

Computer Vision Group Prof. Daniel Cremers

Technische Universität München

9. Probabilistic Graphical Models Directed Models

The Bayes Filter (Rep.)

$$\begin{split} & \operatorname{Bel}(x_t) = p(x_t \mid u_1, z_1, \dots, u_t, z_t) \\ & \text{(Bayes)} &= \eta \; p(z_t \mid x_t, u_1, z_1, \dots, u_t) p(x_t \mid u_1, z_1, \dots, u_t) \\ & \text{(Markov)} &= \eta \; p(z_t \mid x_t) p(x_t \mid u_1, z_1, \dots, u_t) \\ & \text{(Tot. prob.)} &= \eta \; p(z_t \mid x_t) \int p(x_t \mid u_1, z_1, \dots, u_t, x_{t-1}) \\ & \quad p(x_{t-1} \mid u_1, z_1, \dots, u_t) dx_{t-1} \\ & \text{(Markov)} &= \eta \; p(z_t \mid x_t) \int p(x_t \mid u_t, x_{t-1}) p(x_{t-1} \mid u_1, z_1, \dots, u_t) dx_{t-1} \\ & \text{(Markov)} &= \eta \; p(z_t \mid x_t) \int p(x_t \mid u_t, x_{t-1}) p(x_{t-1} \mid u_1, z_1, \dots, z_{t-1}) dx_{t-1} \\ &= \eta \; p(z_t \mid x_t) \int p(x_t \mid u_t, x_{t-1}) \operatorname{Bel}(x_{t-1}) dx_{t-1} \end{split}$$



Graphical Representation (Rep.)

We can describe the overall process using a Dynamic Bayes Network:



• This incorporates the following Markov assumptions: $p(z_t \mid x_{0:t}, u_{1:t}, z_{1:t}) = p(z_t \mid x_t) \text{ (measurement)}$ $p(x_t \mid x_{0:t-1}, u_{1:t}, z_{1:t}) = p(x_t \mid x_{t-1}, u_t) \text{ (state)}$



Definition

A Probabilistic Graphical Model is a diagrammatic representation of a probability distribution.

- In a Graphical Model, random variables are represented as **nodes**, and statistical dependencies are represented using **edges** between the nodes.
- The resulting graph can have the following properties:
- Cyclic / acyclic
- Directed / undirected
- The simplest graphs are Directed Acyclig Graphs (DAG).



Simple Example

- Given: 3 random variables a, b, and c
- Joint prob: p(a, b, c) = p(c|a, b)p(a, b) = p(c|a, b)p(b|a)p(a)



A Graphical Model based on a DAG is called a **Bayesian Network**

Simple Example

- In general: K random variables x_1, x_2, \ldots, x_K
- Joint prob:

 $p(x_1,\ldots,x_K) = p(x_K|x_1,\ldots,x_{K-1})\ldots p(x_2|x_1)p(x_1)$

- This leads to a fully connected graph.
- Note: The ordering of the nodes in such a fully connected graph is arbitrary. They all represent the joint probability distribution:

$$p(a, b, c) = p(a|b, c)p(b|c)p(c)$$
$$p(a, b, c) = p(b|a, c)p(a|c)p(c)$$

Bayesian Networks

Statistical independence can be represented by the **absence** of edges. This makes the computation efficient.

 $p(x_1, \dots, x_7) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3)$ $p(x_5|x_1, x_3)p(x_6|x_4)p(x_7|x_4, x_5)$

Intuitively: only x_1 and x_3 have an influence on x_5

Bayesian Networks

We can now define a one-to-one mapping from graphical models to probabilistic formulations:

General Factorization:

$$p(\mathbf{x}) = \prod_{k=1}^{K} p(x_k | \mathrm{pa}_k)$$

where

 $pa_k \triangleq ancestors of p_k$

and

$$p(\mathbf{x}) = p(x_1, \ldots, x_K)$$

Elements of Graphical Models

In case of a series of random variables with equal dependencies, we can subsume them using a **plate:**

$$p(\mathbf{t}, \mathbf{w}) = p(\mathbf{w}) \prod_{n=1}^{N} p(t_n | \mathbf{w})$$

Elements of Graphical Models (2)

We distinguish between **input** variables and explicit **hyper-parameters**:

3 7

$$p(\mathbf{t}, \mathbf{w} | \mathbf{x}, \alpha, \sigma^2) = p(\mathbf{w} | \alpha) \prod_{n=1}^{N} p(t_n | \mathbf{w}, x_n, \sigma^2).$$

