

A Study on Comparison of Bayesian Network Structure Learning Algorithm for Selecting Appropriate Model

유재성

Dept. of Statistics

November 21, 2014

Title

1 서론

- 목표

2 비교 기준

- 비교 기준
- 어떻게 비교할 것인가?
- 전제

3 Real Data를 이용한 구조 학습 결과 비교

- Asia DataSet by Lauritzen and Spiegelhalter
- Insurance Evaluation Network DataSet
- ALARM Monitoring System DataSet
- The HailFinder Weather Forecast System DataSet

4 Data Generator 성능 평가

5 Topology에 따른 비교

6 결론

- 결론
- Discussion

목표

변수에 대한 정보만 존재하고 그들의 구조를 잘 모를 때 베이지안 네트워크 구조 학습을 진행하게 된다.

다양한 분야에서의 각 이론적 배경에 맞는 "적합한" 모형 선택을 돋기 위해,
베이지안 네트워크의 구조 학습 방법을 선택하기 위한 여러가지 방법을 비교하고자 한다.

서론	비교 기준
Real Data를 이용한 구조 학습 결과 비교	어떻게 비교할 것인가?
Data Generator 성능 평가	전제
Topology에 따른 비교	
결론	

Outline

1 서론

- 목표

2 비교 기준

- 비교 기준
- 어떻게 비교할 것인가?
- 전제

3 Real Data를 이용한 구조 학습 결과 비교

- Asia DataSet by Lauritzen and Spiegelhalter
- Insurance Evaluation Network DataSet
- ALARM Monitoring System DataSet
- The HailFinder Weather Forecast System DataSet

4 Data Generator 성능 평가

5 Topology에 따른 비교

6 결론

- 결론
- Discussion

비교 기준	서론
Real Data를 이용한 구조 학습 결과 비교	
Data Generator 성능 평가	
Topology에 따른 비교 결론	

비교 기준
어떻게 비교할 것인가? 전제

비교 기준

점수를 이용한 비교 BDe, Log Likelihood, AIC, BIC

주의! 점수가 "클수록" 좋다.(*1)

목표 네트워크와 학습된 네트워크와의 직접 비교(*2) 학습된 모형 자체가 실제 기대한 모형과 일치한지 자체에 주목한다.

- C (Correct Arcs) : 목표 네트워크 O, 학습 네트워크 O, 방향 일치
- M (Missing Arcs) : 목표 네트워크 O, 학습 네트워크에 X
- WO (Wrongly Oriented Arcs) : 목표 네트워크 O, 학습 네트워크 O, 방향 불일치
- WC (Wrolgly Connected Arcs) : 목표 네트워크 X, 학습 네트워크 O

(*1) Silvia Acid 외 5명,

"A comparison of learning algorithms for Bayesian networks: a case study based on data from an emergency medical service", Artificial Intelligence in Medicine 30 (2004) 215–232.

(*2) Reference : Fadhl, M. Al-Akwa, Mohammed M. Ikhawlani, (2012),

"Comparison of the Bayesian Network Structure Learning Algorithms",

International Journal of Advanced Research in Computer Science and Software Engineering Volume 2, Issue 3.

어떻게 비교할 것인가?

- 어떤 Algorithm을 적용했는가에 따라
- Topology의 유형에 따라
- Node 개수가 증가함에 따라
- Sample Size가 증가함에 따라

전제

- 각 Node에 대한 Cardinality는 2로 제한한다.
(즉, 이 모든 Node가 Binary Data인 경우만 다룬다.)
- 이 연구에서 사용하는 구조 학습 알고리즘은 R의 bnlearn 패키지에서 제공하는 것으로 한정한다.
- 모든 실험은 100번 반복하였다. 그 결과들을 받아 종합하였다.
- Synthetic Data로 분석할 때, 확률 관계를 정의할 때, 확률값을 $U(0,1)$ 사이의 값에서 임의로 주었다.

Outline

1 서론

- 목표

2 비교 기준

- 비교 기준
- 어떻게 비교할 것인가?
- 전제

3 Real Data를 이용한 구조 학습 결과 비교

- Asia DataSet by Lauritzen and Spiegelhalter
- Insurance Evaluation Network DataSet
- ALARM Monitoring System DataSet
- The HailFinder Weather Forecast System DataSet

4 Data Generator 성능 평가

5 Topology에 따른 비교

6 결론

- 결론
- Discussion

Asia DataSet

Description Small synthetic data set from Lauritzen and Spiegelhalter (1988) about lung diseases (tuberculosis, lung cancer or bronchitis) and visits to Asia.

Number of nodes 8

Number of arcs 8

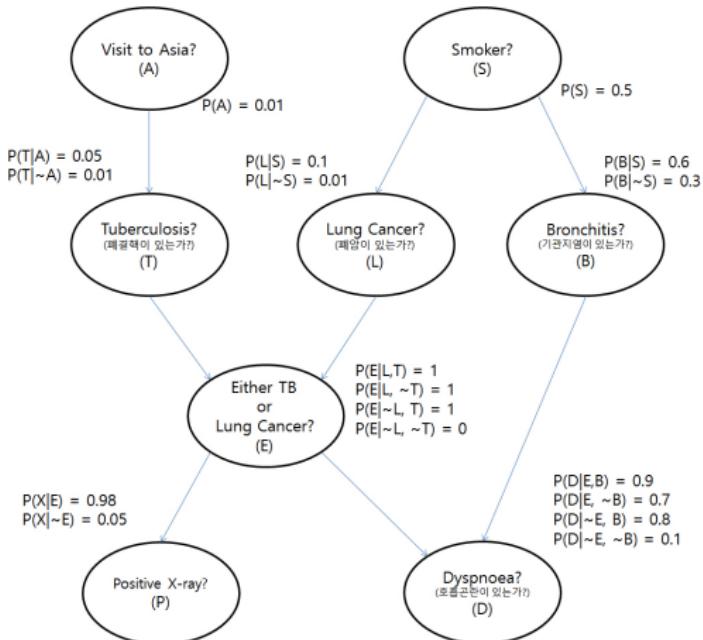
Number of parameters 18

Source Lauritzen S, Spiegelhalter D (1988).

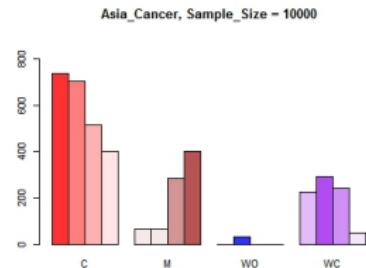
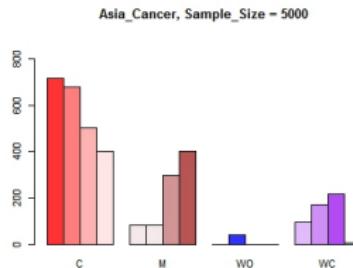
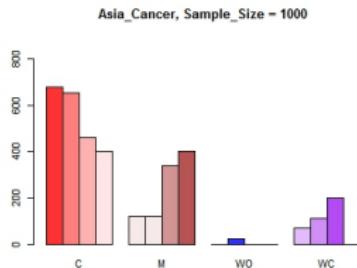
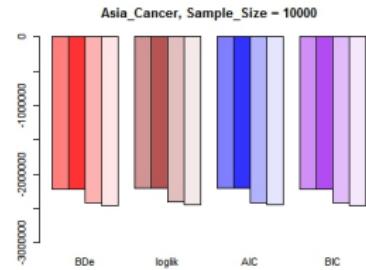
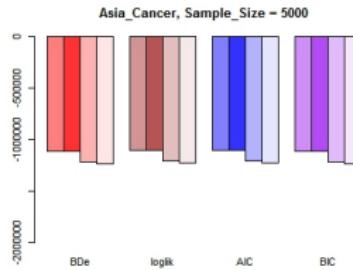
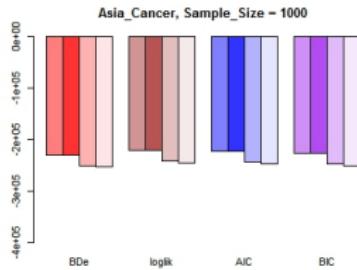
"Local Computation with Probabilities on Graphical Structures and their Application to Expert Systems (with discussion)".

