

Advanced Data Analysis and Machine Learning

Lecture: Bayesian Networks

Lasse Lensu

2015-10-26

Outline

1 Introduction

2 Theory

Bayesian inference

- Simple Bayesian inference is unsuitable for complex modelling problems with a large number of variables.
- High dimensionality \Rightarrow a huge amount of data would be required to learn the joint probability distributions of several random variables.
- Independence of variables would be handy, but this is uncommon.
- In real-world modeling problems it is typical that the majority of the model variables depend on each other to some extent.

Conditional probability and chain rule

- **A priori probability** $P(A)$ is the probability of event A (similarly event B and its a priori probability $P(B)$).
- **Conditional probability** $P(B|A)$ means the probability of observing event B when A has happened and it can be defined

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad (1)$$

where $A \cap B$ says that both A and B happen.

- The chain rule follows from Eq. 1 directly

$$P(A \cap B) = P(A)P(B|A) \quad (2)$$

$$\begin{aligned} \Rightarrow P(A_1 \cap \dots \cap A_k) &= P(A_1)P(A_2|A_1) \dots \\ &\quad P(A_k|A_1 \cap \dots \cap A_{k-1}). \end{aligned} \quad (3)$$

Conditional probability and Bayes

- Bayes' rule is a consequence of the conditional probability and chain rule:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}. \quad (4)$$

- In the case of a background event C affecting the occurrence of other events, Bayes' rule can be written in a more general form

$$P(A|B \cap C) = \frac{P(B|A \cap C)P(A|C)}{P(B|C)}. \quad (5)$$

- (Bayes' rule allows to compute the conditional probability $P(A|B)$ from the “inverse” conditional probability $P(B|A)$.)

Random variables and distributions

- A random variable is defined by a function associating a value with each outcome of a random experiment.
- Depending on the experiment outcome space and the events a random variable can describe, the values arise from a continuous or discrete distribution.
- When several random variables are involved, the interesting problem is that what is the joint distribution which represents the probabilities of co-occurring events.
- The problem can be solved by marginalising the joint distribution as needed (for example, summing up the rows or columns of the joint distribution in the case of two random variables).

Independence and conditional independence



LUT
Lappeenranta
University of Technology

- An event A is independent of event B in P (can be denoted as $P \models (A \perp B)$), if

$$P(A|B) = P(A), \text{ or if} \quad (6)$$

$$P(B) = 0. \quad (7)$$

- Alternatively, a distribution P satisfies $(A \perp B)$ if and only if

$$P(A \cap B) = P(A)P(B) \quad (8)$$

$$\Rightarrow P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2) \cdots P(A_n). \quad (9)$$

Conditional independence in practice

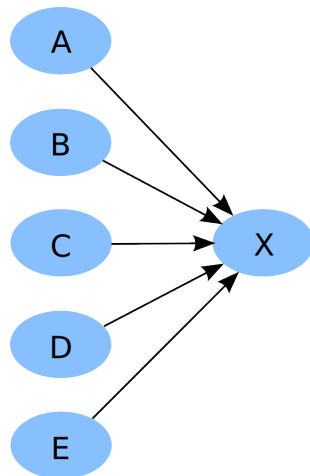
- Two events A and B are conditionally independent of a third event C if and only if, given knowledge of whether C occurs, knowledge of whether A occurs provides no information on the likelihood of B occurring and vice versa.
- Example:
 - John and Jeff usually use the same train to go to work.
 - Event A: “John is late”, B: “Jeff is late”, C: “the train is late”.
 - A and B are conditionally independent (given C) since if we already know that the train is late (C), knowing that John is late (A) does not give any new knowledge whether Jeff is late (B).
- Independence of variables would be handy, but this is uncommon. It is more common that two events are independent given a background event.

Bayesian networks

- A probabilistic graphical model.
- Representation of a set of random variables $\mathcal{X} = X_1, \dots, X_n$ as nodes and their conditional dependences as edges \mathcal{E} by using a directed acyclic graph (DAG) $\mathcal{G} = (\mathcal{X}, \mathcal{E})$.
- If there is a directed edge from X_i to X_j , X_i is X_j 's parent.

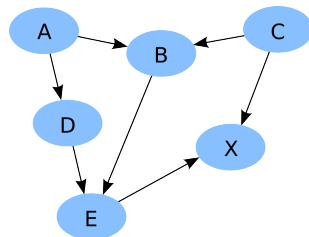
Bayesian networks

- Instead of defining $P(X|A, B, C, D, E) \dots$
- The network can be seen as i) a data structure representing a joint distribution in a factorised way or ii) a representation of the conditional independence assumptions of a distribution.
- Once network is constructed and parameters estimated, it can be used to answer queries related to the joint distribution, for example, compute the posterior probability distribution for any node.

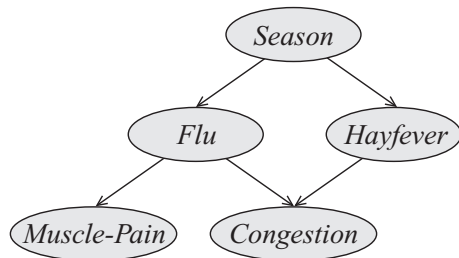


Bayesian networks

- Instead of defining $P(X|A, B, C, D, E)$ we need $P(A), P(B|A, C), P(C), P(D|A), P(E|B, D), P(X|C, E)$.
- The network can be seen as i) a data structure representing a joint distribution in a factorised way or ii) a representation of the conditional independence assumptions of a distribution.
- Once network is constructed and parameters estimated, it can be used to answer queries related to the joint distribution, for example, compute the posterior probability distribution for any node.

