

# Learning the Bayesian Network Structure: Dirichlet Prior versus Data

Harald Steck

IKM CKS Siemens Medical Solutions, USA Harald.Steck@siemens.com

Page 1/22



#### Where does the le



&



# DATA: Independence

# DIRICHLET PRIOR: -> RESULT: Independence Dependence

Image and Knowledge Solutions CKS-HP /USA / MED



## **Outline**

# Formal Problem Statement and Notation

## Explanation, Part I

# Explanation, Part II

# **Review: Dirchlet Prior and BDeu score**

- random variables with multinomial distribution
- Bayesian network model:

$$p(X_1, \dots, X_n) = \prod_{i=1}^n \theta_{X_i \mid \prod_i}$$

- Dirichlet prior over model parameters:
  - … is conjugate prior
  - … ensures likelihood equivalence [Heckerman et al., '95]

$$p( heta_{X_i|\pi_i}) \propto \prod_{x_i} heta_{x_i|\pi_i}^{lpha_{x_i,\pi_i}-1}$$

where:

:  $\alpha_{x_i,\pi_i} = \alpha \cdot q_{x_i,\pi_i}$   $\alpha > 0 \dots$  equivalent sample size (ESS)  $q \dots$  prior distribution

Page 4/22

#### SIEMENS Review: Dirchlet Prior and BDeu score (cont'd)

Scoring function for graph G: marginal likelihood

$$p(D|\alpha, G) = \prod_{i=1}^{n} \prod_{\pi_i} \frac{\Gamma(\alpha_{\pi_i})}{\Gamma(N_{\pi_i} + \alpha_{\pi_i})} \prod_{x_i} \frac{\Gamma(N_{x_i, \pi_i} + \alpha_{x_i, \pi_i})}{\Gamma(\alpha_{x_i, \pi_i})}$$

#### BDeu Score [Buntine, 1991]:

- choose q uniform
- ESS is the only free parameter in scoring function:

$$\alpha_{x_i,\pi_i} = \alpha \cdot q_{x_i,\pi_i} = \frac{\alpha}{|X_i| \cdot |\Pi_i|}$$

Image and Knowledge Solutions CKS-HP /USA / MED

Page 5/22

#### Result: How many Edges ?

Experiments [Silander et al., UAI 2007]:





## **Outline**

## Formal Problem Statement and Notation

# Part I: Large Equivalent Sample Size

# Part II: 'Optimal' Equivalent Sample Size

#### Absolute Score -> Relative Score



log Bayes Factor (log BF):

$$\log \frac{p(D|\alpha, G^+)}{p(D|\alpha, G^-)} = \sum_{a,b,\pi} \log \frac{\Gamma(N_{a,b,\pi} + \alpha_{a,b,\pi})}{\Gamma(\alpha_{a,b,\pi})} + \sum_{\pi} \log \frac{\Gamma(N_{\pi} + \alpha_{\pi})}{\Gamma(\alpha_{\pi})}$$
$$- \sum_{a,\pi} \log \frac{\Gamma(N_{a,\pi} + \alpha_{a,\pi})}{\Gamma(\alpha_{a,\pi})} - \sum_{b,\pi} \log \frac{\Gamma(N_{b,\pi} + \alpha_{b,\pi})}{\Gamma(\alpha_{b,\pi})}$$

Image and Knowledge Solutions CKS-HP /USA / MED

Page 8/22

#### Absolute Score -> Relative Score



log Bayes Factor (log BF):

log BF> 0  $\Leftrightarrow$  edge present  $\Leftrightarrow$  A,B dependent cond. on  $\Pi$ log BF< 0  $\Leftrightarrow$  edge absent  $\Leftrightarrow$  A,B independent cond. on  $\Pi$ 

# **Bayes Factor for large & finite ESS**

Leading-order Taylor-expansion for  $\alpha >> N$  :

$$\log \frac{p(D|\alpha, G^+)}{p(D|\alpha, G^-)} = \frac{N}{2\alpha} \left\{ N \cdot U(\hat{p}(A, B|\Pi)) - d_{\mathsf{F}} \right\} + \mathcal{O}\left(\frac{N^2}{\alpha^2}\right)$$

where 
$$\widehat{p}(a, b|\pi) = N_{a,b,\pi}/N_{\pi}$$
  
 $d_{\mathsf{F}} = |\Pi|(|A|-1)(|B|-1)$ 

$$N \cdot U(\hat{p}(A, B|\Pi)) - d_{\mathsf{F}} > 0$$
  
 $\Leftrightarrow$  edge  $A \leftarrow B$  present for suff. large  $\alpha$ 

# new Uniformity Measure U

 $U(n(A \mid B \mid \Pi))$ 

$$= \sum_{a,b,\pi} p(a,b,\pi) \Big( |A,B,\Pi| \cdot p(a,b,\pi) - |A,\Pi| \cdot p(a,\pi) \Big)$$

$$-|B, \Pi| \cdot p(b, \pi) + |\Pi| \cdot p(\pi)$$

- symmetry:  $U(p(A, B|\Pi)) = U(p(B, A|\Pi))$
- non-negativity:  $U(p(A, B|\Pi)) \ge 0$
- minimality:  $U(p(A, B|\Pi)) = 0$  if and only if
  - \* (conditional) independence:  $p(A, B|\Pi) = p(A|\Pi)p(B|\Pi)$
  - \* and, for each state  $\pi$  with  $p(\pi) > 0$ , at least one of the marginal distributions is uniform:

 $p(A|\pi) = 1/|A|$  or  $p(B|\pi) = 1/|B|$ .

