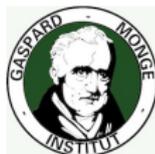


Inference and the sum-product algorithm



Guillaume Obozinski

Ecole des Ponts - ParisTech



INIT/AERFAI Summer school on Machine Learning
Benicàssim, June 26th 2017

Probabilistic inference

Given a discrete Gibbs model of the form:

$$p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where \mathcal{C} is the set of cliques of the graph, inference is the problem of computation of several related quantities of interest:

- Computation of the marginal $p(x_i)$ or more generally, $p(x_C)$.

Probabilistic inference

Given a discrete Gibbs model of the form:

$$p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where \mathcal{C} is the set of cliques of the graph, inference is the problem of computation of several related quantities of interest:

- Computation of the marginal $p(x_i)$ or more generally, $p(x_C)$.
- Computation of the partition function Z

Probabilistic inference

Given a discrete Gibbs model of the form:

$$p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where \mathcal{C} is the set of cliques of the graph, inference is the problem of computation of several related quantities of interest:

- Computation of the marginal $p(x_i)$ or more generally, $p(x_C)$.
- Computation of the partition function Z
- Computation of the conditional marginal $p(x_i | X_j = x_j, X_k = x_k)$

Probabilistic inference

Given a discrete Gibbs model of the form:

$$p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C),$$

where \mathcal{C} is the set of cliques of the graph, inference is the problem of computation of several related quantities of interest:

- Computation of the marginal $p(x_i)$ or more generally, $p(x_C)$.
- Computation of the partition function Z
- Computation of the conditional marginal $p(x_i | X_j = x_j, X_k = x_k)$

Inference is actually necessary

- For the computation of the gradient of the likelihood of a model
- Computation of the expected value of the log-likelihood of an exponential family at step E of the EM algorithm (for example for the HMM)

Hardness of inference

 Problem: In general, the inference problem is NP-hard.

Hardness of inference

⚠ Problem: In general, the inference problem is NP-hard.

For trees

the inference problem is efficient: its cost is linear in n .

Hardness of inference

⚠ Problem: In general, the inference problem is NP-hard.

For trees

the inference problem is efficient: its cost is linear in n .

For “tree-like” graphs

we use the *Junction Tree Algorithm* which enables us to bring the situation back to that of a tree.

Hardness of inference

⚠ Problem: In general, the inference problem is NP-hard.

For trees

the inference problem is efficient: its cost is linear in n .

For “tree-like” graphs

we use the *Junction Tree Algorithm* which enables us to bring the situation back to that of a tree.

In the general case

There are several methods for approximate/inexact inference.

Hardness of inference

⚠ Problem: In general, the inference problem is NP-hard.

For trees

the inference problem is efficient: its cost is linear in n .

For “tree-like” graphs

we use the *Junction Tree Algorithm* which enables us to bring the situation back to that of a tree.

In the general case

There are several methods for approximate/inexact inference.

Rest of this part of the lecture

We will focus in the rest of this part to inference for undirected graphical models that have cliques of size 1 and 2.

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{i,j}(x_i, x_j)$$

Rest of this part of the lecture

We will focus in the rest of this part to inference for undirected graphical models that have cliques of size 1 and 2.

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{i,j}(x_i, x_j)$$

- The general case is not much more difficult but requires the concept of *factor graph*.

Rest of this part of the lecture

We will focus in the rest of this part to inference for undirected graphical models that have cliques of size 1 and 2.

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{i,j}(x_i, x_j)$$

- The general case is not much more difficult but requires the concept of *factor graph*.
- To perform inference on a directed graph, we will simply moralize the graph and apply the inference algorithm for undirected graphs

Rest of this part of the lecture

We will focus in the rest of this part to inference for undirected graphical models that have cliques of size 1 and 2.

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{i,j}(x_i, x_j)$$

- The general case is not much more difficult but requires the concept of *factor graph*.
- To perform inference on a directed graph, we will simply moralize the graph and apply the inference algorithm for undirected graphs
- We will focus on the cases where inference can be done efficiently using *dynamic programming*, that is:
 - in the chain case
 - in the tree case

Inference on a chain

We define X_i a random variable, taking value in $\{1, \dots, K\}$, $i \in V = \{1, \dots, n\}$ with joint distribution

$$p(x) = \frac{1}{Z} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i)$$

Inference on a chain

We define X_i a random variable, taking value in $\{1, \dots, K\}$, $i \in V = \{1, \dots, n\}$ with joint distribution

$$p(x) = \frac{1}{Z} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i)$$

Computing

$$p(x_j) = \sum_{x_{V \setminus \{j\}}} p(x_1, \dots, x_n)$$

has a priori a complexity of $O(K^n)$.

