

# Large Sample Robustness Bayes Nets with Incomplete Information

Jim Smith and Ali Daneshkhan

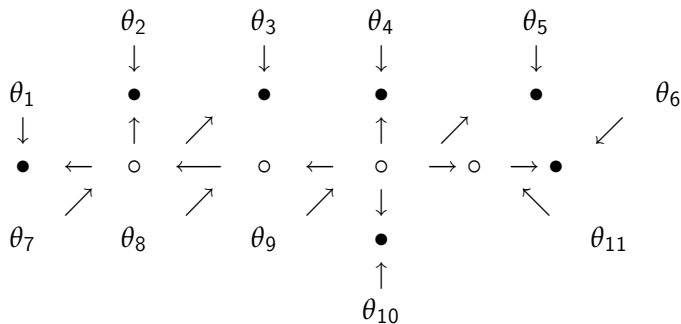
Universities of Warwick and Strathclyde

Durham July 2008

- We often worry about convergence of samplers etc. in a Bayesian analysis. But how much does the precise specification of the prior matter in an analysis of a BN?
- In particular what is the overall effect of local and global independence assumptions on a given model?
- What are the overall inferential implications of using standard priors like product Dirichlets or product logistics?
- In general how hard do I need to think about these issues a priori when I know I will collect a large sample?

# Messy Analyses

- We have a large BN with some expert knowledge incorporated.
- Many of the nodes in our graph are systematically missing, and the sample not random. So even taking account of aliasing, some features of the model may be unidentified even as data size increases



# The Problems

- We only usually have a **numerical** or **algebraic approximation** of our posterior density with respect to our chosen prior. So essentially we see various **approximate summary statistics** (e.g. means, variances, sampled low dimensional margins, ...)
- Even for standard complete sampling on identified systems, there are still robustness issues. **Variation distance**  $d_V(f, g) = \int |f - g|$  between two posteriors can **diverge quickly as sample size increases**, especially when the parameter space is large. This phenomenon might occur whenever our data set contains an outlier (Dawid, 1973) but also more generally (Gustafson and Wasserman, 1995) for local priors.
- In the complex scenario above it is therefore far from clear when posterior inferences are strongly influenced by how we specify our prior.
- **Local De Robertis** separations the key to addressing this issue!

- Local De Robertis (LDR) separations are **easy to calculate** and extend natural parametrizations in exponential families.
- They have an intriguing **prior to posterior invariance** property.
- BN factorization of a density implies **linear** relationships between **clique** marginal separations and joint.
- **Bounds on the variation distance** between two posterior distributions associated with different priors can be **calculated explicitly** as a function of prior LDR bounds and statistics associated with the calculated posterior associated with the functioning prior.
- The **bounds apply** posterior to an observed likelihood, even **when the sample density is misspecified**.

- De Robertis local Separations
- Some Properties of Local De Robertis Separations
- Some useful Theorems concerning LDR and BNs.
- What this means for the robustness of BN's

# The Setting

- Let  $g_0, (g_n)$  our **genuine** prior (posterior) density :  $f_0, (f_n)$  our **functioning** prior (posterior) density
- Default for Bayes  $f_0$  often products of Dirichlets
- $\mathbf{x}_n = (x_1, x_2, \dots, x_n)$ ,  $n \geq 1$ . with *observed* sample densities  $\{p_n(\mathbf{x}_n|\boldsymbol{\theta})\}_{n \geq 1}$ ,
- With missing data, typically these sample densities are typically  $\{p_n(\mathbf{x}_n|\boldsymbol{\theta})\}_{n \geq 1}$  (and hence  $f_n$  and  $g_n$ ) intractable
- $f_n$  therefore approximated either by drawing samples or algebraically.

