

Bayesian multimodeling: graphical models

MIPT

2023

Graphical models

Conditional independence

Events X, Y are conditionally independent w.r.t. Z : $X \perp Y|Z$, if

$$P(X|Y, Z) = P(X|Z).$$

Conditional dependence

Events X, Y are conditionally dependent w.r.t. \mathcal{G} : $X, Y \in \mathcal{G}$, if

$$X \not\perp Y | \mathcal{G} \setminus \{X, Y\}.$$

Graphical models

A probability model is graphical, if it can be represented as a graph, where the edges correspond to conditionally dependent events.

Non-graphical models

- MLP, decision trees, etc.
- Models with complex behaviour.

Types of graphical models

- Directed models (aka Bayesian networks)
 - ▶ Easy to design
- Undirected (Markov models)
- Factor-graphs
 - ▶ Easy to infer and optimize

Plate notation

Plate notation is an alternative visualization for graphical models.

Elements:

- White circles (random variables);
- Grey circles (observed variables);
- Small circles (deterministic values);
- Plates (batching).

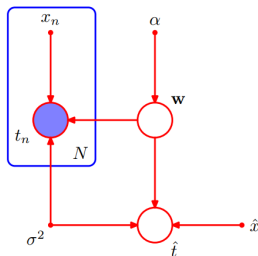


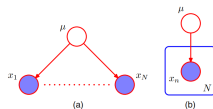
Plate notation for linear regression (Bishop)

Bayesian networks

- Models are set using directed acyclic graphs
- Joint distribution for the graph with K vertices:

$$p(v_1, \dots, v_K) = \prod_{i=1}^K p(v_i | \text{parent}(v_i))$$

- Example: linear regression



DAG and Plate notation (Bishop)

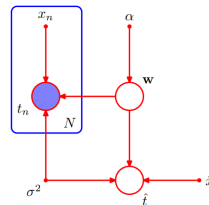


Plate notation for regression model (Bishop)

Causality graph elements

$$X \rightarrow Y \rightarrow Z - \text{chain}$$

Example:

- X — school budget
- Y — average student score
- Z — university acceptance ratio

Properties:

- ① X and Y , Y and Z are dependent:
 $\exists x, y : \mathbf{P}(Y = y | X = x) \neq p(Y = y)$
 $\exists y, z : \mathbf{P}(Z = z | Y = y) \neq p(Z = z)$
- ② Z and X : are (probably) dependent
- ③ $Z \perp X | Y$: are conditionally independent: $\forall x, y, z$

$$\mathbf{P}(Z = z | X = x, Y = y) = \mathbf{P}(Z = z | Y = y)$$

(if Y is fixed, then X and Z are independent)

Causality graph elements

$$X \leftarrow Y \rightarrow Z \text{ — fork}$$

Example:

- X — ice cream sells
- Y — average temperature
- Z — crime ratio

Properties:

- ① X and Y , Y and Z are dependent
- ② X and Z are (probably) dependent
- ③ $X \perp Z | Y$ are conditionally independent

Causality graph elements

$$Y \rightarrow X \leftarrow Z \text{ — collider}$$

Example (illness):

- X — bad symptoms
- Y — age
- Z — chronic diseases

Properties:

- ① Y and X , Z and X are dependent
- ② Y and Z are independent
- ③ $Y \not\perp Z | X$ are conditionally dependent

d-separation

The path P is **blocked** by Z , if:

- ① P contains $A \rightarrow B \rightarrow C$, $A \leftarrow B \rightarrow C$, $B \in Z$
- ② P contains $A \rightarrow B \leftarrow C$, $B \notin Z$ and all children of $B \notin Z$

If Z blocks all the paths from X to Y , then X and Y are **d-separated**:

$$X \perp Y | Z.$$

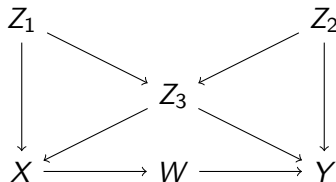
d-separation

The path P is blocked by Z , if:

- 1 P contains $A \rightarrow B \rightarrow C$, $A \leftarrow B \rightarrow C$, $B \in Z$
- 2 P contains $A \rightarrow B \leftarrow C$, $B \notin Z$ and all children of $B \notin Z$

If Z blocks all the paths from X to Y , then X and Y are d-separated.

