Chapter 17 Learning: The Maximum Likelihood Approach

Adnan Darwiche¹

¹Lecture slides for *Modeling and Reasoning with Bayesian Networks*, Adnan Darwiche, Cambridge University Press, 2009.

We will discuss in this chapter the process of learning Bayesian networks from data. The learning process will be studied under different conditions, which relate to the nature of available data and the amount of prior knowledge we have on the Bayesian network.

Introduction



Case	Cold?	Flu?	Tonsillitis?	Chilling?	Bodyache?	Sorethroat?	Fever?
1	true	false	?	true	false	false	false
2	false	true	false	true	true	false	true
3	?	?	true	false	?	true	false
•	•	•	•	•	•	•	•

Terminology

data (complete, incomplete); sample; case/observation/unit.

Adnan Darwiche Chapter 17 Learning: The Maximum Likelihood Approach

A key objective of this chapter is to provide techniques for:

- Estimating the parameters of a network structure given both complete and incomplete data sets.
- Learning the network structure itself, although our focus here will be on complete data sets for reasons that we state later.

One can distinguish between three general approaches to the learning problem.

The first approach is based on the likelihood principle

which favors those estimates that have a maximal likelihood, i.e., ones that maximize the probability of observing the given data set. This approach is therefore known as the maximum likelihood approach to learning.

This is the approach treated in this chapter.

向下 イヨト イヨト

The second approach requires more input to the learning process

as it demands one to define a meta distribution over network structures and parameters. It then reduces the problem of learning to a problem of classical inference in which the data set is viewed as evidence. In particular, it first conditions the meta distribution on the given data set, and then uses the posterior meta distribution as a criterion for defining estimates.

This approach is known as the Bayesian approach to learning and will be treated in Chapter 18.

向下 イヨト イヨト

A third approach for learning Bayesian networks is known as the constraint-based approach

and applies mostly to learning network structures. According to this approach, one seeks structures that respect the conditional independencies exhibited by the given data set.

We do not treat this approach in this book, but provide some references in the bibliographical remarks section.

Previous approaches induce Bayesian networks that are meaningful independent of the tasks for which they are intended.

Consider a network which models diseases and symptoms:

- Diagnostic tasks: inferring most likely disease given symptoms.
- Prediction tasks: inferring most likely symptom given diseases.
- If we focus on diagnostics, we can use a more specialized learning principle that optimizes the diagnostic performance of the learned network.
- This leads to a discriminative model: used to discriminate among patients according to a predefined set of classes (e.g., has cancer or not).
- A generative model: evaluated based on its ability to generate the given data set, regardless of how it performs on any particular task.

We do not cover discriminative approaches to learning in this book.



(a) network structure

(b) complete data

Т Т Т 2/16 F 0/16 Т Т т F Т 9/16 Т F F 1/16F Т Т 0/16 F Т F 1/16F F Т 2/16 F F F 1/16

Е

H S

 $\Pr_{\mathcal{D}}(.)$

(c) empirical distribution

э

Assumption

Data simulated from the true Bayesian network: the cases generated independently according to their true probabilities.

Empirical distribution summarizes data set.

Н	S	Ε	$\Pr_{\mathcal{D}}(.)$
Т	Т	Т	2/16
Т	Т	F	0/16
Т	F	Т	9/16
Т	F	F	1/16
F	Т	Т	0/16
F	Т	F	1/16
F	F	Т	2/16
F	F	F	1/16

The empirical probability of instantiation h, s, e

is its frequency of occurrence in the data set:

$$\operatorname{Pr}_{\mathcal{D}}(h, s, e) = \frac{\mathcal{D}\#(h, s, e)}{N},$$

where $\mathfrak{D}\#(h, s, e)$ is the number of cases in the data set \mathfrak{D} that satisfy instantiation h, s, e, and N is the data set size.

Estimate parameters based on the empirical distribution

Consider the parameter $\theta_{s|h}$ for example, which corresponds to the probability that a person will smoke given that they are health aware, $\Pr(s|h)$. Our estimate for this parameter is now given by:

$$\Pr_{\mathcal{D}}(s|h) = \frac{\Pr_{\mathcal{D}}(s,h)}{\Pr_{\mathcal{D}}(h)} = \frac{2/16}{12/16} = 1/6$$

This corresponds to the simplest estimation technique we discussed in Chapter 15: the method of direct sampling.

・同ト くきト くきり

Basic definitions

A data set \mathcal{D} for variables **X** is a vector $\mathbf{d}_1, \ldots, \mathbf{d}_N$, where each \mathbf{d}_i is called a case and represents a partial instantiation of variables **X**. The data set is complete if each case is a complete instantiation of variables **X**; otherwise, the data set is incomplete. The empirical distribution for a complete data set \mathcal{D} is defined as follows:

$$\Pr_{\mathcal{D}}(\alpha) \stackrel{\text{def}}{=} \frac{\mathcal{D}\#(\alpha)}{N},$$

where $\mathcal{D}\#(\alpha)$ is the number of cases \mathbf{d}_i in the data set \mathcal{D} that satisfy event α , that is, $\mathbf{d}_i \models \alpha$.

 $\mathfrak{D}\#(\alpha) = N$ when α is a valid event (α satisfied by every case \mathbf{d}_i)

・ 同 ト ・ ヨ ト ・ ヨ ト

Case	Н	S	Е
1	Т	F	Т
2	Т	F	Т
3	F	Т	F
4	F	F	Т
5	Т	F	F
6	Т	F	Т
7	F	F	F
8	Т	F	Т
9	Т	F	Т
10	F	F	Т
11	Т	F	Т
12	Т	Т	Т
13	Т	F	Т
14	Т	Т	Т
15	Т	F	Т
16	T	F	Т

$$\begin{aligned} \mathcal{D}\#(\alpha) &= 9, \text{ when } \alpha \text{ is } (H=T) \land (S=F) \land (E=T); \\ \mathcal{D}\#(\alpha) &= 12, \text{ when } \alpha \text{ is } (H=T); \\ \mathcal{D}\#(\alpha) &= 14, \text{ when } \alpha \text{ is } (H=T) \lor (E=T). \end{aligned}$$

< ≣ →

æ

We estimate the parameter $\theta_{x|\mathbf{u}}$ by the empirical probability

$$\theta_{x|\mathbf{u}}^{ml} \stackrel{def}{=} \Pr_{\mathcal{D}}(x|\mathbf{u}) = \frac{\mathcal{D}\#(x,\mathbf{u})}{\mathcal{D}\#(\mathbf{u})}$$

The count $\mathcal{D}\#(x, \mathbf{u})$ is called a sufficient statistic in this case.

More generally though, any function of the data is called a statistic. Moreover, a sufficient statistic is a statistic that contains all of the information in the data set that is needed for a particular estimation task.



(a) network structure

(b) complete data

Т Т Т 2/16 F 0/16 Т Т т F Т 9/16 Т F F 1/16F Т Т 0/16 F Т F 1/16F F Т 2/16 F F F 1/16

Е

H S

 $\Pr_{\mathcal{D}}(.)$

(c) empirical distribution

э



We have the following parameter estimates:

	Н	S	$\theta_{S H}^{ml}$	Н	Ε	$\theta_{E H}^{ml}$
$H \mid \theta_H^{ml}$	h	5	1/6	h	е	11/12
h 3/4	h	5	5/6	h	ē	1/12
\bar{h} 1/4	\bar{h}	5	1/4	\bar{h}	е	1/2
	ĥ	5	3/4	ĥ	ē	1/2

4 3 4 3 4 3 4

- Estimate $\theta_{x|\mathbf{u}}^{ml}$ will have different values depending on the given data set.
- The variance of this estimate will decrease as the data set increases in size.

