

Computer Vision Group Prof. Daniel Cremers

Technische Universität München

4a. Inference in Graphical Models

Inference on a Chain (Rep.) $\mu_{lpha}(x_{n-1}) = \mu_{lpha}(x_n) = \mu_{eta}(x_n) = \mu_{eta}(x_{n+1})$

- x_n x_{n-1} x_N x_{n+1} x_1
- The first values of μ_{α} and μ_{β} are:

$$\mu_{\alpha}(x_2) = \sum_{x_1} \psi_{1,2}(x_1, x_2) \qquad \qquad \mu_{\beta}(x_{N-1}) = \sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N)$$

The partition function can be computed at any node:

$$Z = \sum_{x_n} \mu_{\alpha}(x_n) \mu_{\beta}(x_n)$$

• Overall, we have $O(NK^2)$ operations to compute the marginal $p(x_n)$



The message-passing algorithm can be extended to more general graphs:



It is then known as the sum-product algorithm. A special case of this is belief propagation.



The message-passing algorithm can be extended to more general graphs:



An **undirected tree** is defined as a graph that has exactly one path between any two nodes



The message-passing algorithm can be extended to more general graphs:

A directed tree has only one node without parents and all other nodes have exactly one parent



Conversion from a directed to an undirected tree is no problem, because no links are inserted

The same is true for the conversion back to a directed tree



The message-passing algorithm can be extended to more general graphs:

Polytrees can contain nodes with several parents, therefore moralization can remove independence relations Polytree





- The Sum-product algorithm can be used to do inference on undirected and directed graphs.
- A representation that generalizes directed and undirected models is the factor graph.





 $f(x_1, x_2, x_3) = p(x_1)p(x_2)p(x_3 \mid x_1, x_2)$ Factor graph



- The Sum-product algorithm can be used to do inference on undirected and directed graphs.
- A representation that generalizes directed and undirected models is the factor graph.





Factor graphs

- can contain multiple factors for the same nodes
- are more general than undirected graphs
- are bipartite, i.e. they consist of two kinds of nodes and all edges connect nodes of different kind





- Directed trees convert to tree-structured factor graphs
- The same holds for undirected trees
- Also: directed polytrees convert to tree-structured factor graphs
- And: Local cycles in a directed graph can be removed by converting to a factor graph





Assumptions:

- all variables are discrete
- the factor graph has a tree structure

The factor graph represents the joint distribution as a product of factor nodes:

$$p(\mathbf{x}) = \prod_{s} f_s(\mathbf{x}_s)$$

The marginal distribution at a given node x is

$$p(x) = \sum_{\mathbf{x} \setminus x} p(\mathbf{x})$$









 $F_s(x, X_s) = f_s(x, x_1, \dots, x_M) G_1(x_1, X_{s_1}) \dots G_M(x_M, X_{s_M})$

The messages can then be computed as

$$\mu_{f_s \to x}(x) = \sum_{x_1} \cdots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \operatorname{ne}(f_s) \setminus x} \sum_{X_{s_m}} G_m(x_m, X_{s_m})$$
$$= \sum_{x_1} \cdots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m \in \operatorname{ne}(f_s) \setminus x} \mu_{x_m \to f_s}(x_m)$$
"Messages from nodes to factors"

 x_M

 x_m

 $G_m(x_m, X_{sm})$

 $\mu_{x_M \to f_s}(x_M)$

 $\mu_{f_s \to x}(x)$

 J_s





The factors *G* of the neighboring nodes can again be factorized further:

$$G_M(x_m, X_{s_m}) = \prod_{l \in \operatorname{ne}(x_m) \setminus f_s} F_l(x_m, X_{m_l})$$

This results in the exact same situation as above! We can now recursively apply the derived rules:

$$\mu_{x_m \to f_s}(x_m) = \prod_{l \in \operatorname{ne}(x_m) \setminus f_s} \sum_{X_{m_l}} F_l(x_m, X_{m_l})$$
$$= \prod_{l \in \operatorname{ne}(x_m) \setminus f_s} \mu_{f_l \to x_m}(x_m)$$



