

Dirichlet Bayesian Network Scores and the Maximum Entropy Principle



UNIVERSITY OF
OXFORD

Marco Scutari

scutari@stats.ox.ac.uk

Department of Statistics
University of Oxford

September 21, 2017

Bayesian Network Structure Learning

Learning a BN $\mathcal{B} = (\mathcal{G}, \Theta)$ from a data set \mathcal{D} is performed in two steps:

$$\underbrace{P(\mathcal{B} | \mathcal{D}) = P(\mathcal{G}, \Theta | \mathcal{D})}_{\text{learning}} = \underbrace{P(\mathcal{G} | \mathcal{D})}_{\text{structure learning}} \cdot \underbrace{P(\Theta | \mathcal{G}, \mathcal{D})}_{\text{parameter learning}} .$$

In a Bayesian setting **structure learning** consists in finding the DAG with the best $P(\mathcal{G} | \mathcal{D})$ (BIC [6] is a common alternative) with some heuristic search algorithm. We can decompose $P(\mathcal{G} | \mathcal{D})$ into

$$P(\mathcal{G} | \mathcal{D}) \propto P(\mathcal{G}) P(\mathcal{D} | \mathcal{G}) = P(\mathcal{G}) \int P(\mathcal{D} | \mathcal{G}, \Theta) P(\Theta | \mathcal{G}) d\Theta$$

where $P(\mathcal{G})$ is the **prior distribution over the space of the DAGs** and $P(\mathcal{D} | \mathcal{G})$ is the **marginal likelihood** of the data given \mathcal{G} averaged over all possible parameter sets Θ ; and then

$$P(\mathcal{D} | \mathcal{G}) = \prod_{i=1}^N \left[\int P(X_i | \Pi_{X_i}, \Theta_{X_i}) P(\Theta_{X_i} | \Pi_{X_i}) d\Theta_{X_i} \right]$$

where Π_{X_i} are the parents of X_i in \mathcal{G} .

The Bayesian Dirichlet Marginal Likelihood

If \mathcal{D} contains no missing values and assuming:

- a **Dirichlet conjugate prior** ($X_i | \Pi_{X_i} \sim \text{Mult}(\Theta_{X_i} | \Pi_{X_i})$ and $\Theta_{X_i} | \Pi_{X_i} \sim \text{Dir}(\alpha_{ijk})$, $\sum_{jk} \alpha_{ijk} = \alpha_i$ the imaginary sample size);
- **positivity** (all conditional probabilities $\pi_{ijk} > 0$);
- **parameter independence** (π_{ijk} for different parent configurations are independent) and **modularity** (π_{ijk} in different nodes are independent);

Heckerman *et al.* [4] derived a closed form expression for $P(\mathcal{D} | \mathcal{G})$:

$$\begin{aligned} \text{BD}(\mathcal{G}, \mathcal{D}; \boldsymbol{\alpha}) &= \prod_{i=1}^N \text{BD}(X_i, \Pi_{X_i}; \boldsymbol{\alpha}_i) = \\ &= \prod_{i=1}^N \prod_{j=1}^{q_i} \left[\frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + n_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(\alpha_{ijk} + n_{ijk})}{\Gamma(\alpha_{ijk})} \right] \end{aligned}$$

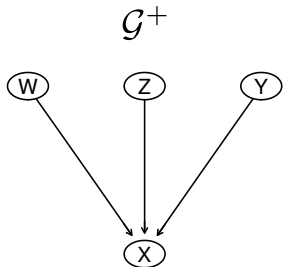
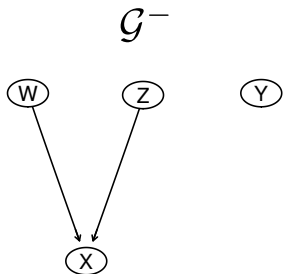
where r_i is the number of states of X_i ; q_i is the number of configurations of Π_{X_i} ; $n_{ij} = \sum_k n_{ijk}$; and $\alpha_{ij} = \sum_k \alpha_{ijk}$.