# A Bayes Rule Identity

Let  $\Theta(n) = \{\theta \in \Theta : p(\mathbf{x}_n|\theta) > 0\}$  For all  $\theta \in \Theta(n)$  then

$$\log g_n(\theta) = \log g_0(\theta) + \log p_n(\mathbf{x}_n|\theta) - \log p_g(\mathbf{x}_n)$$

$$\log f_n(\theta) = \log f_0(\theta) + \log p_n(\mathbf{x}_n|\theta) - \log p_f(\mathbf{x}_n)$$

where

$$p_g(\mathbf{x}_n) = \int_{\theta \in \Theta(n)} p(\mathbf{x}_n|\theta) g_0(\theta) d\theta, \quad p_f(\mathbf{x}_n) = \int_{\theta \in \Theta(n)} p(\mathbf{x}_n|\theta) f_0(\theta) d\theta,$$

(When  $\theta \in \Theta \setminus \Theta(n)$  set  $g_n(\theta) = f_n(\theta) = 0$ )

So

$$\log f_n(\theta) - \log g_n(\theta) = \log f_0(\theta) - \log g_0(\theta) + \log p_g(\mathbf{x}_n) - \log p_f(\mathbf{x}_n)$$



# From Bayes Rule to LDR

For any subset  $A \subseteq \Theta(n)$  let

$$d_A^L(f, g) \triangleq \sup_{\theta \in A} (\log f(\theta) - \log g(\theta)) - \inf_{\phi \in A} (\log f(\phi) - \log g(\phi))$$

Then since

$$\log f_n(\theta) - \log g_n(\theta) = \log f_0(\theta) - \log g_0(\theta) + \log p_g(\mathbf{x}_n) - \log p_f(\mathbf{x}_n)$$

for any sequence  $\{p(\mathbf{x}_n|\theta)\}_{n \geq 1}$  - however complicated -

$$d_A^L(f_n, g_n) = d_A^L(f_0, g_0)$$

$$d_A^L(f_n, g_n) = d_A^L(f_0, g_0)$$

- So for  $A \subseteq \Theta(n)$  the posterior quality of the approximation of  $f_n$  to  $g_n$  is identical to that of  $f_0$  to  $g_0$ .

$$d_A^L(f_n, g_n) = d_A^L(f_0, g_0)$$

- So for  $A \subseteq \Theta(n)$  the posterior quality of the approximation of  $f_n$  to  $g_n$  is identical to that of  $f_0$  to  $g_0$ .
- When  $A = \Theta(n)$  this property (De Robertis, 1978) used for density ratio metrics and the specification of neighbourhoods.

$$d_A^L(f_n, g_n) = d_A^L(f_0, g_0)$$

- So for  $A \subseteq \Theta(n)$  the posterior quality of the approximation of  $f_n$  to  $g_n$  is identical to that of  $f_0$  to  $g_0$ .
- When  $A = \Theta(n)$  this property (De Robertis, 1978) used for density ratio metrics and the specification of neighbourhoods.
- Trivially posterior distances between densities can be calculated effortlessly from priors.

$$d_A^L(f_n, g_n) = d_A^L(f_0, g_0)$$

- So for  $A \subseteq \Theta(n)$  the posterior quality of the approximation of  $f_n$  to  $g_n$  is identical to that of  $f_0$  to  $g_0$ .
- When  $A = \Theta(n)$  this property (De Robertis, 1978) used for density ratio metrics and the specification of neighbourhoods.
- Trivially posterior distances between densities can be calculated effortlessly from priors.
- Separation of two priors lying in standard families can usually be expressed explicitly and always explicitly bounded.

# Some notation

We will be especially interested in small sets  $A$ .

- Let  $B(\boldsymbol{\mu}; \rho)$  denote the open ball centred at  $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_k)$  and of radius  $\rho$
- Let

$$d_{\boldsymbol{\mu}; \rho}^L(f, g) \triangleq d_{B(\boldsymbol{\mu}; \rho)}^L(f, g)$$

- For any subset  $\Theta_0 \subseteq \Theta$ , let

$$d_{\Theta_0; \rho}^L(f, g) = \sup_{\boldsymbol{\mu} \in \Theta_0} d_{\boldsymbol{\mu}; \rho}^L(f, g)$$