Summary marginalization:

- 1.Consider the node *x* as a root note
- 2.Initialize the recursion at the leaf nodes as:

 $\mu_{f \to x}(x) = 1$ (var) or $\mu_{x \to f}(x) = f(x)$ (fac)

- 3.Propagate the messages from the leaves to the root *x*
- 4.Propagate the messages back from the root to the leaves
- 5.We can get the marginals at every node in the graph by multiplying all incoming messages





Sum-product is used to find the marginal distributions at every node, but:

How can we find the setting of all variables that **maximizes** the joint probability? And what is the value of that maximal probability?

Idea: use sum-product to find all marginals and then report the value for each node *x* that maximizes the marginal p(x)

However: this does not give the **overall** maximum of the joint probability



Observation: the max-operator is distributive, just like the multiplication used in sum-product:

 $\max(ab, ac) = a \max(b, c)$ if $a \ge 0$

Idea: use max instead of sum as above and exchange it with the product

Chain example: $\max_{\mathbf{x}} p(\mathbf{x}) = \frac{1}{Z} \max_{x_1} \dots \max[\psi_{1,2}(x_1, x_2) \dots \psi_{N-1,N}(x_{N-1}, x_N)]$ $= \frac{1}{Z} \max_{x_1} [\psi_{1,2}(x_1, x_2) [\dots \max \psi_{N-1,N}(x_{N-1}, x_N)]]$

Message passing can be used as above!



To find the maximum value of $p(\mathbf{x})$, we start again at the leaves and propagate to the root.

Two problems:

- no summation, but many multiplications; this leads to numerical instability (very small values)
- when propagating back, multiple configurations of x can maximize *p*(x), leading to wrong assignments of the overall maximum

Solution to the first:

Transform everything into log-space and use sums



Solution to the second problem:

Keep track of the arg max in the forward step, i.e. store at each node which value was responsible for the maximum:

$$\phi(x_n) = \arg\max_{x_{n-1}} \left[\ln f_{n-1,n}(x_{n-1}, x_n) + \mu_{x_{n-1} \to f_{n-1,n}}(x_n) \right]$$

Then, in the back-tracking step we can recover the arg max by recursive substitution of:

$$x_{n-1}^{\max} = \phi(x_n^{\max})$$



Other Inference Algorithms

Junction Tree Algorithm:

- Provides exact inference on general graphs.
- Works by turning the initial graph into a junction
 tree and then running a sum-product-like algorithm
- A junction tree is obtained from an undirected model by triangulation and mapping cliques to nodes and connections of cliques to edges
- It is the maximal spanning tree of cliques

Problem: Intractable on graphs with large cliques.

Cost grows exponentially with the number of variables in the largest clique ("tree width").





Other Inference Algorithms

Loopy Belief Propagation:

- Performs Sum-Product on general graphs, particularly when they have loops
- Propagation has to be done several times, until a convergence criterion is met
- No guarantee of convergence and no global optimum
- Messages have to be scheduled
- Initially, unit messages passed across all edges
- Approximate, but tractable for large graphs



Conditional Random Fields

- Another kind of undirected graphical model is known as Conditional Random Field (CRF).
- CRFs are used for classification where labels are represented as discrete random variables y and features as continuous random variables x
- A CRF represents the conditional probability

$$p(\mathbf{y} \mid \mathbf{x}, \mathbf{w}) = \frac{\prod_{C} \phi_{C}(\mathbf{x}_{C}, \mathbf{y}_{C}; \mathbf{w})}{\sum_{\mathbf{y}'} \prod_{C} \phi_{C}(\mathbf{x}_{C}, \mathbf{y}'_{C}; \mathbf{w})}$$

where w are parameters learned from training data.