Journal of the Royal Statistical Society: Series B (Statistical Methodology), 50(2), 157-224.

Asia DataSet



선행연구와 결과 비교 : Asia DataSet



Insurance DataSet

Description Insurance is a network for evaluating car insurance risks.

Number of nodes 27

Number of arcs 52

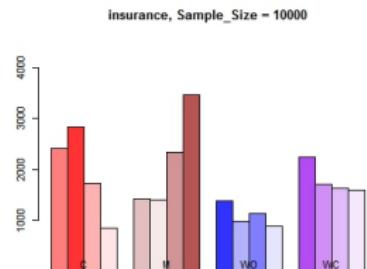
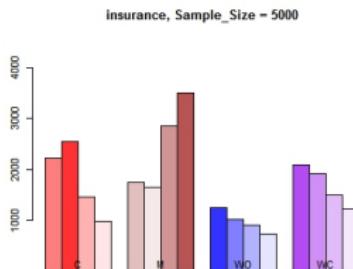
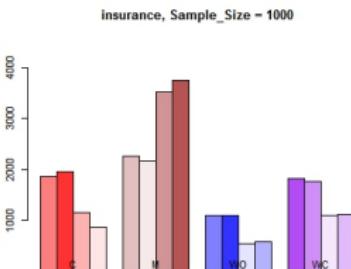
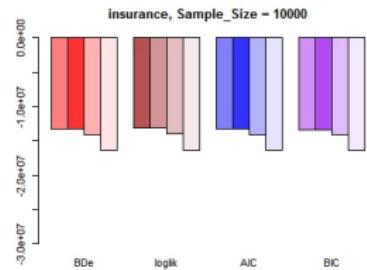
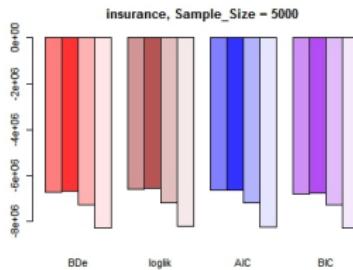
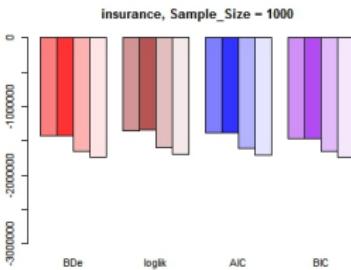
Number of parameters 984

Source Binder J, Koller D, Russell S, Kanazawa K (1997).
"Adaptive Probabilistic Networks with Hidden Variables".
Machine Learning, 29(2-3), 213-244.

Insurance DataSet



Insurance DataSet



Alarm DataSet

Description The ALARM ("A Logical Alarm Reduction Mechanism") is a Bayesian network designed to provide an alarm message system for patient monitoring.

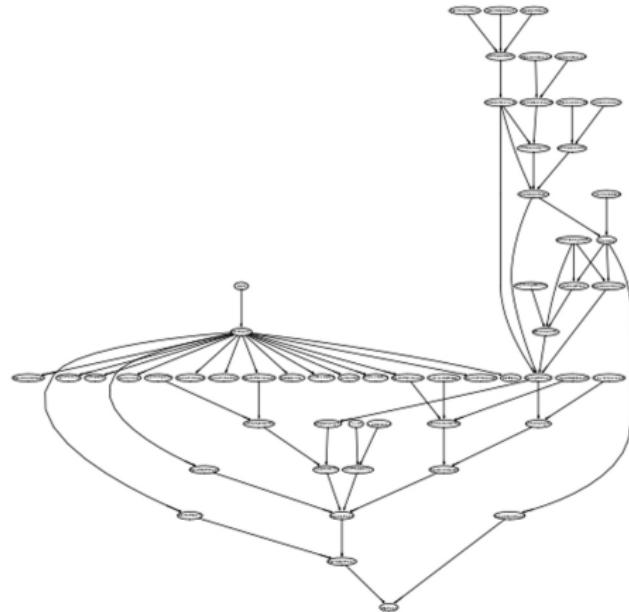
Number of nodes 37

Number of arcs 46

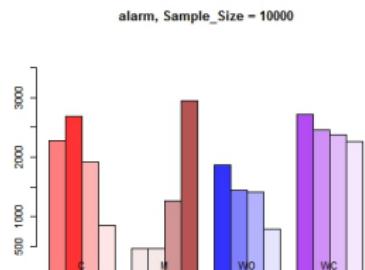
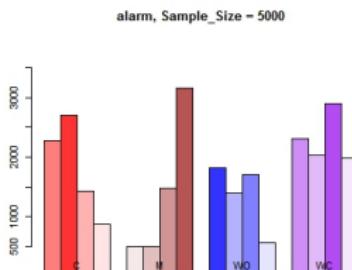
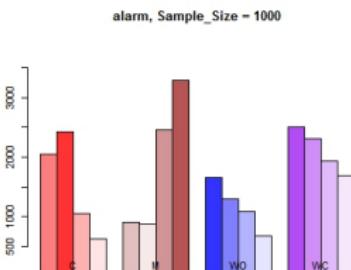
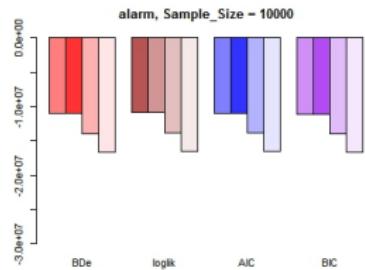
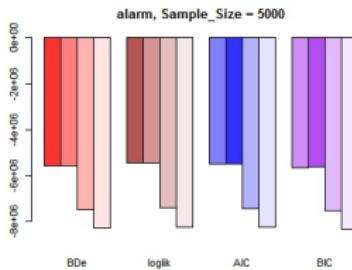
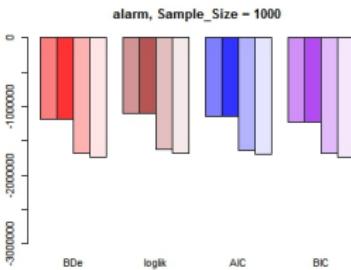
Number of parameters 509

Source Beinlich I, Suermondt HJ, Chavez RM, Cooper GF (1989).
"The ALARM Monitoring System: A Case Study with Two Probabilistic Inference Techniques for Belief Networks."
In "Proceedings of the 2nd European Conference on Artificial Intelligence in Medicine", pp. 247-256. Springer-Verlag.

Alarm DataSet



Alarm DataSet



HailFinder DataSet

Description Hailfinder is a Bayesian network designed to forecast severe summer hail in northeastern Colorado.

