

Computer Vision Group Prof. Daniel Cremers

Technische Universität München

# 4. Probabilistic Graphical Models Directed Models

# **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

- ${\scriptstyle \bullet}$  Given: 3 random variables a , b , and
- Joint probp(a, b, c) = p(c|a, b)p(a, b) = p(c|a, b)p(b|a)p(a)



Random variables can be discrete or continuous

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



### Simple Example

In general: K random variables x<sub>1</sub>, x<sub>2</sub>,..., x<sub>K</sub>
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  $pa_k \triangleq \text{ancestors of } p_k$ and  $p(\mathbf{x}) = p(x_1, \dots, x_K)$ 



# **Elements of Graphical Models**

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





### **Elements of Graphical Models (2)**

 We distinguish between input variables and explicit hyperparameters:

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





### **Elements of Graphical Models (3)**

 We distinguish between observed variables and hidden variables:

 $p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w}) \prod p(t_n|\mathbf{w})$ (deterministic paran=1meters omitted)  $x_n$ w  $t_n$ 



#### **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  $\mu_1, \ldots, \mu_K$  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<sup>M</sup> -1 parameters

If the distribution can be described by this graph:



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.



### **Gaussian Variables**

Assume all random variables are Gaussian and we define

 $p(x_i \mid pa_i) = \mathcal{N}\left(x_i; \sum_{j \in pa_i} w_{ij}x_j + b_i, v_i\right)$ Then the joint probability p(x) is a multivariate Gaussian where

- Mean and covariance can be calculated recursively
- 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:  $x \perp \!\!\!\perp y \mid z$ 





• This graph represents the probability distribution:

p(a, b, c) = p(a|c)p(b|c)p(c)Marginalizing out *c* on both sides gives

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

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



 $\bullet$  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:  $a \not\perp b \mid \emptyset$ 



 ${\scriptstyle \bullet}$  As before, now we condition on c :



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





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





#### ullet Again, we condition on $_c$





### **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:



### **D-Separation**

•Say: A, B, and C are non-intersecting subsets of nodes in a directed graph. **D-Separation is a**  A path ; a node : property of graphs a)the r tail-toand not of tail a b)the neither probability the I et C. •If all p aid to distributions be d-s



### **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:



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

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. dist. 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

ПП

# **Hidden Markov Models**

# **Graphical Representation (Rep.)**

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



• This incorporates the following Markov assumptions:

$$p(z_t \mid x_{0:t}, u_{1:t}, z_{1:t}) = p(z_t \mid x_t) \text{ (measurement)}$$

$$p(x_t \mid x_{0:t-1}, u_{1:t}, z_{1:t}) = p(x_t \mid x_{t-1}, u_t) \text{ (state)}$$



# **Graphical Representation**

We can describe the overall process using a *Markov* chain of latent variables:



• This incorporates the following Markov assumptions:

Computer Vision

$$p(z_t \mid x_{0:t}, \qquad z_{1:t}) = p(z_t \mid x_t) \text{ (measurement)}$$

$$p(x_t \mid x_{0:t-1}, \qquad z_{1:t}) = p(x_t \mid x_{t-1} \quad ) \quad \text{(state)}$$
Machine Learning for Dr. Rudolph Triebel

Computer Vision Group

# Example

- Place recognition for mobile robots
- 3 different states: corridor, room, doorway
- Problem: misclassifications
- Idea: use information from previous time step





# Formulation as HMM

- 1.Discrete random variables
  - Observation variables:  $\{z_n\}, n = 1..N$
  - State variables (unobservable):  $\{x_n\}, n = 1..N$
  - Number of states *K*:  $x_n \in \{1..K\}$
- 2. Transition model  $p(x_i | x_{i-1})$ 
  - Markov assumption ( $x_i$  only depends on  $x_i$ )
  - Represented as a *K*×*K* transition matrix *A*
  - Initial probability:  $p(x_0)$  repr. as  $\pi_1, \pi_2, \pi_3$

3. Observation model  $p(z_i|x_i)$  with parameters  $\varphi$ 

- Observation only depends on the current state
- Example: output of a "local" place classifier



Model Parameters

θ

### **The Trellis Representation**





# **Application Example (1)**

- Given an observation sequence  $z_1, z_2, z_3...$
- Assume that the model parameters  $\theta = (A, \pi, \phi)$  are known
- What is the probability that the given observation sequence is actually observed under this model,
   i.e. *p*(*Z*| *θ*)?
- If we are given several different models, we can choose the one with highest probability
- Expressed as a supervised learning problem, this can be interpreted as the inference step (classification step)



# **Application Example (2)**

- Given an observation sequence  $z_1, z_2, z_3...$
- Assume that the model parameters  $\theta = (A, \pi, \varphi)$  are known
- What is the state sequence  $x_1, x_2, x_3...$  that explains best the given observation sequence?
- In the case of place recognition: which is the sequence of truly visited places that explains best the sequence of obtained place labels (classifications)?