If data set ${\mathcal D}$ is a sample of size N simulated from distribution \Pr

The distribution of estimate $\theta_{x|\mathbf{u}}^{ml}$ is asymptotically Normal and can be approximated by a Normal distribution with mean $\Pr(x|\mathbf{u})$ and variance:

$$\frac{\Pr(x|\mathbf{u})(1 - \Pr(x|\mathbf{u}))}{N \cdot \Pr(\mathbf{u})}$$

If data set ${\mathfrak D}$ is a sample of size N simulated from distribution \Pr

The distribution of estimate $\theta_{x|\mathbf{u}}^{ml}$ is asymptotically Normal and can be approximated by a Normal distribution with mean $\Pr(x|\mathbf{u})$ and variance:

$$\frac{\Pr(x|\mathbf{u})(1 - \Pr(x|\mathbf{u}))}{N \cdot \Pr(\mathbf{u})}$$

If probability $Pr(\mathbf{u})$ is too small, and the data set is not large enough, it is not uncommon for the empirical probability $Pr_{\mathcal{D}}(\mathbf{u})$ to be zero. Under these conditions, the estimate $Pr_{\mathcal{D}}(x|\mathbf{u})$ is not well defined, leading to what is known as the problem of zero counts.

Likelihood of parameter estimates

Let θ be the set of all parameter estimates for network structure G, and let $Pr_{\theta}(.)$ be the probability distribution induced by structure G and estimates θ . The likelihood of these estimates is:

$$L(\theta|\mathcal{D}) \stackrel{def}{=} \prod_{i=1}^{N} Pr_{\theta}(\mathbf{d}_{i})$$

Likelihood of estimates θ is the probability of observing the data set $\mathcal D$ under these estimates.

Let ${\mathfrak D}$ be a complete data set

The parameter estimates defined earlier are the only estimates that maximize the likelihood function:^{*a*}

$$heta^{\star} = rgmax \operatorname{L}(heta | \mathcal{D}) ext{ iff } heta^{\star}_{x \mid \mathsf{u}} = \operatorname{Pr}_{\mathcal{D}}(x \mid \mathsf{u})$$

aAssumes $\Pr_{\mathfrak{D}}(u)>0$ for every instantiation u of every parent set U

It is for this reason that these estimates are called maximum likelihood (ML) estimates and are denoted by θ^{ml} :

$$\theta^{ml} = \operatorname*{argmax}_{\theta} \operatorname{L}(\theta | \mathfrak{D})$$

伺下 イヨト イヨト

We defined these estimates based on the empirical distribution and then showed that they maximize the likelihood function.

Yet, it is quite common to start with the goal of maximizing the likelihood function and then derive these estimates accordingly.

This alternative approach is justified by some strong, desirable, properties that are satisfied by estimates that maximize the likelihood function.

We will indeed follow this approach when dealing with incomplete data in the next section.

伺い イヨト イヨト

Another property of our ML estimates is that they minimize the KL–divergence between the learned Bayesian network and the empirical distribution.

Let
$$\mathcal{D}$$
 be a complete data set over variables X
 $\operatorname{argmax}_{\theta} L(\theta | \mathcal{D}) = \operatorname{argmin}_{\theta} KL(\operatorname{Pr}_{\mathcal{D}}(X), \operatorname{Pr}_{\theta}(X))$

Since ML estimates are unique for a given structure G and complete data set \mathcal{D}

the likelihood of these parameters is then a function of the structure ${\it G}$ and data set ${\cal D}$

We will therefore define the likelihood of structure G given data set \mathcal{D} as follows:

$$L(G|\mathcal{D}) \stackrel{def}{=} L(\theta^{ml}|\mathcal{D}),$$

where θ^{ml} are the ML estimates for structure G and data set ${\cal D}$

More convenient to work with the logarithm of likelihood

$$\operatorname{LL}(\theta | \mathcal{D}) \stackrel{def}{=} \log \operatorname{L}(\theta | \mathcal{D}) = \sum_{i=1}^{N} \log \operatorname{Pr}_{\theta}(\mathbf{d}_i)$$

The log-likelihood of structure G is defined similarly:

$$LL(G|\mathcal{D}) \stackrel{def}{=} \log L(G|\mathcal{D})$$

Likelihood is ≥ 0 while log-likelihood is ≤ 0

Maximizing likelihood is equivalent to maximizing log-likelihood.

We will use \log_2 but write \log_2 .

・ 同下 ・ ヨト ・ ヨト

Log-likelihood decomposes into family-based components

Let G be a network structure and \mathfrak{D} be a complete data set of size N. If XU ranges over the families of structure G, then

$$LL(G|\mathcal{D}) = -N \sum_{X\mathbf{U}} ENT_{\mathcal{D}}(X|\mathbf{U}),$$

where $ENT_{\mathcal{D}}(X|\mathbf{U})$ is the conditional entropy defined as follows:

$$\operatorname{ENT}_{\mathcal{D}}(X|\mathbf{U}) = -\sum_{x\mathbf{u}} \operatorname{Pr}_{\mathcal{D}}(x\mathbf{u}) \log_2 \operatorname{Pr}_{\mathcal{D}}(x|\mathbf{u})$$

Decomposition is critical when learning network structure.

伺下 イヨト イヨト

The parameter estimates we considered are unique, asymptotically Normal, and maximize the probability of data. These estimates are easily computable by performing a single pass on the data set.

We have proven some of these properties independently, yet some of them follow from the others under more general conditions.

For example, maximum likelihood estimates are known to be asymptotically Normal for a large class of models that include but is not limited to Bayesian networks.

It is therefore common to seek maximum likelihood estimates for incomplete data sets. The properties of these estimates, however, will depend on the nature of incompleteness we have.

Consider a network structure $C \rightarrow T$, where C represents a medical condition and T represents a test for detecting this condition:

Note: Pr(T = +ve) = Pr(T = -ve) = 1/2

Consider now the following data sets, all of which are incomplete:

\mathcal{D}^1	C	Т	\mathfrak{D}^2	С	Т	\mathcal{D}^3	С	Т
1	?	+ve	1	yes	+ve	1	yes	+ve
2	?	+ve	2	yes	+ve	2	yes	+ve
3	?	–ve	3	yes	–ve	3	?	-ve
4	?	–ve	4	no	?	4	no	?
5	?	—ve	5	yes	-ve	5	yes	—ve
6	?	+ve	6	yes	+ve	6	?	+ve
7	?	+ve	7	no	?	7	no	?
8	?	–ve	8	no	-ve	8	no	–ve

\mathfrak{D}^1	С	Т
1	?	+ve
2	?	+ve
3	?	–ve
4	?	-ve
5	?	–ve
6	?	+ve
7	?	+ve
8	?	–ve

Values of variable *C* are missing in all cases of the first data set, perhaps because we can never determine this condition directly. We will say in this situation that variable *C* is hidden or latent.

\mathfrak{D}^2	С	Т
1	yes	+ve
2	yes	+ve
3	yes	–ve
4	no	?
5	yes	–ve
6	yes	+ve
7	no	?
8	no	–ve

Variable C is always observed, while variable T has some missing values, but is not hidden.

\mathcal{D}^3	С	Т
1	yes	+ve
2	yes	+ve
3	?	–ve
4	no	?
5	yes	–ve
6	?	+ve
7	no	?
8	no	–ve

Both variables have some missing values, but neither is hidden.

\mathfrak{D}^1	С	Т
1	?	+ve
2	?	+ve
3	?	—ve
4	?	—ve
5	?	–ve
6	?	+ve
7	?	+ve
8	?	-ve

Cases are split equally between the +ve and -ve values of T. We expect this to be true in the limit, given the distribution generating this data.

