

Confident Bayesian Learning of Graphical Models

MIKKO KOIVISTO

Machine Learning Coffee Seminar
22 January 2018

AGENDA

I Graphical Models

II Confident Bayesian Structure Learning

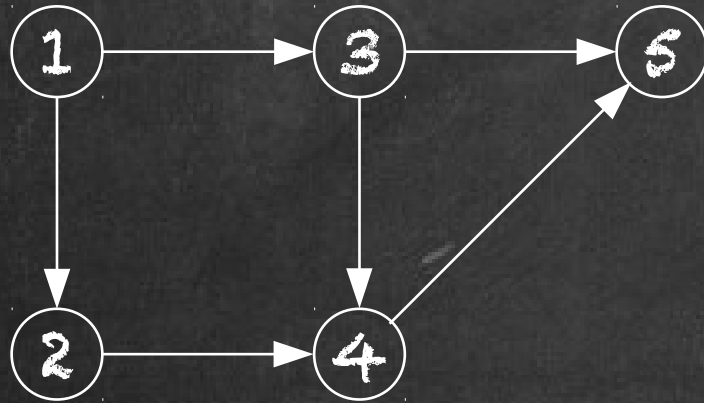
III State of the Art - Summary

IV Computational Techniques & Sample Results

Part I

Graphical Models

BAYESIAN NETWORK (BN)



Structure

G : a DAG

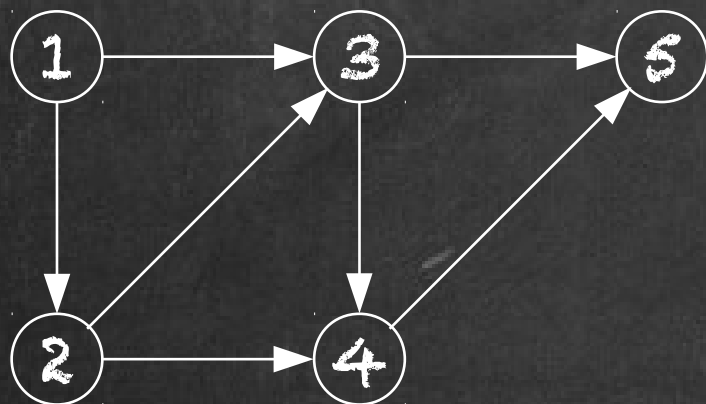
G_v : parents of v

$$G_4 = \{2, 3\}$$

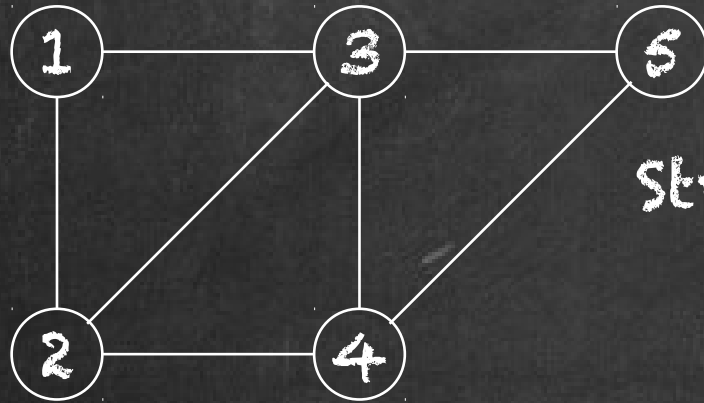
Local conditional distributions

$$p(\mathbf{x}) = \prod_{v=1}^d p(x_v | x_{G_v})$$

CHORDAL MARKOV NETWORK (CMN)



CHORDAL MARKOV NETWORK (CMN)



Structure

G : a chordal graph

C : a clique in G

Cliques: $\{1, 2, 3\}$,
 $\{2, 3, 4\}$, $\{3, 4, 5\}$

Local marginal distributions

Separators: $\{2, 3\}$,
 $\{3, 4\}$

$$p(\mathbf{x}) = \frac{\prod_{\text{clique } C} p(\mathbf{x}_C)}{\prod_{\text{separator } S} p(\mathbf{x}_S)}$$

CHORDAL MARKOV NETWORK (CMN)