#### Presence of Edge A <- B



# **Skewness ---- An Intuitive Explanation**

Average as regularized parameter estimate:

$$\bar{\theta}_{x_i|\pi_i} \equiv E_{p(\theta_{x_i|\pi_i}|D,G)}[\theta_{x_i|\pi_i}] = \frac{N_{x_i,\pi_i} + \alpha_{x_i,\pi_i}}{N_{\pi_i} + \alpha_{\pi_i}}$$

... weighted sum of uniform prior *q* and (skewed) empirical distribution

 $N \cdot \hat{p}(A)\hat{p}(B) + \alpha \cdot q(A)q(B) \neq c \cdot p(A)p(B)$ 

=> In general: dependence



#### Outline

## Formal Problem Statement and Notation

#### Part I: Large Equivalent Sample Size

Part II: 'Optimal' Equivalent Sample Size: What properties of the data determine the value of the optimal ESS?

# 'Optimal' ESS

- ESS treated as an additional parameter to be learned
- Objective [Silander et al., UAI 2007]:

$$(\alpha^*, G^*) = \arg \max_{(\alpha, G)} p(D|\alpha, G)$$

# Coordinate-wise ascent:

 optimize the graph for a fixed ESS-value: G<sup>\*</sup><sub>k</sub> = arg max<sub>G</sub> p(D|α<sup>\*</sup><sub>k-1</sub>,G),

 optimize the ESS-value for a fixed graph:

 (D| G<sup>\*</sup>)

$$\alpha_k^* = \arg \max_{\alpha} p(D|\alpha, G_k^*).$$

# **Approximation of optimal ESS**



Page 16/22

Image and Knowledge Solutions CKS-HP /USA / MED

# Approximation of optimal ESS

Approximation for optimal ESS given graph G:  $\alpha^* \approx \frac{d_G^{\text{eff}}}{E_{\hat{p}(X)}[\log \hat{p}(X|G)] - E_{q(X)}[\log \hat{p}(X|G)]} + O\left(\frac{\alpha^{*2}}{N}\right)$ 

 $d_G^{\text{eff}}$  ... effective number of parameters in BN

Denominator is positive for uniform *q* :

 $E_{\hat{p}(X)}[\log \hat{p}(X|G)] - E_{q(X)}[\log \hat{p}(X|G)]$ =  $H(q(X|G)) - H(\hat{p}(X|G)) + \mathsf{KL}(q(X|G)||\hat{p}(X|G))$ 

# **Properties of optimal ESS**

Approximation for optimal ESS given graph G:  $\alpha^* \approx \frac{d_G^{\text{eff}}}{E_{\hat{p}(X)}[\log \hat{p}(X|G)] - E_{q(X)}[\log \hat{p}(X|G)]} + O\left(\frac{\alpha^{*2}}{N}\right)$ 

 $d_G^{\text{eff}}$  ... effective number of parameters in BN

Large skewness or dependences in data:

- $\Rightarrow$  large ML
- $\Rightarrow$  optimal ESS is small

# **Properties of optimal ESS**

Approximation for optimal ESS given graph G:  $\alpha^* \approx \frac{d_G^{\text{eff}}}{E_{\hat{p}(X)}[\log \hat{p}(X|G)] - E_{q(X)}[\log \hat{p}(X|G)]} + O\left(\frac{\alpha^{*2}}{N}\right)$ 

 $d_G^{\text{eff}}$  ... effective number of parameters in BN

Sample size *N* : has no explicit impact  $\Rightarrow$  for large N, we have  $\alpha^* \approx const$ 

# **Properties of optimal ESS**

Approximation for optimal ESS given graph G:  $\alpha^* \approx \frac{d_G^{\text{eff}}}{E_{\hat{p}(X)}[\log \hat{p}(X|G)] - E_{q(X)}[\log \hat{p}(X|G)]} + O\left(\frac{\alpha^{*2}}{N}\right)$ 

 $d_G^{\text{eff}}$  ... effective number of parameters in BN

## Number *n* of nodes in graph:

enumerator and denominator are additive in the number of nodes

 $\Rightarrow$  optimal ESS approx. unaffected by *n* 

k

## **Experimental Validation of Approximations**

	Data	N	n	$\alpha^{I}$	$\alpha^{M}$	$\alpha^*$
	balance	625	5	1100	48	44
	iris	150	5	13	2	2
	thyroid	215	6	22	2	3
Comparison to exact	liver	345	7	36	4	3
results from [Silander	ecoli	336	8	710	8	8
	abalone	4177	9	66	6	7
et al., UAI 2007]:	diabetes	768	9	35	4	3
	post op	90	9	35	3	3
	yeast	1484	9	16	6	6
->	cancer	286	10	610	8	7
_/	shuttle	58000	10	13	3	3
	tictac	958	10	5162	51	60
	bc wisc	699	11	715	8	5
$\alpha^*$ has correct order of	glass	214	11	56	6	6
	page	5473	11	33	3	3
maanitude	heart cl	303	14	1316	13	9
<b>9</b> • • • • • • • • • • • • • • • • • • •	heart hu	294	14	56	5	5
	heart st	270	14	710	10	10
	wine	178	14	88	8	7
	adult	32561	15	4858	50	49

Page 21/22

Image and Knowledge Solutions CKS-HP /USA / MED



#### Conclusions

# Iarge ESS:

even if data and prior imply independence, the **complete** graph may have the highest marginal likelihood.

# • optimal' ESS:

- approx. independent of sample size and number of nodes
- in data
- is large if small skewness and weak dependencies in data