

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

An exploration of methods for verification of probability models based on Bayesian networks.

Pavel Novikov
Supervised by: Oleg Senko

Moscow, 2015

Table of contents

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

Results

- 1 Problem overview
- 2 Extracting subset distribution
- 3 Performing hypothesis testing
- 4 Generating random Bayesian net
- 5 Results

Bayesian networks

Problem
overview

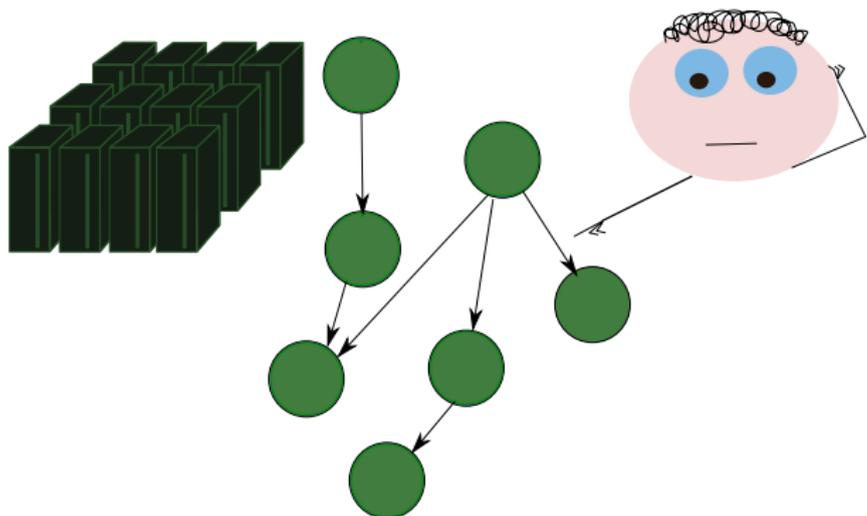
Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- Bayesian networks are good for encoding distributions.
- Problem: learning from scratch requires huge amounts of computational resources and doesn't guarantee good result.
- Alternative: expert knowledge (still no guarantee, though).



Verification

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

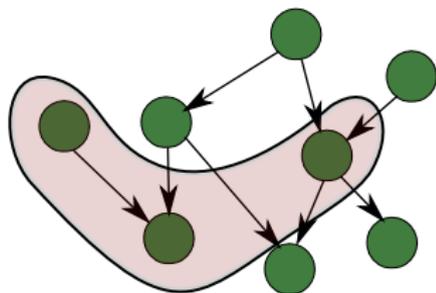
Results

Problem:

- Learned network may be of low quality.
- Expert knowledge can be flawed.

Proposed verification procedure:

- Look at various marginal distributions.
- Use statistical testing to check if they fit to data.



Exploration procedure

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

Results

- 1 Generate random Bayes net
- 2 Generate sample from this net
- 3 Extract subset distributions
- 4 Perform statistical testing

Repeat many times and examine false rejection portion.

Table of contents

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- 1 Problem overview
- 2 Extracting subset distribution**
- 3 Performing hypothesis testing
- 4 Generating random Bayesian net
- 5 Results

Where it all starts

Problem
overview

Extracting
subset
distribution

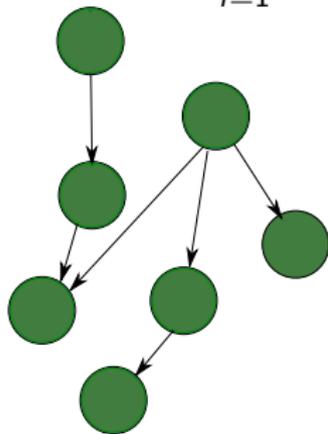
Performing
hypothesis
testing

Generating
random
Bayesian net

Results

$$P(X) = P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1}) = \\ = \prod_{i=1}^n P(x_i | \text{Par}(x_i))$$

- Some scopes of each factor may be dummy.
- Graph is used to show which ones are.