Clique tree



Structure

G : a chordal graph

C : a clique in G

Cliques: $\{1, 2, 3\}$,
 $\{2, 3, 4\}$, $\{3, 4, 5\}$

Local marginal distributions

Separators: $\{2, 3\}$,
 $\{3, 4\}$

$$p(\mathbf{x}) = \frac{\prod_{\text{clique } C} p(\mathbf{x}_C)}{\prod_{\text{separator } S} p(\mathbf{x}_S)}$$

Part II

Confident Bayesian Structure Learning

BAYESIAN STRUCTURE LEARNING

Given a data set $X = (x^1, \dots, x^N) = (X_1, \dots, X_d)$, summarize the posterior of G .

Mode

$$\operatorname{argmax} P(G) P(X | G)$$

Expectations

$$E[f(G) | X]$$

Example: $f(G) = 1\{u \text{ is a parent of } v\}$

Normalizing constant

$$\sum_G P(G) P(X | G) = P(X)$$

CONFIDENT BAYESIAN LEARNING

Attach the estimate of the target quantity Q with a characterization of accuracy.

Example (confidence interval):

Output an interval I and a number p such that Q belongs to I with probability at least p

Example (exact deterministic):

$$I = \{Q\} \text{ and } p = 1$$

Characterization
required only
a posteriori!

Part III

State of the Art - Summary

BAYESIAN NETWORKS

NP-hard [Chickering 1996, Chickering et al. 2004]

Exact algorithms, time $2^d \text{poly}(d)$ or $3^d \text{poly}(d)$
[Ott et al. 2004, K. & Sood 2004, Singh & Moore 2005, Silander & Myllymäki 2006, Tian & He 2009]

Complete solvers based on B&B, A*, ILP, constraint programming [de Campos & Ji 2011, Yuan & Malone 2013, Bartlett & Cussens 2015, van Beek & Hoffmann]

Tractable subclasses

[Gaspers et al. 2014, Korhonen & Parviainen 2015]

Markov chain Monte Carlo (MCMC) [Madigan & York 1995, Friedman & Koller 2003, Niinimäki et al. 2011, 2015, 2016]

Mode only!

Mode only!

Mostly
unconfident!

CHORDAL MARKOV NETWORKS

NP-hard [Srebro 2003]

Exact algorithms, time $4^d \text{poly}(d)$
[Kangas et al. 2014, 2015]

Complete solvers based on constraint programming, ILP, B&B [Corander et al. 2013, Bartlett & Cussens 2015, Rantanen et al. 2017]

For bounded-treewidth networks [Korkkari & Parviainen 2013, Berg et al. 2014, Parviainen et al. 2014]

Markov chain Monte Carlo (MCMC)
[Giudici & Green 1999, Green & Thomas 2013]

Mode only!

Mode only!

Unconfident!

Part IV

Computational Techniques & Sample Results

THE COVERING TECHNIQUE

Cover the **structure space** by "nice-behaving" sets.

Example (MCMC for BNs):

$$\{\text{DAGs}\} = \bigcup_{\text{constraint } \mathcal{D}} \{\text{DAGs compatible with } \mathcal{D}\}$$

Constraint families

DAGs

[Madigan & York 1995]

Linear orders

[Friedman & Koller 2003]