10

Elements of Graphical Models (3)

We distinguish between **observed** variables and **hidden** variables:

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w}) \prod_{n=1}^{N} p(t_n|\mathbf{w})$$

(deterministic parameters omitted in formula)

Regression as a Graphical Model

Regression: Prediction of a new target value \hat{t}

Two Special Cases

- We consider two special cases:
- All random variables are discrete; i.e. Each x_i is represented by values μ_1, \ldots, μ_K where

All random variables are Gaussian

Discrete Variables: Example

• Two dependent variables: $K^2 - 1$ parameters Here: K = 2

• Independent joint distribution: 2(K-1) parameters

$$K - 1 + K - 1 = 2(K - 1)$$

Discrete Variables: General Case

In a general joint distribution with M variables we need to store K^M -1 parameters

If the distribution can be described by this graph:

then we have only K-1 + (M-1) K(K-1) parameters. This graph is called a **Markov chain** with M nodes. The number of parameters grows only **linearly** with the number of variables.

Gaussian Variables

Assume all random variables are Gaussian and we define

$$p(x_i \mid \text{pa}_i) = \mathcal{N}\left(x_i; \sum_{j \in \text{pa}_i} w_{ij}x_j + b_i, v_i\right)$$

Then one can show that the joint probability p(x) is a multivariate Gaussian. Furthermore:

$$x_i = \sum_{j \in pa_i} w_{ij} x_j + b_j + \sqrt{v_i} \epsilon_i \qquad \epsilon_i \sim \mathcal{N}(0, 1)$$

Thus:

$$E[x_i] = \sum_{j \in pa_i} w_{ij} E[x_j] + b_i$$

i.e., we can compute the mean values recursively.

Machine Learning for Computer Vision PD Dr. Rudolph Triebel Computer Vision Group

Gaussian Variables

Assume all random variables are Gaussian and we define

$$p(x_i \mid \mathrm{pa}_i) = \mathcal{N}\left(x_i; \sum_{i \in \mathrm{pa}_i} w_{ij}x_j + b_i, v_i\right)$$

The same can be shown for the covariance. Thus:

Mean and covariance can be calculated recursively

Furthermore it can be shown that:

- The fully connected graph corresponds to a Gaussian with a general symmetric covariance matrix
- The non-connected graph corresponds to a diagonal covariance matrix

Independence (Rep.)

Definition 1.4: Two random variables X and Y are *independent* iff: p(x, y) = p(x)p(y)

For independent random variables X and Y we have:

$$p(x \mid y) = \frac{p(x, y)}{p(y)} = \frac{p(x)p(y)}{p(y)} = p(x)$$

Notation:
$$x \perp \!\!\!\perp y \mid \emptyset$$

Independence does **not** imply conditional independence! The same is true for the opposite case.

Conditional Independence (Rep.)

Definition 1.5: Two random variables X and Y are conditional independent given a third random variable Z iff:

$$p(x, y \mid z) = p(x \mid z)p(y \mid z)$$

This is equivalent to:

$$p(x \mid z) = p(x \mid y, z) \text{ and}$$
$$p(y \mid z) = p(y \mid x, z)$$

Notation:
$$x \perp \!\!\!\perp y \mid z$$

p(a, b, c) = p(a|c)p(b|c)p(c)

Marginalizing out *c* on both sides gives

$$p(a,b) = \sum_{c} p(a|c)p(b|c)p(c)$$

This is in general not equal to p(a)p(b).

Thus: *a* and *b* are not independent: $a \not\!\!\perp b \mid \emptyset$

a

Now, we condition on c (it is assumed to be known):

Thus: *a* and *b* are conditionally independent given *c*: $a \perp b \mid c$ We say that the node at *c* is a **tail-to-tail node** on the path between *a* and *b*

This graph represents the distribution:

p(a, b, c) = p(a)p(c|a)p(b|c)

Again, we marginalize over c:

$$p(a,b) = p(a) \sum_{c} p(c|a)p(b|c) = p(a) \sum_{c} p(c|a)p(b|c,a)$$
$$= p(a) \sum_{c} \frac{p(c,a)p(b,c,a)}{p(a)p(c,a)} = p(a) \sum_{c} p(b,c \mid a)$$
$$= p(a)p(b|a)$$

And we obtain: $a \not\perp b \mid \emptyset$

As before, now we condition on c:

And we obtain: $a \perp b \mid c$

We say that the node at c is a head-to-tail node on the path between a and b.