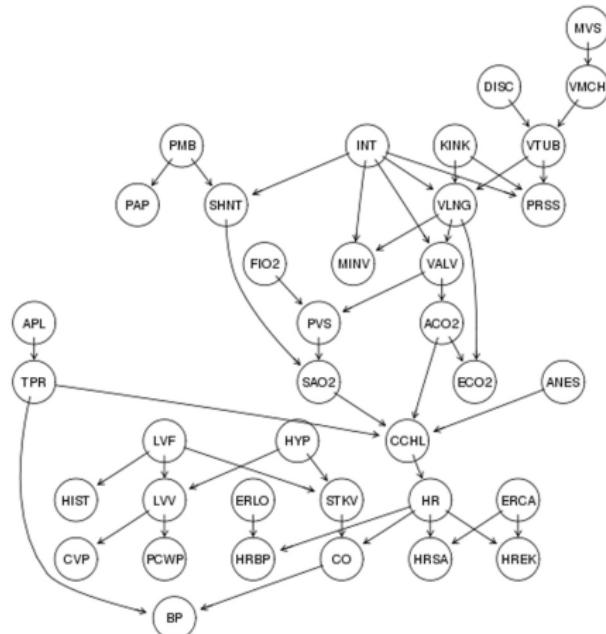
Number of nodes 56

Number of arcs 66

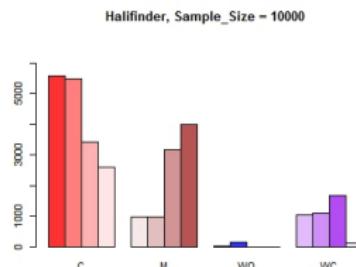
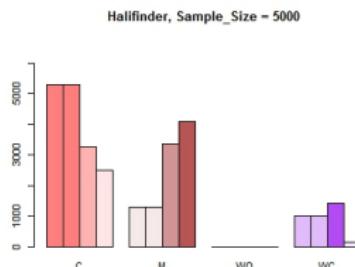
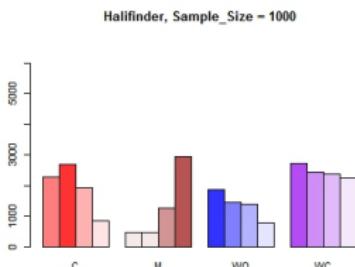
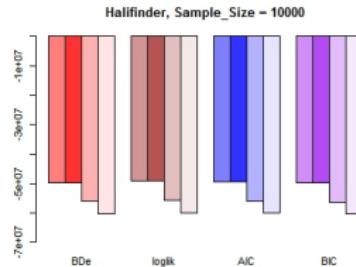
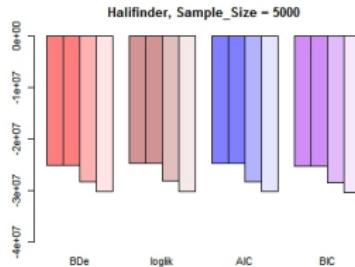
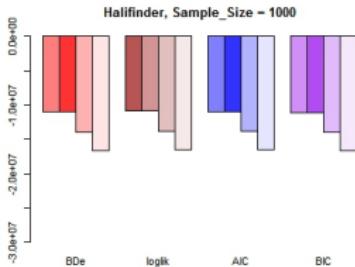
Number of parameters 2656

Source Abramson B, Brown J, Edwards W, Murphy A, Winkler RL (1996).
"Hailfinder: A Bayesian system for forecasting severe weather".
International Journal of Forecasting, 12(1), 57-71.

HailFinder DataSet



HailFinder DataSet



Summary

Sample Size 1000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
Asia	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	3	2	1	4
Insurance	2	1	3	4	2	1	3	4	3	4	2	1	2	1	4	3	1	2	4	3
Alarm	2	1	3	4	2	1	3	4	3	4	2	1	1	2	3	4	1	2	3	4
HailFinder	2	1	3	4	2	1	3	4	4	4	2	1	1	2	3	4	1	2	3	4
Sample Size 5000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
Asia	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	3	2	1	4
Insurance	2	1	3	4	2	1	3	4	3	4	2	1	1	2	3	4	1	2	3	4
Alarm	1	2	3	4	2	1	3	4	4	3	2	1	1	3	2	4	2	3	1	4
HailFinder	1	1	3	4	1	1	3	4	4	4	2	1	4	4	4	4	2	2	1	4
Sample Size 10000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
Asia	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	3	1	2	4
Insurance	2	1	3	4	2	1	3	4	3	4	2	1	1	3	2	4	1	2	3	4
Alarm	2	1	3	4	2	1	3	4	4	4	2	1	1	2	3	4	1	2	3	4
HailFinder	2	1	3	4	1	2	3	4	4	3	2	1	2	1	4	4	3	2	1	4

Outline

1 서론

- 목표

2 비교 기준

- 비교 기준
- 어떻게 비교할 것인가?
- 전제

3 Real Data를 이용한 구조 학습 결과 비교

- Asia DataSet by Lauritzen and Spiegelhalter
- Insurance Evaluation Network DataSet
- ALARM Monitoring System DataSet
- The HailFinder Weather Forecast System DataSet

4 Data Generator 성능 평가

5 Topology에 따른 비교

6 결론

- 결론
- Discussion

Data Generator 성능 평가

BN_Data_Generator {User-Defined Function}

베이지안 네트워크 데이터 생성기

Description 베이지안 네트워크 모형을 기반으로 Synthetic Data를 생성하여 준다.

Usage BN_Data_Generator (arcs, input_Probs, n, node_names)

Arguments

arcs	(matrix)	모형 내 arc 유무와 방향을 결정하는 행렬
input_Probs	(list)	arc가 이어져있는 node간 관계의 조건부 확률
n	(constant)	Sample Size
node_names	(vector)	Node의 이름

Data Generator 성능 평가

```

[File] run_main.R [File] arach.R [File] BN_Data_Generator.R [Run]
1  BN_Data_Generator = function(arcs, input_Probs, n, node_names)
2  {
3    # Node 개수
4    num_of_nodes = dim(arcs)[1];
5    # 각 Node의 Parent Node 개수
6    num_of_parent_nodes = apply(arcs, 2, sum);
7    list_parent_nodes = list();
8    for(i in 1:num_of_nodes)
9    {
10      if (length(which(arcs[,i] == 1)) == 0)
11      {
12        list_parent_nodes[[i]] = NULL;
13      } else {
14        list_parent_nodes[[i]] = which(arcs[,i] == 1);
15      }
16    }
17  }
18
19  # Root node의 개수
20  root_nodes = sum(num_of_parent_nodes == 0);
21
22  # 결과는 여기서 정정이 된다.
23  result_mat = matrix(0, n, num_of_nodes);
24  dimnames(result_mat)[[1]] = node_names;
25  # result_mat
26
27  # 지정해야 할 조건부 확률 개수
28  num_of_probs = (as.matrix(result_mat[,1])) * num_of_parent_nodes;
29  dimnames(num_of_probs)[[2]] = node_names;
30  num_of_probs
31
32
33
34
35  # 지정해야 할 조건부 확률 개수만큼 input_prob에 맞는지 확인. 만약 false이면 프로그램 종료
36  input_prob_len = length(input_Probs);
37  for (i in 1:input_prob_len)
38  {
39    if (length(input_Probs[[i]]) != num_of_probs[i])
40    {
41      cat("Error");
42      return(break);
43    }
44  }
45
46
47
48  # Root Node Initialization
49  for(i in root_nodes)
50  {
51    p = input_Probs[[i]];
52    result_mat[[i]] = sample(c("Y", "N"), 1, prob=c(p, 1-p), rep=T);
53  }
54
55
56
57  # Generator
58  i=1

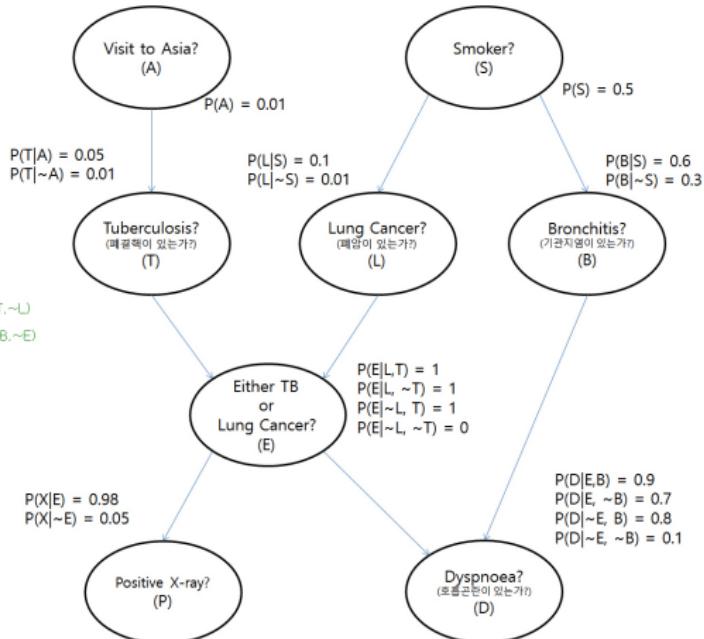
```