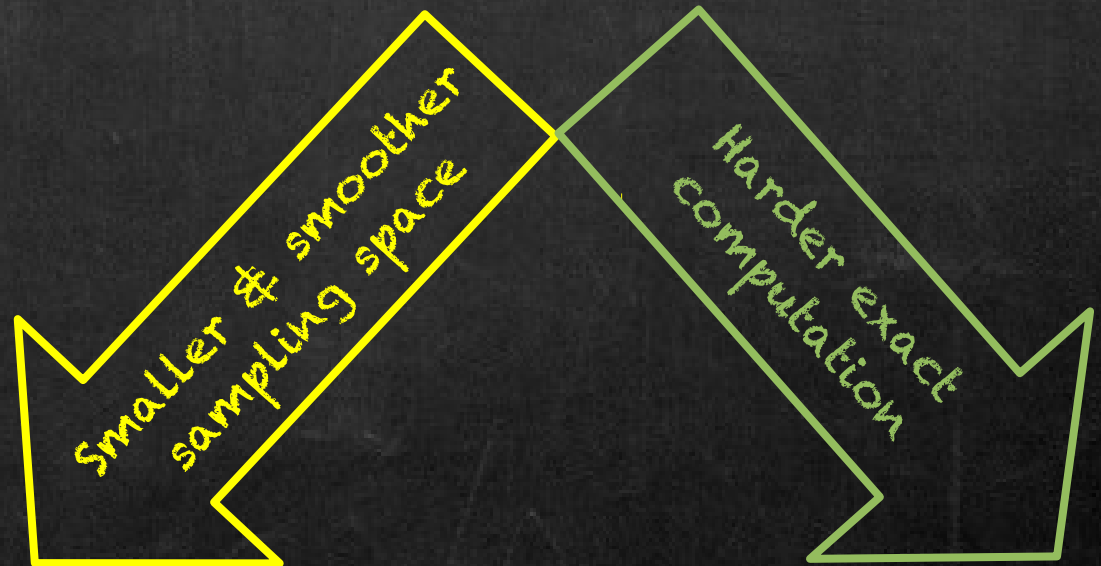
Partial orders

[Niinimäki et al. 2011-16]

No constraints

[K. & Sood 2004]

