Probabilistic Model Building GP with Belief Propagation

Hiroyuki Sato Yoshihiko Hasegawa,

Danushka Bollegala and Hitoshi Iba



June. 12, 2012

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

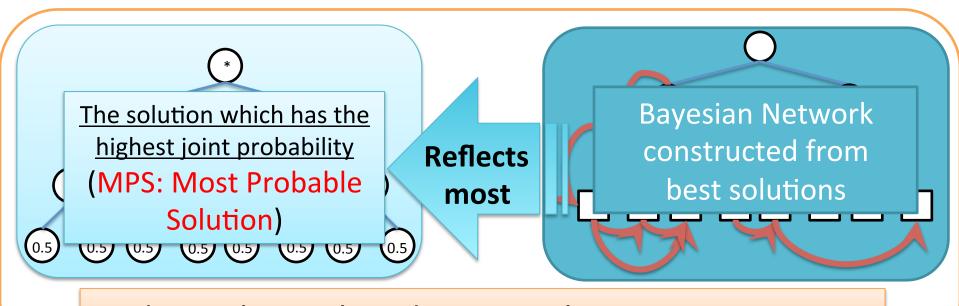
Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

PMBGP: Probabilistic Model Building Genetic Programing (GP)

- EDAs deal with tree structures, like GP
- PCFG (Probabilistic Context Free Grammar) based methods
 - Learn production rule probabilities
- 2. Prototype tree based methods (current research)
 - Convert full α-ary tree structures to 1 dimensional arrays in breadth-first order and apply GA-type EDAs

The weak point of prototype tree based PMBGP



traditional sampling does not always generate MPS

The same problem is pointed out in GA-type EDAs, but sampling using Loopy Belief Propagation is known as a way to overcome it

Motivation

Improvement of sampling to generate always MPS will enhance PMBGP

Loopy Belief Propagation, which effects EDAs well, is the appropriate way to improve sampling of prototype tree based PMBGP

We propose POLE-BP, and compare its performance against POLE and Simple GP in 3 benchmark experiments

POLE-BP

POLE (a prototype tree based PMBGP) + Loopy Belief Propagation

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

POLE: Program Optimization with Linkage Estimation

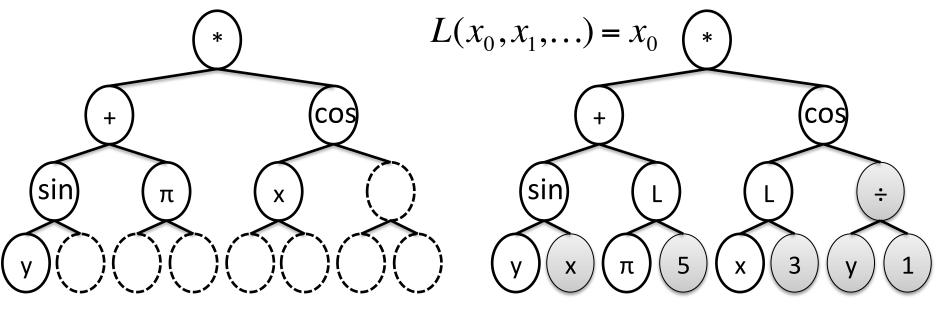
- Prototype tree based PMBGP
- Use Bayesian Network as probabilistic model
- Use EPT (Expanded Parse Tree) as chromosome
 - Guarantee syntactic correctness
 - Reduce symbols on trunk and conditional probability table size, and estimate probabilities more accurately

EPT: Expanded Parse Tree

$$f(x,y) = (\sin(y) + \pi)\cos(x)$$

Normal GP tree

EPT (gray nodes are introns)



- 1. Initialization
- 2. Evaluation
- 3. Selection
- 4. Construction of Bayesian Network by K2 algorithm
 - Bayesian information criteria (BIC) is used to evaluate networks
- 5. Generation of new individuals
 - Sampling ALL individuals by Probabilistic Logic
 Sampling (PLS)

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

LBP: Loopy Belief Propagation (Loopy max-sum)

- Infer MPS approximately, but succeeded in many applications: image processing, multi-agent system and EDAs
- Repeatedly update Messages, which are locally calculated joint probabilities
- Must be run on factor graph

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

POLE-BP: The proposed method

- Initialization
- 2. Evaluation
- 3. Selection
- 4. Construction of Bayesian Network
- 5. Generation of new individuals by sampling
 - Sample M-1 individuals (M: population size)
- 6. Graph transformation
 - The Bayesian Network constructed at step 4 is transformed to the equivalent factor graph
- 7. Generate MPS by LBP on the factor graph

