

ECE 175B: Probabilistic Reasoning and Graphical Models

Lecture 4: Bayesian Belief Networks

Florian Meyer & Ken Kreutz-Delgado

*Electrical and Computer Engineering Department
University of California San Diego*

Factorizing Joint Distributions

- **Complexity of Joint Probability Distributions:** Let $\mathcal{X} = \{x_1, \dots, x_N\}$ be a set of N random variables with joint pmf $P(\mathcal{X}) = P(x_1, \dots, x_N)$
 - assume each random variable x_i is categorically taking $k_i \geq 2$ values; how many possible values does \mathcal{X} have?
 - imagine $k_i = 2$, and $N = 300$, then at least $2^{300} - 1 \approx 10^{90}$ values (v.s. number of atoms in the universe $10^{78} \sim 10^{82}$)
- To address this issue we exploit the structure of $P(\mathcal{X})$
 - factorize $P(\mathcal{X})$ based on the **product rule (a.k.a. chain rule)** and make use of the **conditional independencies** that exists among the random variables x_1, \dots, x_N
 - encode the resulting factorization of $P(\mathcal{X})$ by a directed acyclic graph (DAG) called a **Bayes Network** or a **Belief Network (BN)**

The Chain Rule

- We can factorize the joint probability distribution $P(\mathcal{X})$ as product of conditional distributions

$$\begin{aligned}P(\mathcal{X}) &= P(x_1, \dots, x_N) \\&= P(x_N | x_{N-1}, \dots, x_1) \cdots P(x_j | x_{j-1}, \dots, x_1) \cdots P(x_1) \\&= \prod_{j=1}^N P(x_j | x_{j-1}, \dots, x_1)\end{aligned}$$

- Note that there are $N!$ different ways to factorizing the joint distribution (each corresponding to a permutation of x_1, \dots, x_N)
- To simplify the notation we just relabel any permutation x_{i_1}, \dots, x_{i_N} as $x_1 \leftarrow x_{i_1}, \dots, x_N \leftarrow x_{i_N}$ such that we have ancestral ordering, i.e., the last distribution in the chain is always $P(x_1)$

Markov Property of Conditional Probability Distributions

- It is still equally challenging to represent the conditional probability distribution $P(x_j | x_{j-1}, \dots, x_1)$ for a large j
- **Question:** Do we really need to include the whole prior set $\mathbf{pr}(x_j) = \{x_{j-1}, \dots, x_1\}$?
- **Answer:** If there is conditional statistical independence, we only need to consider the smaller parents set of x_j , denoted as $\mathbf{pa}(x_j)$

Markov Property of Conditional Probability Distributions

- Define the parents of x_j , denoted $\mathbf{pa}(x_j)$, to be the **smallest subset** of $\mathbf{pr}(x_j)$, i.e.,

$$\mathbf{pa}(x_j) \subset \mathbf{pr}(x_j) \subset \mathcal{X}$$

such that

$$P(x_j | \mathbf{pa}(x_j)) = P(x_j | \mathbf{pr}(x_j))$$

- This implies that

$$x_j \perp\!\!\!\perp (\mathbf{pr}(x_j) \setminus \mathbf{pa}(x_j)) \mid \mathbf{pa}(x_j)$$

- Now we can write $P(\mathcal{X}) = \prod_{j=1}^N P(x_j | \mathbf{pa}(x_j))$ and encode this factorization by a DAG

Ancestral Order of a BN

- Once a BN $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is constructed corresponding to the factorization $P(\mathcal{X}) = \prod_{j=1}^N P(x_j | \mathbf{pa}(x_j))$, we can define \mathcal{X} -subsets of interest in terms of the corresponding subsets of \mathcal{V}

- $h \in \mathbf{anc}(j) \iff x_h \in \mathbf{anc}(x_j) \quad (h < j)$
- $i \in \mathbf{pa}(j) \iff x_i \in \mathbf{pa}(x_j) \quad (i < j)$
- $k \in \mathbf{child}(j) \iff x_k \in \mathbf{child}(x_j) \quad (k > j)$
- $l \in \mathbf{desc}(j) \iff x_l \in \mathbf{desc}(x_j) \quad (l > j)$

- Note that for x_h, x_i, x_k, x_l related as defined above, we have $h \leq i < k \leq l$, i.e., we have **ancestral ordering** of the nodes of the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$
- Note that $\mathbf{pa}(x_j) \subseteq \mathbf{anc}(x_j) \subseteq \mathbf{pr}(x_j) \subset \mathcal{X}$ (why?) and thus $x_j \perp\!\!\!\perp (\mathbf{anc}(x_j) - \mathbf{pa}(x_j)) \mid \mathbf{pa}(x_j)$

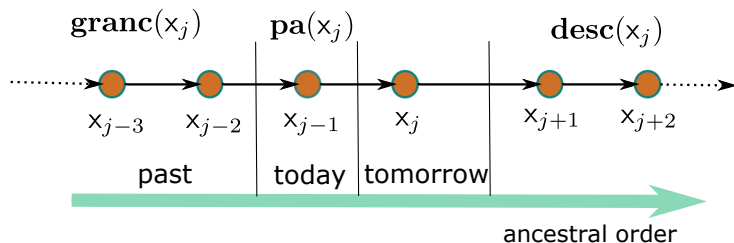
Ancestral Order of a BN (cont.)

