# MC<sup>4</sup>: a tempering algorithm for large-sample network inference



D. Barker<sup>1</sup>, S. Hill<sup>1</sup> and S. Mukherjee<sup>2</sup>

<sup>1</sup>Complexity Science DTC, <sup>2</sup> Dept. of Statistics University of Warwick, England CV4 7AL

PRIB Sept. 2010

- Networks (e.g. genes, proteins, metabolites) important notion in current biology.
- Statistical models called Probabilistic Graphical Models (PGM) are a key approach
- Example of signalling network; RTK.



(Weinberg 2007, Yarden & Sliwkowski 2001)

# Graphical Models

- Stochastic Models where a graph is used to describe probability relationships between components.
- Graph specifies the form of **conditional independence statements**.
- Graph must be directed and acyclic (DAG).

Special cases include HMMs, Bayesian Networks (BN), Dynamics BNs.



Interested in the probability of a certain graph G given some observed data **X**.

For certain models it is possible to obtain closed form for **posterior probability** upto a constant.

 $P(G|\mathbf{X}) \propto P(\mathbf{X}|G)P(G)$ 

Maximising  $P(G|\mathbf{X})$  can have robustness problems; If posterior has several highly scoring graphs how do we choose between them?

• For this reason we use **model averaging**.

# Model Averaging

Probability E(e) of seeing an edge e averaged over all graphs G is more robust.

• Edges which repeatedly appear in likely graphs have high E(e).

Knowledge of proportionality constant requires *enumeration* of whole p-node DAG space G.

•  $\mathcal{G}$  grows super-exponentially with p.

Thus we must use MCMC to estimate the posterior probabilities  $P(G|\mathbf{X})$ .

# Monte Carlo

- Move around *G* by performing elementary moves on current graph *G*.
- Accept or reject new graphs G' based on MH acceptance probability;

$$\alpha = \frac{P(\mathbf{X}|G')|\eta(G)|}{P(\mathbf{X}|G)|\eta(G')|}$$

(for uniform priors)

Called MC<sup>3</sup> (Madigan & York 1995)

Addition

Neighbourhood  $\eta(G)$  is all graphs reachable from *G*.

Estimate of posterior probability given by

$$\hat{P}(G|\mathbf{X}) = rac{1}{t_{\max}}\sum_{t=1}^{t_{\max}} I(g^{(t)} = G)$$

# Sample Size





Having more data is clearly a good thing.

• High Throughput exprements, FACS, etc...

# Catuion!

In certain situations large sample size N can cause problems.

 $MC^3$  guaranteed to converge given enough time but can be slow.

#### Motivation

- Posterior for p=4 node system with two different sample sizes N = 5 and N = 10
- Posterior mass concentrates on a few highly likely graphs.
- If these are hard to get between Markov chain mixing is slow.

Note: As  $N \to \infty$  we pick out all graphs from the correct data generating class.



# MC<sup>4</sup> Scheme

- Aim is to allow Markov chains (MC) to move between high scoring graphs.
- Utilise physics approach of **parallel tempering**.
- Couple high temperature MCs to one with desired posterior.



Temperature analogy acheived by raising posterior score to  $\beta = \frac{1}{\tau}$ :

 $P(G|\mathbf{X})^{eta} \propto (P(\mathbf{X}|G)P(G))^{eta}$ 

Set up *m* MCs at temperatures  $T_1, ..., T_m$ .

MCs at higher temperature can explore the space more freely.

• Each chain simulated using often used MH scheme.

Every iteration randomly swap graphs between neighbouring chains i and j with probability  $p_{swap}$ 

• Accept the swap with probability  $\rho$ .

Swapping probability

$$\rho = \frac{(P(\mathbf{X}|G_j)P(G_j))^{\beta_i} (P(\mathbf{X}|G_i)P(G_i))^{\beta_j}}{(P(\mathbf{X}|G_i)P(G_i))^{\beta_i} (P(\mathbf{X}|G_j)P(G_j))^{\beta_j}}$$

# Simulation

First we examine performance on synthetic data generated from the known network shown earlier.

Data is generated using

- $A \sim N(0, \sigma)$  for parent nodes.
- $C \sim N(A + B + \gamma AB, \sigma)$  for child nodes. (with parents A and B)

Since we know the underlying graph from which the data were generated we can draw ROC curves...

# **ROC Curves**

Curves paramterised by threshold t; keep in output graph all edges with E(e) > t.



 $MC^4$  has picked up fewer false positive edges compared to  $MC^3$  for the same number of true edges.

(Xie & Geng 2008)

# **ROC Curves**



13 of 18

# Proteomic Data

Such methods are only useful if they provide a benefit in practical problems.

We examine here the application to inferring the underlying DBN from a set of proteomic data.

Due to certain factorisation for DBNs we can calculate **exact edge probabilities**.

• Gives us gold standard comparison!

We examine;

- Correlation  $\rho$  between the exact and MC estimated edge probabilities.
- Normalised **sum difference** *s* between the exact and MC estimated edge probs.

14 of 18

# Edge Probabilities



T = 1.0, 1.25, 1.5, 1.75, 2.0 and  $p_{swap} = 0.1$ , averaged over 4 runs.

#### Edge Probabilities

If we look at the individual edge probabilities we see better performace (closer to x = y) for MC<sup>4</sup>:



Toughest edges to infer are significantly better estimated by MC<sup>4</sup>.

16 of 18

# Conclusions

- As sample size increases posterior mass can concentrate around several hard to move between graphs.
- Widely employed MCMC schemes can fail to estimate edges properly in these increasingly common situations.
- Counter this by using higher temperature chains coupled to desired posterior: MC<sup>4</sup>/PT.
- Important to draw robust conclusions from biological data.



# Acknowledgements

I would like to thank my co-authors

- Sach Mukherjee,
- Steven Hill,

as well as

- R. Goudie
- M. Nicodemi
- N. Podd

for useful discussions, EPSRC for funding ...

**EPSRC** Pioneering research

and skills

THE UNIVERSITY OF

And finally, thank you for listening.