Now consider this graph:

And the result is: $a \perp b \mid \emptyset$

Again, we condition on_c

We say that the node at c is a head-to-head node on the path between a and b.

To Summarize

When does the graph represent (conditional) independence?

Tail-to-tail case: if we condition on the tail-to-tail node Head-to-tail case: if we cond. on the head-to-tail node Head-to-head case: if we do not condition on the headto-head node (and neither on any of its descendants)

In general, this leads to the notion of **D-separation** for directed graphical models.

D-Separation

Say: A, B, and C are non-intersecting subsets of nodes in a directed graph.

A path from A to B is **blocked** by C if it contains a node such that either

 a) the arrows on the path meet either head-to-tail or tail-totail at the node, and the node is in the set C, or

b) the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, are in the set C.
If all paths from A to B are blocked, A is said to be d-separated from B by C.

Notation: dsep(A, B|C)

D-Separation

Say: A, B, and C are non-intersecting subsets of **D-Separation is a** nodes A path ntains property of graphs a nod a) the a ^r tail-toand not of tail at t probability b) the a neither the noc **J**. distributions If all p aid to be d-separated from B by C. Notation: dsep(A, B|C)

D-Separation: Example

$\neg \operatorname{dsep}(a, b|c)$

We condition on a descendant of e, i.e. it does not block the path from a to b.

$\operatorname{dsep}(a, b|f)$

We condition on a tail-to-tail node on the only path from a to b, i.e f blocks the path.

I-Map

Definition 4.1: A graph G is called an I-map for a distribution p if every D-separation of G corresponds to a conditional independence relation satisfied by p:

$\forall A, B, C : \operatorname{dsep}(A, B, C) \Rightarrow A \perp\!\!\!\perp B \mid C$

Example: The fully connected graph is an I-map for any distribution, as there are no D-separations in that graph.

D-Map

Definition 4.2: A graph G is called an **D-map** for a distribution p if for every conditional independence relation satisfied by p there is a D-separation in G :

$\forall A, B, C : A \perp \!\!\!\perp B \mid C \Rightarrow \operatorname{dsep}(A, B, C)$

Example: The graph without any edges is a D-map for any distribution, as all pairs of subsets of nodes are D-separated in that graph.

Perfect Map

Definition 4.3: A graph G is called a perfect map for a distribution p if it is a D-map and an I-map of p.

$\forall A, B, C : A \perp\!\!\!\perp B \mid C \Leftrightarrow \operatorname{dsep}(A, B, C)$

A perfect map uniquely defines a probability distribution.

The Markov Blanket

Consider a distribution of a node x_i conditioned on all other nodes:

$$|\mathbf{x}_{\{j\neq i\}}\rangle = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_M)}{\int p(\mathbf{x}_1, \dots, \mathbf{x}_M) d\mathbf{x}_i}$$
$$= \frac{\prod_k p(\mathbf{x}_k | \mathbf{pa}_k)}{\int \prod_k p(\mathbf{x}_k | \mathbf{pa}_k) d\mathbf{x}_i}$$
$$= p(\mathbf{x}_i | \mathbf{x}_{\mathcal{M}_i})$$

Factors independent of \mathbf{x}_i cancel between numerator and denominator.

Summary

- Graphical models represent joint probability distributions using nodes for the random variables and edges to express (conditional) (in)dependence
- A prob. distribution can always be represented using a fully connected graph, but this is inefficient
- In a directed acyclic graph, conditional independence is determined using **D-separation**
- A perfect map implies a one-to-one mapping between c.i. relations and D-separations
- The Markov blanket is the minimal set of observed nodes to obtain conditional independence

Computer Vision Group Prof. Daniel Cremers

Technische Universität München

4. Probabilistic Graphical Models Undirected Models

Repetition: Bayesian Networks

Directed graphical models can be used to represent **probability distributions** This is useful to do **inference** and to **generate samples** from the distribution efficiently

$$p(x_1, \dots, x_7) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3)$$
$$p(x_5|x_1, x_3)p(x_6|x_4)p(x_7|x_4, x_5)$$

Repetition: D-Separation

- D-separation is a property of graphs that can be easily determined
- An I-map assigns every d-separation a c.i. rel
- A D-map assigns every c.i. rel a d-separation
- Every Bayes net determines a unique prob. dist.