ML estimates are characterized by

$$\theta_{T=+ve|C=yes} \cdot \theta_{C=yes} + \theta_{T=+ve|C=no} \cdot \theta_{C=no} = \frac{1}{2}$$

ML estimates are characterized by

$$\theta_{T=+ve|C=ves} \cdot \theta_{C=ves} + \theta_{T=+ve|C=vo} \cdot \theta_{C=vo} = \frac{1}{2}$$

The true parameter values satisfy the above equation. But the following estimates do as well:

$$\theta_{\textit{C}=\textit{yes}} = 1, \quad \theta_{\textit{T}=+\textit{ve}|\textit{C}=\textit{yes}} = 1/2,$$

with $\theta_{T=+ve|C=no}$ taking any value.

ML estimates are not unique.

\mathfrak{D}^2	С	Т
1	yes	+ve
2	yes	+ve
3	yes	-ve
4	no	?
5	yes	-ve
6	yes	+ve
7	no	?
8	no	—ve

Consider the following two scenarios:

People who do not suffer from the condition tend not to take the test. That is, the data is missing because the test is not performed to start with.

People who test negative tend not to report the result. That is, the test is performed, but its value is not recorded.

In the second scenario, the fact that a value is missing does provide some evidence that this value must be negative.

The ML approach will give the intended results when applied under the first scenario, but will give unintended results under the second scenario as it does not integrate all of the information we have about this scenario.

The ML approach can still be applied under the second scenario, but that requires some explication of the mechanism that causes the data to be missing.

イロト イポト イヨト イヨト

We will next present two methods that search for ML estimates under incomplete data.

Both methods are based on local search, which start with some initial estimates, and then iteratively improve on them until some stopping condition is met.

Both methods are generally more expensive than the method for complete data, yet neither is generally guaranteed to find ML estimates.

Our first local search method, called Expectation Maximization (EM), is based on the method of complete data we discussed earlier.

This method will first complete the data set, inducing an empirical distribution, and then use it to estimate parameters as we did earlier.

The new set of parameters are guaranteed to have no less likelihood than the initial parameters, so this process can be repeated until some convergence condition is met.

向下 イヨト イヨト


Our goal is to find ML estimates for the given data set.

Expectation Maximization

					(
			Α	В	$\theta_{b a}^{0}$	Α	С	$\theta_{c a}^{0}$	В	D	$\theta^0_{d b}$
A	θ_a^0	_	a_1	b_1	.75	a_1	c_1	.50	b_1	d_1	.20
a_1	.20		a_1	b_2	.25	a_1	<i>c</i> ₂	.50	b_1	d_2	.80
a_2	.80		a ₂	b_1	.10	a 2	c_1	.25	b_2	d_1	.70
			a_2	b_2	.90	a_2	<i>c</i> ₂	.75	b_2	d_2	.30

A Bayesian network inducing a probability distribution $Pr_{\theta^0}(.)$

(< ≥) < ≥)</p>

э

The initial estimates θ^0 have the following likelihood:

$$L(\theta^{0}|\mathcal{D}) = \prod_{i=1}^{5} \Pr_{\theta^{0}}(\mathbf{d}_{i})$$

= $\Pr_{\theta^{0}}(b_{1}, c_{2})\Pr_{\theta^{0}}(b_{1}, d_{2})\Pr_{\theta^{0}}(b_{2}, c_{1}, d_{1})\Pr_{\theta^{0}}(b_{2}, c_{1}, d_{1})\Pr_{\theta^{0}}(b_{1}, d_{2})$
= $(.135)(.184)(.144)(.144)(.184)$
= 9.5×10^{-5}

Evaluating the terms in the above product would generally require inference on the Bayesian network.

To illustrate the process of completing a data set, consider again the data set:

D	Α	В	С	D
\mathbf{d}_1	?	b_1	с2	?
d ₂	?	b_1	?	d_2
d ₃	?	b_2	<i>c</i> ₁	d_1
d ₄	?	b_2	<i>c</i> ₁	d_1
\mathbf{d}_5	?	b_1	?	d_2

The first case in this data set has two variables with missing values, A and D. Hence, there are four possible completions for this case. Although we do not know which one of these completions is the correct one, we can compute the probability of each completion based on the initial set of parameters we have.

Expectation Maximization

D	Α	В	С	D	$\Pr_{\theta^0}(\mathbf{C}_i \mathbf{d}_i)$
d_1	?	b_1	<i>c</i> ₂	?	
	a ₁	b_1	<i>c</i> ₂	d_1	$.111 = \Pr_{\theta 0}(a_1, d_1 b_1, c_2)$
	a ₁	b_1	<i>c</i> ₂	<i>d</i> ₂	.444
	a2	b_1	<i>c</i> ₂	d_1	.089
	a2	b_1	<i>c</i> ₂	d_2	.356
d ₂	?	b_1	?	<i>d</i> ₂	
	a_1	b_1	c_1	d_2	$.326 = \Pr_{\theta^0}(a_1, c_1 b_1, d_2)$
	a_1	b_1	<i>c</i> ₂	<i>d</i> ₂	.326
	a2	b_1	c_1	d_2	.087
	a2	b_1	<i>c</i> ₂	d_2	.261
d ₃	?	b ₂	c_1	d_1	
	a ₁	<i>b</i> ₂	c_1	d_1	$.122 = \Pr_{\theta^0}(a_1 b_2, c_1, d_1)$
	а ₁ а ₂	b ₂ b ₂	с ₁ с ₁	d_1 d_1	$.122 = \Pr_{\theta^0}(a_1 b_2, c_1, d_1)$.878
d ₄	a ₁ a ₂ ?	b ₂ b ₂ b ₂	c ₁ c ₁ c ₁	d_1 d_1 d_1	$.122 = \Pr_{\theta^0}(a_1 b_2, c_1, d_1)$.878
d ₄	a ₁ a ₂ ? a ₁	b ₂ b ₂ b ₂ b ₂	c ₁ c ₁ c ₁ c ₁	$\begin{array}{c} d_1 \\ d_1 \\ \hline d_1 \\ \hline d_1 \\ \hline d_1 \end{array}$	$.122 = \Pr_{\theta^0}(a_1 b_2, c_1, d_1)$.878 $.122 = \Pr_{\theta^0}(a_1 b_2, c_1, d_1)$
d ₄	a1 a2 ? a1 a2	b ₂ b ₂ b ₂ b ₂ b ₂ b ₂	c ₁ c ₁ c ₁ c ₁ c ₁	$\begin{array}{c} d_1 \\ d_1 \\ \hline d_1 \\ \hline d_1 \\ d_1 \\ d_1 \end{array}$	$\begin{array}{c} .122 = \Pr_{\theta 0}(a_1 b_2, c_1, d_1) \\ .878 \\ \hline \\ .122 = \Pr_{\theta 0}(a_1 b_2, c_1, d_1) \\ .878 \end{array}$
d ₄	a ₁ a ₂ ? a ₁ a ₂ ?	$\begin{array}{c} b_2\\ b_2\\ b_2\\ b_2\\ b_2\\ b_2\\ b_2\\ b_1 \end{array}$	c ₁ c ₁ c ₁ c ₁ c ₁ ?	$ \begin{array}{c} d_1\\ d_1\\ d_1\\ d_1\\ d_1\\ d_1\\ d_2\\ \end{array} $	$\begin{array}{c} .122 = \Pr_{\theta 0} \left(a_1 b_2, c_1, d_1 \right) \\ .878 \\ \hline \\ .122 = \Pr_{\theta 0} \left(a_1 b_2, c_1, d_1 \right) \\ .878 \end{array}$
d ₄ d ₅	a ₁ a ₂ ? a ₁ a ₂ ? a ₁ a ₂	$b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_1 \\ b_1 \\ b_1$	c ₁ c ₁ c ₁ c ₁ c ₁ ? c ₁	$ \begin{array}{c} d_1\\ d_1\\ d_1\\ d_1\\ d_1\\ d_2\\ d_2\\ d_2 \end{array} $	$\begin{array}{l} .122 = \Pr_{\theta 0}(a_1 b_2,c_1,d_1)\\ .878\\ \hline\\ .122 = \Pr_{\theta 0}(a_1 b_2,c_1,d_1)\\ .878\\ \hline\\ .326 = \Pr_{\theta 0}(a_1,c_1 b_1,d_2) \end{array}$
d ₄	 a1 a2 ? a1 a2 ? a1 a2 ? a1 a1 a2 a1 a1 a1 	$\begin{array}{c} b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_1 \\ b_1 \\ b_1 \\ b_1 \end{array}$	C1 C1 C1 C1 C1 C1 C1 C1 C2	$ \begin{array}{c} d_1\\ d_1\\ d_1\\ d_1\\ d_1\\ d_2\\ d_2\\ d_2\\ d_2\\ d_2 \end{array} $	$\begin{array}{l} .122 = \Pr_{\theta^0}(a_1 b_2,c_1,d_1)\\ .878\\ \hline\\ .122 = \Pr_{\theta^0}(a_1 b_2,c_1,d_1)\\ .878\\ \hline\\ .326 = \Pr_{\theta^0}(a_1,c_1 b_1,d_2)\\ .326 \end{array}$
d ₄	a1 a2 ? a1 a2 ? a1 a1 a1 a2	$\begin{array}{c} b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_2 \\ b_1 \\ b_1 \\ b_1 \\ b_1 \end{array}$	C1 C1 C1 C1 C1 C1 C1 C1 C2 C2 C1	$ \begin{array}{c} d_1 \\ d_1 \\ d_1 \\ d_1 \\ d_2 \\ d_2 \\ d_2 \\ d_2 \\ d_2 \\ d_2 \end{array} $	$\begin{array}{c} .122 = \Pr_{\theta 0} \left(a_{1} \mid b_{2}, c_{1}, d_{1} \right) \\ .878 \\ \hline \\ .122 = \Pr_{\theta 0} \left(a_{1} \mid b_{2}, c_{1}, d_{1} \right) \\ .878 \\ \hline \\ .326 = \Pr_{\theta 0} \left(a_{1}, c_{1} \mid b_{1}, d_{2} \right) \\ .326 \\ .087 \end{array}$