Bayesian Dirichlet Equivalent Uniform (BDeu)

The most common implementation of BD assumes $\alpha_{ijk} = \alpha / (r_i q_i)$, $\alpha = \alpha_i$ and is known from [4] as the **Bayesian Dirichlet equivalent uniform** (BDeu) marginal likelihood. However, there is evidence that assuming a flat prior over the parameters can be problematic:

- The prior is **actually not uninformative** [5].
- MAP DAGs selected using BDeu are **highly sensitive to the choice of α** and can have markedly different number of arcs even for reasonable α [8].
- In the limits $\alpha \rightarrow 0$ and $\alpha \rightarrow \infty$ it is possible to obtain both very simple and very complex DAGs, and **model comparison may be inconsistent** for small \mathcal{D} and small α [8, 10].
- The sparseness of the MAP network is determined by a **complex interaction between α and \mathcal{D}** [10, 12].
- There are formal proofs of all this in [11, 12].

Exhibits A and B


 \mathcal{D}_1

X	Z	W	Y
1	0	0	0
0	0	0	0
0	0	0	0
0	0	1	0
1	0	1	0
1	0	1	0
0	1	0	0
1	1	0	0
1	1	0	0
1	1	1	1
0	1	1	1
0	1	1	1

 \mathcal{D}_2

X	Z	W	Y
0	0	0	0
0	0	0	0
0	0	0	0
1	0	1	0
1	0	1	0
1	0	1	0
1	1	0	0
1	1	0	0
1	1	0	0
0	1	1	1
0	1	1	1
0	1	1	1

Exhibit A

The sample frequencies (n_{ijk}) for $X | \Pi_X$ are:

		Z, W			
		0,0	1,0	0,1	1,1
X	0	2	1	1	2
	1	1	2	2	1

and those for $X | \Pi_X \cup Y$ are as follows.

		Z, W, Y							
		0,0,0	1,0,0	0,1,0	1,1,0	0,0,1	1,0,1	0,1,1	1,1,1
X	0	2	1	1	0	0	0	0	2
	1	1	2	2	0	0	0	0	1

Even though $X | \Pi_X$ and $X | \Pi_X \cup Y$ have the **same empirical entropy**,

$$H(X | \Pi_X) = H(X | \Pi_X \cup Y) = 4 \left[-\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} \right] = 2.546 \dots$$

Exhibit A

... \mathcal{G}^- has a **higher entropy** than \mathcal{G}^+ *a posteriori* with $\alpha = 1$...

$$\begin{aligned} H(X | \Pi_X; \alpha) &= 4 \left[-\frac{1 + 1/8}{3 + 1/4} \log \frac{1 + 1/8}{3 + 1/4} - \frac{2 + 1/8}{3 + 1/4} \log \frac{2 + 1/8}{3 + 1/4} \right] \\ &= 2.580, \end{aligned}$$

$$\begin{aligned} H(X | \Pi_X \cup Y; \alpha) &= 4 \left[-\frac{1 + 1/16}{3 + 1/8} \log \frac{1 + 1/16}{3 + 1/8} - \frac{2 + 1/16}{3 + 1/8} \log \frac{2 + 1/16}{3 + 1/8} \right] \\ &= 2.564 \end{aligned}$$

... and BDeu with $\alpha = 1$ chooses accordingly, so things fortunately work out:

$$\begin{aligned} \text{BDeu}(X | \Pi_X) &= \left(\frac{\Gamma(1/4)}{\Gamma(1/4 + 3)} \left[\frac{\Gamma(1/8 + 2)}{\Gamma(1/8)} \cdot \frac{\Gamma(1/8 + 1)}{\Gamma(1/8)} \right] \right)^4 \\ &= 3.906 \times 10^{-7}, \end{aligned}$$

$$\begin{aligned} \text{BDeu}(X | \Pi_X \cup Y) &= \left(\frac{\Gamma(1/8)}{\Gamma(1/8 + 3)} \left[\frac{\Gamma(1/16 + 2)}{\Gamma(1/16)} \cdot \frac{\Gamma(1/16 + 1)}{\Gamma(1/16)} \right] \right)^4 \\ &= 3.721 \times 10^{-8}. \end{aligned}$$

Exhibit B

The sample frequencies for $X \mid \Pi_X$ are:

		Z, W			
		0,0	1,0	0,1	1,1
X	0	3	0	0	3
	1	0	3	3	0

and those for $X \mid \Pi_X \cup Y$ are as follows.