Inference on a chain: a dynamic program

$$p(x_j)$$

Inference on a chain: a dynamic program

$$p(x_j) = \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i)$$

Inference on a chain: a dynamic program

$$\begin{aligned} & p(x_j) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n) \end{aligned}$$

Inference on a chain: a dynamic program

$$\begin{aligned} & p(x_j) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \sum_{x_n} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n) \end{aligned}$$

Inference on a chain: a dynamic program

$$\begin{aligned} & p(x_j) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^n \psi_i(x_i) \prod_{i=2}^n \psi_{i-1,i}(x_{i-1}, x_i) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \sum_{x_n} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n) \\ = & \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \underbrace{\sum_{x_n} \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n)}_{\mu_{n \rightarrow n-1}(x_{n-1})} \end{aligned}$$

Inference on a chain: a dynamic program

$$= \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \underbrace{\sum_{x_n} \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n)}_{\mu_{n \rightarrow n-1}(x_{n-1})}$$

Inference on a chain: a dynamic program

$$\begin{aligned} &= \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \underbrace{\sum_{x_n} \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n)}_{\mu_{n \rightarrow n-1}(x_{n-1})} \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{j,n,n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1,i}(x_{i-1}, x_i) \times \\ &\quad \times \underbrace{\sum_{x_{n-1}} \psi_{n-1}(x_{n-1}) \psi_{n-2,n-1}(x_{n-2}, x_{n-1}) \mu_{n \rightarrow n-1}(x_{n-1})}_{\mu_{n-1 \rightarrow n-2}(x_{n-2})} \end{aligned}$$

Inference on a chain: a dynamic program

$$\begin{aligned} &= \frac{1}{Z} \sum_{x_{V \setminus \{j,n\}}} \prod_{i=1}^{n-1} \psi_i(x_i) \prod_{i=2}^{n-1} \psi_{i-1,i}(x_{i-1}, x_i) \underbrace{\sum_{x_n} \psi_n(x_n) \psi_{n-1,n}(x_{n-1}, x_n)}_{\mu_{n \rightarrow n-1}(x_{n-1})} \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{j,n,n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1,i}(x_{i-1}, x_i) \times \\ &\quad \times \underbrace{\sum_{x_{n-1}} \psi_{n-1}(x_{n-1}) \psi_{n-2,n-1}(x_{n-2}, x_{n-1}) \mu_{n \rightarrow n-1}(x_{n-1})}_{\mu_{n-1 \rightarrow n-2}(x_{n-2})} \end{aligned}$$

Descending messages:

$$\mu_{i \rightarrow i-1}(x_{i-1}) = \sum_{x_i} \psi_{i-1,i}(x_{i-1}, x_i) \psi_i(x_i) \mu_{i+1 \rightarrow i}(x_i)$$

Inference on a chain: ascending messages

Now marginalizing from the left, $p(x_j) =$

$$= \frac{1}{Z} \sum_{x_{V \setminus \{j, n, n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2})$$

Inference on a chain: ascending messages

Now marginalizing from the left, $p(x_j) =$

$$\begin{aligned} &= \frac{1}{Z} \sum_{x_{V \setminus \{j, n, n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2}) \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{1, j, n, n-1\}}} \underbrace{\sum_{x_1} \psi_1(x_1) \psi_{1,2}(x_1, x_2)}_{\mu_{1 \rightarrow 2}(x_2)} \times \prod_{i=2}^{n-2} \psi_i(x_i) \prod_{i=3}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \times \\ &\hspace{20em} \times \mu_{n-1 \rightarrow n-2}(x_{n-2}) \end{aligned}$$

Inference on a chain: ascending messages

Now marginalizing from the left, $p(x_j) =$

$$\begin{aligned} &= \frac{1}{Z} \sum_{x_{V \setminus \{j, n, n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2}) \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{1, j, n, n-1\}}} \underbrace{\sum_{x_1} \psi_1(x_1) \psi_{1,2}(x_1, x_2)}_{\mu_{1 \rightarrow 2}(x_2)} \times \prod_{i=2}^{n-2} \psi_i(x_i) \prod_{i=3}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \times \\ & \hspace{15em} \times \mu_{n-1 \rightarrow n-2}(x_{n-2}) \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{1, j, n, n-1\}}} \mu_{1 \rightarrow 2}(x_2) \prod_{i=2}^{n-2} \psi_i(x_i) \prod_{i=3}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2}) \end{aligned}$$

Inference on a chain: ascending messages

Now marginalizing from the left, $p(x_j) =$

$$\begin{aligned} &= \frac{1}{Z} \sum_{x_{V \setminus \{j, n, n-1\}}} \prod_{i=1}^{n-2} \psi_i(x_i) \prod_{i=2}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2}) \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{1, j, n, n-1\}}} \underbrace{\sum_{x_1} \psi_1(x_1) \psi_{1,2}(x_1, x_2)}_{\mu_{1 \rightarrow 2}(x_2)} \times \prod_{i=2}^{n-2} \psi_i(x_i) \prod_{i=3}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \times \\ & \hspace{20em} \times \mu_{n-1 \rightarrow n-2}(x_{n-2}) \\ &= \frac{1}{Z} \sum_{x_{V \setminus \{1, j, n, n-1\}}} \mu_{1 \rightarrow 2}(x_2) \prod_{i=2}^{n-2} \psi_i(x_i) \prod_{i=3}^{n-2} \psi_{i-1, i}(x_{i-1}, x_i) \mu_{n-1 \rightarrow n-2}(x_{n-2}) \end{aligned}$$

Ascending messages:

$$\mu_{i \rightarrow i+1}(x_{i+1}) = \sum_{x_i} \mu_{i-1 \rightarrow i}(x_i) \psi_i(x_i) \psi_{i, i+1}(x_i, x_{i+1})$$

Ascending and descending messages

$$\mu_{i \rightarrow i-1}(x_{i-1}) = \sum_{x_i} \psi_{i-1,i}(x_{i-1}, x_i) \psi_i(x_i) \mu_{i+1 \rightarrow i}(x_i)$$

$$\mu_{i \rightarrow i+1}(x_{i+1}) = \sum_{x_i} \mu_{i-1 \rightarrow i}(x_i) \psi_i(x_i) \psi_{i,i+1}(x_i, x_{i+1})$$

Ascending and descending messages

$$\mu_{i \rightarrow i-1}(x_{i-1}) = \sum_{x_i} \psi_{i-1,i}(x_{i-1}, x_i) \psi_i(x_i) \mu_{i+1 \rightarrow i}(x_i)$$

$$\mu_{i \rightarrow i+1}(x_{i+1}) = \sum_{x_i} \mu_{i-1 \rightarrow i}(x_i) \psi_i(x_i) \psi_{i,i+1}(x_i, x_{i+1})$$

Note that node i always receives a message which is a function of x_i .