Α	В	С	D	$\Pr_{\mathfrak{D},\theta^0}(.)$
a1	b_1	c_1	d_1	0
a1	b_1	c_1	d2	.130
a1	b_1	c ₂	d_1	.022
a1	b_1	<i>c</i> ₂	d2	.219
a_1	<i>b</i> ₂	c_1	d_1	.049
a1	<i>b</i> ₂	c_1	d2	0
a_1	b_2	c_2	d_1	0
a1	<i>b</i> ₂	<i>c</i> ₂	d2	0
a2	b_1	c_1	d_1	0
a2	b_1	c_1	d2	.035
a2	b_1	<i>c</i> ₂	d_1	.018
a2	b_1	c_2	d_2	.176
a2	b_2	c_1	d_1	.351
a2	b_2	c_1	d2	0
a2	b_2	<i>c</i> ₂	d_1	0
a ₂	b2	C2	d ₂	0

・ロト ・回ト ・ヨト ・ヨト

(a) completed data set, with expected values of completed cases

(b) expected empirical distribution

There are three occurrences of the instantiation a_1, b_1, c_2, d_2 in the completed data set, which result from completing the cases $\mathbf{d}_1, \mathbf{d}_2$ and \mathbf{d}_5 .

The probability of seeing these completions is given by:

$$\begin{aligned} \Pr_{\mathcal{D},\theta^{0}}(a_{1}, b_{1}, c_{2}, d_{2}) &= \frac{\Pr_{\theta^{0}}(a_{1}, d_{2}|\mathbf{d}_{1}) + \Pr_{\theta^{0}}(a_{1}, c_{2}|\mathbf{d}_{2}) + \Pr_{\theta^{0}}(a_{1}, c_{2}|\mathbf{d}_{5})}{N} \\ &= \frac{.444 + .326 + .326}{5} \\ &= .219 \end{aligned}$$

Note here that we are using $Pr_{\mathcal{D},\theta^0}(.)$ to denote the expected empirical distribution based on parameters θ^0

The expected empirical distribution of data set \mathcal{D} under parameters θ^k is defined as follows

$$\Pr_{\mathcal{D},\theta^{k}}(\alpha) \stackrel{def}{=} \frac{1}{N} \sum_{\mathbf{d}_{i},\mathbf{c}_{i}\models\alpha} \Pr_{\theta^{k}}(\mathbf{c}_{i}|\mathbf{d}_{i}),$$

where α is an event and C_i are the variables with missing values in case d_i .

Recall that $\mathbf{d}_i, \mathbf{c}_i \models \alpha$ means that event α is satisfied by complete case $\mathbf{d}_i, \mathbf{c}_i$. Hence, we are summing $\Pr_{\theta^k}(\mathbf{c}_i | \mathbf{d}_i)$ for all cases \mathbf{d}_i and their completions \mathbf{c}_i that satisfy event α .

・ 同 ト ・ ヨ ト ・ ヨ ト

When the data set is complete, $\Pr_{\mathcal{D},\theta^k}(.)$ reduces to the empirical distribution $\Pr_{\mathcal{D}}(.)$ which is independent of parameters θ^k .

Moreover, $N \cdot \Pr_{\mathcal{D}, \theta^k}(\mathbf{x})$ is called the expected count of instantiation \mathbf{x} in data set \mathcal{D} , just as $N \cdot \Pr_{\mathcal{D}}(\mathbf{x})$ represents the count of instantiation \mathbf{x} in a complete data set \mathcal{D} .

We can now use this expected empirical distribution to estimate parameters, just as we did for complete data.

For example, we have the following estimate for parameter $\theta_{c_1|a_2}$:

$$\theta_{c_1|a_2}^1 = \operatorname{Pr}_{\mathcal{D},\theta^0}(c_1|a_2) = \frac{\operatorname{Pr}_{\mathcal{D},\theta^0}(c_1,a_2)}{\operatorname{Pr}_{\mathcal{D},\theta^0}(a_2)} = .666$$

Expectation Maximization

		Α	В	$\theta^1_{b a}$	Α	С	$\theta_{c a}^1$	В	D	$\theta^1_{d b}$
A	θ_a^1	a_1	b_1	.883	a_1	c_1	.426	b_1	d_1	.067
a_1	.420	a_1	b_2	.117	a_1	<i>c</i> ₂	.574	b_1	d_2	.933
<i>a</i> ₂	.580	a_2	b_1	.395	a 2	c_1	.666	b_2	d_1	1.00
		a 2	b_2	.605	<i>a</i> ₂	<i>c</i> ₂	.334	b_2	d_2	0.00

A Bayesian network inducing a probability distribution $Pr_{\theta^1}(.)$

御 とくきとくきとうき

The new estimates θ^1 have likelihood:

$$\begin{split} \mathrm{L}(\theta^{1}|\mathcal{D}) &= \prod_{i=1}^{5} \mathrm{Pr}_{\theta^{1}}(\mathbf{d}_{i}) \\ &= (.290)(.560)(.255)(.255)(.560) \\ &= 5.9 \times 10^{-3} \\ &> \mathrm{L}(\theta^{0}|\mathcal{D}) \end{split}$$

The new estimates have a higher likelihood than the initial ones we started with. This holds more generally as we will now show.