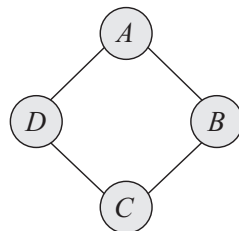


Bayesian and Markov networks



Independences

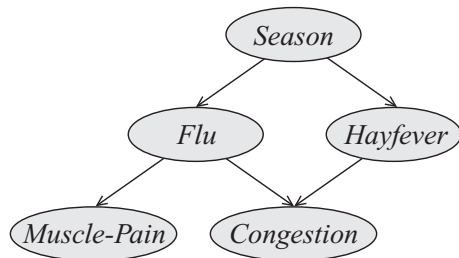
$$\begin{aligned}
 &(F \perp H | S) \\
 &(C \perp S | F, H) \\
 &(M \perp H, C | F) \\
 &(M \perp C | F)
 \end{aligned}$$



Independences

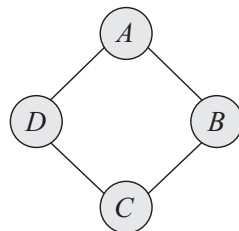
$$\begin{aligned}
 &(A \perp C | B, D) \\
 &(B \perp D | A, C)
 \end{aligned}$$

Bayesian and Markov networks



Factorisation

$$P(S, F, H, C, M) = P(S)P(F|S) \\ P(H|S)P(C|F, H)P(M|F)$$



Factorisation

$$P(A, B, C, D) = \frac{1}{Z} \phi_1(A, B) \\ \phi_2(B, C) \phi_3(C, D) \phi_4(A, D)$$

Network interpretation

■ Causality

- Arrows and their directionality implies causal relationship.
- However, this does not need to be the case.

Network interpretation



LUT
Lappeenranta
University of Technology

■ Causality

- Arrows and their directionality implies causal relationship.
- However, this does not need to be the case.



Inference

- Answering queries of interest based on a joint probability distribution over multiple random variables.
- Inference can be realised by using the following $[\mathcal{X} = X_1, \dots, X_n]$ is a set of random variables whose joint distribution is $P(X_1, \dots, X_n)$]:
 - Computing the marginal distribution of one or more nodes $\mathbf{Y} \subset \mathcal{X}$ (model output) based on evidence $\mathbf{E} \subset \mathcal{X}$ (model inputs): $P(\mathbf{Y}|\mathbf{E} = \mathbf{e})$, that is, the posterior probability distribution of $\mathbf{y} \in \mathbf{Y}$ conditioned on the fact that $\mathbf{E} = \mathbf{e}$.
 - Finding the most probable explanations, that is the maximum a posteriori (MAP) assignment of the values for several nodes:

$$\text{MAP}(\mathbf{W}|\mathbf{e}) = \arg \max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) \quad (10)$$

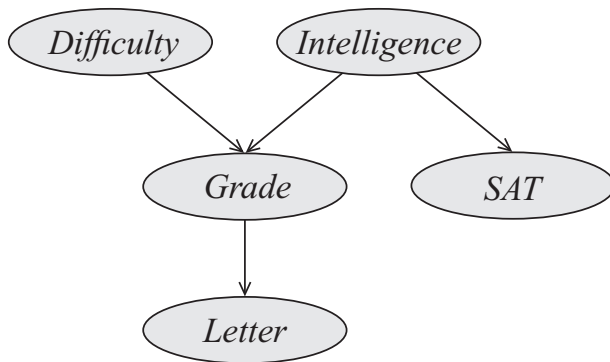
where $\mathbf{W} = \mathcal{X} - \mathbf{E}$.

Learning Bayesian networks

- Parameter learning
 - Learning probability distributions for nodes conditional upon their parents.
 - Efficient algorithms do exist.
- Structure learning
 - Learning conditional dependencies between variables.
 - NP-complete problem.

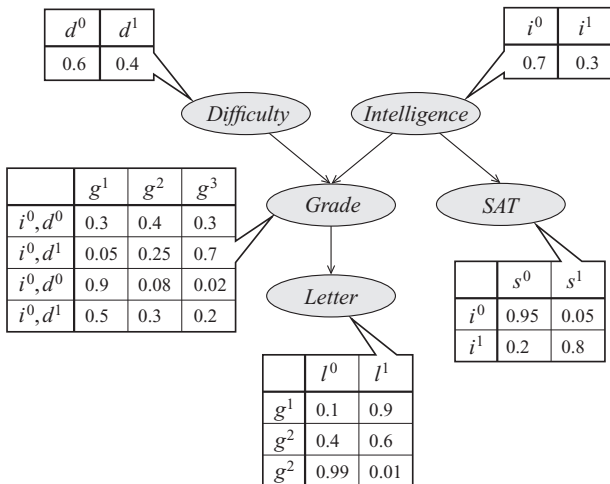
Bayesian network example

- Student example [1] :



Bayesian network example

- Student example [1] (Can you find the misspelled symbols?):



Summary

- Simple Bayesian inference is unsuitable for complex modelling problems with a large number of variables.
- Independence of variables would be handy, but in real-world modeling problems it is typical that the model variables depend on each other to some extent.
- Bayesian networks are probabilistic graphical models that represent the random variable dependencies (and independences) as a directed acyclic graph (DAG).
- The conditional joint distributions can be parameterised and learnt to perform inference.

References



Daphne Koller and Nir Friedman.

Probabilistic Graphical Models: Principles and Techniques.

The MIT Press, Cambridge, Massachusetts, 2009.