Ascending and descending messages

$$\mu_{i \rightarrow i-1}(\mathbf{x}_{i-1}) = \sum_{\mathbf{x}_i} \psi_{i-1,i}(\mathbf{x}_{i-1}, \mathbf{x}_i) \psi_i(\mathbf{x}_i) \mu_{i+1 \rightarrow i}(\mathbf{x}_i)$$

$$\mu_{i \rightarrow i+1}(\mathbf{x}_{i+1}) = \sum_{\mathbf{x}_i} \mu_{i-1 \rightarrow i}(\mathbf{x}_i) \psi_i(\mathbf{x}_i) \psi_{i,i+1}(\mathbf{x}_i, \mathbf{x}_{i+1})$$

Note that node i always receives a message which is a function of \mathbf{x}_i .

After propagating all the messages, we have

$$p(\mathbf{x}_j) = \frac{1}{Z} \mu_{j-1 \rightarrow j}(\mathbf{x}_j) \psi_j(\mathbf{x}_j) \mu_{j+1 \rightarrow j}(\mathbf{x}_j).$$

Ascending and descending messages

$$\mu_{i \rightarrow i-1}(x_{i-1}) = \sum_{x_i} \psi_{i-1,i}(x_{i-1}, x_i) \psi_i(x_i) \mu_{i+1 \rightarrow i}(x_i)$$

$$\mu_{i \rightarrow i+1}(x_{i+1}) = \sum_{x_i} \mu_{i-1 \rightarrow i}(x_i) \psi_i(x_i) \psi_{i,i+1}(x_i, x_{i+1})$$

Note that node i always receives a message which is a function of x_i .

After propagating all the messages, we have

$$p(x_j) = \frac{1}{Z} \mu_{j-1 \rightarrow j}(x_j) \psi_j(x_j) \mu_{j+1 \rightarrow j}(x_j).$$

How do we compute Z ?

Ascending and descending messages

$$\mu_{i \rightarrow i-1}(x_{i-1}) = \sum_{x_i} \psi_{i-1,i}(x_{i-1}, x_i) \psi_i(x_i) \mu_{i+1 \rightarrow i}(x_i)$$

$$\mu_{i \rightarrow i+1}(x_{i+1}) = \sum_{x_i} \mu_{i-1 \rightarrow i}(x_i) \psi_i(x_i) \psi_{i,i+1}(x_i, x_{i+1})$$

Note that node i always receives a message which is a function of x_i .

After propagating all the messages, we have

$$p(x_j) = \frac{1}{Z} \mu_{j-1 \rightarrow j}(x_j) \psi_j(x_j) \mu_{j+1 \rightarrow j}(x_j).$$

How do we compute Z ?

$$1 = \sum_{x_j} p(x_j) \quad \Rightarrow \quad Z = \sum_{x_j} \mu_{j-1 \rightarrow j}(x_j) \psi_j(x_j) \mu_{j+1 \rightarrow j}(x_j)$$

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$
$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$
$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

How do we compute $p(x_{j'})$ for $j' \neq j$?

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$

$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

How do we compute $p(x_{j'})$ for $j' \neq j$? Can use the same messages !

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$

$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

How do we compute $p(x_{j'})$ for $j' \neq j$? Can use the same messages !

If we compute all $n - 1$ ascending messages

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \rightarrow \dots \rightarrow \boxed{\mu_{n-2 \rightarrow n-1}} \rightarrow \boxed{\mu_{n-1 \rightarrow n}}$$

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$

$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

How do we compute $p(x_{j'})$ for $j' \neq j$? Can use the same messages !

If we compute all $n - 1$ ascending messages

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \rightarrow \dots \rightarrow \boxed{\mu_{n-2 \rightarrow n-1}} \rightarrow \boxed{\mu_{n-1 \rightarrow n}}$$

... and all $n - 1$ descending messages

$$\boxed{\mu_{2 \rightarrow 1}} \leftarrow \boxed{\mu_{3 \rightarrow 1}} \leftarrow \dots \leftarrow \boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

Computing all marginals $p(x_1), p(x_2), \dots, p(x_n)$

Computing $p(x_j)$ requires to compute

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \quad \text{and}$$
$$\boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

How do we compute $p(x_{j'})$ for $j' \neq j$? Can use the same messages !

If we compute all $n - 1$ ascending messages

$$\boxed{\mu_{1 \rightarrow 2}} \rightarrow \boxed{\mu_{2 \rightarrow 3}} \rightarrow \dots \rightarrow \boxed{\mu_{j-1 \rightarrow j}} \rightarrow \dots \rightarrow \boxed{\mu_{n-2 \rightarrow n-1}} \rightarrow \boxed{\mu_{n-1 \rightarrow n}}$$

... and all $n - 1$ descending messages

$$\boxed{\mu_{2 \rightarrow 1}} \leftarrow \boxed{\mu_{3 \rightarrow 1}} \leftarrow \dots \leftarrow \boxed{\mu_{j+1 \rightarrow j}} \leftarrow \dots \leftarrow \boxed{\mu_{n-1 \rightarrow n-2}} \leftarrow \boxed{\mu_{n \rightarrow n-1}}$$

then all marginals are computed with an additional $O(n)$ computation.