Example:



| Pair | d-separation set |
|------------|------------------|
| (Z_1, W) | X |

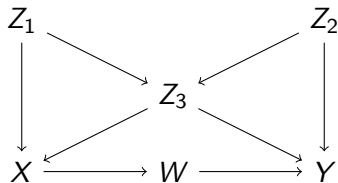
d-separation

The path P is blocked by Z , if:

- 1 P contains $A \rightarrow B \rightarrow C$, $A \leftarrow B \rightarrow C$, $B \in Z$
- 2 P contains $A \rightarrow B \leftarrow C$, $B \notin Z$ and all children of $B \notin Z$

If Z blocks all the paths from X to Y , then X and Y are d-separated.

Example:



| Pair | d-separation set |
|------------|------------------------------------|
| (Z_1, W) | X |
| (Z_1, Y) | $\{Z_3, X, Z_2\}, \{Z_3, W, Z_2\}$ |

Model selection for Bayesian networks

- Generally, NP-hard problem
- Reduces to optimization problem with predefined search space or sampling problem
- Independence determination:
 - ▶ ML and MAP
 - ▶ Evidence
 - ▶ Information criteria

Simple algorithm for predefined vertices

- ① $\forall A, B \in V$ search a set $S_{AB}: A \perp B | S_{AB}, A, B \notin S_{AB}$. If S_{AB} does not exist, make an edge AB .
- ② $\forall A, B$, not connected by edge and having a common neighbor C , check: $C \in S_{AB}$? If not, replace $A - C, C - B$ by $A \rightarrow C, C \leftarrow B$
- ③ Recursively:
 - ▶ if there is an oriented path from A to B $A \rightarrow \dots \rightarrow B$, then replace $A - B$ by $A \rightarrow B$;
 - ▶ if A and B are not connected, $A \rightarrow C, C - B$, then replace $C - B$ by $C \rightarrow B$.

Markov random fields

Models are represented as undirected graphs.

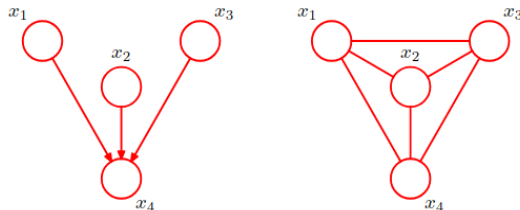
Difference from Bayesian networks:

- No direction \rightarrow cannot infer causality.
- The likelihood is factorized as follows:

$$p(\mathbf{x}) = \frac{1}{Z} \prod_C \psi(\mathbf{x}_C),$$

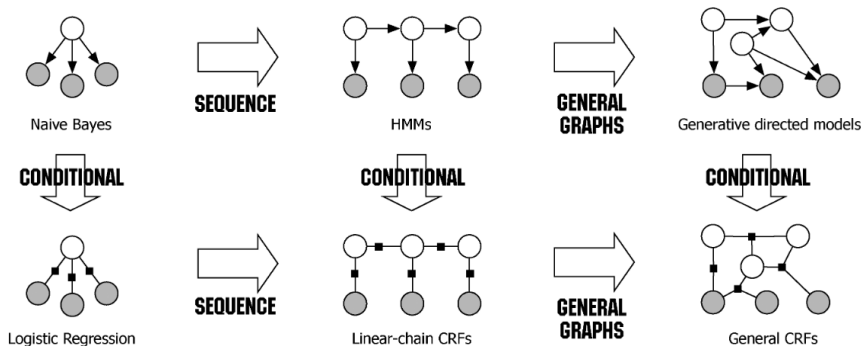
where \mathbf{x}_C is a maximal clique, $\psi \geq 0$ is a potential function.

- Conditional independence: if all the paths from A to B go through C , then $A \perp B | C$.