In-depth: The Head-to-Head Node

$$p(a) = 0.9 \qquad p(b) = 0.9$$

$$a \qquad b \qquad p(c)$$

$$1 \qquad 1 \qquad 0.8$$

$$1 \qquad 0 \qquad 0.2$$

$$0 \qquad 1 \qquad 0.2$$

$$0 \qquad 0 \qquad 0.1$$

Example:

- a: Battery charged (0 or 1)
- b: Fuel tank full (0 or 1)
- c: Fuel gauge says full (0 or 1)
- We can compute $p(\neg c) = 0.315$

and $p(\neg c \mid \neg b) = 0.81$

and obtain $p(\neg b \mid \neg c) \approx 0.257$

similarly: $p(\neg b \mid \neg c, \neg a) \approx 0.111$

"*a* explains *c* away"

Repetition: D-Separation

Directed vs. Undirected Graphs

Using D-separation we can identify conditional independencies in directed graphical models, but:

- Is there a simpler, more intuitive way to express conditional independence in a graph?
- Can we find a representation for cases where an "ordering" of the random variables is inappropriate (e.g. the pixels in a camera image)?

Yes, we can: by removing the directions of the edges we obtain an Undirected Graphical Model, also known as a Markov Random Field

Example: Camera Image

- directions are counter-intuitive for images
- Markov blanket is not just the direct neighbors when using a directed model

Markov Random Fields

All paths from *A* to *B* go through *C*, i.e. *C* blocks all paths.

Markov Blanket

We only need to condition on the direct neighbors of

x to get c.i., because these already block every path from x to any other node.

Factorization of MRFs

Any two nodes x_i and x_j that are not connected in an MRF are conditionally independent given all other nodes:

 $p(x_i, x_j \mid \mathbf{x}_{\backslash \{i,j\}}) = p(x_i \mid \mathbf{x}_{\backslash \{i,j\}}) p(x_j \mid \mathbf{x}_{\backslash \{i,j\}})$

In turn: each factor contains only nodes that are connected

This motivates the consideration of cliques in the graph:

- A clique is a fully connected subgraph.
- A maximal clique can not be extended with another node without loosing the property of full connectivity.

Maximal Clique

Factorization of MRFs

In general, a Markov Random Field is factorized as

$$p(\mathbf{x}) = \frac{\prod_C \phi_C(\mathbf{x}_C)}{\sum_{\mathbf{x}'} \prod_C \phi_C(\mathbf{x}'_C)} = \frac{1}{Z} \prod_C \phi_C(\mathbf{x}_C)$$
(4.1)

where *C* is the set of all (maximal) cliques and Φ_C is a positive function of a given clique \mathbf{x}_C of nodes, called the **clique potential**. *Z* is called the **partition function**. **Theorem (Hammersley/Clifford):** Any undirected model with associated clique potentials Φ_C is a perfect map for the probability distribution defined by Equation (4.1).

As a conclusion, all probability distributions that can be factorized as in (4.1), can be represented as an MRF.

Converting Directed to Undirected Graphs (1)

Converting Directed to Undirected Graphs (2)

$$p(\mathbf{x}) = p(x_1)p(x_2)p(x_2)p(x_4 \mid x_1, x_2, x_3)$$

In general: conditional distributions in the directed graph are mapped to cliques in the undirected graph

However: the variables are **not** conditionally independent given the head-to-head node

Therefore: Connect all parents of head-to-head nodes with each other (moralization)

Converting Directed to Undirected Graphs (2)

 $p(\mathbf{x}) = p(x_1)p(x_2)p(x_2)p(x_4 \mid x_1, x_2, x_3)$

 $p(\mathbf{x}) = \phi(x_1, x_2, x_3, x_4)$

Problem: This process can remove conditional independence relations (inefficient)

Generally: There is no one-to-one mapping between the distributions represented by directed and by undirected graphs.

Representability

- As for DAGs, we can define an I-map, a D-map and a perfect map for MRFs.
- The set of all distributions for which a DAG exists that is a perfect map is different from that for MRFs.

Using Graphical Models

We can use a graphical model to do inference:

- Some nodes in the graph are observed, for others we want to find the posterior distribution
- Also, computing the local marginal distribution p(x_n) at any node x_n can be done using inference.

Question: How can inference be done with a graphical model?

We will see that when exploiting conditional independences we can do efficient inference.