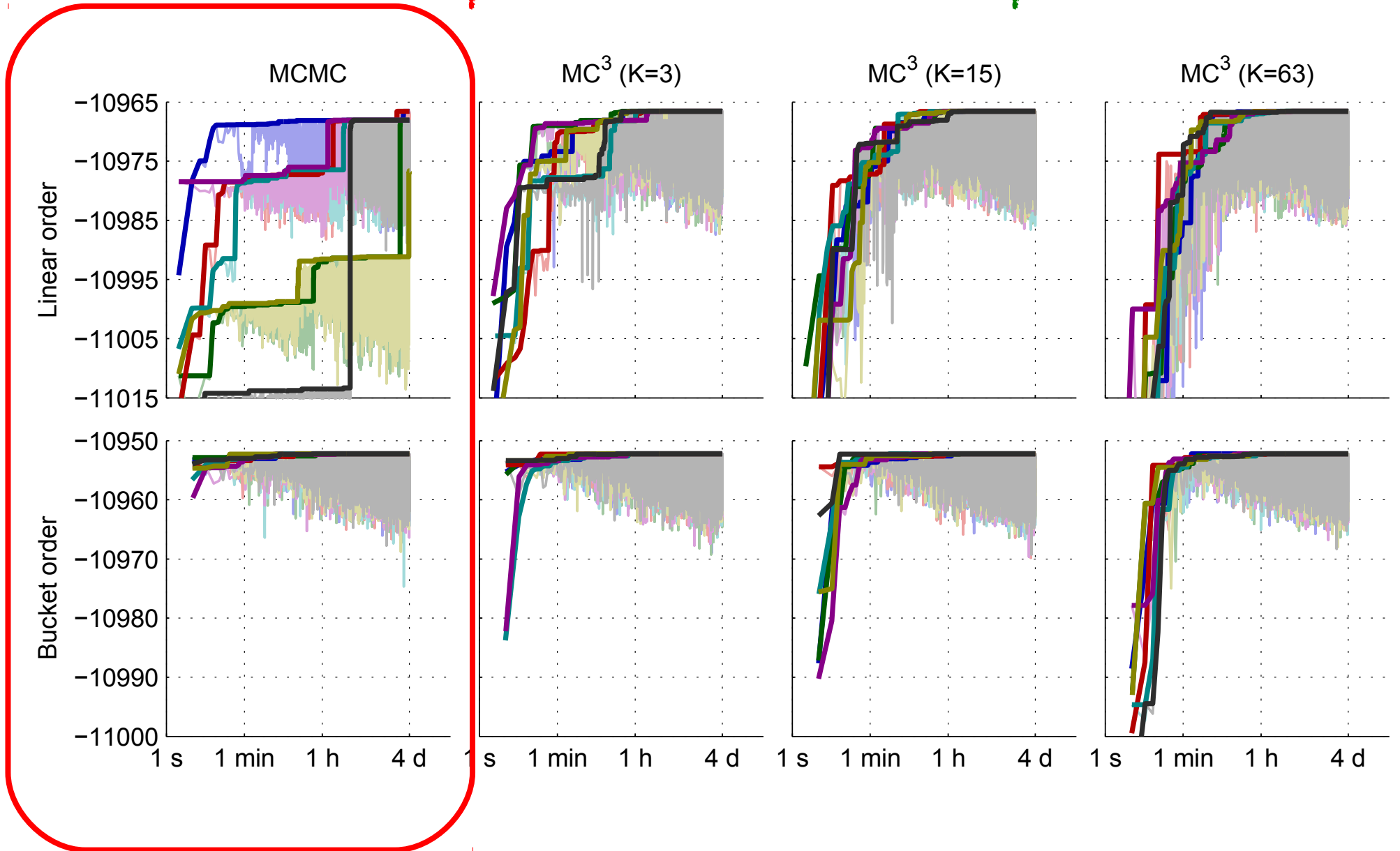
Tradeoffs



MCMC CONVERGENCE [Niinimäki et al. 2016]

Dataset: Mushroom, 1000 points, $d = 22$ variables

Bucket orders superior, even with tempered MCMC



DYNAMIC PROGRAMMING

Exploit a **factorization** to memorize repeated subproblems.

Chordal Markov networks:

- Call a function **f** decomposable if