Data Generator 성능 평가

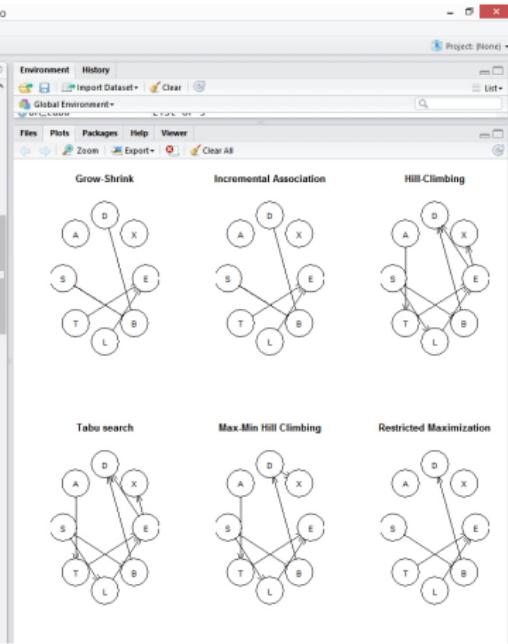
```

# Asia
arcs = rbind(
  | # A S T L B E X D
  | c(0, 0, 1, 0, 0, 0, 0, 0), #A
  | c(0, 0, 0, 1, 0, 0, 0, 0), #S
  | c(0, 0, 0, 0, 1, 0, 0, 0), #T
  | c(0, 0, 0, 0, 0, 1, 0, 0), #L
  | c(0, 0, 0, 0, 0, 0, 1, 0), #B
  | c(0, 0, 0, 0, 0, 0, 0, 1), #E
  | c(0, 0, 0, 0, 0, 0, 0, 0), #X
  | c(0, 0, 0, 0, 0, 0, 0, 0), #D
)
arc_name = c("A", "S", "T", "L", "B", "E", "X", "D")
dimnames(arcs)[[1]] = arc_name
dimnames(arcs)[[2]] = arc_name

Probs = list(
  | c(0.01),      # P(A)
  | c(0.5),       # P(S)
  | c(0.05, 0.01), # P(T|A), P(T|~A)
  | c(0.1, 0.01), # P(L|S), P(L|~S)
  | c(0.6, 0.5),  # P(B|S), P(B|~S)
  | c(1, 1, 0),   # P(E|T,L), P(E|~T,L), P(E|T,~L), P(E|~T,~L)
  | c(0.98, 0.05), # P(X|E), P(X|~E)
  | c(0.9, 0.7, 0.8, 0.1) # P(D|E,E), P(D|~E,E), P(D|E,~E), P(D|~E,~E)
)
  )
  
```



Data Generator 성능 평가



```

RStudio
File Edit Code View Plots Session Build Debug Tools Help
Project: (None)
Console (bytempho_제작자면 배드민턴/Scores/R) 
> ##### 시작
> # set.seed(1234)
> require(bnlearn)
 필요한 키워드를 포함한입니다: bnlearn
>
> n = 1000
>
> # Asia
> arcs = rbind(
+   c("A", "S", "T", "B", "E", "X", "D"),
+   c(0, 0, 1, 0, 0, 0, 0, 0), #A
+   c(0, 0, 0, 1, 1, 0, 0, 0), #S
+   c(0, 0, 0, 0, 0, 1, 0, 0), #T
+   c(0, 0, 0, 0, 0, 1, 0, 0), #B
+   c(0, 0, 0, 0, 0, 0, 1, 0), #E
+   c(0, 0, 0, 0, 0, 0, 0, 1), #X
+   c(0, 0, 0, 0, 0, 0, 0, 0) #D
+ )
> arc_name = c("A", "S", "T", "B", "E", "X", "D")
> dimnames(arcs)[[1]] = arc_name
> dimnames(arcs)[[2]] = arc_name
>
> Probs = list(
+   c(0.01),      # P(A)
+   c(0.5),       # P(S)
+   c(0.05, 0.01), # P(T|A), P(T|~A)
+   c(0.1, 0.01), # P(L|S), P(L|~S)
+   c(0.6, 0.3),  # P(B|S), P(B|~S)
+   c(0.9, 0.01), # P(E|L), P(E|~L)
+   c(0.98, 0.02),# P(X|L), P(X|~L)
+   c(0.9, 0.7, 0.8, 0.1) # P(D|B,E), P(D|B,~E), P(D|~B,E), P(D|~B,~E)
+ )
>
> res = BN_data_generator(arcs, Probs, n, arc_name)
> data = res$data
> head(data)
A S T L B E X D
1 N N N N N N N N
2 N N N N N N N N
3 N Y N N Y N N Y
4 N N N N N N N N
5 N Y N N Y N N N
6 N Y N N N N N N
> str(data)
'data.frame': 1000 obs. of 8 variables:
 $ A: num 0 0 0 0 0 0 0 0 0 0 ...
 $ S: num 0 0 0 0 0 0 0 0 0 0 ...
 $ T: num 0 0 0 0 0 0 0 0 0 0 ...
 $ L: num 0 0 0 0 0 0 0 0 0 0 ...
 $ B: num 0 0 0 0 0 0 0 0 0 0 ...
 $ E: num 0 0 0 0 0 0 0 0 0 0 ...
 $ X: num 0 0 0 0 0 0 0 0 0 0 ...
 $ D: num 0 0 0 0 0 0 0 0 0 0 ...
> # constraint-based algorithms
> bn_gs = gs(data)      # the Grow-Shrink(GS)
> bn_mx
  
```

성능 평가를 위한 선행 연구

Pekka Parviainen, Hossein Shahrabi Farahani, and Jens Lagergren (2014).

"Learning Bounded Tree-width Bayesian Networks using Integer Linear Programming"

Proceedings of the 17th International Conference on Artificial Intelligence and Statistics
(AISTATS)

Data Generator 성능 평가 via Hill Climbing (HC)

Dataset	Num. of Nodes	Sample Size	Score via HC				
			ILP	BDe	BDe	loglik	AIC
Asia real dataset	100	-245.64	-251.07	-194.08	-210.08	-230.92	
		1000	-2317.41	-2281.98	-2188.63	-2205.63	-2247.35
		10000	-22466.40	-21937.24	-21812.37	-21832.37	-21904.47
	1000		-271.30	-209.69	-222.69	-239.62	
		1000		-2350.21	-2262.16	-2279.16	-2320.87
		10000		-22529.57	-22419.24	-22437.24	-22502.14
	10000	100		-	-	-	-
		1000		-2289.26	-2197.43	-2214.43	-2256.15
		10000		-22521.67	-22403.53	-22420.53	-22481.82

Data Generator 성능 평가 via TABU Search

Dataset	Num. of Nodes	Sample Size	Score via TABU				
			ILP	BDe	BDe	loglik	AIC
Asia real dataset	100	-245.64		-270.99	-206.14	-219.14	-236.08
		1000	-2317.41	-2342.17	-2249.24	-2266.24	-2307.96
		10000	-22466.40	-21996.98	-21870.55	-21889.55	-21958.05
	1000			-270.74	-208.81	-221.81	-238.74
		100		-2350.21	-2262.16	-2279.16	-2320.87
		10000		-22529.57	-22419.24	-22437.24	-22502.14
	8	100		-	-	-	-
		1000		-2289.26	-2197.43	-2214.43	-2256.15
		10000		-22521.67	-22403.53	-22420.53	-22481.82
Gen. by WEKA							
Gen. by Generator							

