus on the learning problem. In particular, we consider three different cases:

- **Parameter learning**  
Given a graph, the parameters are learned from the observation of the entire set of realizations of all random variables.
- **Graphical model learning**  
Both the parameters and the graph are learned from the observations of the entire set of realizations of all random variables.
- **Latent graphical model learning**  
The parameters and the graph are learned from partial observations: some of the random variables are assumed to be hidden.

### 4.1 Parameter learning

#### 4.1.1 Single parameter learning

We consider the following simple setting where  $x_i$  is a Bernoulli random variable with parameter  $\theta$ :

$$P_X(x_i, \theta) = \begin{cases} \theta & \text{if } x_i = 1, \\ 1 - \theta & \text{if } x_i = 0. \end{cases} \quad (89)$$

Given observations  $\{x_1, \dots, x_S\}$ , we are interested in the MAP estimation of the parameter  $\theta$ :

$$\begin{aligned} \hat{\theta}^{MAP} &= \arg \max_{\theta \in [0,1]} P(\theta | x_1, \dots, x_S), \\ &= \arg \max_{\theta \in [0,1]} P(x_1, \dots, x_S | \theta) p(\theta), \end{aligned} \quad (90)$$

where maximizing  $P(x_1, \dots, x_S | \theta)$  leads to the maximum likelihood (ML) estimator  $\hat{\theta}^{ML}$  of  $\theta$ .

Denoting  $\mathcal{D} \triangleq \{x_1, \dots, x_S\}$  the observed set of realizations, we define the empirical likelihood as follows:

$$\begin{aligned} \ell(\mathcal{D}; \theta) &= \frac{1}{S} \log P(x_1, \dots, x_S | \theta), \\ &= \frac{1}{S} \sum_i \log P(x_i | \theta), \\ &= \hat{P}(1) \log \theta + \hat{P}(0) \log(1 - \theta), \end{aligned} \quad (91)$$

with  $\hat{P}(1) = \frac{1}{S} \sum_i \mathbb{1}_{\{x_i=1\}}$ . Derivating (91) and setting the result to zero, we obtain the **maximal likelihood estimator**  $\hat{\theta}^{ML}$ :

$$\begin{aligned} \frac{\partial}{\partial \theta} \ell(\mathcal{D}; \theta) &= \frac{\hat{P}(1)}{\theta} - \frac{\hat{P}(0)}{1 - \theta} = 0, \\ \Rightarrow \quad &\boxed{\hat{\theta}^{ML} = \hat{P}(1)} \end{aligned} \quad (92)$$

What is the amount of samples  $S$  needed to achieve  $\hat{\theta}^{ML}(S) \approx (1 \pm \epsilon)\theta$ ? Considering the binomial variable  $B(S, \theta)$  (which is the sum of  $S$  independently drawn Bernoulli variables from (89)), we can write

$$\begin{aligned} P(|B(S, \theta) - S\theta| > \epsilon S\theta) &\sim \exp(-\epsilon^2 S\theta) \leq \delta, \\ \Rightarrow \quad &\boxed{S \geq \frac{1}{\theta} \frac{1}{\epsilon^2} \log \frac{1}{\delta}} \end{aligned} \quad (93)$$

#### 4.1.2 Directed graphs

We consider the following setting in which we have not one, but many random variables to learn on a directed graph:

$$P_{\underline{X}}(\underline{x}) \propto \prod_i P_{X_i | X_{\Pi_i}}(x_i | x_{\Pi_i}), \quad (94)$$

where  $\Pi_i$  stands for the parents of node  $i$ , and  $P_{X_i | X_{\Pi_i}}(x_i | x_{\Pi_i}) \triangleq \theta_{x_i, x_{\Pi_i}}$ . Again, we look at the empirical likelihood

$$\begin{aligned} \ell(\mathcal{D}; \underline{\theta}) &= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i, x_{\Pi_i}) \log \theta_{x_i, x_{\Pi_i}}, \\ &= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i | x_{\Pi_i}) \hat{P}(x_{\Pi_i}) \left[ \log \frac{\theta_{x_i, x_{\Pi_i}}}{\hat{P}(x_i | x_{\Pi_i})} + \log \hat{P}(x_i | x_{\Pi_i}) \right], \\ &= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i | x_{\Pi_i}) \hat{P}(x_{\Pi_i}) \log \frac{\theta_{x_i, x_{\Pi_i}}}{\hat{P}(x_i | x_{\Pi_i})}, \end{aligned} \quad (95)$$

and set the derivative to zero in order to obtain the ML estimation of  $\underline{\theta}$ , resulting in

$$\begin{aligned} \sum_{x_i} \hat{P}(x_i | x_{\Pi_i}) \log \frac{\theta_{x_i, x_{\Pi_i}}}{\hat{P}(x_i | x_{\Pi_i})} &= \mathbb{E}_{\hat{P}} \left[ \log \frac{\theta_{x_i, x_{\Pi_i}}}{\hat{P}(x_i | x_{\Pi_i})} \right], \\ \Rightarrow \quad &\boxed{\hat{\theta}_{x_i, x_{\Pi_i}}^{ML} = \hat{P}(x_i | x_{\Pi_i})} \end{aligned} \quad (96)$$