・ 同 ト ・ ヨ ト ・ ヨ ト …

The EM estimates for data set \mathfrak{D} and parameters θ^k

$$\theta_{x|\mathbf{u}}^{k+1} \stackrel{def}{=} \operatorname{Pr}_{\mathcal{D},\theta^k}(x|\mathbf{u})$$

EM estimates are based on the expected empirical distribution, just as our estimates for complete data were based on the empirical distribution. We now have the following key result.

EM estimates satisfy the following property

 $\operatorname{LL}(\theta^{k+1}|\mathfrak{D}) \geq \operatorname{LL}(\theta^k|\mathfrak{D})$

This is a corollary of Theorems to be discussed later, which characterize the EM algorithm and also explain its name

通 と く ヨ と く ヨ と

EM estimates can be computed without constructing the expected empirical distribution.

The expected empirical distribution of data set ${\mathcal D}$ given parameters θ^k can be computed as follows

$$\operatorname{Pr}_{\mathfrak{D},\theta^{k}}(\alpha) = \frac{1}{N} \sum_{i=1}^{N} \operatorname{Pr}_{\theta^{k}}(\alpha | \mathbf{d}_{i})$$

That is, we simply iterate over the data set cases, while computing the probability of α given each case (i.e., no need to explicitly consider the completion of each case).

The EM estimates for data set \mathfrak{D} and parameters θ^k can now be computed as follows:

$$\theta_{x|\mathbf{u}}^{k+1} = \frac{\sum_{i=1}^{N} \Pr_{\theta^{k}}(x\mathbf{u}|\mathbf{d}_{i})}{\sum_{i=1}^{N} \Pr_{\theta^{k}}(\mathbf{u}|\mathbf{d}_{i})}$$

Does not reference the expected empirical distribution.

Equation computes EM estimates by performing inference on a Bayesian network parameterized by the previous parameter estimates θ^k

For example,

$$\theta_{c_1|a_2}^1 = \frac{\sum_{i=1}^5 \Pr_{\theta^0}(c_1, a_2 | \mathbf{d}_i)}{\sum_{i=1}^5 \Pr_{\theta^0}(a_2 | \mathbf{d}_i)} = \frac{0 + .087 + .878 + .878 + .087}{.444 + .348 + .878 + .878 + .348} = .666$$

Expectation Maximization

ML_EM(
$$G, \theta^0, \mathfrak{D}$$
)

input:

- G: Bayesian network structure with families XU
- θ^0 : parametrization of structure G
- \mathcal{D} : data set of size N

output: ML/EM parameter estimates for structure G.

main: 1: $k \leftarrow 0$ 2: while $\theta^k \neq \theta^{k-1}$ do {this test is different in practice} $c_{xu} \leftarrow 0$ for each family instantiation xufor i = 1 to N do for each family instantiation xu do $c_{xu} \leftarrow c_{xu} + \Pr_{\rho k}(xu|\mathbf{d}_i)$ {requires inference on network (G, θ^k) } end for end for compute parameter estimates θ^{k+1} using $\theta_{x|u}^{k+1} = c_{xu} / \sum_{x^*} c_{x^*u}$ 10: $k \leftarrow k + 1$ 11: end while 12; return θ^k

イロン イ理と イヨン ・

EM may converge to different parameters, with different likelihoods, depending on the initial estimates θ^0 that it starts with.

Each iteration of the EM algorithm will have to perform inference on a Bayesian network.

In each iteration, the algorithm computes the probability of each instantiation x**u** given each case **d**_i as evidence.

All of these computations correspond to posterior marginals over network families.

伺い イヨト イヨト

Recall the log-likelihood function:

$$\operatorname{LL}(\theta | \mathcal{D}) = \sum_{i=1}^{N} \log \operatorname{Pr}_{\theta}(\mathbf{d}_i)$$

We have seen earlier how one can maximize this function for a complete data set by choosing parameter estimates based on the empirical distribution:

$$\theta_{x|\mathbf{u}} = \Pr_{\mathcal{D}}(x|\mathbf{u})$$

Consider now the following new function of parameters, called the expected log-likelihood, which computes the log-likelihood of parameters but with respect to a completed data set:

$$\mathrm{ELL}(\theta|\boldsymbol{\mathcal{D}},\theta^{k}) \stackrel{def}{=} \sum_{i=1}^{N} \sum_{\mathbf{c}_{i}} \Big[\log \mathrm{Pr}_{\theta}(\mathbf{c}_{i},\mathbf{d}_{i}) \Big] \mathrm{Pr}_{\theta^{k}}(\mathbf{c}_{i}|\mathbf{d}_{i})$$

As before, C_i are the variables with missing values in case d_i .

Recall again the EM estimates based on the expected empirical distribution:

$$\theta_{x|\mathbf{u}}^{k+1} = \Pr_{\mathcal{D},\theta^k}(x|\mathbf{u})$$

向下 イヨト イヨト

Following result draws a parallel between the two cases of log-likelihood and expected log-likelihood.

EM parameter estimates are the only estimates that maximize the expected log-likelihood function

$$\theta^{k+1} = \operatorname*{argmax}_{\theta} \operatorname{ELL}(\theta | \mathcal{D}, \theta^k) \text{ iff } \theta^{k+1}_{x | \mathbf{u}} = \operatorname{Pr}_{\mathcal{D}, \theta^k}(x | \mathbf{u})$$

Hence, EM is indeed searching for estimates that maximize the expected log-likelihood function, which also explains its name.

Parameters that maximize the expected log-likelihood function cannot decrease the log-likelihood function

If
$$\theta^{k+1} = \underset{\theta}{\operatorname{argmax}} \operatorname{ELL}(\theta | \mathcal{D}, \theta^k)$$
, then $\operatorname{LL}(\theta^{k+1} | \mathcal{D}) \ge \operatorname{LL}(\theta^k | \mathcal{D})$

EM is capable of converging to every local maxima of the log-likelihood function

The fixed points of EM are precisely the stationary points of the log-likelihood function.

The EM algorithm is known to converge very slowly if the fraction of missing data is quite large.

Another approach for maximizing the log-likelihood function is to view the problem as one of optimizing a continuous nonlinear function.

This is a widely studied problem, where most of the solutions are based on local search, which starts by assuming some initial value $\theta_{x|u}^0$ for each parameter $\theta_{x|u}$, and then move through the parameter space in steps of the form $\theta_{x|u}^{k+1} = \theta_{x|u}^k + \delta_{x|u}^k$

Different algorithms will use different values for the increment $\delta_{x|u}^k$, yet most of them will use gradient information for determining this increment.

Recall that for a function $f(v_1, \ldots, v_n)$, the gradient is the vector of partial derivatives $\partial f/\partial v_1, \ldots, \partial f/\partial v_n$

When evaluated at a particular point (v_1, \ldots, v_n) , the gradient gives the direction of the greatest increase in the value of f

Hence, a direct use of the gradient, called gradient ascent, suggests that we move in the direction of the gradient by incrementing each variable v_i with $\eta \frac{\partial f}{\partial v_i}(v_1, \ldots, v_n)$, where η is a constant known as the learning rate.

イロト イヨト イヨト イヨト

Variable *I* indicates whether the test result is missing in the data.



The missing data depends on the condition (e.g., people who do not suffer from the condition tend not to take the test). Variable *I* indicates whether the test result is missing in the data.



The missing of data depends on the test result (e.g., individuals who test negative tend not to report the result).

The two structures explicate different missing-data mechanisms.

Our goal is to discuss ML estimates that one would obtain with respect to structures that explicate missing-data mechanisms, and to compare these estimates with the ones obtained when ignoring such mechanisms (e.g., using the simpler structure $C \rightarrow T$ as we did earlier).