Data Generator 성능 평가 via MMHC

Dataset	Num. of Nodes	Sample Size	ILP	Score via MMHC			
			BDe	BDe	loglik	AIC	BIC
Asia real dataset	100	-245.64	-301.29	-246.32	-257.32	-271.65	
		1000	-2317.41	-2504.74	-2421.80	-2436.80	-2473.61
		10000	-22466.40	-24346.08	-24232.89	-24249.89	-24311.18
	1000		-272.31	-213.08	-225.08	-240.71	
		1000		-2508.05	-2423.51	-2439.51	-2478.77
		10000		-22815.07	-22709.96	-22725.96	-22783.64
	Gen. by WEKA	100		-	-	-	-
		1000		-2446.05	-2360.82	-2376.82	-2416.08
		10000		-24137.75	-24026.34	-24042.34	-24100.02
Gen. by Generator							

Data Generator 성능 평가 via RSMAX2

Dataset	Num. of Nodes	Sample Size	ILP	Score via RSMAX2			
			BDe	BDe	loglik	AIC	BIC
Asia real dataset	100	-245.64	-299.97	-246.07	-260.07	-278.31	
		1000	-2317.41	-2531.66	-2451.91	-2464.91	-2496.81
		10000	-22466.40	-24295.92	-24194.06	-24207.06	-24253.93
	1000		-272.31	-213.08	-225.08	-240.71	
		1000		-2534.36	-2456.36	-2469.36	-2501.26
		10000		-22823.55	-22715.80	-22730.80	-22784.88
	10000	100		-	-	-	-
		1000		-2479.98	-2401.17	-2414.17	-2446.07
		10000		-24504.43	-24403.62	-24416.62	-24463.49

Outline

1 서론

- 목표

2 비교 기준

- 비교 기준
- 어떻게 비교할 것인가?
- 전제

3 Real Data를 이용한 구조 학습 결과 비교

- Asia DataSet by Lauritzen and Spiegelhalter
- Insurance Evaluation Network DataSet
- ALARM Monitoring System DataSet
- The HailFinder Weather Forecast System DataSet

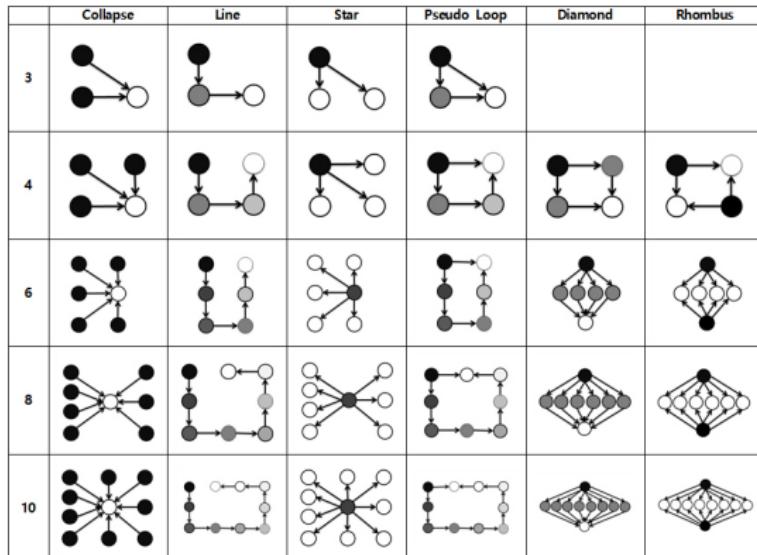
4 Data Generator 성능 평가

5 Topology에 따른 비교

6 결론

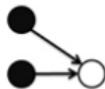
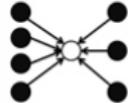
- 결론
- Discussion

Topology에 따른 비교

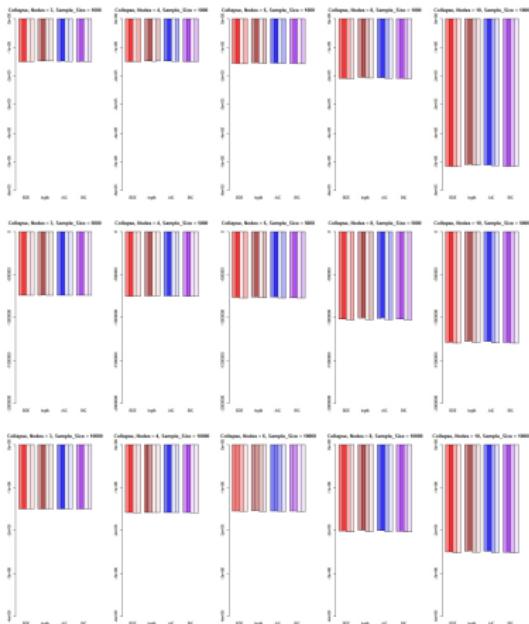


Eitel J. M. Lauría,
 "An Information-Geometric Approach to Learning Bayesian Network Topologies from Data",
 Innovations in Bayesian Networks Studies in Computational Intelligence Volume 156, 2008, pp 187-217

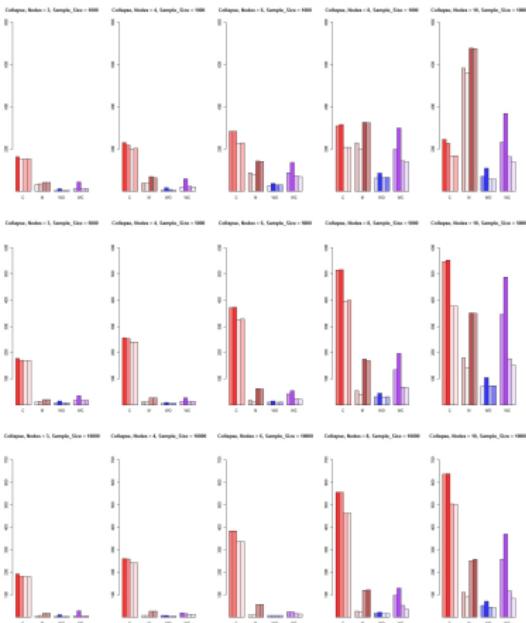
Topology에 따른 비교 분석 : Collapse

	3	4	6	8	10
Collapse					

Topology에 따른 비교 : Collapse (Score)



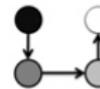
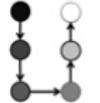
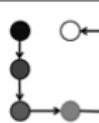
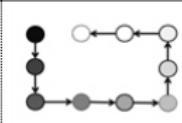
Topology에 따른 비교 : Collapse (Arcs)