$$f(G) = \frac{\prod_{\text{clique } C} f'(C)}{\prod_{\text{separator } S} f''(S)}$$

- Require a decomposable structure prior $P(G)$.

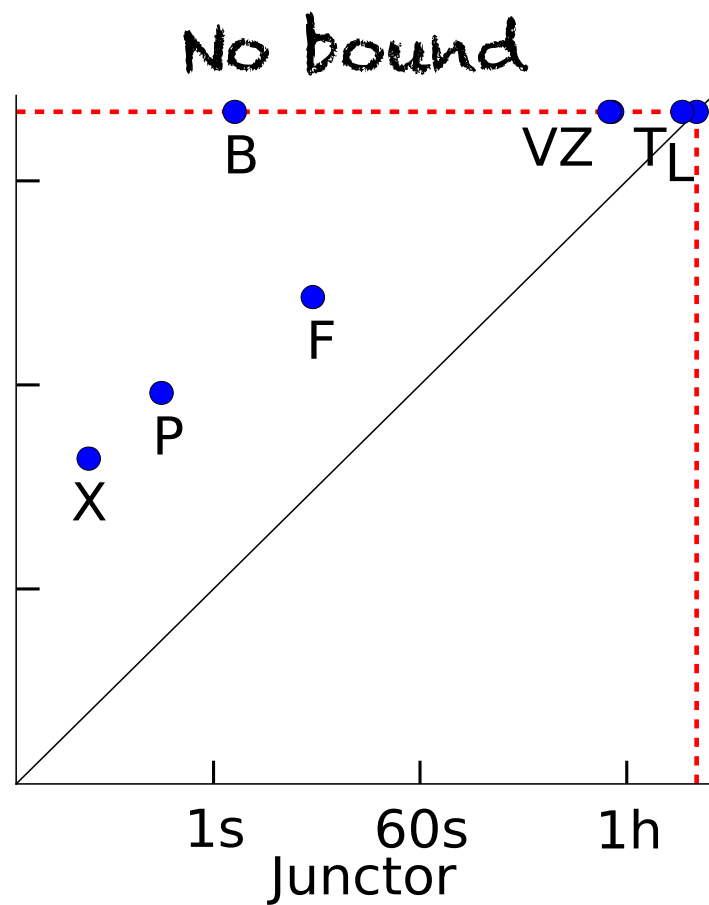
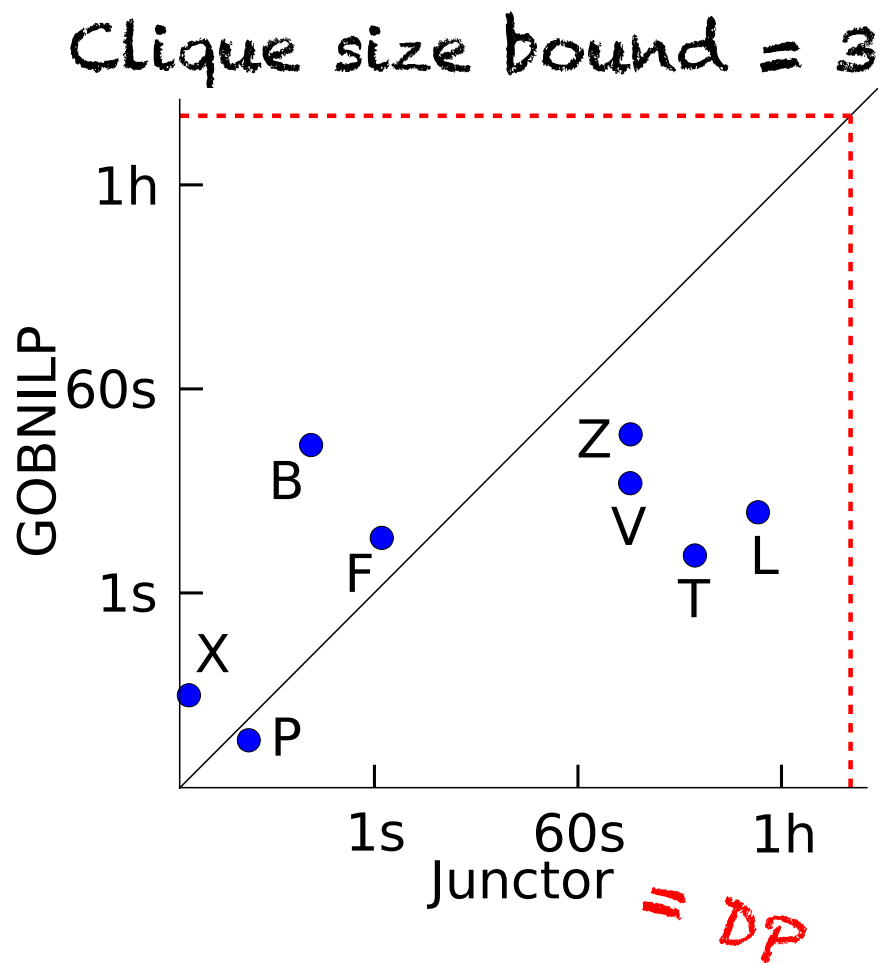
- Require a "nice" parameter prior of $p(x_c)$ for each clique **C** to get a decomposable $P(X|G)$.