Let G be a network structure, \mathcal{D} be a corresponding data set, and let **M** be the variables of G that have missing values in the data set. Let **I** be a set of variables, called missing-data indicators, that are in one-to-one correspondence with variables **M**. A network structure that results from adding variables **I** as leaf nodes to G is said to explicate the missing-data mechanism, and will be denoted by $G_{\mathbf{I}}$ We will generally use \mathcal{D}_I to denote an extension of the data set $\mathcal{D},$ which includes missing-data indicators.

\mathfrak{D}	С	Т	\mathfrak{D}_{I}	С	Т	- 1
1	yes	+ve	1	yes	+ve	no
2	yes	+ve	2	yes	+ve	no
3	yes	-ve	3	yes	-ve	no
4	no	?	4	no	?	yes
5	yes	-ve	5	yes	-ve	no
6	yes	+ve	6	yes	+ve	no
7	no	?	7	no	?	yes
8	no	-ve	8	no	-ve	no

Now that we have two different missing-data mechanisms, we can apply the ML approach in three different ways:

- ${\small \bigcirc} \ \ {\rm To \ the \ original \ structure \ } C \to \ T \ {\rm and \ the \ data \ set \ } {\frak D}$
- **2** To the first extended structure G_I and the data set \mathcal{D}_I
- **③** To the second extended structure G_I and the data set \mathcal{D}_I

All three approaches will yield estimates for variables C and T as they are shared among all three structures.

The question we now face is whether ignoring the missing-data mechanism will change the ML estimates for these variables.

・ 同 ト ・ ヨ ト ・ ヨ ト

The Missing-Data Mechanism



ignorable mechanism non-ignorable mechanism

As it turns out, the first and second approaches will indeed yield identical estimates, which are different from the ones obtained by the third approach.

・ 同 ト ・ ヨ ト ・ ヨ ト

Let G_{I} be a network structure that explicates the missing-data mechanism of structure G and data set \mathcal{D}

Let **O** be variables that are always observed in the data set \mathcal{D} , and **M** be the variables that have missing values in the data set. We will say that G_I satisfies the missing at random (MAR) assumption if I and M are d-separated by **O** in structure G_I

Intuitively, G_{I} satisfies the MAR assumption if once we know the values of variables **O**, the specific values of variables **M** become irrelevant to whether these values will be missing in the data set.

向下 イヨト イヨト

The Missing-Data Mechanism





ignorable mechanism

satisfies MAR

non-ignorable mechanism

does not

イロン イ理と イヨン ・

э

If the MAR assumption holds, the missing-data mechanism can be ignored

Let G_{I} and \mathcal{D}_{I} be a structure and a data set that explicate the missing-data mechanism of G and \mathcal{D} . Let θ be the parameters of structure G, and θ_{I} be the parameters of indicator variables I in structure G_{I} . If G_{I} satisfies the MAR assumption, then:

$$\operatorname*{argmax}_{\theta} \mathrm{LL}(\theta | \mathcal{D}) = \operatorname*{argmax}_{\theta} \max_{\theta_{\mathbf{l}}} \mathrm{LL}(\theta, \theta_{\mathbf{l}} | \mathcal{D}_{\mathbf{l}})$$

Hence, under the MAR assumption, we obtain the same ML estimates θ whether we include or ignore the missing-data mechanism.

・ 同 ト ・ ヨ ト ・ ヨ ト

Consider a data set with a single case:

$$C = \operatorname{no}, \quad T = ?$$

If we ignore the missing-data mechanism, i.e., compute ML estimates with respect to structure $C \rightarrow T$, we get the following:

- No one has the condition $C: \theta_{C=no}$ is 1
- Nothing is learned about the reliability of test *T* (its parameters are unconstrained)

If we now compute ML estimates with respect to the single case:

$$C = no, T = ?, I = yes,$$

and the second missing-data mechanism, we also get $\theta_{C=no}$ is 1

But we also get the following additional constraint.

If the missing-data mechanism is not trivial—that is, if it is not the case that $\theta_{I=yes|T=-ve} = \theta_{I=yes|T=+ve} = 1$ —then either:

- The test has a true negative rate of 100% and negative test results are always missing. That is, $\theta_{T=\rightarrow ve|C=no} = 1$ and $\theta_{l=yes|T=\rightarrow ve} = 1$
- **2** The test has a false positive rate of 100% and positive test results are always missing. That is, $\theta_{T=+ve|C=no} = 1$ and $\theta_{\modelsves|T=+ve} = 1$

Our main approach for estimating network parameters has been to search for ML estimates, that is, ones that maximize the probability of observing the given data set.

We will now assume that the structure itself is unknown, and suggest methods for learning it from the given data set.

It is natural here to adopt the same approach we adopted for parameter estimation, that is, search for network structures that maximize the probability of observing the given data set.

We will indeed start with this approach first, and then show that it needs some further refinements, leading to a general class of scoring functions for network structures.

イボト イラト イラト
Consider the ML estimates for the following structure and data set.



The log-likelihood of this network structure is given by:

 $LL(G|\mathcal{D}) = -13.3$

Learning Network Structure



A network structure with its maximum likelihood parameters. The log-likelihood of this structure is -13.3

글 > 글

Consider the ML estimates for the following structure and data set.



The log-likelihood of this network structure is given by:

$$\mathrm{LL}(G^{\star}|\mathcal{D}) = -14.1,$$

which is smaller than the likelihood for structure G.

Learning Network Structure



		Α	В	$\theta_{b a}^{ml}$	Α	С	$\theta_{c a}^{ml}$	Α	D	$\theta_{d a}^{ml}$
A	θ_a^{ml}	a_1	b_1	3/4	a_1	c_1	1/4	a_1	d_1	1/2
a_1	4/5	a_1	b_2	1/4	a_1	<i>c</i> ₂	3/4	a_1	d_2	1/2
a 2	1/5	a 2	b_1	1	a 2	c_1	1	a 2	d_1	0
		a_2	b_2	0	a_2	<i>c</i> ₂	0	a_2	d_2	1

A network structure with its maximum likelihood parameters. The log-likelihood of this structure is $-14.1\,$

We will next present an algorithm for finding ML tree structures in time and space that are quadratic in the number of nodes in the structure.

Consider the mutual information between two variables in the empirical distribution:

$$\mathrm{MI}_{\mathcal{D}}(X, U) \stackrel{def}{=} \sum_{x, u} \mathrm{Pr}_{\mathcal{D}}(x, u) \log \frac{\mathrm{Pr}_{\mathcal{D}}(x, u)}{\mathrm{Pr}_{\mathcal{D}}(x) \mathrm{Pr}_{\mathcal{D}}(u)}$$

Given a tree structure G with edges $U \rightarrow X$, its score is given by

$$\operatorname{tScore}(G|\mathcal{D}) \stackrel{def}{=} \sum_{U \to X} \operatorname{MI}_{\mathcal{D}}(X, U)$$

Trees having a maximal likelihood are precisely those trees that maximize the above score.

If G is a tree structure, and \mathcal{D} is a complete data set, then $\operatorname{argmax}_{G} \operatorname{tScore}(G|\mathcal{D}) = \operatorname{argmax}_{G} \operatorname{LL}(G|\mathcal{D})$

Learning Tree Structures







(b) maximum spanning tree



(d) maximum likelihood tree

・ロト ・四ト ・ヨト ・ヨトー

э

Learning Tree Structures

We can obtain log-likelihood by computing the probability of each case in the data set using any of these tree structures (and its corresponding ML estimates).

We can also use an earlier result, which shows that the log-likelihood corresponds to a sum of terms, one term for each family in the network.