### 4.1.3 Undirected graphs

Let us now consider the case of undirected graphs. To reduce the amount of indices, we will write  $i$  instead of  $x_i$  in the following.

On a tree, 
$$P_X = \prod_i P_i \prod_{ij} \frac{P_{ij}}{P_i P_j} \rightarrow \text{possible estimator: } \hat{P}_i \frac{\hat{P}_{ij}}{\hat{P}_i \hat{P}_j}$$

on a chordal graph, 
$$P_X \propto \frac{\prod_C \phi_C(x_C)}{\prod_S \phi_S(x_S)} \rightarrow \text{possible estimator: } \frac{\hat{P}_C}{\hat{P}_S}$$

on a triangle-free graph, 
$$P_X \propto \prod_i \phi_i \prod_{ij} \psi_{ij}$$

For the last case, let us use the Hammersley-Clifford theorem. Let  $\mathcal{X} = \{0, 1\}$ . On a triangle-free graph, the maximal clique size is 2, and therefore we can write

$$P_{\underline{X}}(\underline{x}) \propto \exp \left( \sum_i U_i(x_i) + \sum_{ij} V_{ij}(x_i, x_j) \right). \quad (97)$$

Using the fact that we have a MRF, we get

$$\frac{P(X_i = 1, X_{rest} = 0)}{P(X_i = 0, X_{rest} = 0)} \propto \exp(Q(i)). \quad (98)$$

Also, because of the fact that on a MRF, a variable conditioned on its neighbours is independent of all the others, we can write

$$\frac{P(X_i = 1, X_{rest} = 0)}{P(X_i = 0, X_{rest} = 0)} = \frac{P(X_i = 1, X_{N(i)} = 0)}{P(X_i = 0, X_{N(i)} = 0)} \quad (99)$$

and therefore this quantity can be calculated with  $2^{|N(i)|+1}$  operations.

## 4.2 Graphical model learning

What can we learn from a set of realizations of variables when the underlying graph is not known? We focus now in the following maximisation

$$\max_{\mathcal{G}, \theta_{\mathcal{G}}} \ell(\mathcal{D}; \mathcal{G}, \theta_{\mathcal{G}}) = \max_{\mathcal{G}} \underbrace{\max_{\theta_{\mathcal{G}}} \ell(\mathcal{D}; \mathcal{G}, \theta_{\mathcal{G}})}_{\hat{\ell}(\mathcal{D}; \mathcal{G}) \triangleq \ell(\mathcal{D}; \mathcal{G}, \hat{\theta}_{\mathcal{G}}^{ML})}. \quad (100)$$

From the previous subsection, we have  $\hat{\theta}_{\mathcal{G}}^{ML}$ , and therefore we only need to find a way to evaluate the maximization on the possible graphs.



### 4.2.1 Directed graphs

On a directed graph  $\mathcal{G} \rightarrow (i, \Pi_i)$ , the empirical likelihood reads

$$\begin{aligned}
\hat{\ell}(\mathcal{D}; \mathcal{G}) &= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i, x_{\Pi_i}) \log \hat{P}(x_i | x_{\Pi_i}), \\
&= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i, x_{\Pi_i}) \log \left[ \frac{\hat{P}(x_i, x_{\Pi_i})}{\hat{P}(x_i) \hat{P}(x_{\Pi_i})} \hat{P}(x_i) \right], \\
&= \sum_i \sum_{x_i, x_{\Pi_i}} \hat{P}(x_i, x_{\Pi_i}) \log \frac{\hat{P}(x_i, x_{\Pi_i})}{\hat{P}(x_i) \hat{P}(x_{\Pi_i})} + \sum_{x_i} \hat{P}(x_i) \log \hat{P}(x_i), \\
&= \sum_i I(\hat{X}_i; \hat{X}_{\Pi_i}) - H(\hat{X}_i). \tag{101}
\end{aligned}$$

Looking for the graph maximizing the empirical likelihood thus consists in maximising the mutual information:  $\max_{\mathcal{G}} \sum_i I(\hat{X}_i; \hat{X}_{\Pi_i})$ . In a general setting, this is not easy. **Reducing the search space to trees** however, some methods exist, like the Chow-Liu algorithm [1], which relies on the procedure used to get the maximum weighted spanning tree (cf. section 2).