Basic variable elimination method

Problem
overview

Extracting
subset
distribution

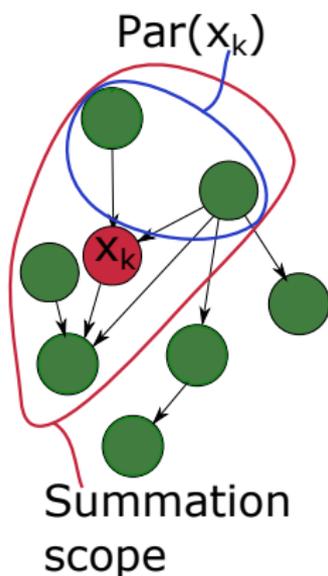
Performing
hypothesis
testing

Generating
random
Bayesian net

Results

We can simply sum variables out one by one.

$$\sum_{x_k \in \text{Dom}(x_k)} \prod_i P(x_i | \text{Par}(x_i)) =$$
$$\left(\prod_{\{i: i \neq k, x_k \notin \text{Par}(x_i)\}} \left[P(x_i | \text{Par}(x_i)) \right] \right)$$
$$\left(\sum_{\{x_k \in \text{Dom}(x_k)\}} \left[P(x_k | \text{Par}(x_k)) \right] \right)$$
$$\prod_{j: x_k \in \text{Par}(x_j)} \left(P(x_j | \text{Par}(x_j)) \right)$$



Problems of variable elimination method

Problem
overview

Extracting
subset
distribution

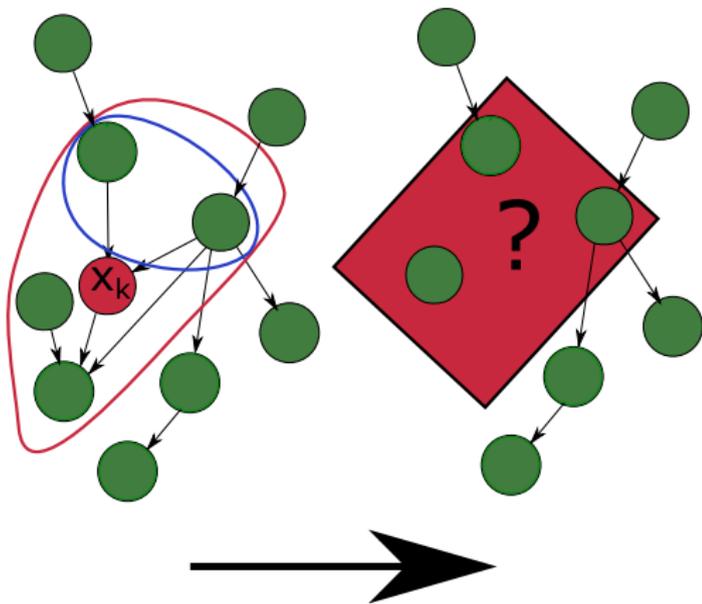
Performing
hypothesis
testing

Generating
random
Bayesian net

Results

There are some drawbacks:

- it has exponential complexity
- we lose bayesian network



Recovering Bayesian network structure

Problem
overview

Extracting
subset
distribution

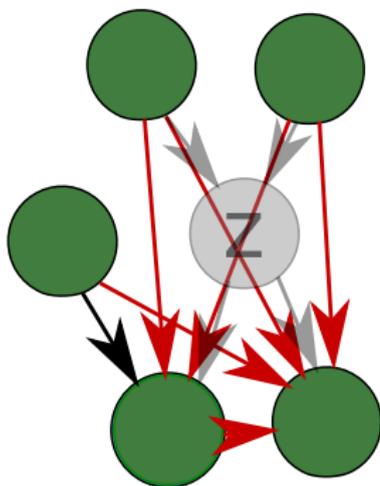
Performing
hypothesis
testing

Generating
random
Bayesian net

Results

Assertion

Distribution $P(X \setminus \{z\})$ factorizes over a graph G' , produced from graph G by connecting every child c of removed vertex z with all vertices of the summation scope, preceding c in some fixed topological order.