- Obviously for any  $A \subseteq B(\boldsymbol{\mu}; \rho)$ ,  $\boldsymbol{\mu} \in \Theta_0 \subseteq \Theta$ ,

$$d_A^L(f, g) \leq d_{\Theta_0; \rho}^L(f, g)$$

# Separation of two Dirichlets

- Let  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$   $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_k)$ ,  $\theta_i, \alpha_i > 0$ ,  $\sum_{i=1}^k \theta_i = 1$
- Let  $f_0(\boldsymbol{\theta}|\boldsymbol{\alpha}_f)$  and  $g_0(\boldsymbol{\theta}|\boldsymbol{\alpha}_g)$  be Dirichlet( $\boldsymbol{\alpha}$ ) so that

$$f_0(\boldsymbol{\theta}|\boldsymbol{\alpha}_f) \propto \prod_{i=1}^k \theta_i^{\alpha_{i,f}-1}, \quad g_0(\boldsymbol{\theta}|\boldsymbol{\alpha}_g) \propto \prod_{i=1}^k \theta_i^{\alpha_{i,g}-1}$$

- Let  $\boldsymbol{\mu}_n = (\mu_{1,n}, \mu_{2,n}, \dots, \mu_{k,n})$  be the mean of  $f_n$  if  $\rho_n < \mu_n^0 = \min \{ \mu_n : 1 \leq i \leq k \}$

$$d_{\mu_n^0, \rho_n}^L(f_0, g_0) \leq 2k\rho_n (\mu_n^0 - \rho_n)^{-1} \bar{\alpha}(f_0, g_0)$$

where

$$\bar{\alpha}(f_0, g_0) = k^{-1} \sum_{i=1}^k |\alpha_{i,f} - \alpha_{i,g}|$$

is the average distance between hyperparameters of  $f_0$  and  $g_0$ .

## Where Separations might be large

$$d_{\mu; \rho_n}^L(f_0, g_0) \leq 2\rho_n (\mu_n^0 - \rho_n)^{-1} \sum_{i=1}^k |\alpha_{i,f} - \alpha_{i,g}|$$

- So  $d_{\mu; \rho_n}^L(f_0, g_0)$  is uniformly bounded whenever  $\mu_n$  all away from 0 and converging approximately linearly in  $n$ .
- OTOH if  $f_n$  tends to mass near a zero probability, then even when  $\bar{\alpha}(f, g)$  is small, it can be shown that at least some likelihoods will force the variation distance between the posterior densities to stay large for increasing  $n$ : Smith(2007). The smaller the smallest probability the slower any convergence can be.



# BN's with local and global independence

From defn. if functioning prior  $f(\boldsymbol{\theta})$  and genuine prior  $g(\boldsymbol{\theta})$  factorize on subvectors  $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_k\}$  so that

$$f(\boldsymbol{\theta}) = \prod_{i=1}^k f_i(\boldsymbol{\theta}_i), \quad g(\boldsymbol{\theta}) = \prod_{i=1}^k g_i(\boldsymbol{\theta}_i)$$

where  $f_i(\boldsymbol{\theta}_i)$  ( $g_i(\boldsymbol{\theta}_i)$ ) are the functioning (genuine) margin on  $\boldsymbol{\theta}_i$ ,  $1 \leq i \leq k$ , then (like K-L separations)

$$d_A^L(f, g) = \sum_{i=1}^k d_{A_i}^L(f_i, g_i)$$

So local prior distances grow linearly with no. of defining conditional probability vectors.

# Some conclusions

- BN's with smaller nos of edges intrinsically **more stable** and the effects of possibly erroneous prior information will endure longer than more complex models encoding less conditional independences.
- However - as with K-L - marginal densities are never more separated than their joint densities - so if a utility is only on a particular margin then these distances may be much less.
- Bayes Factors automatically select simpler models but note also inferences of a more complex model tends to be more sensitive to wrongly specified priors.

- There are certain features in the prior which will always endure.
- If there is a point where locally LDR diverges - in a sense which violates the condition above then it is possible to construct a "regular" likelihood such that the variation distance between posteriors remains bounded away from zero as  $n \rightarrow \infty$ .
- However if the mass is converging on to a small set because then we can focus on a small set  $A$
- Usually  $d_A^L(f_0, g_0)$  is small when  $A$  lies in a small ball.