Computing pairwise marginals $p(x_j, x_{j+1})$

Passing

- ascending messages until node j
- descending messages until node $j + 1$

one gets:

$$p(x_j, x_{j+1}) = \frac{1}{Z} \mu_{j-1 \rightarrow j}(x_j) \psi_j(x_j) \psi_{j+1}(x_{j+1}) \psi_{j,j+1}(x_j, x_{j+1}) \mu_{j+1 \rightarrow j}(x_j)$$

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root
- Let \mathcal{C}_j be the set of children of node j

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root
- Let \mathcal{C}_j be the set of children of node j
- Let \mathcal{D}_j be the set of descendants of node j

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root
- Let \mathcal{C}_j be the set of children of node j
- Let \mathcal{D}_j be the set of descendants of node j
- Let π_j be the parent of node j

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root
- Let \mathcal{C}_j be the set of children of node j
- Let \mathcal{D}_j be the set of descendants of node j
- Let π_j be the parent of node j

Let

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

Inference in undirected trees

Let (V, E) be a tree and assume a tree graphical model

$$p(x) = \frac{1}{Z} \prod_{i \in V} \psi_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Again we wish to compute $p(x_r)$ for some node r

- Orient the tree by setting node r to be the root
- Let \mathcal{C}_j be the set of children of node j
- Let \mathcal{D}_j be the set of descendants of node j
- Let π_j be the parent of node j

Let

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

Then we have

$$p(x) = \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

$$F(x_i, x_j, x_{\mathcal{D}_j}) = \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k})$$

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

$$F(x_i, x_j, x_{\mathcal{D}_j}) = \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k})$$

As a consequence

$$\sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j}) = \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \sum_{x_{\mathcal{D}_j}} \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k}) \right]$$

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

$$F(x_i, x_j, x_{\mathcal{D}_j}) = \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k})$$

As a consequence

$$\begin{aligned} \sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j}) &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \sum_{x_{\mathcal{D}_j}} \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \\ &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \sum_{x_k, x_{\mathcal{D}_k}} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \end{aligned}$$

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

$$F(x_i, x_j, x_{\mathcal{D}_j}) = \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k})$$

As a consequence

$$\begin{aligned} \sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j}) &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \sum_{x_{\mathcal{D}_j}} \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \\ &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \sum_{x_k, x_{\mathcal{D}_k}} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \end{aligned}$$

So if we define $\mu_{j \rightarrow i}(x_i) := \sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j})$,

Exploiting a recursion on the tree

$$F(x_i, x_j, x_{\mathcal{D}_j}) := \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{D}_j} [\psi_k(x_k) \psi_{\pi_k, k}(x_{\pi_k}, x_k)]$$

$$F(x_i, x_j, x_{\mathcal{D}_j}) = \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k})$$

As a consequence

$$\begin{aligned} \sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j}) &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \sum_{x_{\mathcal{D}_j}} \prod_{k \in \mathcal{C}_j} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \\ &= \sum_{x_j} \left[\psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \sum_{x_k, x_{\mathcal{D}_k}} F(x_j, x_k, x_{\mathcal{D}_k}) \right] \end{aligned}$$

So if we define $\mu_{j \rightarrow i}(x_i) := \sum_{x_j, x_{\mathcal{D}_j}} F(x_i, x_j, x_{\mathcal{D}_j})$, then we have

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \mu_{k \rightarrow j}(x_j)$$

Rooted sum-product algorithm: ascending scheme

The recursion

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \mu_{k \rightarrow j}(x_j)$$

defines an algorithm starting at the leaves.

Rooted sum-product algorithm: ascending scheme

The recursion

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \mu_{k \rightarrow j}(x_j)$$

defines an algorithm starting at the leaves.

- 1 For any leaf j with parent $i = \pi_j$

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j)$$

Rooted sum-product algorithm: ascending scheme

The recursion

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \mu_{k \rightarrow j}(x_j)$$

defines an algorithm starting at the leaves.

- 1 For any leaf j with parent $i = \pi_j$

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j)$$

- 2 Then for any node that only has leaves, we can compute the message to their parents
- 3 etc

Rooted sum-product algorithm: ascending scheme

The recursion

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{C}_j} \mu_{k \rightarrow j}(x_j)$$

defines an algorithm starting at the leaves.

- 1 For any leaf j with parent $i = \pi_j$

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j)$$

- 2 Then for any node that only has leaves, we can compute the message to their parents
- 3 etc

At the end of the algorithm:

Each node i has sent a message $\mu_{i \rightarrow \pi_i}(x_{\pi_i})$ to its parent.