Algorithm for recomputing CPD

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- 1 initialize $\mathcal{P}_0 \leftarrow P(z|Par(z))$, $i \leftarrow 1$
- 2 take next variable c in topological order from $Child(z)$
- 3 $\mathcal{P}_i = \mathcal{P}_{i-1}P(c|Par(c))$,
- 4 $P(x_c|\tilde{P}ar(x_c)) = \frac{\sum_z \mathcal{P}_i}{\sum_z \mathcal{P}_{i-1}}$.
- 5 $i \leftarrow i + 1$
- 6 if there are still variables in $Child(z)$ go to step 2;

Table of contents

Problem
overview

Extracting
subset
distribution

**Performing
hypothesis
testing**

Generating
random
Bayesian net

Results

- 1 Problem overview
- 2 Extracting subset distribution
- 3 Performing hypothesis testing**
- 4 Generating random Bayesian net
- 5 Results

The idea of hypothesis testing

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

H_0 : zero hypothesis, "default".

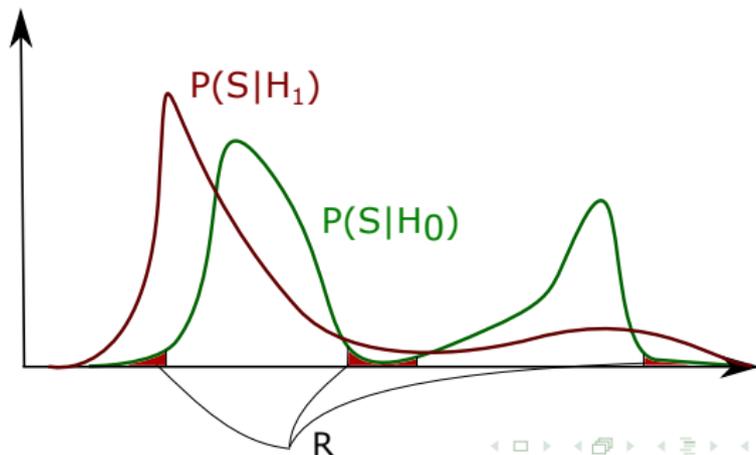
H_1 : alternative.

S : statistic.

R : rejection region.

Idea: pick R so that $P(S \in R | H_0) < \alpha$ (significance level).

p-value : minimal significance level that allows to reject
particular hypothesis.



Pearson χ^2 and g-Test

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

Results

P - probability given H_0

F - observed counts

	P	F
e_1	P_{e_1}	n_{e_1}
e_2	P_{e_2}	n_{e_2}
		...
e_k	P_{e_k}	n_{e_k}
\sum	1	n

$$\chi^2 = \sum_{i=1}^k \frac{(n_{e_i} - nP_{e_i})^2}{nP_{e_i}}$$
$$G = 2 \sum_{i=1}^k nP_{e_i} \log\left(\frac{nP_{e_i}}{n_{e_i}}\right)$$

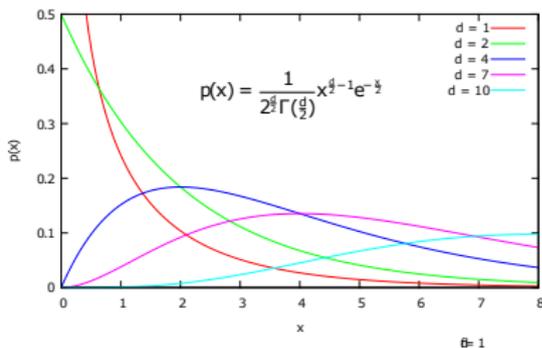
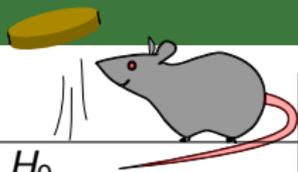


Figure: χ^2 distribution

Multiple hypothesis testing



H_0	Coin tossing Coin is fair	Food admixture Admixture affects mice
H_1	Coin is not fair	Admixture does not affect mice
Rejection	Equal results of all tosses	Specialized test pro- cedure for some vital signs in 2 groups of mice
Procedure 1 OK	Throw N times.	Measure one type of vital signs
Procedure 2 Not OK	Make M people throw N times each	Measure M types

Multiple hypothesis testing really makes a difference!