If we consider the tree structure G in (c) above, this result gives:

 $\operatorname{LL}(\mathcal{G}|\mathcal{D})$

- $= -N \times (\operatorname{ENT}_{\mathcal{D}}(A|C) + \operatorname{ENT}_{\mathcal{D}}(B) + \operatorname{ENT}_{\mathcal{D}}(C|B) + \operatorname{ENT}_{\mathcal{D}}(D|B))$
- $= -5 \times (.400 + .722 + .649 + .649)$
- = -12.1

The terms correspond to the families of given tree structure: AC, B, CB and DB.

Suppose now that our goal is to find a maximum likelihood structure, but without restricting ourselves to tree structures.



Consider the DAG structure in (b) earlier, which is obtained by adding an edge $D \rightarrow A$ to the tree structure in (a).

The log-likelihood of this DAG is given by:

$$LL(G|\mathcal{D})$$

$$= -N \times (ENT_{\mathcal{D}}(A|C, D) + ENT_{\mathcal{D}}(B) + ENT_{\mathcal{D}}(C|B) + ENT_{\mathcal{D}}(D|B))$$

$$= -5 \times (0 + .722 + .649 + .649)$$

$$= -10.1$$

which is larger than the log-likelihood of the tree in (a).

伺下 イヨト イヨト

Learning DAG Structures

Only difference between two likelihoods is the entropy term for variable A, since this is the only variable with different families.

The family of A is AC in the tree, and it is ACD in the DAG. Moreover,

 $\operatorname{ENT}_{\mathcal{D}}(A|C, D) < \operatorname{ENT}_{\mathcal{D}}(A|C),$

and, hence,

$$-\mathrm{ENT}_{\mathcal{D}}(A|C,D)>-\mathrm{ENT}_{\mathcal{D}}(A|C),$$

which is why the DAG has a larger log-likelihood than the tree.

More generally

If $\mathbf{U} \subseteq \mathbf{U}^{\star}$, then $\operatorname{ENT}(X|\mathbf{U}) \geq \operatorname{ENT}(X|\mathbf{U}^{\star})$

By adding more parents to a variable, we will never increase its entropy term and, hence, will never decrease the log-likelihood of resulting structure.

If DAG G^{\star} is the result of adding edges to DAG G, then

 $LL(G^*|\mathcal{D}) \ge LL(G|\mathcal{D}).$

If we simply search for a network structure with maximal likelihood, we will end up choosing a complete network structure; that is, a DAG to which no more edges can be added (without introducing directed cycles).²

²Recall that there are *n*! complete DAGs over *n* variables. Each of these DAGs corresponds to a total variable ordering X_1, \ldots, X_n in which variable X_i has X_1, \ldots, X_{i-1} as its parents.

Complete DAGs are undesirable for a number of reasons:

- They make no assertions of conditional independence and, hence, their topology does not reveal any properties of the distribution they induce.
- 2 A complete DAG over n variables has a treewidth of n-1 and is therefore impossible to work with practically.
- Complete DAGs suffer from the problem of overfitting, which refers to the use of a model that has too many parameters compared to the available data.

個 と く ヨ と く ヨ と

A classical example for illustrating the problem of overfitting is that of finding a polynomial that fits a given set of data points $(x_1, y_1), \ldots, (x_n, y_n)$.

Consider the following data points for an example:

X	y
1	1.1
5	4.5
10	11
15	14.5
20	22

If we insist on a perfect fit, however, we can use a fourth degree polynomial, which is guaranteed to fit the data perfectly.

伺下 イヨト イヨト

Learning DAG Structures



(a) straight line (b) 4th degree polynomial

The problem of overfitting: Even though the fit in (b) is perfect, the polynomial does not appear to provide a good generalization of the data beyond the range of the observed data points.

In summary, the problem of overfitting materializes when one focuses on learning a model that fits the data well, without constraining enough the number of free model parameters.

The result is that one ends up adopting models that are more complex than necessary.

Moreover, such models tend to provide poor generalizations of the data and will, therefore, perform poorly on cases that are not part of the given data set.

Even though there is no agreed upon solution to the problem of overfitting, all available solutions tend to be based on a common principle known as Occam's razor, which says that one should prefer simpler models over more complex models, others things being equal.

To realize this principle, one needs a measure of model complexity, and a method for balancing the complexity of a model with its data fit.

For Bayesian networks (and many other modeling frameworks), model complexity is measured using the number of independent parameters in the model.

(4月) イヨト イヨト

Model complexity

Let *G* be a DAG over variables X_1, \ldots, X_n with corresponding parents U_1, \ldots, U_n , and let $\mathbf{Y}^{\#}$ denote the number of instantiations for variables \mathbf{Y} . The dimension of DAG *G* is defined as follows:

$$||G|| \stackrel{def}{=} \sum_{i=1}^{n} ||X_i \mathbf{U}_i||$$
$$X_i \mathbf{U}_i|| \stackrel{def}{=} (X_i^{\#} - 1) \mathbf{U}_i^{\#}$$

Dimension is number of independent parameters in CPTs.

・ 回 ト ・ ヨ ト ・ ヨ ト

Learning DAG Structures



(a) dimension 7 (b) dimension 9

(c) dimension 15

- * 聞 > - * ヨ > - * ヨ >

æ

Scoring measures for structure G and data set \mathfrak{D} of size N:

Score(
$$G|\mathcal{D}$$
) $\stackrel{def}{=}$ LL($G|\mathcal{D}$) – $\psi(N) \cdot ||G||$

Note: Score is ≤ 0

The first component of this score, $LL(G|\mathcal{D})$, is the log-likelihood function we considered before.

The second component, $\psi(N) \cdot ||G||$, is a penalty term that favors simpler models, i.e., ones with a smaller number of independent parameters.

Penalty term has a weight, $\psi(N) \ge 0$, which is a function of the data set size N

When the penalty weight $\psi(N)$ is a constant that is independent of N, one gets score in which model complexity is a secondary issue.

Log-likelihood function $LL(G|\mathcal{D})$ grows linearly in the data set size N and will quickly dominate the penalty term.

Model complexity will only be used to distinguish between models that have relatively equal log-likelihood terms.

Scoring measure is known as the Akaike Information Criterion (AIC).

向下 イヨト イヨト

Another, yet more common, choice of the penalty weight is $\psi(N) = \frac{\log_2 N}{2}$, which leads to a more influential penalty term.

This term grows logarithmically in N, while the log-likelihood term grows linearly in N.

The influence of model complexity will decrease as N grows, allowing the log-likelihood term to eventually dominate the score.

This penalty weight gives rise to the Minimum Description Length (MDL) score:

$$\mathrm{MDL}(G|\mathcal{D}) \stackrel{def}{=} \mathrm{LL}(G|\mathcal{D}) - \left(\frac{\log_2 N}{2}\right) ||G||$$

伺い イヨト イヨト

Learning DAG Structures



MDL prefers first structure even though it has smaller log-likelihood.

伺下 イヨト イヨト

The MDL score is also known as the Bayesian Information Criterion (BIC).

It is sometimes expressed as the negative of the given score, where the goal is to minimize the score instead of maximizing it. Searching for a network structure that optimizes a particular score can be quite expensive due to the very large number of structures one may need to consider.

Greedy algorithms tend to be of more practical use when learning network structures.

Systematic search algorithms can also be practical, but only under some conditions.

Both classes of algorithms rely for their efficient implementation on a property that most scoring functions have.

decomposability or modularity: allows one to decompose the score into an aggregate of local scores, one for each network family.