		Z, W, Y							
		0,0,0	1,0,0	0,1,0	1,1,0	0,0,1	1,0,1	0,1,1	1,1,1
X	0	3	0	0	0	0	0	0	3
	1	0	3	3	0	0	0	0	0

The empirical entropy of X is equal to zero for both \mathcal{G}^+ and \mathcal{G}^- , since the value of X is completely determined by the configurations of its parents in both cases.

Exhibit B

Again, the posterior entropies for \mathcal{G}^+ and \mathcal{G}^- differ:

$$\begin{aligned} H(X \mid \Pi_X; \alpha) &= 4 \left[-\frac{0 + 1/8}{3 + 1/4} \log \frac{0 + 1/8}{3 + 1/4} - \frac{3 + 1/8}{3 + 1/4} \log \frac{3 + 1/8}{3 + 1/4} \right] \\ &= 0.652, \end{aligned}$$

$$\begin{aligned} H(X \mid \Pi_X \cup Y; \alpha) &= 4 \left[-\frac{0 + 1/16}{3 + 1/8} \log \frac{0 + 1/16}{3 + 1/8} - \frac{3 + 1/16}{3 + 1/8} \log \frac{3 + 1/16}{3 + 1/8} \right] \\ &= 0.392. \end{aligned}$$

However, BDeu with $\alpha = 1$ yields

$$\text{BDeu}(X \mid \Pi_X) = \left(\frac{\Gamma(1/4)}{\Gamma(1/4 + 3)} \left[\frac{\Gamma(1/8 + 3)}{\Gamma(1/8)} \cdot \frac{\Gamma(1/8)}{\Gamma(1/8)} \right] \right)^4 = 0.032,$$

$$\text{BDeu}(X \mid \Pi_X \cup Y) = \left(\frac{\Gamma(1/8)}{\Gamma(1/8 + 3)} \left[\frac{\Gamma(1/16 + 3)}{\Gamma(1/16)} \cdot \frac{\Gamma(1/16)}{\Gamma(1/16)} \right] \right)^4 = 0.044,$$

preferring \mathcal{G}^+ over \mathcal{G}^- even though **the additional arc $Y \rightarrow X$ does not provide any additional information** on the distribution of X , and even though **4 out of 8 conditional distributions in $X \mid \Pi_X \cup Y$** are not observed at all in the data.

Better Than BDeu: Bayesian Dirichlet Sparse (BDs)

If the positivity assumption is violated or the sample size n is small, there may be configurations of some Π_{X_i} that are not observed in \mathcal{D} . And then

$$\text{BDeu}(X_i, \Pi_{X_i}; \alpha) = \prod_{j:n_{ij}=0} \left[\frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(\alpha_{ijk})}{\Gamma(\alpha_{ijk})} \right] \prod_{j:n_{ij}>0} \left[\frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + n_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(\alpha_{ijk} + n_{ijk})}{\Gamma(\alpha_{ijk})} \right],$$

so the **effective imaginary sample size decreases as the number of unobserved parents configurations increases**. We can prevent that by replacing α_{ijk} with

$$\tilde{\alpha}_{ijk} = \begin{cases} \alpha / (r_i \tilde{q}_i) & \text{if } n_{ij} > 0 \\ 0 & \text{otherwise} \end{cases}, \quad \tilde{q}_i = \{\text{number of } \Pi_{X_i} \text{ such that } n_{ij} > 0\}$$

and plugging it in BD instead of $\alpha_{ijk} = \alpha / (r_i q_i)$ to obtain BDs.