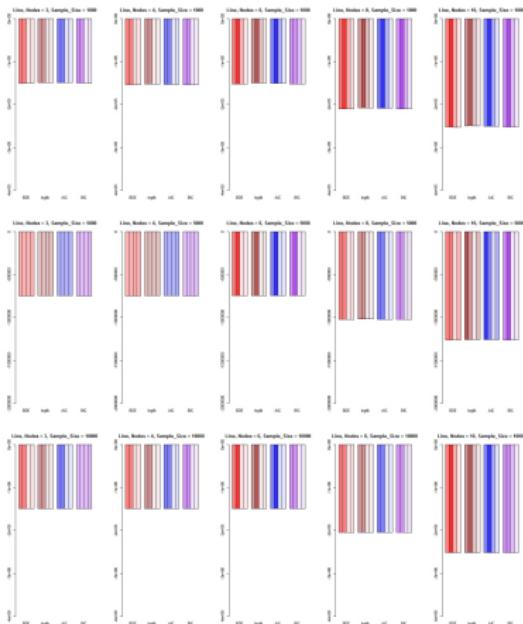
Topology에 따른 비교 : Collapse

Sample Size 1000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	4	2	2	4	3	1	1	4	1	4	4	4	1	4	4
4	2	1	4	3	1	2	4	3	3	4	1	2	4	1	2	4	4	1	2	4
6	2	1	4	3	1	1	4	3	3	4	1	2	4	1	2	3	2	1	3	4
8	2	1	4	3	2	1	4	4	3	4	1	2	4	1	3	2	2	1	3	4
10	2	1	3	4	1	2	3	4	3	4	1	2	2	1	4	4	2	1	3	4
Sample Size 5000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	2	4	4	4	3	1	1	4	1	4	4	4	1	4	4
4	2	1	4	4	1	2	4	4	3	4	1	1	4	1	2	2	4	1	4	4
6	2	1	4	3	2	1	4	3	3	4	1	2	2	1	4	2	2	1	3	4
8	2	1	4	3	2	1	4	3	3	4	1	2	2	1	4	2	2	1	3	4
10	2	1	3	4	2	1	4	4	3	4	1	2	2	1	4	2	2	1	3	4
Sample Size 10000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	2	4	4	4	3	1	1	4	1	4	4	4	1	4	4
4	2	1	4	4	1	2	4	4	4	4	1	1	4	1	4	4	1	2	4	4
6	2	1	3	4	1	1	4	4	4	4	1	1	1	1	1	1	1	1	3	4
8	2	1	3	4	1	2	3	4	3	4	2	1	2	1	4	4	2	1	3	4
10	2	1	3	4	2	1	3	4	3	4	2	1	2	1	3	4	2	1	3	4

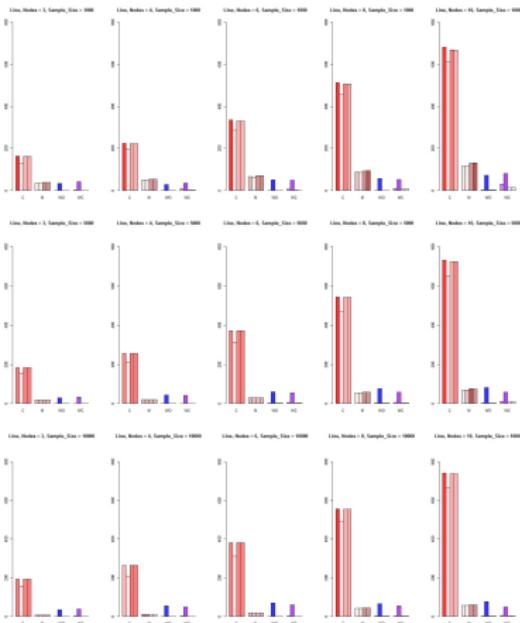
Topology에 따른 비교 분석 : Line

	3	4	6	8	10
Line					

Topology에 따른 비교 : Line (Score)



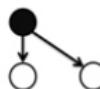
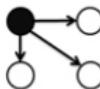
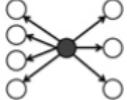
Topology에 따른 비교 : Line (Arcs)



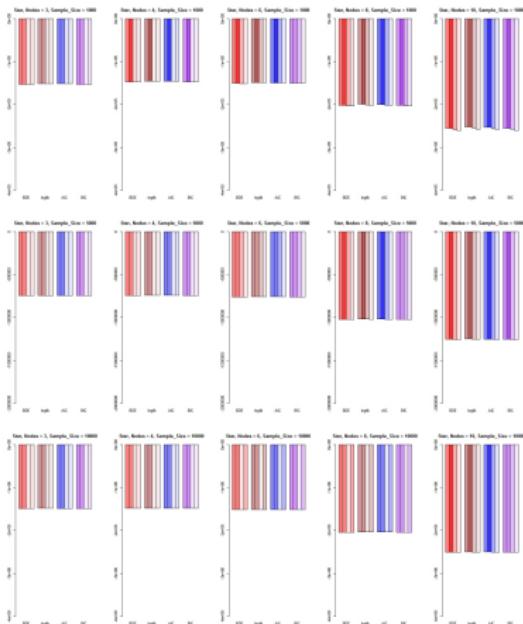
Topology에 따른 비교 : Line

Sample Size 1000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	4	4
4	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	3	4
6	2	1	3	4	1	4	2	2	3	4	1	1	4	1	4	4	2	1	3	4
8	2	1	3	4	1	4	2	3	3	4	2	1	4	1	4	4	2	1	4	4
10	2	1	3	4	1	4	2	3	4	4	2	1	4	1	4	4	2	1	4	4
Sample Size 5000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	1	1	1	4	1	1	1	1	1	1	4	1	4	4	4	1	4	4
4	1	1	1	1	1	4	1	1	1	1	1	1	4	1	4	4	4	1	4	4
6	2	1	4	4	1	4	1	1	1	1	1	1	4	1	4	4	2	1	4	4
8	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	4	4
10	1	1	4	3	1	4	3	2	4	3	1	2	4	1	4	4	2	1	4	4
Sample Size 10000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	4	4	1	4	1	1	1	1	1	1	4	1	4	4	2	1	4	4
4	1	1	4	4	3	4	1	1	1	1	4	4	4	1	4	4	2	1	4	4
6	2	1	4	4	1	4	1	1	1	1	1	1	4	1	4	4	2	1	4	4
8	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
10	2	1	3	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	3	4

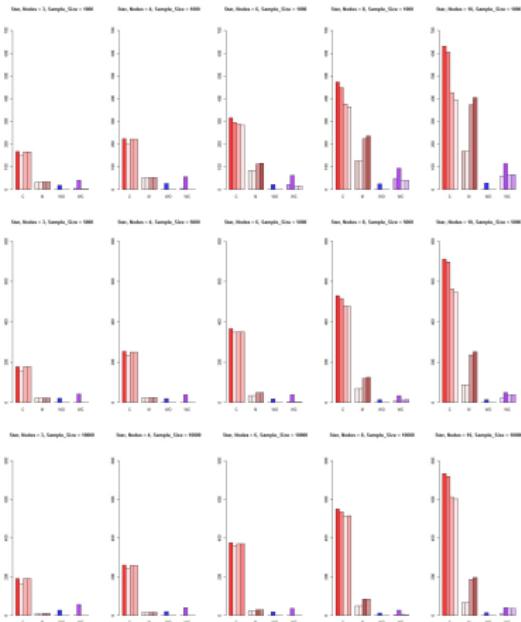
Topology에 따른 비교 분석 : Star

	3	4	6	8	10
Star					

Topology에 따른 비교 : Star (Score)



Topology에 따른 비교 : Star (Arcs)



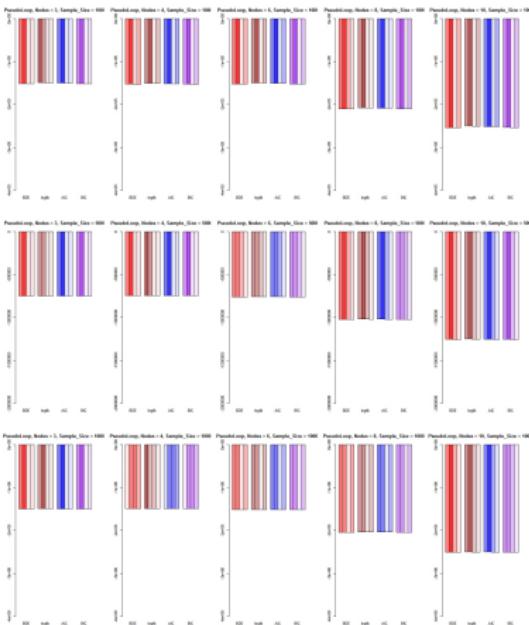
Topology에 따른 비교 분석 : Star

Sample Size 1000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	4	4
4	2	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	2	1	4	4
6	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	2	1	4	4
8	2	1	3	4	1	2	3	4	3	4	2	1	4	1	4	4	2	1	4	3
10	2	1	3	4	1	2	3	4	3	4	2	1	4	1	4	4	4	1	3	2
Sample Size 5000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
4	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
6	1	1	3	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
8	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	3	2
10	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	3	2
Sample Size 10000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
4	1	1	4	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
6	1	1	4	3	1	2	4	3	4	4	1	2	4	1	4	4	4	1	4	4
8	1	1	4	3	1	2	4	3	4	4	1	2	4	1	4	4	4	1	2	3
10	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	2	2