# Salvation!

- When  $n$  is large  $A$  will lie in a small ball with high probability
- it is usually reasonable to assume that  $f_0$  and  $g_0$  for  $A$  lying in a small ball  $d_A^L(f_0, g_0)$  is small.
- Can usually assume for open balls  $B(\boldsymbol{\mu}; \rho)$  centred at  $\boldsymbol{\mu}$  and of radius  $\rho$ ,  $f_0, g_0 \in \mathcal{F}(\Theta_0, M(\Theta_0), \rho(\Theta_0))$  meaning

$$\sup_{\boldsymbol{\theta}, \boldsymbol{\phi} \in B(\boldsymbol{\mu}; \rho)} |\log f_0(\boldsymbol{\theta}) - \log f_0(\boldsymbol{\phi})| \leq M(\Theta_0) \rho^{0.5\rho(\Theta_0)}$$

$$\sup_{\boldsymbol{\theta}, \boldsymbol{\phi} \in B(\boldsymbol{\mu}; \rho)} |\log g_0(\boldsymbol{\theta}) - \log g_0(\boldsymbol{\phi})| \leq M(\Theta_0) \rho^{0.5\rho(\Theta_0)}$$

# A simple smoothness/roughness condition

- When  $p(\Theta_0) = 2$  just demands that  $\log f_0$  and  $\log g_0$  both have bounded derivatives within the set  $\Theta_0$  - used to determine where  $f_n$  concentrates its mass. Then it is easily shown (see Smith and Rigat, 2008) that

$$d_{\Theta_0, \rho}^L(f, g) \leq 2M(\Theta_0)\rho^{1/2p(\Theta_0)}$$

- So **rate** of convergence to zero of  $d_{\Theta_0, \rho}^L(f, g)$  governed by the "roughness" parameter  $p(\Theta_0)$ .
- This is always true for densities with inverse polynomial tails like the **Student  $t$  density**. If densities have tighter tails than this then provided they are continuously differentiable on a closed bounded interval  $\Theta_0$ .
- For continuous  $f, g$  when  $\Theta_0$  closed and bounded  $d_{\Theta_0, \rho}^L(f, g)$  converges to zero (this boundedness prevents divergence due to outliers).

Consider the typical hierarchical models used in e.g. BUGS

$$\begin{array}{ccc} \mathbf{X}_1 & & \mathbf{X}_2 \\ \uparrow & & \uparrow \\ \theta_1 & \leftarrow \theta & \rightarrow \theta_2 \end{array}$$

e.g.  $i = 1, 2$ ,  $\theta_i = \theta \div \varepsilon_i$  where  $\varepsilon_i$  is an independent error term, (Gaussian, Student  $t$ ) etc. provided the error term is smooth then this automatically forces the prior margin  $g_0(\theta_1, \theta_2)$  to be smooth (even if  $\theta$  is discrete) regardless of the smoothness of  $\theta$ .

Moral: nearly all conventional hierarchical BN's with enough depth have implicit priors on parameters of the likelihood are smooth in the sense above (making them robust in the sense below).

## But why worry about LDR separation?

- Without the LDR condition above large sample variation convergence cannot hold in general.
- Conversely with a regularity condition and a technical device convergence will happen. .
- **Regularity Condition.** Call a genuine prior *c-rejectable* if the ratio of marginal likelihood  $\frac{p_g(\mathbf{x})}{p_f(\mathbf{x})} < c$ .

If  $f_0$  does not explain the data much better than  $g_0$  we would expect this ratio to be small - certainly not *c-rejectable* for a moderately large values of  $c \geq 1$ .

## A Second Tail convergence condition

- Say density  $f$   $\Lambda$ -tail dominates a density  $g$  if

$$\sup_{\theta \in \Theta} \frac{g(\theta)}{f(\theta)} = \Lambda < \infty$$

When  $g(\theta)$  is bounded then this condition requires that the tail convergence of  $g$  is no slower than  $f$ .