- (Only expectations of decomposable functions.)

FINDING A MODE [Kangas et al. 2014]

Eight benchmark datasets, $d = 10$ to 19 variables

DP superior, especially if allowing large cliques



IMPORTANCE SAMPLING & BIAS CORRECTION

Sample from a "nice" proxy distribution q .
Correct **bias** by importance **weighting**.

Draw G^1, \dots, G^T independently from

$$q(G) = P(G|X) b(G) \cdot \text{Constant}$$

Put $w^t = 1/b(G^t)$ and estimate

$$E[f(G)|X] \approx \sum_t f(G^t) w^t / \sum_t w^t$$

Examples:

BNS: $b(G) = \#$ topological sorts of G

CMNs: $b(G) = \#$ rooted clique trees of G

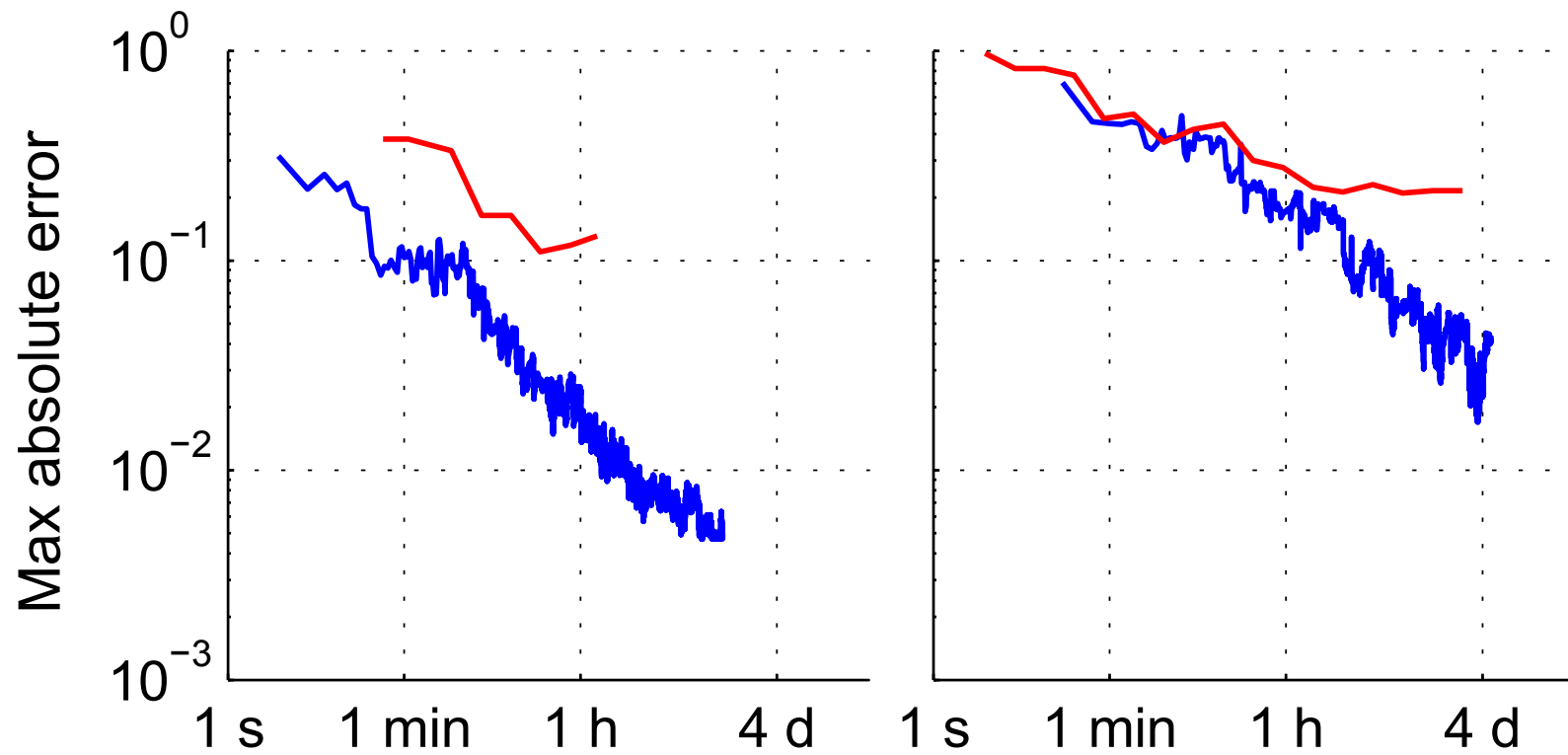
ARC POSTERIOR [Niinimäki et al. 2016]

Importance weighting is superior to

a heuristic by Ellis & Wong [JASA 2008]

Datasets: Flare, $d = 13$

German, $d = 20$



ANNEALING & MARKOV'S INEQUALITY

Construct a sampling distribution q close to the target. Independent samples yield lower bounds.

Annealed importance sampling [Neal 2001]

Generate ϑ_k along a Markov chain:

$$q(\vartheta_1, \dots, \vartheta_k) = q_1(\vartheta_1) q_2(\vartheta_2 | \vartheta_1) \dots q_k(\vartheta_k | \vartheta_{k-1})$$

Theorem (Markov lower bound) [Gomes et al. 2007]