Topology에 따른 비교 분석 : PseudoLoop

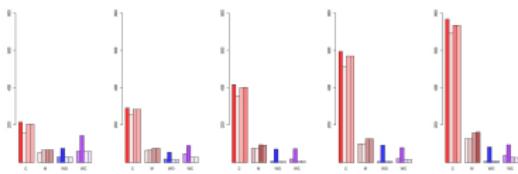
	3	4	6	8	10
Pseudo Loop					

Topology에 따른 비교 분석 : PseudoLoop (Score)

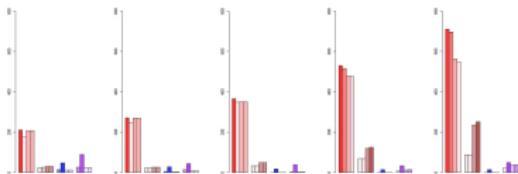


Topology에 따른 비교 분석 : PseudoLoop (Arc)

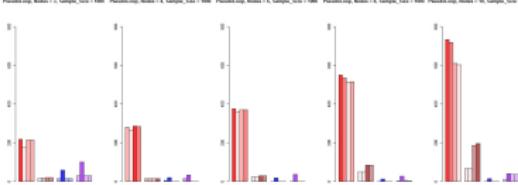
PseudoLoop_Nodes = 1, Sample_Size = 1000 PseudoLoop_Nodes = 4, Sample_Size = 1000 PseudoLoop_Nodes = 5, Sample_Size = 1000 PseudoLoop_Nodes = 6, Sample_Size = 1000 PseudoLoop_Nodes = 10, Sample_Size = 1000



PseudoLoop_Nodes = 1, Sample_Size = 1000 PseudoLoop_Nodes = 4, Sample_Size = 1000 PseudoLoop_Nodes = 5, Sample_Size = 1000 PseudoLoop_Nodes = 6, Sample_Size = 1000 PseudoLoop_Nodes = 10, Sample_Size = 1000



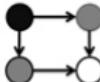
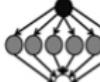
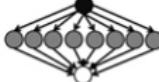
PseudoLoop_Nodes = 1, Sample_Size = 1000 PseudoLoop_Nodes = 4, Sample_Size = 1000 PseudoLoop_Nodes = 5, Sample_Size = 1000 PseudoLoop_Nodes = 6, Sample_Size = 1000 PseudoLoop_Nodes = 10, Sample_Size = 1000



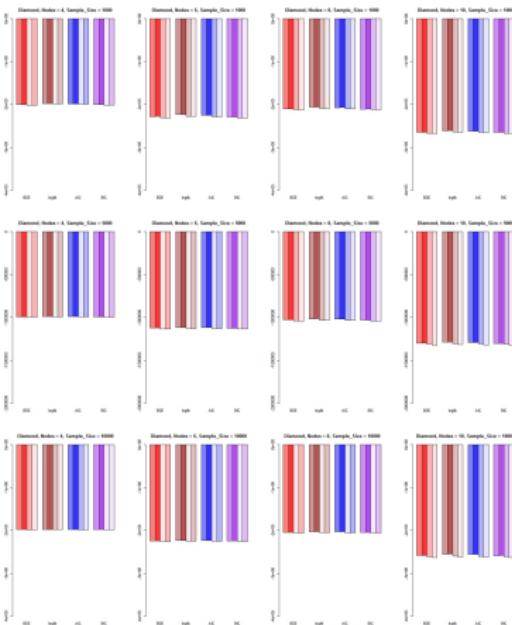
Topology에 따른 비교 분석 : PseudoLoop

Sample Size 1000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	4	2	2	4	3	1	1	2	1	4	4	2	1	4	4
4	2	1	4	3	1	4	2	2	4	3	1	1	2	1	4	4	2	1	4	4
6	2	1	4	3	1	4	3	2	3	4	1	2	2	1	4	4	2	1	4	4
8	2	1	3	4	1	4	2	2	3	4	1	1	4	1	4	4	2	1	4	4
10	2	1	4	3	1	4	2	3	4	3	2	1	4	1	4	4	2	1	3	4
Sample Size 5000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	4	2	2	3	2	1	1	2	1	4	4	2	1	4	4
4	2	1	4	4	1	4	2	2	4	4	1	1	2	1	4	4	2	1	4	4
6	1	1	3	4	1	4	2	2	4	4	1	1	4	1	4	4	4	1	4	4
8	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	3	2
10	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	3	2
Sample Size 10000	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
3	2	1	4	4	1	4	2	2	4	3	1	1	4	1	4	4	4	1	4	4
4	4	1	1	3	2	4	1	2	4	4	4	1	2	1	4	4	2	1	4	4
6	1	1	4	3	1	4	3	2	4	4	1	2	4	1	4	4	4	1	4	4
8	1	1	4	3	1	2	4	3	4	4	1	2	4	1	4	4	4	1	2	3
10	2	1	3	4	1	2	3	4	4	4	2	1	4	1	4	4	4	1	2	2

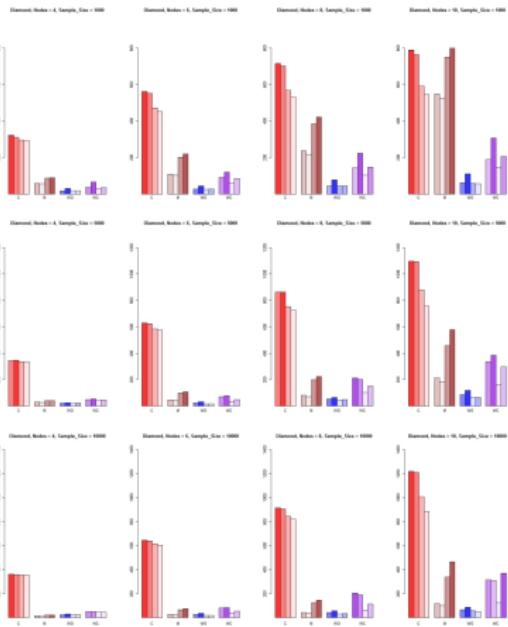
Topology에 따른 비교 분석 : Diamond

	3	4	6	8	10
Diamond					

Topology에 따른 비교 분석 : Diamond (Score)



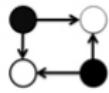
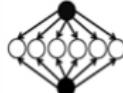
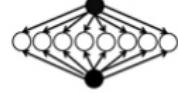
Topology에 따른 비교 분석 : Diamond (Arc)