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

$$p(x_r) = \sum_{x_{\mathcal{D}_r}} p(x)$$

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

$$p(x_r) = \sum_{x_{\mathcal{D}_r}} p(x) = \frac{1}{Z} \psi_r(x_r) \sum_{x_{\mathcal{D}_r}} \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

$$\begin{aligned} p(x_r) &= \sum_{x_{\mathcal{D}_r}} p(x) = \frac{1}{Z} \psi_r(x_r) \sum_{x_{\mathcal{D}_r}} \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i}) \\ &= \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \sum_{x_i, x_{\mathcal{D}_i}} F(x_r, x_i, x_{\mathcal{D}_i}) \end{aligned}$$

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

$$\begin{aligned} p(x_r) = \sum_{x_{\mathcal{D}_r}} p(x) &= \frac{1}{Z} \psi_r(x_r) \sum_{x_{\mathcal{D}_r}} \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i}) \\ &= \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \sum_{x_i, x_{\mathcal{D}_i}} F(x_r, x_i, x_{\mathcal{D}_i}) \\ &= \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \mu_{i \rightarrow r}(x_r) \end{aligned}$$

Finally computing $p(x_r)$

Remember

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i})$$

$$\begin{aligned} p(x_r) = \sum_{x_{\mathcal{D}_r}} p(x) &= \frac{1}{Z} \psi_r(x_r) \sum_{x_{\mathcal{D}_r}} \prod_{i \in \mathcal{C}(r)} F(x_r, x_i, x_{\mathcal{D}_i}) \\ &= \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \sum_{x_i, x_{\mathcal{D}_i}} F(x_r, x_i, x_{\mathcal{D}_i}) \\ &= \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \mu_{i \rightarrow r}(x_r) \end{aligned}$$

So

$$p(x_r) = \frac{1}{Z} \psi_r(x_r) \prod_{i \in \mathcal{C}(r)} \mu_{i \rightarrow r}(x_r)$$

Towards the sum-product without root

Towards the sum-product without root

Notice

- if we change the root, the edges whose orientation does not change still transmit the same message !

Towards the sum-product without root

Notice

- if we change the root, the edges whose orientation does not change still transmit the same message !
- the edges whose orientation has changed now send a message in the other direction...

Towards the sum-product without root

Notice

- if we change the root, the edges whose orientation does not change still transmit the same message !
- the edges whose orientation has changed now send a message in the other direction...

This leads to a general scheme in which each edge will at some point be traversed by a message in each direction.

Towards the sum-product without root

Notice

- if we change the root, the edges whose orientation does not change still transmit the same message !
- the edges whose orientation has changed now send a message in the other direction...

This leads to a general scheme in which each edge will at some point be traversed by a message in each direction.

Message passing without orientation of the tree

Let \mathcal{N}_i denote the set of neighbors of node i .

Consider the update:

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \mu_{k \rightarrow j}(x_j)$$

Towards the sum-product without root

Notice

- if we change the root, the edges whose orientation does not change still transmit the same message !
- the edges whose orientation has changed now send a message in the other direction...

This leads to a general scheme in which each edge will at some point be traversed by a message in each direction.

Message passing without orientation of the tree

Let \mathcal{N}_i denote the set of neighbors of node i .

Consider the update:

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \mu_{k \rightarrow j}(x_j)$$

Key remark: a node can only send a message if it has received messages from **all of its neighbors except one**.

Sum-product algorithm: sequential version

At each iteration:

Any node j that has received messages from **all but one** – call it i – of its neighbors sends a message $\mu_{j \rightarrow i}(x_i)$ to i with

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \mu_{k \rightarrow j}(x_j)$$

Flooding sum-product algorithm (aka parallel SP)

Initialize all messages at random

At each iteration:

For each node j

For each neighbor $i \in \mathcal{N}_j$

- j sends to i the message $\mu_{j \rightarrow i}(x_i)$ with

$$\mu_{j \rightarrow i}(x_i) = \sum_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \mu_{k \rightarrow j}(x_j)$$

The algorithm converges after a number of iterations equal to the diameter of the tree.

Computing the marginals from the messages

At each node $i \in V$

$$p(x_i) = \frac{1}{Z} \psi_i(x_i) \prod_{k \in \mathcal{N}_i} \mu_{k \rightarrow i}(x_i).$$

For each edge,

$$p(x_i, x_j) = \frac{1}{Z} \psi_i(x_i) \psi_j(x_j) \psi_{ij}(x_i, x_j) \left(\prod_{\substack{k \in \mathcal{N}_i \\ k \neq j}} \mu_{k \rightarrow i}(x_i) \right) \left(\prod_{\substack{\ell \in \mathcal{N}_j \\ \ell \neq i}} \mu_{\ell \rightarrow j}(x_j) \right)$$

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Example

$$p(x_i \mid x_5 = 3, x_{10} = 2) \propto p(x_i, x_5 = 3, x_{10} = 2)$$

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Example

$$p(x_i \mid x_5 = 3, x_{10} = 2) \propto p(x_i, x_5 = 3, x_{10} = 2)$$

Idea: redefine potentials

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Example

$$p(x_i \mid x_5 = 3, x_{10} = 2) \propto p(x_i, x_5 = 3, x_{10} = 2)$$

Idea: redefine potentials

$$\tilde{\psi}_5(x_5) = \psi_5(x_5) \delta(x_5, 3)$$

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Given observations $X_j = x_{j0}$ for $j \in B$, define modified potentials:

$$\tilde{\psi}_j(x_j) = \psi_j(x_j) \delta(x_j, x_{j0})$$

so that

$$p(x \mid X_B = x_{B0}) = \frac{1}{\tilde{Z}} \prod_{i \in V} \tilde{\psi}_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Given observations $X_j = x_{j0}$ for $j \in B$, define modified potentials:

$$\tilde{\psi}_j(x_j) = \psi_j(x_j) \delta(x_j, x_{j0})$$

so that

$$p(x \mid X_B = x_{B0}) = \frac{1}{\tilde{Z}} \prod_{i \in V} \tilde{\psi}_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Indeed we have

$$p(x \mid X_B = x_{B0}) p(X_B = x_{B0}) = p(x) \prod_{j \in B} \delta(x_j, x_{j0})$$