イロン 不得と 不良と 不良とう

Score for structure G and data set \mathfrak{D} of size N

Score(
$$G|\mathcal{D}$$
) $\stackrel{def}{=}$ LL($G|\mathcal{D}$) – $\psi(N) \cdot ||G||$

Let $X\mathbf{U}$ range over the families of DAG G

This score can be decomposed as follows:

$$\operatorname{Score}(G|\mathcal{D}) = \sum_{X\mathbf{U}} \operatorname{Score}(X,\mathbf{U}|\mathcal{D}),$$

where

$$\operatorname{Score}(X, \mathbf{U} | \mathcal{D}) \stackrel{def}{=} - N \cdot \operatorname{ENT}_{\mathcal{D}}(X | \mathbf{U}) - \psi(N) \cdot ||X\mathbf{U}||$$

(A) < (A)

Local Search



Adding or removing an edge will change only one family, while reversing an edge will change only two families. Hence, the score can always be updated locally as a result of the local network change induced by adding, removing or reversing an edge. The local modifications to the structure are then constrained to: adding an edge, removing an edge, or reversing an edge, while ensuring that the structure remains a DAG.

These local changes to the network structure will also change the score, possibly increasing or decreasing it.

The goal, however, is to commit to the change that will increase the score the most.

If none of the local changes can increase the score, the algorithm will terminate and return the current structure.

向下 イヨト イヨト

Local search is not guaranteed to return an optimal network structure, i.e., one that has the largest score.

The only guarantee provided by the algorithm is that the structure it returns will be locally optimal in that no local change can improve its score.

This sub-optimal behavior of local search can usually be improved by techniques such as random restarts.

According to this technique, one would repeat the local search multiple times, each time starting with a different initial network, and then return the network with the best score across all repetitions.

通 と く ヨ と く ヨ と

A common technique for reducing the search space size is to assume a total ordering on network variables and then search only among network structures that are consistent with the chosen order.

If we use the variable order X_1, \ldots, X_n , the search process can by viewed as trying to find, for each variable X_i , a set of parents $U_i \subseteq X_1, \ldots, X_{i-1}$

Not only does this technique reduce the size of search space, but it also allows one to decompose the search problem into n independent problems, each concerned with finding a set of parents for some network variable.

イロン 不得と 不良と 不良とう

One of the more effective heuristic algorithms for optimizing a family score is known as ${\rm K3.}^3$

This algorithm starts with an empty set of parents, successively adding variables to the set, one at a time, until such additions will no longer increase the score.

³The name K3 refers to the version of this algorithm that optimizes the MDL score, although it also applies to other scores as well. Another version of this heuristic algorithm, called K2, works similarly but using a different score to be discussed in Chapter 18

Constraining the Search Space



Greedy search for a parent set for variable X_5

Greedy Search

Suppose the goal is to find a set of parents for X_5 from the set of variables X_1, \ldots, X_4 . The K3 algorithm will start by setting \mathbf{U}_5 to the empty set, and then find a variable X_i (if any), $i = 1, \ldots, 4$, that will maximize

 $\operatorname{Score}(X_5, X_i | \mathcal{D}) \geq \operatorname{Score}(X_5 | \mathcal{D})$

Suppose that X_3 happens to be such a variable. The algorithm will then set $U_5 = \{X_3\}$ and search for another variable X_i in X_1, X_2, X_4 that will maximize

$$\operatorname{Score}(X_5, X_3X_i|\mathcal{D}) \geq \operatorname{Score}(X_5, X_3|\mathcal{D})$$

Suppose again that X_2 happens to be such a variable, leading to the new set of parents $U_5 = \{X_2, X_3\}$

It may happen that adding X_1 to this set will not increase the score, and neither will adding X_4

In this case, K3 will terminate, returning $\mathbf{U}_5 = \{X_2, X_3\}$ as the parent set for X_5

K3 is a greedy algorithm that is not guaranteed to identify the optimal set of parents \mathbf{U}_i , i.e., the one that maximizes $\text{Score}(X_i, \mathbf{U}_i | \mathcal{D})$

Therefore, it is not uncommon to use the structure obtained by this algorithm as a starting point for other algorithms, such as the local search algorithm discussed earlier, or the optimal search algorithm we shall discuss next. We will next discuss an optimal search algorithm for network structures, which is based on branch-and-bound depth-first search.

Similar to K3, the algorithm will assume a total order of network variables, X_1, \ldots, X_n and search only among network structures that are consistent with this order.

As mentioned earlier, this allows one to decompose the search process into n independent search problems.

Optimal Search



Tree nodes are in one-to-one correspondence with parent sets for X_5 . A search tree for variable X_i will have a total of 2^{i-1} nodes, corresponding to the number of subsets one can choose from variables X_1, \ldots, X_{i-1}
Optimal Search



Tree nodes are in one-to-one correspondence with parent sets for X_5 . A search tree for variable X_i will have a total of 2^{i-1} nodes, corresponding to the number of subsets one can choose from variables X_1, \ldots, X_{i-1}

One can search the tree using depth-first search, while maintaining the score s of the best parent set visited thus far.

The complexity of this algorithm can be improved on average if one can compute for each search node U_i an upper bound on $\text{Score}(X_i, U_i^* | \mathcal{D})$, where $U_i \subseteq U_i^*$

If the computed upper bound at node U_i is not better than the best score *s* obtained thus far, then one can prune U_i and all nodes below it in the search tree, since none of these parent sets can be better than the one found thus far.

This pruning allows one to escape the exponential complexity in some cases.

The extent of pruning depends on the quality of upper bound used.

(本間) (本語) (本語)

Upper bound for MDL score

Let \mathbf{U}_i be a parent set, and let \mathbf{U}_i^+ be the largest parent set appearing below \mathbf{U}_i in the search tree. If \mathbf{U}_i^* is a parent set in the tree rooted at \mathbf{U}_i , then

$$\mathrm{MDL}(X_i, \mathbf{U}_i^* | \mathcal{D}) \leq -N \cdot \mathrm{ENT}_{\mathcal{D}}(X_i | \mathbf{U}_i^+) - \psi(N) \cdot ||X_i \mathbf{U}_i||$$

Consider tree in (a). At the search node $U_5 = \{X_2\}$, we get $U_5^+ = \{X_2, X_3, X_4\}$. Moreover, U_5^* ranges over parent sets $\{X_2\}$, $\{X_2, X_3\}$, $\{X_2, X_4\}$ and $\{X_2, X_3, X_4\}$

・ 同 ト ・ ヨ ト ・ ヨ ト



Formulation: Search space: Variable subsets

[Yuan, Malone, Wu, IJCAI-11]

1/39 -



Formulation: Search space: Variable subsets Start node: Empty set

[Yuan, Malone, Wu, IJCAI-11]



Formulation: Search space: Variable subsets Start node: Empty set Goal node: Complete set

[Yuan, Malone, Wu, IJCAI-11]



Formulation:Search space:Variable subsetsStart node:Empty setGoal node:Complete setEdges:Select parentsEdge cost:BestMDL(X,U) for
edge $U \rightarrow U \cup \{X\}$

[Yuan, Malone, Wu, IJCAI-11]





Formulation:Search space: Variable subsetsStart node:Empty setGoal node:Complete setEdges:Select parentsEdge cost:BestMDL(X,U) for
edge $U \rightarrow U \cup \{X\}$

Task: find the shortest path between start and goal nodes

[Yuan, Malone, Wu, IJCAI-11]



BestMDL(*X*,*U*) for edge $U \rightarrow U \cup \{X\}$ Task: find the shortest path between start and goal nodes

Empty set

Complete set

Select parents

[Yuan, Malone, Wu, IJCAI-11]

Our discussion on the search for network structures has been restricted to complete data sets.

The main reason for this is computational.

The likelihood of a network structure does not admit a closed form when the data set is incomplete and does not decompose into components.

Algorithms for learning structures with incomplete data will typically involve two searches: an outer search in the space of network structures, and an inner search in the space of network parameters.

向下 イヨト イヨト