Bonferroni correction

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

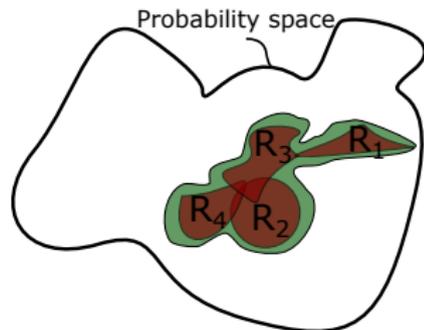
N hypotheses: $H_0^1, H_1^1; \dots; H_0^N, H_1^N$

N significance levels: $\alpha^1; \dots; \alpha^N$

N rejection regions: $R^1; \dots; R^N$

N statistics: $S^1; \dots; S^N$

$$P(\bigvee_{i=1}^N (S_i \in R_i)) \leq \sum_{i=1}^N P((S_i \in R_i))$$



Idea: fix α_i to be equal α/N .

Stepwise correction procedures

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

Results

- sort hypotheses by p-values p_k .
- compare p_k with $\frac{\alpha}{N+1-k}$
- reject all $H_0^i : i \leq r$, where r is:
 - Holm step-down: $r = \min(\{k : p_k > \frac{\alpha}{N+1-k}\}) - 1$
 - Hochberg step-up: $r = \max(\{k : p_k \leq \frac{\alpha}{N+1-k}\})$

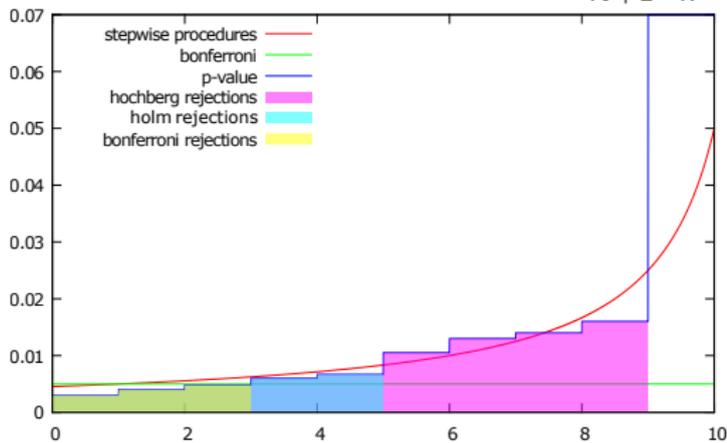


Figure: example of rejection procedures

Table of contents

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- 1 Problem overview
- 2 Extracting subset distribution
- 3 Performing hypothesis testing
- 4 Generating random Bayesian net**
- 5 Results

Random graphs - Erdős-Rényi

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

First formulation: each edge can be added to the graph independently of others with probability p .

Second formulation: random set of k edges is chosen uniformly at random.



Figure: 50 vertices, 100 edges

Random graphs - Barabási–Albert

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

Procedure starts with fully connected graph with n vertices. Each new vertex added to the graph is connected to n old vertices.

Probability to chose a particular old vertex to connect to is proportional to its degree.