Computing *conditional marginals* $p(x_i \mid X_B = x_{B0})$

Given observations $X_j = x_{j0}$ for $j \in B$, define modified potentials:

$$\tilde{\psi}_j(x_j) = \psi_j(x_j) \delta(x_j, x_{j0})$$

so that

$$p(x \mid X_B = x_{B0}) = \frac{1}{Z} \prod_{i \in V} \tilde{\psi}_i(x_i) \prod_{(i,j) \in E} \psi_{i,j}(x_i, x_j)$$

Indeed we have

$$p(x \mid X_B = x_{B0}) p(X_B = x_{B0}) = p(x) \prod_{j \in B} \delta(x_j, x_{j0})$$

We then simply apply the SPA to these new potentials.

Inference beyond trees

- The SPA is also called *belief propagation* or *message passing*.
- On trees, it is an exact inference algorithm.

Inference beyond trees

- The SPA is also called *belief propagation* or *message passing*.
- On trees, it is an exact inference algorithm.
- If G is not a tree, the algorithm can be extended to *loopy belief propagation* but

Inference beyond trees

- The SPA is also called *belief propagation* or *message passing*.
- On trees, it is an exact inference algorithm.
- If G is not a tree, the algorithm can be extended to *loopy belief propagation* but
 - It does not converge in general to the correct marginals
 - it sometimes gives reasonable approximations.

Inference beyond trees

- The SPA is also called *belief propagation* or *message passing*.
- On trees, it is an exact inference algorithm.
- If G is not a tree, the algorithm can be extended to *loopy belief propagation* but
 - It does not converge in general to the correct marginals
 - it sometimes gives reasonable approximations.

Junction trees

- The exact SPA can be extended to graphs that are close to trees, in the sense that they can be viewed as trees if nodes can be grouped together into cliques to form a tree.
- Elegant theory, but rarely used in practice.

Generalized distributive law

We need to compute

$$p(x_{i_0}) = \sum_{x_{V \setminus \{i_0\}}} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

Generalized distributive law

We needed to compute

$$p(x_{i_0}) = \sum_{x_{V \setminus \{i_0\}}} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

The only key property we used is that

The addition is *distributive* over the product

$$(a + b)(c + d) = ac + bc + ad + bd$$

Generalized distributive law

We needed to compute

$$p(x_{i_0}) = \sum_{x_V \setminus \{i_0\}} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

The only key property we used is that

The addition is *distributive* over the product

$$(a + b)(c + d) = ac + bc + ad + bd$$

Max product for decoding (MAP inference)

Decoding consists in computing the most probable configuration

$$p(x^*) = \max_{x_V} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

Remark: this is often the most probable configuration of a conditional model (given some input data).

Generalized distributive law

We needed to compute

$$p(x_{i_0}) = \sum_{x_{V \setminus \{i_0\}}} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

The only key property we used is that

The addition is *distributive* over the product

$$(a + b)(c + d) = ac + bc + ad + bd$$

Max product for decoding (MAP inference)

Decoding consists in computing the most probable configuration

$$p(x^*) = \max_{x_V} \left[\prod_{i \in V} \psi_i(x_i) \prod_{\{i,j\} \in E} \psi_{ij}(x_i, x_j) \right]$$

Remark: this is often the most probable configuration of a conditional model (given some input data).

Decoding (aka MAP inference) with max-product

Apply the same algorithms as for SPA, but replace the update with

$$\check{\mu}_{j \rightarrow i}(x_i) = \max_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \check{\mu}_{k \rightarrow j}(x_j)$$

Decoding (aka MAP inference) with max-product

Apply the same algorithms as for SPA, but replace the update with

$$\check{\mu}_{j \rightarrow i}(x_i) = \max_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \check{\mu}_{k \rightarrow j}(x_j)$$

- Known as the Viterbi algorithm in the case of chains.

Decoding (aka MAP inference) with max-product

Apply the same algorithms as for SPA, but replace the update with

$$\check{\mu}_{j \rightarrow i}(x_i) = \max_{x_j} \psi_{i,j}(x_i, x_j) \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus \{i\}} \check{\mu}_{k \rightarrow j}(x_j)$$

- Known as the Viterbi algorithm in the case of chains.

Max-sum algorithm

Same on log-scale:

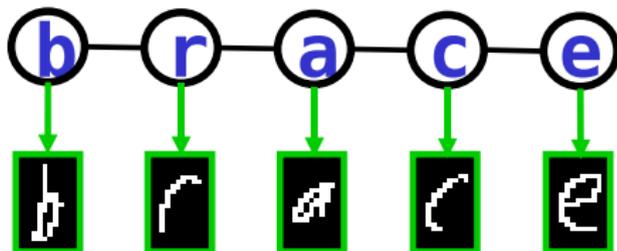
$$\log \check{\mu}_{j \rightarrow i}(x_i) = \max_{x_j} \left[\log \psi_{i,j}(x_i, x_j) + \log \psi_j(x_j) + \sum_{k \in \mathcal{N}_j \setminus \{i\}} \log \check{\mu}_{k \rightarrow j}(x_j) \right]$$