Let Z_1, \dots, Z_T be independent and nonnegative with mean μ . Then with probability at least p :

$$(1-p)^{1/T} \min\{Z_1, \dots, Z_T\} \leq \mu$$

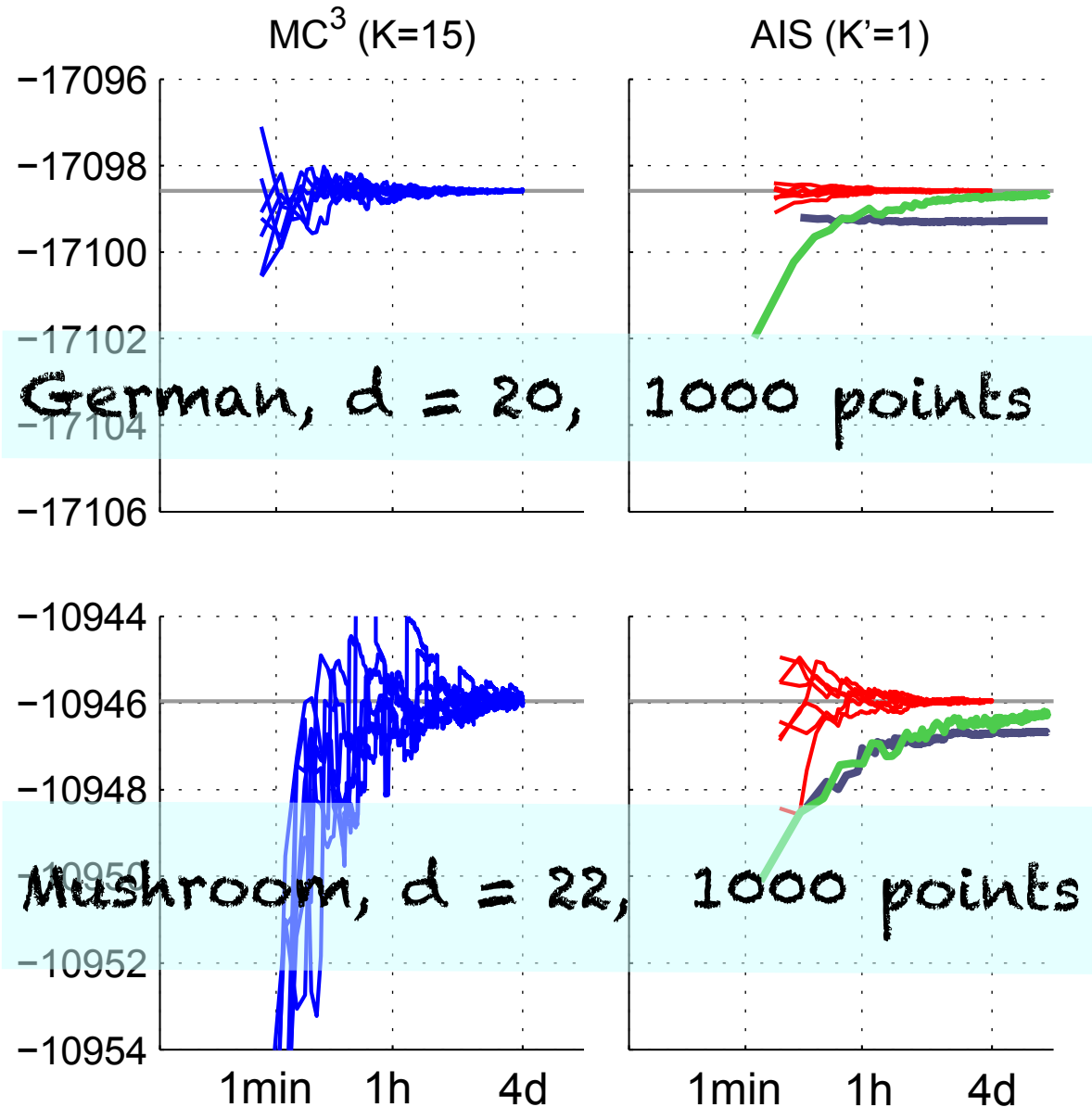
NORMALIZING CONSTANT [Niinimäki et al. 2016]

AIS is accurate and yields good lower bounds

Use B averages, each over T/B samples:

$$B = 5$$

$$B = \sqrt{T}$$



SUMMARY

Confident Bayesian Learning is hard but desirable

Some recent progress in structure Learning

Bidirectional approach: enhance unconfident schemes - expand exact algorithms

For details, see the PhD theses of Pekka P. [2012], Janne K. [2014], Teppo N. [2015], and Kustaa K. [2016]