Then $\text{BDs}(X_i, \Pi_{X_i}; \alpha) = \text{BDeu}(X_i, \Pi_{X_i}; \alpha q_i / \tilde{q}_i)$.

BDeu and BDs Compared

$$\begin{array}{c}
 \underbrace{\hspace{10em}}_{\Pi_{X_i}} \\
 \pi_1 \quad \pi_2 \quad \dots \quad \pi_{q_i} \\
 \left. \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_{r_i} \end{array} \right\} X_i \quad \begin{array}{c} \frac{\alpha}{r_i q_i} \quad \frac{\alpha}{r_i q_i} \quad \dots \quad \frac{\alpha}{r_i q_i} \\ \frac{\alpha}{r_i q_i} \quad \frac{\alpha}{r_i q_i} \quad \dots \quad \frac{\alpha}{r_i q_i} \\ \vdots \quad \vdots \quad \vdots \quad \vdots \\ \frac{\alpha}{r_i q_i} \quad \frac{\alpha}{r_i q_i} \quad \dots \quad \frac{\alpha}{r_i q_i} \end{array}
 \end{array}
 \quad
 \begin{array}{c}
 \underbrace{\hspace{10em}}_{\Pi_{X_i}} \\
 \pi_1 \quad \pi_2 \quad \dots \quad \pi_{q_i} \\
 \left. \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_{r_i} \end{array} \right\} X_i \quad \begin{array}{c} \frac{\alpha}{r_i \tilde{q}_i} \quad 0 \quad \dots \quad \frac{\alpha}{r_i \tilde{q}_i} \\ \frac{\alpha}{r_i \tilde{q}_i} \quad 0 \quad \dots \quad \frac{\alpha}{r_i \tilde{q}_i} \\ \vdots \quad \vdots \quad \vdots \quad \vdots \\ \frac{\alpha}{r_i \tilde{q}_i} \quad 0 \quad \dots \quad \frac{\alpha}{r_i \tilde{q}_i} \end{array}
 \end{array}$$

Cells that correspond to (X_i, Π_{X_i}) combinations that are not observed in the data are in red, observed combinations are in green.

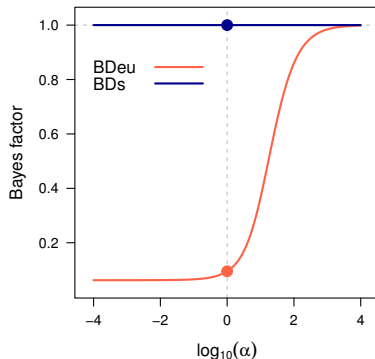
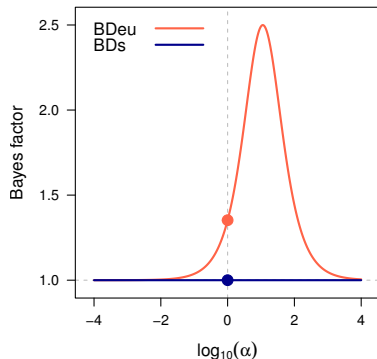
Exhibits A and B, Once More

BDs does not suffer from the bias arising from $\tilde{q}_i < q_i$ and it assigns the same score to \mathcal{G}^- and \mathcal{G}^+ in both examples,

Exhibit A: $\text{BDs}(X | \Pi_X) = \text{BDs}(X | \Pi_X \cup Y) = 3.9 \times 10^{-7}$,

Exhibit B: $\text{BDs}(X | \Pi_X) = \text{BDs}(X | \Pi_X \cup Y) = 0.032$.

It also avoids giving wildly different Bayes factors depending on the value of α .



This Left Me with a Few Questions...

The obvious one being:

1. The behaviour of BDeu is certainly undesirable, but **it is it wrong?**

Followed by:

2. Posterior entropy and BDeu rank \mathcal{G}^- and \mathcal{G}^+ in the same order for Exhibit A, but they do not for Exhibit B. **Why is that?**

And the reason why I found that surprising is that:

3. Maximum (relative) entropy [7, 9, 1] represents a very general approach that **includes Bayesian posterior estimation as a particular case** [3]; it can also be seen as a particular case of MDL [2].