### 4.2.2 Undirected graphs

What can we do in the case of undirected graphs? Let us restrict ourselves to the binary case  $\underline{x} \in \{0, 1\}^N$  and to exponential families:

$$P_{\underline{X}}(\underline{x}) = \exp \left( \sum_i \theta_i x_i + \sum_{i,j} \theta_{ij} x_i x_j - \log Z(\underline{\theta}) \right). \tag{102}$$

Again, we denote  $\mathcal{D} = \{\underline{x}^1, \dots, \underline{x}^S\}$  the observed dataset, and the log-likelihood can be written as

$$\ell(\mathcal{D}; \underline{\theta}) = \underbrace{\sum_i \theta_i \mu_i + \sum_{i,j} \theta_{ij} \mu_{ij}}_{\langle \underline{\theta}, \underline{\mu} \rangle} - \log Z(\underline{\theta}). \tag{103}$$

As  $\ell(\mathcal{D}; \underline{\theta})$  is a concave function of  $\underline{\theta}$ , it can be efficiently solved using a gradient descent algorithm of the form

$$\underline{\theta}^{t+1} = \underline{\theta}^t + \alpha(t) \nabla_{\underline{\theta}} \ell(\mathcal{D}; \underline{\theta})|_{\underline{\theta}=\underline{\theta}^t} \tag{104}$$

The difficulty in this formula is the evaluation of the gradient:

$$\nabla_{\underline{\theta}} \ell(\mathcal{D}; \underline{\theta}) = \underline{\mu} - \frac{\mathbb{E}(\underline{X})}{\underline{\theta}}, \tag{105}$$

whose second term is an expectation that has to be calculated, using the sum-product algorithm or with a Markov chain Monte Carlo method for instance.

Another question is whether we will be learning interesting graphs at all. Graph-learning algorithms tend to link variables that are not linked in the real underlying graph. To avoid this, complicated graphs should be penalized by introducing a regularizer. Unfortunately, this is a highly non-trivial problem, and graphical model learning algorithms do not always perform well to this day.

### 4.3 Latent Graphical Model learning: the Expectation-maximization algorithm

In this last case, we distinguish two different variables:

- $Y$  stands for observed variables,
- $X$  denotes the hidden variables.

The parameter  $\theta$  is estimated from the observations, namely

$$\hat{\theta}^{ML} = \arg \max_{\theta} \log P_Y(y; \theta). \quad (106)$$

The log-likelihood is derived by marginalizing on the hidden variables

$$\begin{aligned} \ell(y; \theta) &= \log P_Y(y; \theta), \\ &= \log \sum_x P_{X,Y}(x, y; \theta), \end{aligned} \quad (107)$$

$$= \log \sum_x q(x|y) \frac{P_{X,Y}(x, y; \theta)}{q(x|y)}, \quad (108)$$

$$= \log \mathbb{E}_q \left[ \frac{P}{q} \right] \geq \mathbb{E}_q \left[ \frac{P}{q} \right] \triangleq \mathcal{L}(q; \theta). \quad (109)$$

This gives rise to the Expectation-Maximisation (EM) algorithm [2].

---

#### EM algorithm

Until convergence, iterate between

- **E-step:** estimation of the distribution  $q$   
 $\theta^t \rightarrow q^{t+1} = \arg \max_q \mathcal{L}(q; \theta^t)$ .
  - **M-step:** estimation of the parameter  $\theta$   
 $q^{t+1} \rightarrow \theta^{t+1} = \arg \max_{\theta} \mathcal{L}(q^{t+1}; \theta)$ .
-

## References

- [1] C. K. Chow and C. N. Liu. Approximating discrete probability distributions with dependence trees. *Information Theory, IEEE Transactions on*, 14(3):462–467, 1968.
- [2] A. P. Dempster, N. M. Laird, and D. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38, 1977.
- [3] G. R. Grimmet. A theorem about random fields. *Bulletin of the London Mathematical Society*, 5(1):81–84, 1973.
- [4] J. M. Hammersley and P. Clifford. Markov fields on finite graphs and lattices. Available online: <http://www.statslab.cam.ac.uk/~grg/books/hammfest/hamm-cliff.pdf>, 1971.