(Bishop)

Example: CRF and HMM



(Sutton, McCallum)

Inference in chains



(Bishop)

Naive likelihood calculation for x_n :

$$p(x_n) = \sum_{x_1} \sum_{x_2} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x}),$$

For N discrete variables with K values the complexity is $O(K^N)$

Inference in chains: regrouping

$$p(x_n) = \sum_{x_1} \sum_{x_2} \dots, \sum_{x_{n-1}} \sum_{x_{n+1}} \dots \sum_{x_N} p(\mathbf{x}),$$

$$p(\mathbf{x}) = \psi(x_1, x_2) \psi(x_2, x_3) \dots \psi(x_{N-1}, x_N).$$

Regroup the sum:

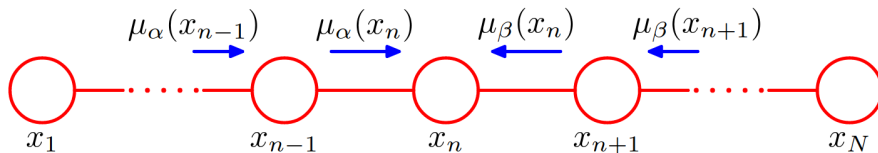
$$\begin{aligned} p(x_n) = & \sum_{x_{n-1}} \psi(x_{n-1}, x_n) \dots \left(\sum_{x_1} \psi(x_1, x_2) \right) \times \\ & \times \left(\sum_{x_{n+1}} \psi(x_n, x_{n+1}) \dots \left(\sum_{x_N} \psi(x_{N-1}, x_N) \right) \right). \end{aligned}$$

Now complexity is $O(NK^2)$.

Message passing

$$p(x_n) = \underbrace{\sum_{x_{n-1}} \psi(x_{n-1}, x_n) \dots \left(\sum_{x_1} \psi(x_1, x_2) \right)}_{\mu_a(x_n)} \times \underbrace{\left(\sum_{x_{n+1}} \psi(x_n, x_{n+1}) \dots \left(\sum_{x_N} \psi(x_{N-1}, x_N) \right) \right)}_{\mu_b(x_n)}.$$

Interpretation: $\mu_a(x_n)$ is a message transferred from x_{n-1} to x_n , $\mu_b(x_n)$ is a backward message from x_{n+1} .



Inference in chains: details

The inference is iterative:

- calculate $\sum_{x_1} \psi(x_1, x_2) = \mu_a(x_2)$, that stores $\mu_a(x_2)$ for each value of x_2 ;
- calculate $\sum_{x_2} \psi(x_2, x_3) (\sum_{x_1} \psi(x_1, x_2)) = \sum_{x_2} \psi(x_2, x_3) \mu_a(x_2) = \mu_a(x_3)$;
- ...
- calculate $\sum_{x_{n+1}} \psi(x_n, x_{n+1}) \mu_b(x_{n+1}) = \mu_b(x_n)$.
- for directed variables, where

$$\psi(x_1, x_2) = p(x_1)p(x_2|x_1), \quad \psi(x_i, x_{i+1}) = p(x_{i+1}|x_i),$$

μ_b should not be calculated:

$$\begin{aligned} \mu_b(x_n) &= \sum_{x_{n+1}} \psi(x_n, x_{n+1}) \dots \left(\sum_{x_N} \psi(x_{N-1}, x_N) \right) = \\ &= \sum_{x_{n+1}} p(x_{n+1}|x_n) \dots \left(\sum_{x_N} p(x_N|x_{N-1}) \right) = 1. \end{aligned}$$

Factor graph

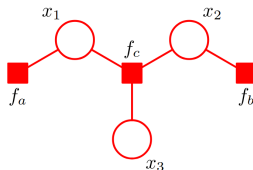
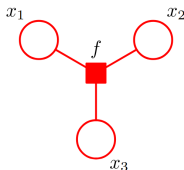
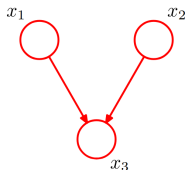
Definition

Factor-graph is a bipartite graph with two types of vertices: variables and factors.
The likelihood is a production of factors:

$$p(\mathbf{x}) = \prod_i f_i.$$

Example: model $p(x_1)p(x_2)p(x_3|x_2, x_1)$ has two variants of factorization:

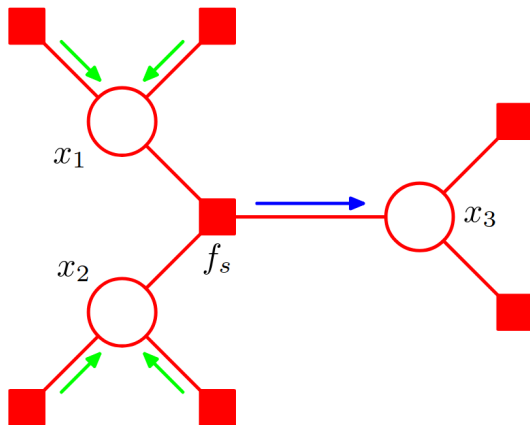
$$f = p(x_1)p(x_2)p(x_3|x_2, x_1), \quad f_a = p(x_1), f_b = p(x_2), f_c = p(x_1)p(x_2)p(x_3|x_2, x_1).$$



(Bishop)

Inference in factor-graphs: example

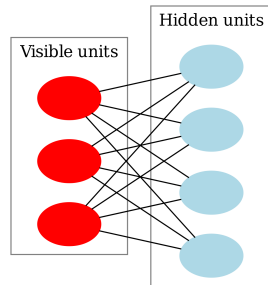
Sum-product: likelihood is a composition of messages from factors to variables.



Model examples: RBM

$$p(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{h})),$$

$$E = -\mathbf{w}_1^T \mathbf{x} - \mathbf{w}_2^T \mathbf{h} - \mathbf{x}^T \mathbf{W}_3 \mathbf{h}.$$



Model examples: Structured VAEs

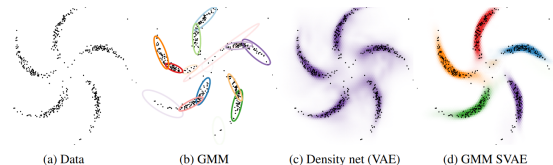
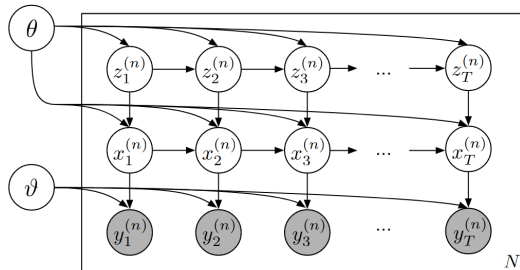
Based on SLDS:

$$z_{t+1}|z_t \sim \pi^{t+1},$$

$$\mathbf{y}_t \sim \mathcal{N}(\text{MLP}^{z_t}(\mathbf{x}_t)).$$

Optimization: optimize ELBO.

Inference: message-passing.



References

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – T. 128. – №. 9.
- Edwards D. Introduction to graphical modelling. – Springer Science & Business Media, 2012.
- Pearl J., Glymour M., Jewell N. P. Causal inference in statistics: A primer. – John Wiley & Sons, 2016.
- Hinton G. E., Salakhutdinov R. R. Reducing the dimensionality of data with neural networks //science. – 2006. – T. 313. – №. 5786. – C. 504-507.
- https://en.wikipedia.org/wiki/Restricted_Boltzmann_machine
- Johnson M. J. et al. Structured VAEs: Composing probabilistic graphical models and variational autoencoders //arXiv preprint arXiv:1603.06277. – 2016. – T. 2. – C. 2016.
- Johnson M. J. et al. Composing graphical models with neural networks for structured representations and fast inference //Advances in neural information processing systems. – 2016. – T. 29. – C. 2946-2954.
- Højsgaard, S., Edwards, D., Lauritzen, S. (2012). Graphical models with R. Springer Science & Business Media.
- Sutton, C, McCallum, A 2012, 'An Introduction to Conditional Random Fields', Foundations and Trends in Machine Learning, vol. 4, no. 4, pp. 267-373

$$\begin{aligned}
\nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[\int_z p_{\theta}(z) f_{\theta}(z) dz \right] \\
&= \int_z \nabla_{\theta} [p_{\theta}(z) f_{\theta}(z)] dz \\
&= \int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz + \int_z p_{\theta}(z) \nabla_{\theta} f_{\theta}(z) dz \\
&= \underbrace{\int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz}_{\text{What about this?}} + \mathbb{E}_{p_{\theta}(z)} [\nabla_{\theta} f_{\theta}(z)]
\end{aligned}$$

<https://gregorygundersen.com/blog/2018/04/29/reparameterization/>