# **Application Example (3)**

- Given an observation sequence  $z_1, z_2, z_3...$
- What are the optimal model parameters  $\theta = (A, \pi, \phi)$ ?
- This can be interpreted as the training step
- Is in general the most difficult problem



# Summary: 3 Operations on HMMs

- 1. Compute data likelihood  $p(Z|\theta)$  from a known model
  - Can be computed with the forward-backward algorithm
- 2. Compute optimal state sequence with a known model
  - Can be computed with the Viterbi-Algorithm
- 1. Learn model parameters for an observation sequence
  - Can be computed using Expectation-Maximization (or Baum-Welch)



# **1. Computing the Data Likelihood**

- Assume: given a state sequence  $x_1, x_2, x_3...$
- Compute  $p(Z|X, \theta) = \prod_{n=1}^{N} p(z_n|x_n, \theta)$
- The probability of this state sequence is  $p(X|\theta) = p(x_0) p(x_1|x_0) p(x_2|x_1)...$
- Thus, we have  $p(Z,X|\theta) = p(Z|X,\theta) p(X|\theta)$
- We need  $p(Z|\theta) = \prod_{all X} p(Z,X|\theta)$ , but this is intractable



• Define  $\alpha(x_n) = p(z_1, z_2, ..., z_n, x_n)$ 



- Define  $\alpha(x_n) = p(z_1, z_2, ..., z_n, x_n)$
- This can be recursively computed:

 $\alpha(x_n) = p(z_n | x_n) \sum_{x_{n-1}} \alpha(x_{n-1}) p(x_n | x_{n-1})$ 





• Define 
$$\alpha(x_n) = p(z_1, z_2, ..., z_n, x_n)$$

• This can be recursively computed:

$$\alpha(x_n) = p(z_n | x_n) \sum_{x_{n-1}} \alpha(x_{n-1}) p(x_n | x_{n-1})$$

• Then we have:

$$\mathbf{p}(Z|\theta) = \sum_{x_N} \alpha(x_N)$$

Solution to first problem!



• Define 
$$\alpha(x_n) = p(z_1, z_2, ..., z_n, x_n)$$

• This can be recursively computed:

$$\alpha(x_n) = p(z_n | x_n) \sum_{x_{n-1}} \alpha(x_{n-1}) p(x_n | x_{n-1})$$

• Then we have:

$$p(Z|\theta) = \sum_{x_N} \alpha(x_N)$$

- Similarly, we define  $\beta(x_n) = p(z_{n+1}, z_{n+2}, \dots, z_N | x_n)$
- This can also be recursively computed:

$$\beta(x_n) = \sum_{x_{n+1}} \beta(x_{n+1}) p(z_{n+1}|x_{n+1}) p(x_{n+1}|x_n)$$

Machine Learning for Computer Vision



this can be used to computey computed:



# 2. Computing the Most Likely States

• Goal: find a state sequence  $x_1, x_2, x_3$ ... that maximizes the probability  $p(X,Z|\theta)$ 

• Define 
$$\delta(x_n) = \max_{x_1, \dots, x_{n-1}} p(x_1, x_2, \dots, x_n, z_1, z_2, \dots, z_n)$$





# 2. Computing the Most Likely States

• Goal: find a state sequence  $x_1, x_2, x_3$ ... that maximizes the probability  $p(X,Z|\theta)$ 

• Define 
$$\delta(x_n) = \max_{x_1, \dots, x_{n-1}} p(x_1, x_2, \dots, x_n, z_1, z_2, \dots, z_n)$$

• This can be recursively computed:  $\delta(x_n) = p(z_n | x_n) \max_{x_{n-1}} [\delta(x_{n-1}) p(x_n | x_{n-1})]$ 



# 2. Computing the Most Likely States

• Goal: find a state sequence  $x_1, x_2, x_3$ ... that maximizes the probability  $p(X,Z|\theta)$ 

• Define 
$$\delta(x_n) = \max_{x_1, \dots, x_{n-1}} p(x_1, x_2, \dots, x_n, z_1, z_2, \dots, z_n)$$

- This can be recursively computed:  $\delta(x_n) = p(z_n | x_n) \max_{x_{n-1}} [\delta(x_{n-1}) p(x_n | x_{n-1})]$
- But we also need the argmax:  $\Psi(x_n) = \operatorname{argmax} \left[\delta(x_{n-1}) p(x_n | x_{n-1})\right]$