Other inference methods

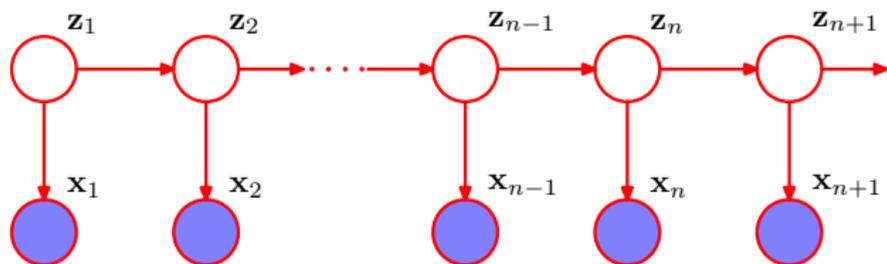
- There are many (exact/approximate/inexact) inference algorithms
- For a general graph exact inference is NP-hard
- The sum-product algorithm (SPA) performs inference on trees, with a complexity linear in the number of nodes.
- SPA can be generalized to graph that are close to trees in some sense using the *junction tree theory*.
- In general, one needs to use approximate or inexact methods
 - Gibbs sampling, and other MCMC sampling methods
 - Variational methods
 - Mean field (/Structured mean field)
 - Loopy belief propagation
 - TRW-entropy based inference

Hidden Markov Model (HMM)

- voice recognition
- natural language processing
- handwritten character recognition
- modelling biological sequence (protein, DNA)



Hidden Markov Model (HMM)



$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) = p(\mathbf{z}_1) \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)$$

Homogeneous Markov chains

- $z_n \in \{0, 1\}^K$ state indicator variable ($1, \dots, K$)
- *homogeneous* Markov chain: $\forall n, p(z_n | z_{n-1}) = p(z_2 | z_1)$
- emitted symbol \mathbf{x}_n ($\{0, 1\}^K$) / observation (\mathbb{R}^d)

Hidden Markov Model (HMM)

Parameterization

distribution of the initial state $p(\mathbf{z}_1; \pi) = \prod_{k=1}^K \pi_k^{z_{1k}}$

Hidden Markov Model (HMM)

Parameterization

distribution of the initial state $p(\mathbf{z}_1; \pi) = \prod_{k=1}^K \pi_k^{z_{1k}}$

transition matrix $p(\mathbf{z}_n | \mathbf{z}_{n-1}; A) = \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$

Hidden Markov Model (HMM)

Parameterization

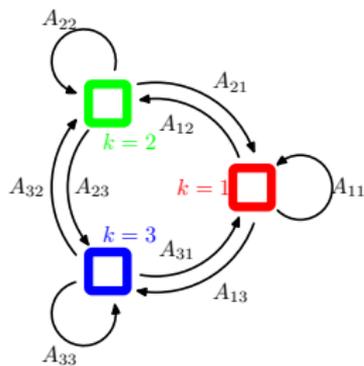
distribution of the initial state	$p(\mathbf{z}_1; \pi) = \prod_{k=1}^K \pi_k^{z_{1k}}$
transition matrix	$p(\mathbf{z}_n \mathbf{z}_{n-1}; A) = \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$
emission probabilities	$p(\mathbf{x}_n \mathbf{z}_n; \phi)$ e.g. Gaussian Mixture

Hidden Markov Model (HMM)

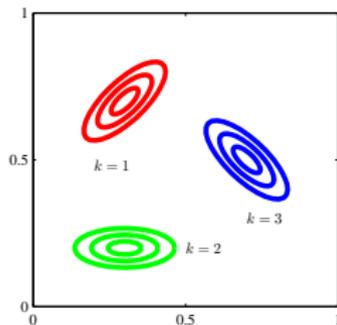
Parameterization

- distribution of the initial state $p(\mathbf{z}_1; \pi) = \prod_{k=1}^K \pi_k^{z_{1k}}$
- transition matrix $p(\mathbf{z}_n | \mathbf{z}_{n-1}; A) = \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$
- emission probabilities $p(\mathbf{x}_n | \mathbf{z}_n; \phi)$ e.g. Gaussian Mixture

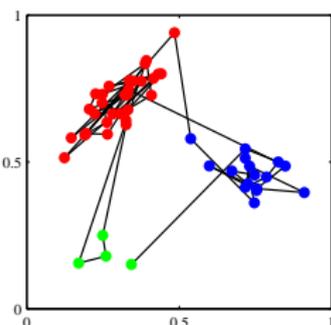
Interpretation



Transitions of \mathbf{z}_n



$p(\mathbf{x}_n | \mathbf{z}_n)$



Path of \mathbf{z}_n

Maximum likelihood for HMMs

Application of the EM algorithm

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t) \quad \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t)$$

Maximum likelihood for HMMs

Application of the EM algorithm

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t) \quad \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t)$$

Expected value of the log-likelihood:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(x_n | \phi_k)$$

Maximum likelihood for HMMs

Application of the EM algorithm

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t) \quad \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t)$$

Expected value of the log-likelihood:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(x_n | \phi_k)$$

Maximizing with respect to the parameters $\{\pi, A\}$, one gets

$$\pi_k^{t+1} = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$

$$A_{jk}^{t+1} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})}$$

Maximum likelihood for HMMs

Application of the EM algorithm

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t) \quad \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t)$$

Expected value of the log-likelihood:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(x_n | \phi_k)$$

Maximizing with respect to the parameters $\{\pi, A\}$, one gets

$$\pi_k^{t+1} = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$

$$A_{jk}^{t+1} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})}$$

If the emissions are Gaussian then we also have:

$$\boldsymbol{\mu}_k^{t+1} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad \boldsymbol{\Sigma}_k^{t+1} = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top}{\sum_{n=1}^N \gamma(z_{nk})}$$

Maximum likelihood for HMMs

Baum-Welch algorithm

The Baum-Welch algorithm is a special instance of the sum-product algorithm. It is also known under the name *forward-backward*.