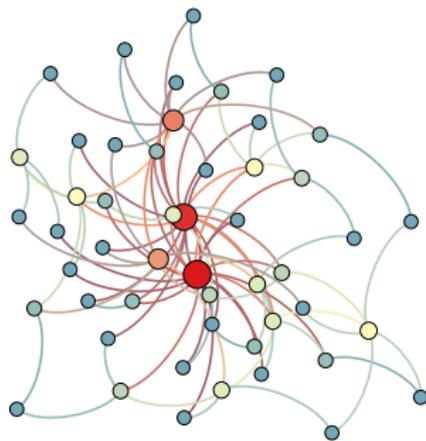


Figure: 50 vertices, 97 edges

Dirichlet and Beta distributions

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

Dirichlet distribution is a distribution over n-dimensional vectors of positive numbers that sum to one - tabular distributions.

$$Dir(x, \alpha) = \frac{1}{B(x)} \prod_{i=1}^K x_i^{\alpha_i - 1}$$

2-dimensional case is referred to as Beta-distribution.

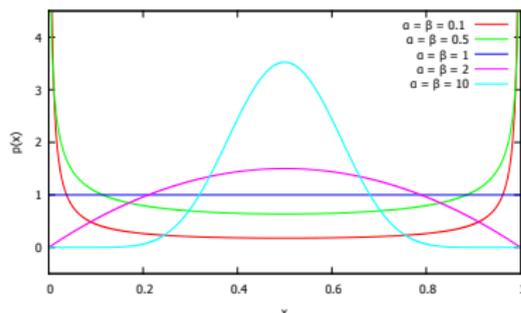


Figure: Beta distribution

Table of contents

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- 1 Problem overview
- 2 Extracting subset distribution
- 3 Performing hypothesis testing
- 4 Generating random Bayesian net
- 5 Results**

χ^2 , 10 vertices, 9 edges, $Dir(x, 1)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

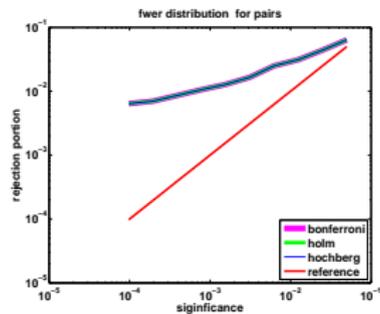
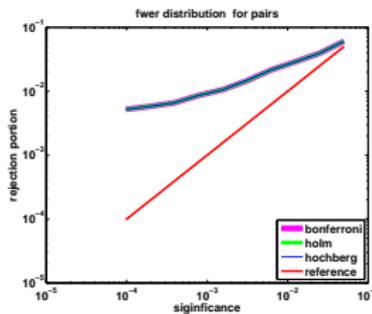
Generating random Bayesian net

Results

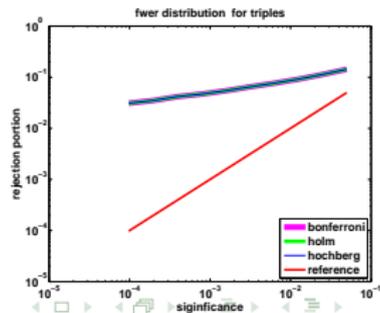
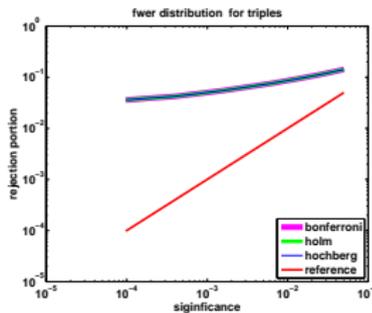
Erdős–Rényi

Barabási–Albert, direct order

Pairs



Triples



χ^2 , 10 vertices, 17-18 edges, $Dir(x, 1)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

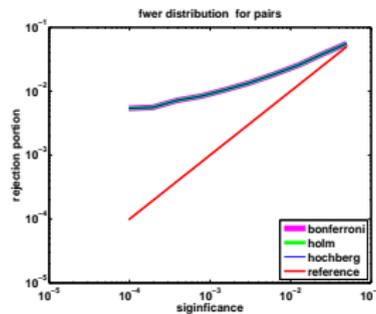
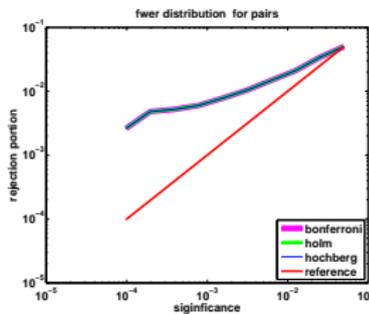
Generating random Bayesian net