The joint probability is given by

$$p(\mathbf{x}) = \frac{1}{Z}\psi_{1,2}(x_1, x_2)\psi_{2,3}(x_2, x_3)\psi_{3,4}(x_3, x_4)\psi_{4,5}(x_4, x_5)$$

The marginal at x_3 is $p(x_3) = \sum_{x_1} \sum_{x_2} \sum_{x_4} \sum_{x_5} p(\mathbf{x})$

In the general case with N nodes we have

$$p(\mathbf{x}) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \cdots \psi_{N-1,N}(x_{N-1}, x_N)$$

and
$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x})$$

 This would mean K^N computations! A more efficient way is obtained by rearranging:

$$p(x_{3}) = \frac{1}{Z} \sum_{x_{1}} \sum_{x_{2}} \sum_{x_{4}} \sum_{x_{5}} \psi_{1,2}(x_{1}, x_{2})\psi_{2,3}(x_{2}, x_{3})\psi_{3,4}(x_{3}, x_{4})\psi_{4,5}(x_{4}, x_{5})$$

$$= \frac{1}{Z} \sum_{x_{2}} \sum_{x_{1}} \sum_{x_{4}} \sum_{x_{5}} \psi_{1,2}(x_{1}, x_{2})\psi_{2,3}(x_{2}, x_{3})\psi_{3,4}(x_{3}, x_{4})\psi_{4,5}(x_{4}, x_{5})$$

$$= \frac{1}{Z} \sum_{x_{2}} \psi_{2,3}(x_{2}, x_{3}) \sum_{x_{1}} \psi_{1,2}(x_{1}, x_{2}) \sum_{x_{4}} \psi_{3,4}(x_{3}, x_{4}) \sum_{x_{5}} \psi_{4,5}(x_{4}, x_{5})$$

$$\mu_{\alpha}(x_{3}) \leftarrow \text{Vectors of size K} \rightarrow \mu_{\beta}(x_{3})$$

JUIIDULEI

VISIOL

In general, we have

$$p(x_n) = \frac{1}{Z} \left[\sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \cdots \left[\sum_{x_1} \psi_{1,2}(x_1, x_2) \right] \cdots \right]$$
$$\mu_{\alpha}(x_n)$$
$$\left[\sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \cdots \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \cdots \right]$$
$$\mu_{\beta}(x_n)$$

The messages μ_{α} and μ_{β} can be computed recursively:

$$\mu_{\alpha}(x_{n}) = \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_{n}) \left[\sum_{x_{n-2}} \cdots \right]$$
$$= \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_{n}) \mu_{\alpha}(x_{n-1}).$$
$$\mu_{\beta}(x_{n}) = \sum_{x_{n+1}} \psi_{n,n+1}(x_{n}, x_{n+1}) \left[\sum_{x_{n+2}} \cdots \right]$$
$$= \sum_{x_{n+1}} \psi_{n,n+1}(x_{n}, x_{n+1}) \mu_{\beta}(x_{n+1}).$$

Computation of μ_{α} starts at the first node and computation of μ_{β} starts at the last node.

• 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 $p(x_n)$

To compute local marginals:

- •Compute and store all forward messages, $\mu_{\alpha}(x_n)$.
- •Compute and store all backward messages, $\mu_{\beta}(x_n)$
- •Compute Z once at a node x_m :

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

•Compute

$$p(x_n) = \frac{1}{Z} \mu_{\alpha}(x_n) \mu_{\beta}(x_n)$$

-1

for all variables required.

More General Graphs

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.

Factor Graphs

- 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

Factor Graphs

- 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.

Sum-Product Inference in General Graphical Models

- 1.Convert graph (directed or undirected) into a factor graph (there are no cycles)
- 2. If the goal is to **marginalize** at node *x*, then consider *x* as a root node
- **3.** Initialize the recursion at the leaf nodes as: $\mu_{f \to x}(x) = 1$ (var) or $\mu_{x \to f}(x) = f(x)$ (fac)
- **4.**Propagate messages from the leaves to *x*
- 5.Propagate messages from *x* to the leaves6.Obtain marginals at every node by multiplying all incoming messages

Further Topics on Graphical Models

Other inference algorithms:

- Max-Sum algorithm: used to maximize the joint probability of all variables (no marginalization)
- Junction Tree algorithm: exact inference for general graphs (even with loops)
- Loopy belief propagation: approximate inference on general graphs (more efficient)

Special kind of undirected GM:

Conditional Random fields (e.g.: classification)

More details: see class of Dr. Domokos http://vision.in.tum.de/teaching/ss2016/lecture_graphical_models

Summary

- Undirected models (aka Markov random fields) provide an intuitive representation of conditional independence
- An MRF is defined as a factorization over clique potentials and normalized globally
- Directed and undirected models have different representative power (no simple "containment")
- Inference on undirected Markov chains is efficient using message passing
- Factor graphs are more general; exact inference can be done efficiently using sum-product