Hence, unless something is wrong with BDeu I would expect the two to agree. Especially because we can use MDL (using BIC), MAP (using BDeu/BDs),

Bayesian Statistics and Information Theory (I)

The derivation of Bayesian posterior as a particular case of maximum (relative) entropy is made clear in Giffin and Caticha [3]. The selected joint posterior $P(X, \Theta)$ is that which **maximises the relative entropy**

$$S(P, P_{old}) = - \int P(X, \Theta) \log \frac{P(X, \Theta)}{P_{old}(X, \Theta)} dX d\Theta.$$

The family of posteriors that reflects the fact that X is now known to take value x' is such that

$$P(X) = \int P(X = x', \Theta) d\Theta = \delta(X - x')$$

which amounts to a (possibly infinite) number of constraints on $P(X, \Theta)$: **for each possible value of X there is one constraint.**

Bayesian Statistics and Information Theory (II)

Maximising $S(P, P_{old})$ subject to those constraints using Lagrange multipliers means solving

$$S(P, P_{old}) + \lambda_0 \underbrace{\left[\int P \, dX \, d\Theta - 1 \right]}_{\text{normalising constraint}} + \int \lambda(x) \underbrace{\left[\int P(X, \Theta) - \delta(X - x') \, d\Theta \right]}_{\text{constraint for each value of } X} dX$$

and yields the familiar Bayesian update rule:

$$P_{new}(X, \Theta) = \frac{P_{old}(X, \Theta)\delta(X - x')}{P_{old}(X)} = P_{old}(\Theta | X)\delta(X - x').$$

Bayesian Statistics and Information Theory (III)

In particular, the updated distribution for Θ is

$$P_{new}(\Theta) = \int P_{new}(X, \Theta) dX = P_{old}(\Theta | X = x')$$

which means that the posterior distribution is that in which we **only update those aspects of our beliefs for which corrective new evidence** (in this case, the data) **has been supplied**. However, we use all the available information (as opposed to just what is in the empirical entropy):

- the information encoded in the distributional assumptions for the prior distribution over Θ ;
- the information encoded in the distributional assumptions for the random variable X ;
- the information encoded in the observed data.

Back to BNs: the Posterior Expected Entropy

Starting from the **Markov property**, for a BN we can write

$$H^{\mathcal{G}}(\mathbf{X}; \Theta) = \sum_{i=1}^N H^{\mathcal{G}}(X_i; \Theta_{X_i}).$$

where $H^{\mathcal{G}}(X_i; \Theta_{X_i})$ is the entropy of X_i given its parents Π_{X_i} in \mathcal{G} .

The **marginal posterior expectation** of $H^{\mathcal{G}}(X_i; \Theta_{X_i})$ with respect to Θ_{X_i} given the data can then be expressed as

$$E(H^{\mathcal{G}}(X_i) | \mathcal{D}) = \int H^{\mathcal{G}}(X_i; \Theta_{X_i}) P(\Theta_{X_i} | \mathcal{D}) d\Theta_{X_i}$$

where we use \mathcal{D} to refer specifically to the observed values for X_i and $\Pi_{X_i}^{\mathcal{G}}$ with a slight abuse of notation.