Results

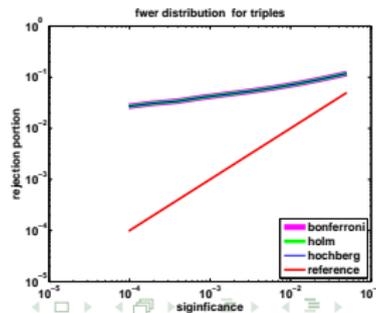
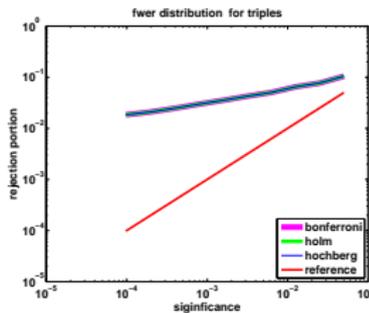
Erdős–Rényi

Barabási–Albert, inverse order

Pairs



Triples



Simple hypothesis testing experiment - 100 000 experiments, 100 samples per experiment

Problem overview

Extracting subset distribution

Performing hypothesis testing

Generating random Bayesian net

Results

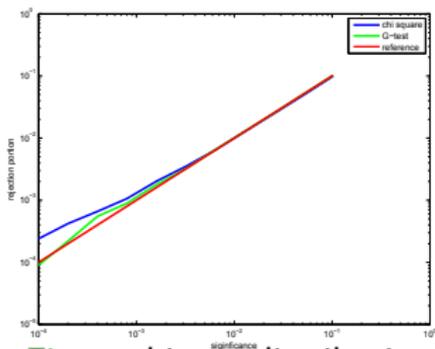


Figure: binary distribution

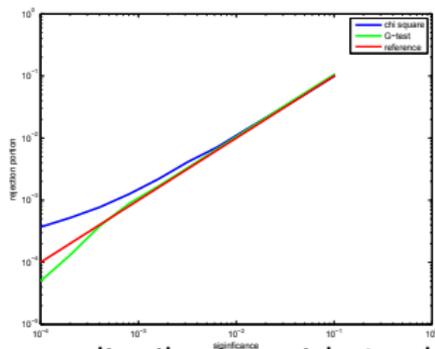


Figure: distribution with 4 values

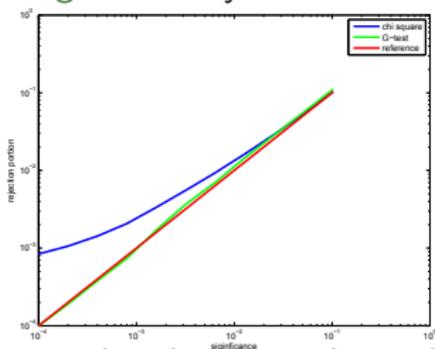


Figure: distribution with 8 values

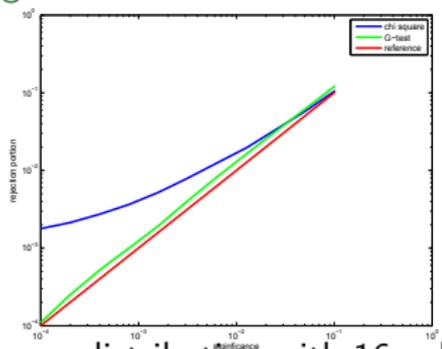


Figure: distribution with 16 values

G – test, 10 vertices, 9 edges, $Dir(x, 1)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