# The Viterbi algorithm

- Initialize:
  - $\delta(x_0) = p(x_0) p(z_0 | x_0)$
  - $\psi(x_0) = 0$
- Compute recursively for n=1...N:
  - $\delta(x_n) = p(z_n | x_{p_n}) \max [\delta(x_{n-1}) p(x_n | x_{n-1})]$
  - $\Psi(x_n)$  = argmax [ $\delta(x_{n-1}) p(x_n | x_{n-1})$ ]
- On termination<sup>\*</sup>:

• 
$$p(Z,X|\theta) = \max \delta(x_N)$$

•  $\mathbf{x}_N^* = \operatorname{argmax}^* \delta(x_N)$ 



# **3. Learning the Model Parameters**

- Given an observation sequence  $z_1, z_2, z_3$ ...
- Find optimal model parameters  $\theta$
- We need to maximize the likelihood  $p(Z|\theta)$
- Can not be solved in closed form
- Iterative algorithm: Expectation Maximization (EM) or for the case of HMMs: Baum-Welch algorithm



# **Expectation maximisation**

- Objective: Find the model parameters knowing the observations:  $\pi$ , A,  $\phi$
- Result:
  - Train the HMM to recognize sequences of input
  - Train the HMM to generate sequences of input
- Technique: Expectation Maximisation
  - E: Find the best state sequence given the parameters
  - M: Find the parameters using the state sequence
  - Maximisation of the log-likelihood:  $\arg \max_{pi,A,\phi} - \log \left( P\left(\{Z_i\}, \pi, A, \phi\right) \right)$



- E-Step (assuming we know  $\pi$ , A,  $\phi$ , i.e.  $\theta^{old}$ )
- Define the posterior probability of being in state i at step k:
- Define  $\gamma(x_n) = p(x_n|Z)$





- E-Step (assuming we know  $\pi$ , A,  $\phi$ , i.e.  $\theta^{old}$ )
- Define the posterior probability of being in state i at step k:
- Define  $\gamma(x_n) = p(x_n|Z)$
- It follows that  $\gamma(x_n) = \alpha(x_n) \beta(x_n) / p(Z)$



- E-Step (assuming we know  $\pi$ ,A, $\phi$ , i.e.  $\theta^{old}$ )
- Define the posterior probability of being in state i at step k:
- Define  $\gamma(x_n) = p(x_n|Z)$
- It follows that  $\gamma(x_n) = \alpha(x_n) \beta(x_n) / p(Z)$
- Define  $\xi(x_{n-1}, x_n) = p(x_{n-1}, x_n | Z)$
- It follows that  $\xi(x_{n-1}, x_n) = \alpha(x_{n-1})p(z_n|x_n)p(x_n|x_{n-1})\beta(x_n) / p(Z)$



- Maximizing Q also maximizes the likelihood:  $p(Z|\theta) \ge p(Z|\theta^{old})$
- M-Step:

• 
$$\pi_k = \gamma(x_1 = k) / \sum_j \gamma(x_1 = j)$$
  
here, we need forward and backward step!

• 
$$A_{ij} = \sum \xi(x_{n-1}=i, x_n=j) / \sum \xi(x_{n-1}=i, x_n=k)$$

- With these new values, Q is recomputed
- This is done until the likelihood does not increase anymore (convergence)



#### The Baum-Welsh algorithm - summary

- Start with an initial estimate of  $\theta = (\pi, A, \phi)$ e.g. uniformly and k-means for  $\phi$
- Compute Q(θ,θ<sup>old</sup>) (E-Step)
- Maximize Q (M-step)
- Iterate E and M until convergence
- In each iteration one full application of the forward-backward algorithm is performed
- Result gives a local optimum
- For other local optima, the algorithm needs to be started again with new initialization



# The Scaling problem

Probability of sequences

 $\prod_{i} P(X_i "...) = 1$ 

Probabilities are very small

- The product of the terms soon is very small
- Usually: convert to log-space works
- But: we have sums of products!
- Solution: Rescale/Normalise the probability during the computation, e.g.:

$$\alpha(x_n) = \alpha(x_n) / p(z_1, z_2, ..., z_n)$$



# Summary

- HMMs are a way to model sequential data
- They assume discrete states
- Three possible operations can be performed with HMMs:
  - Data likelihood, given a model and an observation
  - Most likely state sequence, given a model and an observation
  - Optimal Model parameters, given an observation
- Appropriate scaling solves numerical problems
- HMMs are widely used, e.g. in speech recognition