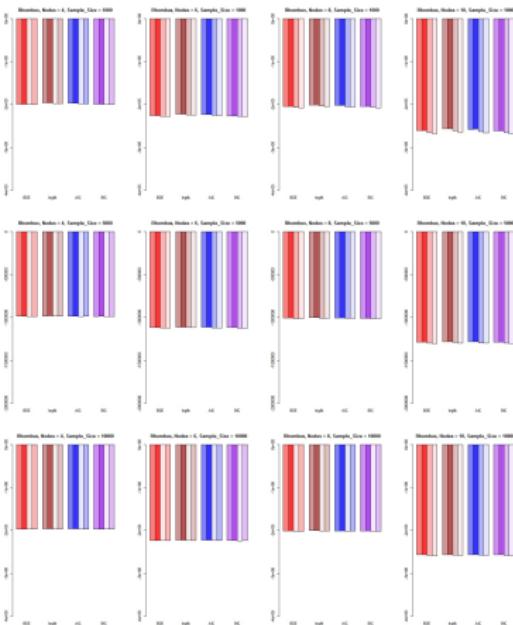
Topology에 따른 비교 분석 : Diamond

Sample Size	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
1000	2	1	4	3	1	2	3	4	3	4	2	1	2	1	4	4	2	1	4	3
4	2	1	3	4	1	2	3	4	3	4	2	1	2	1	4	2	2	1	4	2
6	2	1	3	4	1	2	3	4	3	4	2	1	4	2	1	4	2	2	1	4
8	2	1	3	4	1	2	3	4	3	4	2	1	4	1	4	4	3	1	4	2
10	2	1	3	4	1	2	3	4	3	4	2	1	2	1	3	4	3	1	4	2
<hr/>																				
Sample Size	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
5000	2	1	4	3	2	1	3	4	3	4	1	1	2	1	4	3	2	1	4	3
4	2	1	4	3	1	2	3	4	3	4	2	1	2	1	4	4	2	1	4	3
6	2	1	4	3	2	1	3	4	3	4	2	1	2	1	4	3	2	1	4	3
8	2	1	3	4	1	2	3	4	3	4	2	1	2	1	4	3	1	2	4	3
10	2	1	3	4	1	2	3	4	3	4	2	1	2	1	4	3	2	1	4	3
<hr/>																				
Sample Size	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
10000	2	1	3	4	1	2	3	4	3	4	2	1	2	1	4	4	1	1	4	4
4	2	1	4	3	1	2	3	4	3	4	2	1	2	1	4	4	2	1	4	3
6	2	1	4	3	1	2	3	4	3	4	2	1	2	1	4	3	2	1	4	3
8	2	1	4	3	1	2	3	4	3	4	2	1	2	1	4	3	1	2	4	3
10	2	1	3	4	1	2	3	4	3	4	2	1	2	1	3	4	2	3	4	1

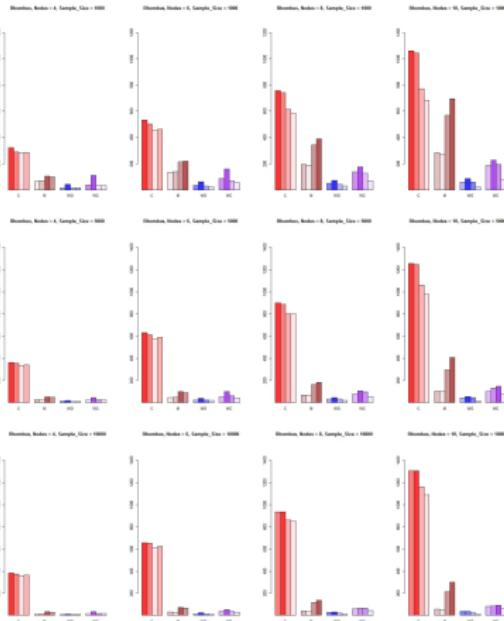
Topology에 따른 비교 분석 : Rhombus

	3	4	6	8	10
Rhombus					

Topology에 따른 비교 분석 : Rhombus (Score)



Topology에 따른 비교 분석 : Rhombus (Arc)



Topology에 따른 비교 분석 : Rhombus

Sample Size	Score				C				M				WO				WC			
	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
1000	2	1	4	3	1	2	4	3	4	3	1	2	4	1	4	4	2	1	4	4
4	2	1	3	4	1	2	4	3	4	3	2	1	2	1	3	4	2	1	3	4
6	2	1	3	4	1	2	3	4	3	4	2	1	2	1	3	4	2	1	3	4
8	2	1	3	4	1	2	3	4	3	4	2	1	2	1	3	4	2	1	3	4
10	2	1	3	4	1	2	3	4	3	4	2	1	3	1	2	4	3	1	2	4
Score				C				M				WO				WC				
5000	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
4	2	1	4	2	1	2	4	3	3	4	1	2	2	1	4	4	4	1	2	4
6	2	1	3	4	1	2	4	3	4	3	1	2	3	1	2	4	3	1	2	4
8	2	1	3	4	1	2	3	4	3	4	2	1	2	1	2	4	3	1	2	4
10	2	1	3	4	1	2	3	4	3	4	2	1	3	1	2	4	3	1	2	4
Score				C				M				WO				WC				
10000	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2	HC	TABU	MMHC	RSMAX2
4	2	1	4	3	1	2	4	3	4	3	1	2	4	1	4	4	4	1	2	4
6	2	1	4	3	1	2	4	3	3	4	1	2	2	1	2	4	2	1	2	4
8	2	1	4	3	2	1	3	4	3	4	2	1	2	1	3	4	3	1	1	4
10	2	1	3	4	2	1	3	4	3	4	2	1	1	1	3	4	3	2	1	4

결론

Algorithm User가 판단해야 할 문제

- Score로 평가할 것인가, Model을 직접 비교하여 평가할 것인가?
- Model로 평가할 때, M, WO, WC에 대한 가중치는?

결론

bnlearn 패키지에서 제공하는, Score-Based Algorithm과 Hybrid Algorithm을 집중 비교해 본 결과 다음과 같은 결론을 내일 수 있었다.

- "Score가 무엇이 더 높은가?", "Correct Arcs가 무엇이 더 많은가?" 어떤 것으로 비교해도 Score-Based 알고리즘(HC, TABU)가 Hybrid 알고리즘(MMHC, RSMAX2)보다 좋게 나타났다.
- 그러나 Score-Based 알고리즘을 적용했을 때, Wrongly Oriented and Connected Arcs도 많이 나타났다.

특히 HC와 TABU만을 놓고 보면, TABU가 score는 더 좋지만, 학습된 모형을 보면 WO, WC 비중이 압도적으로 높았다.

만일 학습된 결과 WO, WC가 많은 것이 C가 적은 것보다 문제가 되는 분야라면, Hybrid Algorithm을 선택하는 것도 방법이 될 수 있다.

결론

bnlearn 패키지에서 제공하는, Score-Based Algorithm과 Hybrid Algorithm을 집중 비교해 본 결과 다음과 같은 결론을 내일 수 있었다.

- Topology별로 ???
- Sample Size별로 ??
- Node 개수별로 ??

Discussion

(컴퓨팅 시간을 기다릴 인내심만 있다면...)

- Topology의 Node 개수를 더 늘려볼 필요도 있다.
- 다른 Algorithm을 적용하여, Sample Size를 더 적게하여, Cardinality를 증가시켜 추가 실험을 진행해 볼 수 있다.
- 서로 다른 Topology를 결합하여 비교 분석 할 수도 있다.
- Bayesian Network Data Generator의 R 패키징 완성
- 확률 관계를 정의할 때 확률값을 $U(0,1)$ 사이의 값에서 임의로 주었는데, 이 확률값을 "순차적"으로 준 것과의 관계를 알아볼 수도 있다.

다른 연구 과제

- Bayesian Network를 이용하여 Missing Value를 컨트롤
(이 때 Bayesian Network Data Generator를 적극 활용할 수 있다.)
 1. 배경 지식을 바탕으로 만든 모형을 이용하여 데이터를 생성하는 방법
 2. Missing이 없는 부분을 이용하여 먼저 구조학습을 한 후, 학습된 모형으로 데이터를 생성하는 방법
 3. 조건부가 Missing인 부분에 대한 조건부 확률관계를 표현하는 Node를 추가로 생성하는 방법