- To make this clearer, define the “grand-ancestors” of x_j as $\mathbf{granc}(x_j) \triangleq \mathbf{anc}(x_j) \setminus \mathbf{pa}(x_j)$ with $\mathbf{anc}(x_j) = \mathbf{granc}(x_j) \cup \mathbf{pa}(x_j)$

- Then we get the **Basic Markov Property** as

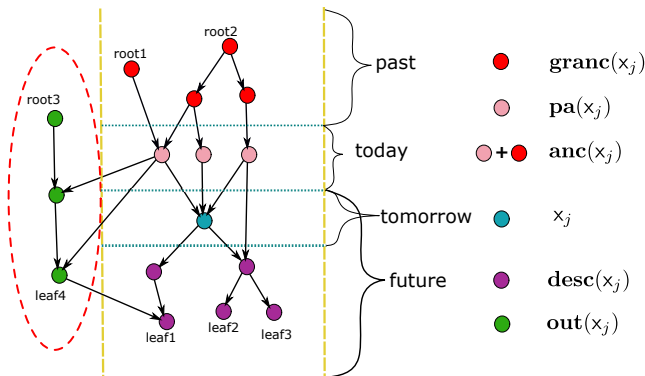
$$P(x_j | \mathbf{pa}(x_j), \mathbf{granc}(x_j)) \iff x_j \perp\!\!\!\perp \mathbf{granc}(x_j) \mid \mathbf{pa}(x_j)$$

- Example:** In a serial chain where the index $i \in \mathcal{V}$ represents time, “tomorrow” is independent of “past” given “today”



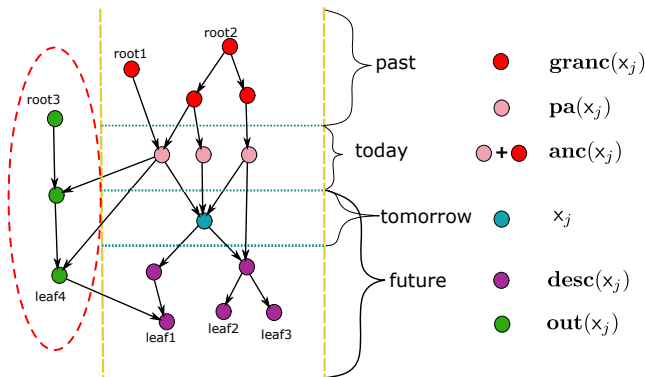
Ancestral Order of a BN (cont.)

- Note that the ancestor/descendent relationship related to x_j restricts a timeline w.r.t. x_j
 - all ancestors of x_j have at least one directed path from it to x_j
 - all descendants of x_j have at least one directed path from x_j to itself
- For nodes that do not share a directed path with x_j , we say they are “out of the timeline” of x_j



Ancestral Order of a BN (cont.)

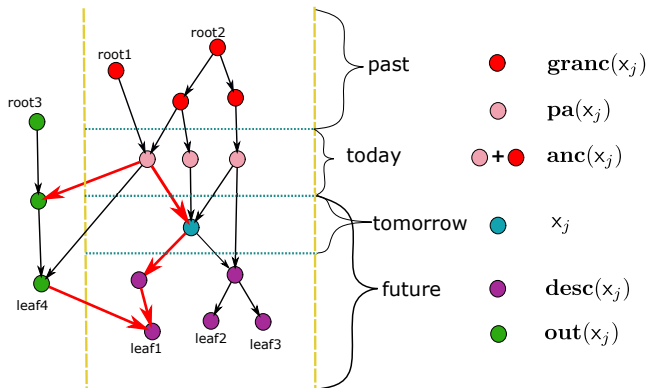
- All nodes that are out of the timeline of x_j are denoted as **out**(x_j)
- The ordering of $x \in \mathbf{out}(x_j)$ w.r.t x_j is unclear



- In a directed tree, every node is “in the timeline” of every other node

Ancestral Order of a BN (cont.)

- Note that in general for $x \in \text{out}(x_j)$, we have $\text{pa}(x) \cap \text{pa}(x_j) \neq \emptyset$
- Non-empty intersections also exists for **anc**, **child**, and **desc**



Finding the Best BN

- Recall that for a certain joint distribution of a set of N random variables
 - there are $N!$ rearrangements of the factorizing order
 - we can build a BN corresponding to a specific factorization
- **Questions:** Do all factorizations represent the same causal relationship between the random variables? Do all factorizations yield the same specification complexity of conditional probability distributions? If not, which factorization should we choose?