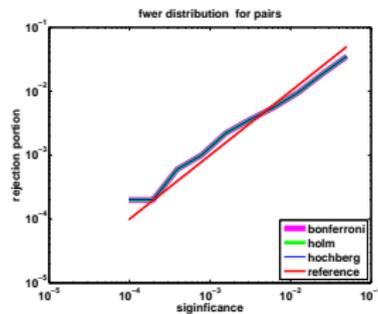
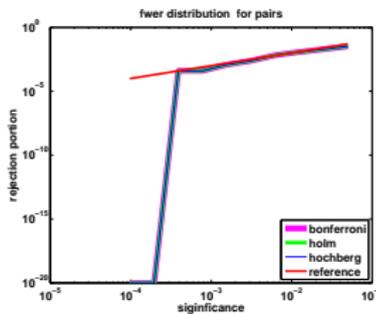
Generating random Bayesian net

Results

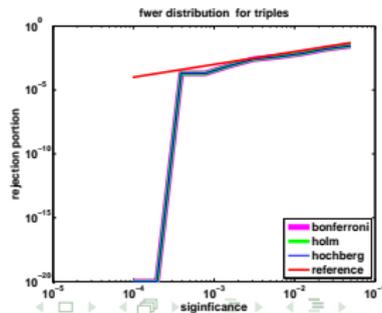
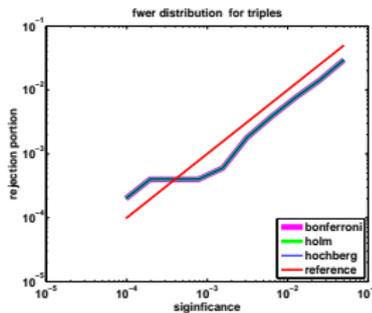
Erdős–Rényi

Barabási–Albert, random order

Pairs



Triples



G – test, 10 vertices, 17-18 edges, $Dir(x, 1)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

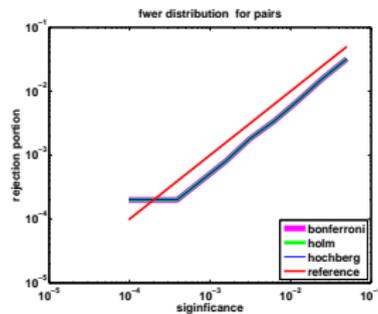
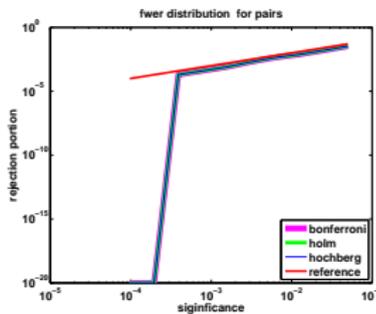
Generating random Bayesian net

Results

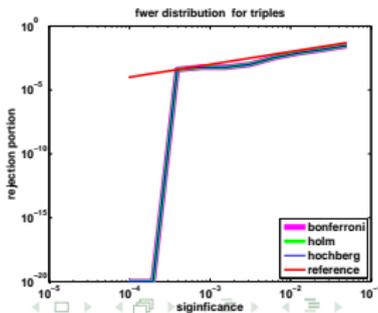
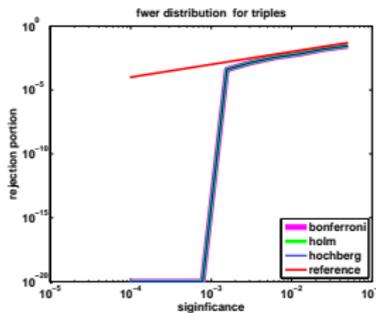
Erdős–Rényi

Barabási–Albert, random order

Pairs



Triples



G – test, 10 vertices, 9 edges, $Dir(x, 0.2)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