One propagates the messages

- forward $\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$

Maximum likelihood for HMMs

Baum-Welch algorithm

The Baum-Welch algorithm is a special instance of the sum-product algorithm. It is also known under the name *forward-backward*.

One propagates the messages

- forward $\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$
- backward $\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$

Maximum likelihood for HMMs

Baum-Welch algorithm

The Baum-Welch algorithm is a special instance of the sum-product algorithm. It is also known under the name *forward-backward*.

One propagates the messages

- forward $\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$
- backward $\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$

that satisfy the following properties:

$$\alpha(\mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) \quad \beta(\mathbf{z}_n) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)$$

Maximum likelihood for HMMs

Baum-Welch algorithm

The Baum-Welch algorithm is a special instance of the sum-product algorithm. It is also known under the name *forward-backward*.

One propagates the messages

- forward $\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$
- backward $\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$

that satisfy the following properties:

$$\alpha(\mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) \quad \beta(\mathbf{z}_n) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)$$

Finally one gets the marginals:

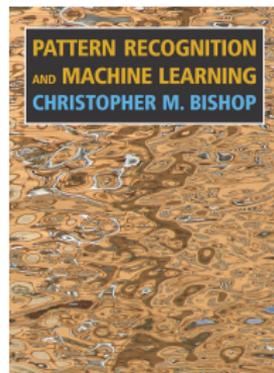
$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^t) = \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X} | \boldsymbol{\theta}^t)}$$

et

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = \frac{\alpha(\mathbf{x}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \beta(\mathbf{x}_n)}{p(\mathbf{X} | \boldsymbol{\theta}^t)}$$

Conclusions

- Graphical models offer a language to design structured probabilistic models
- The graph encodes both
 - A factorization of the probability distribution
 - Implicitly associated conditional independences between the variables
- Three key operations on graphical model are decoding, inference and learning.
- Probabilistic inference and decoding can be done efficiently using sum-product and max-product algorithm if the graph is a tree.
- Learning with the maximum likelihood principle in exponential families is a convex optimization problem requires to solve the probabilistic inference problem.
- Latent variables are easily handled using the Expectation-Maximization algorithm, which again requires to solve the probabilistic inference problem.

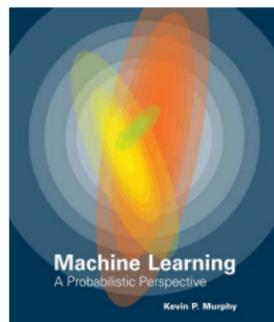


Pattern Recognition and Machine Learning,
Christopher Bishop, Springer (2006).

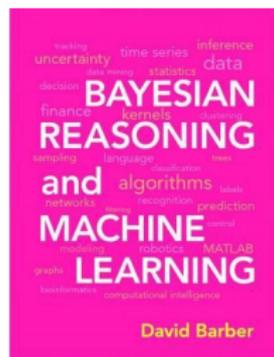
http:

`//research.microsoft.com/~cmbishop/PRML/`

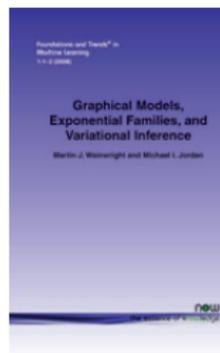
Many of the figures of these lectures are figures from Bishop's book that are kindly made available for teaching purposes.



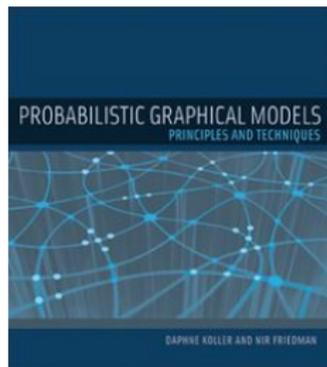
Machine Learning, a probabilistic perspective
Kevin Murphy, MIT Press (2012).



Bayesian reasoning and machine learning,
David Barber,
Cambridge University Press (2012).
<http://www.cs.ucl.ac.uk/staff/d.barber/brml/>



Wainwright, M. J., & Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. Foundations and Trends in Machine Learning.



Daphne Koller and Nir Friedman, Probabilistic Graphical Models - Principles and Techniques, 2009, MIT Press.

Going further

- There are many approximate/inexact variational inference algorithms (Wainwright, 2008)
- Graphical models can be learned in the max-margin setting (Taskar et al., 2003; Tsochantaridis et al., 2005; Pletscher et al., 2010)
- Graphical models techniques are used in deep learning: variational auto-encoders, RBMs, generative adversarial networks, etc.

References I

- Meshi, O., Sontag, D., Globerson, A., and Jaakkola, T. S. (2010). Learning efficiently with approximate inference via dual losses. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 783–790.
- Meshi, O., Srebro, N., and Hazan, T. (2015). Efficient training of structured SVMs via soft constraints. In *AISTATS*.
- Pletscher, P., Ong, C. S., and Buhmann, J. M. (2010). Entropy and margin maximization for structured output learning. In *ECML*.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin markov networks. In *NIPS*.
- Tsochantaridis, I., Joachims, T., Hofmann, T., and Altun, Y. (2005). Large margin methods for structured and interdependent output variables. In *Journal of Machine Learning Research*, pages 1453–1484.
- Wainwright, M. J. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*.