- Condition met provided  $f_0$  is chosen to have a flatter tail than  $g_0$ .
- Note: flat tailed priors recommended for robustness on other grounds e.g. O'Hagan and Forster (2003)



# A typical result (Smith and Rigat (2007))

## Theorem

If the genuine prior  $g_0$  is not  $c$  rejectable with respect to  $f_0$ ,  $f_0$   $\Lambda$ -tail dominates  $g_0$  and  $f_0, g_0 \in \mathcal{F}(\Theta_0, M(\Theta_0), p(\Theta_0))$ . then

$$d_V(f_n, g_n) \leq \inf_{\rho > 0} \{ T_n(1, \rho_n) + 2T_n(2, \rho_n) : B(\mu_n, \rho_n) \subset \Theta_0 \} \quad (1)$$

where

$$T_n(1, \rho_n) = \exp d_{\mu, \rho}^L(f, g) - 1 \leq \exp \left\{ 2M\rho_n^{p/2} \right\} - 1$$

and

$$T_n(2, \rho_n) = (1 + c\Lambda)F_n(\theta \notin B(\mu_n; \rho_n))$$

Easy to bound  $F_n(\theta \notin B(\mu_n; \rho_n))$  in many ways explicitly using Chebychev type inequalities: Smith (2007). Example of bound is given below, specified in terms of the posterior means and variances of the vector of parameters under  $f_n$  routinely approximated.

# An Example of an Explicit Bound

Let  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$  and  $\mu_{j,n}, \sigma_{jj,n}^2$  denote the mean and variance of  $\theta_j$ ,  $1 \leq j \leq k$  under  $f_n$ . Using Chebychev bounds in Tong (1980), p153), writing  $\mu_n = (\mu_{1,n}, \mu_{2,n}, \dots, \mu_{k,n})$

$$F_n(\theta \notin B(\mu_n; \rho_n)) \leq k\rho_n^{-2} \sum_{j=1}^k \sigma_{jj,n}^2$$

where writing  $\sigma_n^2 = k \max_{1 \leq j \leq k} \sigma_{j,n}^2$  this implies

$$T_n(2, \rho_n) \leq c\Lambda \sigma_n^2 \rho_n^{-2}$$

- e.g. if  $\sigma_n^2 \leq n^{-1}\sigma^2$  for some value  $\sigma^2$ ,  $T_n(2, \rho_n) \rightarrow 0$  provided  $\rho_n^2 \leq n^r \rho^2$  where  $0 < r < 1$ .
- In practice for a given data set we just have an approximate value of  $\sigma_n^2$  we can plug in.

# Inference on margins separation

When  $A_1$  is a restriction of  $A$  to  $\theta_1$ ,  $\theta = (\theta_1, \theta_2)$  and  $f_1(\theta_1), g_1(\theta_1)$  contin. margins of  $f(\theta)$  and  $g(\theta)$ , resp. then

$$d_{A_1}^L(f_1, g_1) \leq d_A^L(f, g)$$

- If  $f_n$  converges on a margin, then even if the model is unidentified, provided  $f_0, g_0 \in \mathcal{F}(\Theta_0, M(\Theta_0), p(\Theta_0))$ , then for large  $n$ ,  $f_n$  will be a good surrogate for  $g_n$ .
- BN's with interior systematically hidden variables are unidentified. However if a utility function is only on manifest variables, in standard scenarios under above conditions  $d_V(f_{1,n}, g_{1,n}) \rightarrow 0$  at a rate of at least  $\sqrt[3]{n}$ .
- Instability only on posteriors of functions of probabilities associated with the hidden variables conditional on the manifest variables.