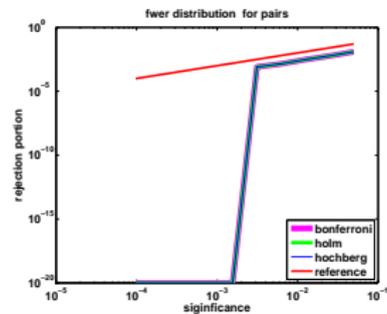
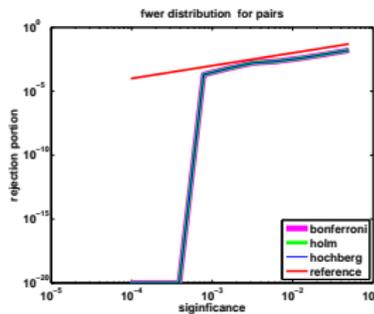
Generating random Bayesian net

Results

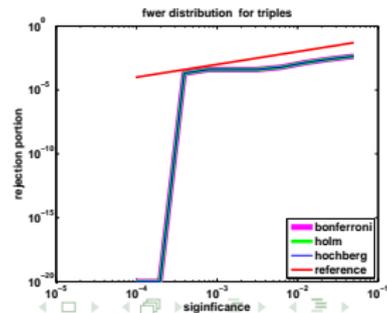
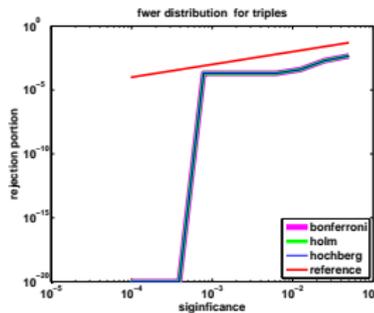
Erdős–Rényi

Barabási–Albert, random order

Pairs



Triples



G – test, 10 vertices, 17-18 edges, $Dir(x, 0.2)$

Problem overview

Extracting subset distribution

Performing hypothesis testing

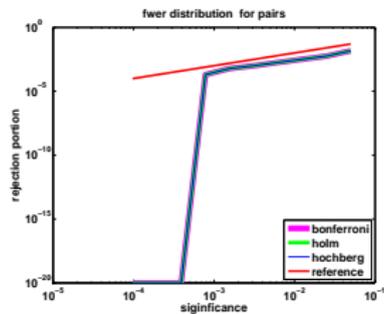
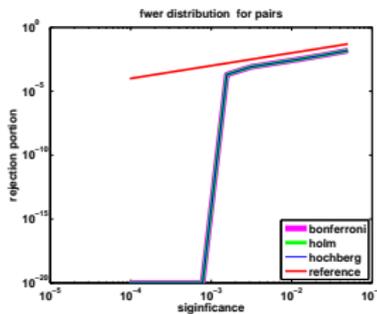
Generating random Bayesian net

Results

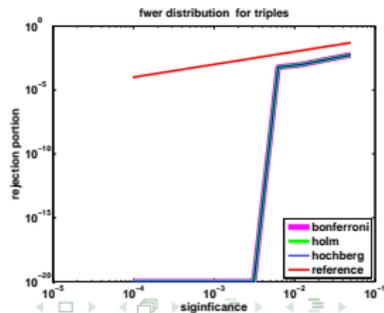
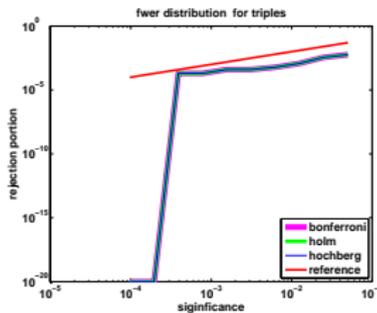
Erdős–Rényi

Barabási–Albert, random order

Pairs



Triples



Conclusions

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

- It's much better to use g-Test than χ^2 test.
- It does not matter which correction procedure we choose.
- Graph does not play a big role.
- Strength of variable interdependence does play an important role.

Problem
overview

Extracting
subset
distribution

Performing
hypothesis
testing

Generating
random
Bayesian net

Results

The End