Adding the Dirichlet Prior

We can then introduce a *Dirichlet*(α_{ijk}) **prior** over Θ_{X_i} with

$$P(\Theta_{X_i} | \mathcal{D}) = \int P(\Theta_{X_i} | \mathcal{D}, \alpha_{ijk}) P(\alpha_{ijk} | \mathcal{D}) d\alpha_{ijk},$$

which leads to

$$\begin{aligned} E(H^{\mathcal{G}}(X_i) | \mathcal{D}) &= \iint H^{\mathcal{G}}(X_i; \Theta_{X_i}) P(\Theta_{X_i} | \mathcal{D}, \alpha_{ijk}) P(\alpha_{ijk} | \mathcal{D}) d\alpha_{ijk} d\Theta_{X_i} \\ &\propto \int E(H^{\mathcal{G}}(X_i) | \mathcal{D}, \alpha_{ijk}) P(\mathcal{D} | \alpha_{ijk}) P(\alpha_{ijk}) d\alpha_{ijk}, \end{aligned}$$

where $P(\alpha_{ijk})$ is a **hyper-prior** distribution over the space of the Dirichlet priors, identified by their parameter sets $\{\alpha_{ijk}\}$.

Components of the Posterior Expected Entropy

$E(H^{\mathcal{G}}(X_i) \mid \mathcal{D}, \alpha_{ijk})$ is the **posterior expected value of the entropy of $X_i \mid \Pi_{X_i}$ given α_{ijk}** , and has closed form

$$E(H^{\mathcal{G}}(X_i) \mid \mathcal{D}, \alpha_{ijk}) = \sum_{j=1}^{q_i} \left[\psi_0(\alpha_{ij} + n_{ij} + 1) - \sum_{k=1}^{r_i} \frac{\alpha_{ijk} + n_{ijk}}{\alpha_{ij} + n_{ij}} \psi_0(\alpha_{ijk} + n_{ijk} + 1) \right].$$

$P(\mathcal{D} \mid \alpha_{ijk})$ follows a Dirichlet-multinomial distribution, so

$$P(\mathcal{D} \mid \alpha_{ijk}) = \left[\prod_{j=1}^{q_i} \frac{n_{ij}!}{\prod_{k=1}^{r_i} n_{ijk}!} \right].$$

$$\left[\prod_{j=1}^{q_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(n_{ij} + \alpha_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(n_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})} \right] \propto \text{BD}(X_i \mid \Pi_{X_i}^{\mathcal{G}}; \alpha_{ijk})$$

making the **link between BD scores and entropy** explicit.

BDeu and the Maximum Entropy Principle

In the case of BDeu, $P(\alpha_{ijk} = \alpha/(r_i q_i)) = 1$ and learning DAGs based on sparse data following the maximum (relative) entropy means

$$\begin{aligned} \mathbb{E} \left(H^{\mathcal{G}^-} (X_i) \mid \mathcal{D}, \alpha_{ijk} \right) \text{BDeu} \left(X_i \mid \Pi_{X_i}^{\mathcal{G}^-}; \alpha_{ijk} \right) \leq \\ \mathbb{E} \left(H^{\mathcal{G}^+} (X_i) \mid \mathcal{D}, \alpha_{ijk} \right) \text{BDeu} \left(X_i \mid \Pi_{X_i}^{\mathcal{G}^+}; \alpha_{ijk} (\tilde{q}_i / q_i) \right) \end{aligned}$$

whereas it should be

$$\begin{aligned} \mathbb{E} \left(H^{\mathcal{G}^-} (X_i) \mid \mathcal{D}, \alpha_{ijk} \right) \text{BDeu} \left(X_i \mid \Pi_{X_i}^{\mathcal{G}^-}; \alpha_{ijk} \right) \leq \\ \mathbb{E} \left(H^{\mathcal{G}^+} (X_i) \mid \mathcal{D}, \alpha_{ijk} \right) \text{BDeu} \left(X_i \mid \Pi_{X_i}^{\mathcal{G}^+}; \alpha_{ijk} \right) \end{aligned}$$

so structure learning with **BDeu may deviate from the maximum (relative) entropy principle when computed from sparse data.**

BDs does not.