CRFs are discriminative and MRFs are generative



Conditional Random Fields

Derivation of the formula for CRFs:

$$p(\mathbf{y} \mid \mathbf{x}, \mathbf{w}) = \frac{p(\mathbf{y}, \mathbf{x} \mid \mathbf{w})}{p(\mathbf{x} \mid \mathbf{w})} = \frac{p(\mathbf{y}, \mathbf{x} \mid \mathbf{w})}{\sum_{y'} p(\mathbf{y}', \mathbf{x} \mid \mathbf{w})} = \frac{\prod_C \phi_C(\mathbf{x}_C, \mathbf{y}_C; \mathbf{w}) \quad Z}{Z \quad \sum_{\mathbf{y}'} \prod_C \phi_C(\mathbf{x}_C, \mathbf{y}'_C; \mathbf{w})}$$

In the training phase, we compute parameters w that maximize the posterior:

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} p(\mathbf{w} \mid \mathbf{x}^*, \mathbf{y}^*) \propto p(\mathbf{y}^* \mid \mathbf{x}^*, \mathbf{w}) p(\mathbf{w})$$

where (x^*, y^*) is the training data and p(w) is a Gaussian prior. In the inference phase we maximize

$$\arg\max_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}, \mathbf{w}^*)$$



Conditional Random Fields



Note: the definition of $x_{i,j}$ and $y_{i,j}$ is different from the one in C.M. Bishop (pg.389)!



CRF Training

We minimize the negative log-posterior:

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \{-\ln p(\mathbf{w} \mid \mathbf{x}^*, \mathbf{y}^*)\} = \arg\min_{\mathbf{w}} \{-\ln p(\mathbf{y}^* \mid \mathbf{x}^*, \mathbf{w}) - \ln p(\mathbf{w})\}$$

Computing the likelihood is intractable, as we have to compute the partition function for each w. We can approximate the likelihood using **pseudo-likelihood**:

where

$$p(\mathbf{y}^* \mid \mathbf{x}^*, \mathbf{w}) \approx \prod_{i} p(y_i^* \mid \mathcal{M}(y_i^*), \mathbf{x}^*, \mathbf{w})$$

$$(\mathbf{y}_i^* \mid \mathbf{M}(y_i^*), \mathbf{x}^*, \mathbf{w}) = \frac{\prod_{C_i} \phi_{C_i}(\mathbf{x}_{C_i}^*, y_i^*, \mathbf{y}_{C_i}^*; \mathbf{w}))}{\sum_{y_i'} \prod_{C_i} \phi_{C}(\mathbf{x}_{C_i}^*, y_i', \mathbf{y}_{C_i}^*; \mathbf{w})}$$



Pseudo Likelihood







Pseudo Likelihood

 $x_{i,j}$

 $y_{i,j}$





Potential Functions

• The only requirement for the potential functions is that they are positive. We achieve that with:

 $\phi_C(\mathbf{x}_C, \mathbf{y}_C, \mathbf{w}) := \exp(\mathbf{w}^T f(\mathbf{x}_C, \mathbf{y}_C))$

Where f is a compatibility function that is large if the labels y_C fit well to the features x_C .

- This is called the log-linear model.
- The function f can be, e.g. a local classifier



CRF Training and Inference

Training:

 Using pseudo-likelihood, training is efficient. We have to minimize:

$$L(\mathbf{w}) = -lpl(\mathbf{y}^* \mid \mathbf{x}^*, \mathbf{w}) + \frac{1}{2\sigma^2} \mathbf{w}^T \mathbf{w}$$

Log-pseudo-likelihood Gaussian prior

This is a convex function that can be minimized using gradient descent

Inference:

 Only approximatively, e.g. using loopy belief propagation





Summary

- Undirected Graphical Models represent conditional independence more intuitively using graph separation
- Their factorization is done based on potential functions The normalizer is called the partition function, which in general is intractable to compute
- Inference in graphical models can be done efficiently using the sum-product algorithm (message passing).
- Another inference algorithm is loopy belief propagation, which is approximate, but tractable
- Conditional Random Fields are a special kind of MRFs and can be used for classification

