

Graphical probabilistic models – learning from data

Jiří Kléma

Department of Cybernetics,
FEE, CTU at Prague



<http://ida.felk.cvut.cz>

Parameter learning from incomplete data – example

- consider a linear connection $A \rightarrow B \rightarrow C$,
- learn network parameters, the samples shown in the table below are available,
- use the EM algorithm to learn with missing values (?).

	s_1	s_2	s_3	s_4
A	F	T	T	T
B	T	F	T	?
C	T	F	T	F

Parameter learning from incomplete data – example

- consider a linear connection $A \rightarrow B \rightarrow C$,
- learn network parameters, the samples shown in the table below are available,
- use the EM algorithm to learn with missing values (?).

	s_1	s_2	s_3	s_4
A	F	T	T	T
B	T	F	T	?
C	T	F	T	F

init: $Pr(a) = \frac{3}{4}$, $Pr(b|a) = \frac{1}{2}$, $Pr(b|\neg a) = 1$, $Pr(c|b) = 1$, $Pr(c|\neg b) = 0$,

E_1 : $Pr(B_4 = T) = Pr(b|a, \neg c) = \frac{Pr(a,b,\neg c)}{Pr(a,\neg c)} = \frac{\frac{3}{4}\frac{1}{2}0}{\frac{3}{4}\frac{1}{2}0 + \frac{3}{4}\frac{1}{2}1} = 0 \rightarrow$ estimated F,

M_1 : $Pr(a) = \frac{3}{4}$, $Pr(b|a) = \frac{1}{3}$, $Pr(b|\neg a) = 1$, $Pr(c|b) = 1$, $Pr(c|\neg b) = 0$,

E_2 : $Pr(B_4 = T) = Pr(b|a, \neg c) = \frac{Pr(a,b,\neg c)}{Pr(a,\neg c)} = \frac{\frac{3}{4}\frac{1}{3}0}{\frac{3}{4}\frac{1}{3}0 + \frac{3}{4}\frac{2}{3}1} = 0 \rightarrow$ estimated F,

M_2 : necessarily the same result as in M_1 , converged, STOP.

Conditional entropy

- information entropy $H(X)$

- a measure of the uncertainty in a random variable,
- the average number of bits per value needed to encode it,
- $H(X) = - \sum_{x \in X} Pr(x) \log_2 Pr(x)$

- conditional (information) entropy $H(Y|X)$

- uncertainty in a random variable Y given that the value of random variable X is known,
- $X \perp\!\!\!\perp Y \Rightarrow H(Y|X) = H(Y)$
- $H(Y|X) = \sum_{x \in X} Pr(x) H(Y|x) = - \sum_{x \in X} Pr(x) \sum_{y \in Y} Pr(y|x) \log_2 Pr(y|x)$

- how to enumerate conditional entropy?

- N_{ij} ... the number of samples, where $parents(P_i)$ take the j -th instantiation of values,
- N_{ijk} ... the number of samples, where P_i takes the k -th value and $parents(P_i)$ the j -th instantiation of values,

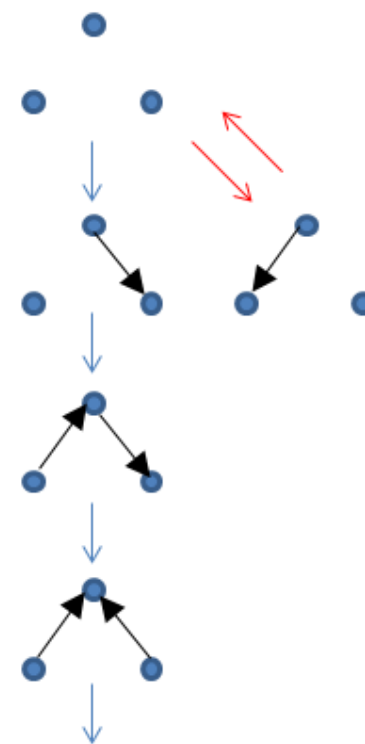
$$H(P_i | parents(P_i)^G) = - \sum_{j=1}^{q_i} \sum_{k=1}^{r_i} \frac{N_{ij}}{M} \frac{N_{ijk}}{N_{ij}} \log_2 \frac{N_{ijk}}{N_{ij}} = - \sum_{j=1}^{q_i} \sum_{k=1}^{r_i} \frac{N_{ijk}}{M} \log_2 \frac{N_{ijk}}{N_{ij}}$$

Score-based structure learning

- however, no evaluation function can be applied to all 2^{n^2} candidate graphs,
- heuristics and metaheuristics known for difficult tasks need to be employed
 - metaheuristic example – **local search**
 - * it starts with a given network (empty, expert's, random),
 - * it constructs all the “near” networks, evaluates them and goes to the best of them,
 - * it repeats the previous step if the local change increases score, otherwise it stops,
 - auxiliary heuristics examples
 - * definition of “near” network,
 - * how to avoid getting stuck in local minima or on plateaux
 - random restarts, simulated annealing, TABU search.

Structure learning – MCMC approach

- **MCMC** = Markov chain Monte-Carlo (for meaning see Gibbs sampling),
- applies **Metropolis-Hastings** (MH) algorithm to search the candidate graph/network space
 1. take an initial graph G
 - user-defined/informed, random, empty with no edges,
 2. evaluate the graph $P(G)$
 - use samples, apply criteria such as BIC or Bayesian,
 3. generate a “neighbor” S of the given graph G
 - insert/remove an edge, change edge direction,
 - check the graph acyclicity constraint,
 - prob of transition from G to S is function of $Q(G, S)$,
 4. evaluate the neighbor graph $P(S)$,
 5. accept or reject the transition to S
 - generate α from $U(0,1)$ (uniform distribution),
 - if $\alpha < \frac{P(S)Q(G,S)}{P(G)Q(S,G)}$ then accept the transition $G \rightarrow S$,
 6. repeat steps 3–5 until convergence or the given number of iterations.



Recommended reading, lecture resources

- Murphy: **A Brief Introduction to Graphical Models and Bayesian Networks.**
 - a practical overview from the author of BN toolbox,
 - <http://www.cs.ubc.ca/~murphyk/Bayes/bayes.html#learn>,
- Friedman, Koller: **Learning Bayesian Networks from Data.**
 - Neural Information Processing Systems conference tutorial, a presentation,
 - <http://www.cs.huji.ac.il/~nirf/Nips01-Tutorial/>,
- Cooper, Herskovits: **A Bayesian Method for the Induction of P.Networks from Data.**
 - theory + K2 algorithm,
 - www.genetics.ucla.edu/labs/sabatti/Stat180/bayesNet.pdf,
- Heckerman: **A Tutorial on Learning With Bayesian Networks.**
 - a theoretical paper, “easy to read”
 - research.microsoft.com/apps/pubs/default.aspx?id=69588,
- Buntine: **Operations for Learning with Graphical Models.**
 - a general, complete and extensive description,
 - <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.52.696&rep=rep1&type=pdf>.