Exhibit A, One Last Time

Combining BDeu with $E(H^{\mathcal{G}}(X_i) | \mathcal{D}, \alpha_{ijk})$ gives

$$E(H^{\mathcal{G}^-}(X) | \mathcal{D}) = 2.066 \cdot 3.906 \times 10^{-7} = 8.071 \times 10^{-7} > \\ 1.514 \times 10^{-7} = 4.069 \cdot 3.731 \times 10^{-8} = E(H^{\mathcal{G}^+}(X) | \mathcal{D})$$

while BDs gives

$$E(H^{\mathcal{G}^-}(X) | \mathcal{D}) = 2.066 \cdot 3.906 \times 10^{-7} = 8.071 \times 10^{-7} = \\ 8.071 \times 10^{-7} = 2.066 \cdot 3.906 \times 10^{-7} = E(H^{\mathcal{G}^+}(X) | \mathcal{D}).$$

Exhibit B, One Last Time

Combining BDeu with $E(H^{\mathcal{G}}(X_i) | \mathcal{D}, \alpha_{ijk})$ gives

$$E(H^{\mathcal{G}^-}(X) | \mathcal{D}) = 0.3931 \cdot 0.0326 = 0.0128 < \\ 0.0252 = 0.5707 \cdot 0.0441 = E(H^{\mathcal{G}^+}(X) | \mathcal{D})$$

while BDs gives

$$E(H^{\mathcal{G}^-}(X) | \mathcal{D}) = 0.3931 \cdot 0.0326 = 0.0128 = \\ 0.0128 = 0.3931 \cdot 0.0326 = E(H^{\mathcal{G}^+}(X) | \mathcal{D}).$$

Summary and Conclusions

- BDeu can be problematic for small/large values of the imaginary sample size; we found that **BDeu can also be problematic regardless if the data are sparse.**
- Then we proposed **BDs as a minimalistic fix** which prevents the imaginary sample size from partially vanishing when there are unobserved parent configurations.
- But is BDeu just not working very well, or **is it methodologically wrong to use it with sparse data?** (Many statistical methods that are methodologically correct but do not work very well on sparse data.)
- One way of looking at this problem is in the context of maximum (relative) entropy. Given the same information in the prior, and the same information from the data, **the assumptions behind BDeu can give a rank more complex, singular BNs over simpler ones.**

References I



A. Caticha.

Relative entropy and inductive inference.

In *Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, pages 75–96, 2004.



M. Feder.

Maximum Entropy as a Special Case of the Minimum Description Length Criterion.

IEEE Transactions on Information Theory, 32(6):847–849, 1986.



A. Giffin and A. Caticha.

Updating Probabilities with Data and Moments.

In *Proceedings of the 27th International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, pages 74–84, 2007.



D. Heckerman, D. Geiger, and D. M. Chickering.

Learning Bayesian Networks: The Combination of Knowledge and Statistical Data.

Machine Learning, 20(3):197–243, 1995.

Available as Technical Report MSR-TR-94-09.



I. Nemenman, F. Shafee, and W. Bialek.

Entropy and Inference, Revisited.

In *Proceedings of the 14th Advances in Neural Information Processing Systems (NIPS) Conference*, pages 471–478, 2002.



G. Schwarz.

Estimating the Dimension of a Model.

The Annals of Statistics, 6(2):461–464, 1978.

References II



J. E. Shore and R. W. Johnson.

Axiomatic Derivation of the Principle of Maximum Entropy and the Principle of Minimum Cross-Entropy.
IEEE Transactions on Information Theory, IT-26(1):26–37, 1980.



T. Silander, P. Kontkanen, and P. Myllymäki.

On Sensitivity of the MAP Bayesian Network Structure to the Equivalent Sample Size Parameter.
In Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence, pages 360–367, 2007.



J. Skilling.

The Axioms of Maximum Entropy.

In Maximum-Entropy and Bayesian Methods in Science and Engineering, pages 173–187, 1988.



H. Steck and T. S. Jaakkola.

On the Dirichlet Prior and Bayesian Regularization.

In Advances in Neural Information Processing Systems 15, pages 713–720. 2003.



M. Ueno.

Learning Networks Determined by the Ratio of Prior and Data.

In Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence, pages 598–605, 2010.



M. Ueno.

Robust Learning of Bayesian Networks for Prior Belief.

In Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence, pages 698–707, 2011.