# Departures from Parameter Independence

$$f(\boldsymbol{\theta}) = f_1(\theta_1) \prod_{i=2}^k f_{i|}.(\theta_i | \boldsymbol{\theta}_{pa_i})$$
$$g(\boldsymbol{\theta}) = g_1(\theta_1) \prod_{i=2}^k g_{i|}.(\theta_i | \boldsymbol{\theta}_{pa_i})$$

we then have the inequality

$$d_A^L(f, g) \leq \sum_{i=2}^k d_{A[i]}^L(f_{[i]}, g_{[i]})$$

where  $f_{[i]}, g_{[i]}$  are respectively the margin of  $f$  and  $g$  on the space  $\Theta[i]$  of the  $i^{\text{th}}$  variable and its parents. So distances bounded by sums on distances on cliques margins.

# Uniformly A Uncertain

Suppose  $g$  is uniformly  $A$  uncertain and factorises as  $f$  and

$$\sup_g \sup_{\theta_i, \phi_i \in A[i]} \{ \log f_{i|}(\theta) - \log g_{i|}(\theta) - \log f_{i|}(\phi) + \log g_{i|}(\phi) \}$$

is not a function of  $\theta_{pa_i}$   $2 \leq i \leq n$ , then we can write

$$d_A^L(f, g) = \sum_{i=1}^k d_{A[i]}^{L*}(f_{i|}, g_{i|})$$

- Separation between the joint densities  $f$  and  $g$  sum of the separation between its component conditionals  $f_{i|}$  and  $g_{i|}$   $1 \leq i \leq k$ .
- Bounds can be calculated *even* when the likelihood destroys the factorisation of the prior. So the critical property we assume here is the fact that we believe a priori that  $f$  respects the same factorisation as  $g$ .

# Conclusions

- Bayesian inference on BN's is most stable to prior settings the simpler the model

# Conclusions

- Bayesian inference on BN's is most stable to prior settings the simpler the model
- For large samples general total variation robustness is lost when posterior masses concentrate near a zero probability.



# Conclusions

- Bayesian inference on BN's is most stable to prior settings the simpler the model
- For large samples general total variation robustness is lost when posterior masses concentrate near a zero probability.
- However robustness can sometimes be retrieved if that probability is not appear in a utility function.

# Conclusions

- Bayesian inference on BN's is most stable to prior settings the simpler the model
- For large samples general total variation robustness is lost when posterior masses concentrate near a zero probability.
- However robustness can sometimes be retrieved if that probability is not appear in a utility function.
- Even for moderate sized samples, explicit bounds on the effects of priors can be calculated on line.

- Bayesian inference on BN's is most stable to prior settings the simpler the model
- For large samples general total variation robustness is lost when posterior masses concentrate near a zero probability.
- However robustness can sometimes be retrieved if that probability is not appear in a utility function.
- Even for moderate sized samples, explicit bounds on the effects of priors can be calculated on line.
- In regular problems, these bounds usually contract surprisingly quickly as data increases.

## A Few references

Daneseshkhan, A (2004) "Estimation in Causal Graphical Models" PhD Thesis University of Warwick.

DeRobertis, L. (1978) "The use of partial prior knowledge in Bayesian inference" Ph.D. dissertation, Yale Univ.

Gustafson, P. and Wasserman, L. (1995) "Local sensitivity diagnostics for Bayesian inference" *Annals Statist*, 23, , 2153 - 2167

French, S. and Rios Insua, D. (2000) "Statistical Decision Theory" Kendall's Library of Statistics Arnold

O'Hagan, A and Forster, J (2004) "Bayesian Inference" Kendall's Advanced Theory of Statistics, Arnold

## A few more References

- Smith, J.Q. "Local Robustness of Bayesian Parametric Inference and Observed Likelihoods" CRiSM Res Rep 07-08
- Smith, J.Q. and Rigat, F.(2008) "Isoseparation and Robustness in Finite Parameter Bayesian Inference" CRiSM Res Rep
- Smith, J.Q. and Croft, J. (2003) "Bayesian networks for discrete multivariate data" J of Multivariate Analysis 84(2), 387 -402
- Tong, Y.L.(1980) "Probability Inequalities in Multivariate Distributions" Academic Press New York
- Wasserman, L.(1992a) "Invariance properties of density ratio priors" Ann Statist, 20, 2177- 2182