Comparison of individual generation

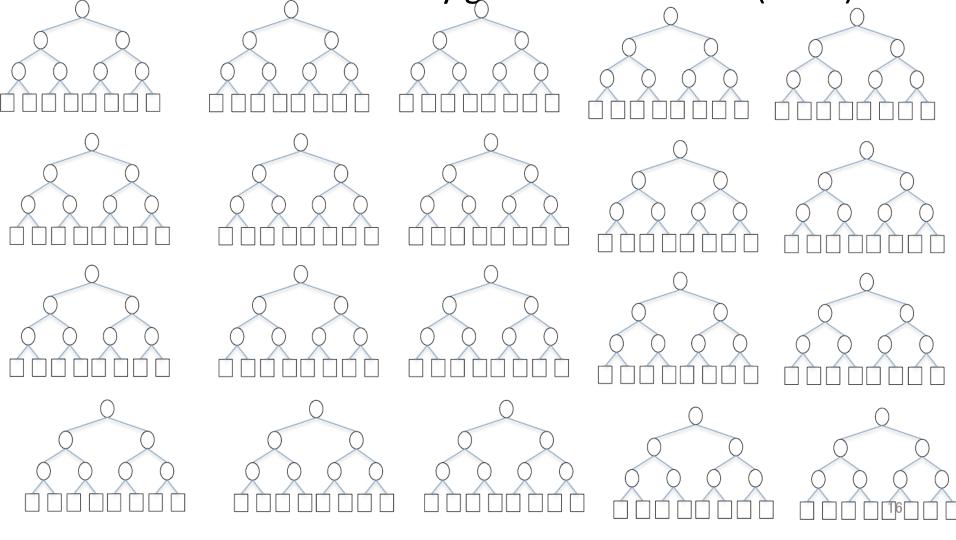
Conventional (POLE)

- 1. Building Bayesian Network
- Sampling M individuals from Bayesian Network
 (M: population size)
- 3. M sampled individuals are carried to the next generation

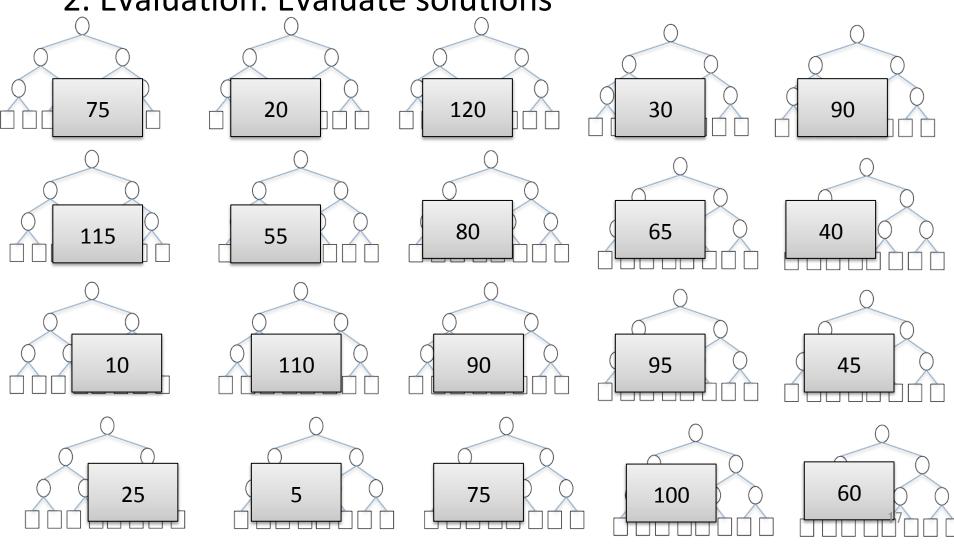
Proposed (POLE-BP)

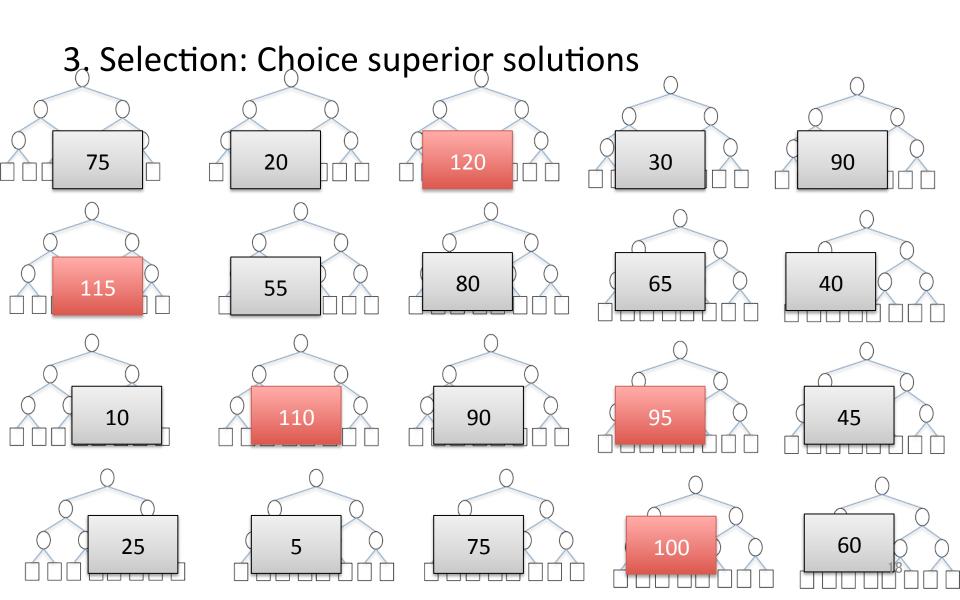
- 1. Building Bayesian Network
- Sampling M-1 individuals from Bayesian Network
- 3. Transforming Bayesian Network to Factor Graph
- 4. Generating the individual having the highest joint probability (MPS) by LBP
- M-1 sampled individuals and MPS generated by LBP are carried to the next generation

1. Initialization: Randomly generate solutions (trees)

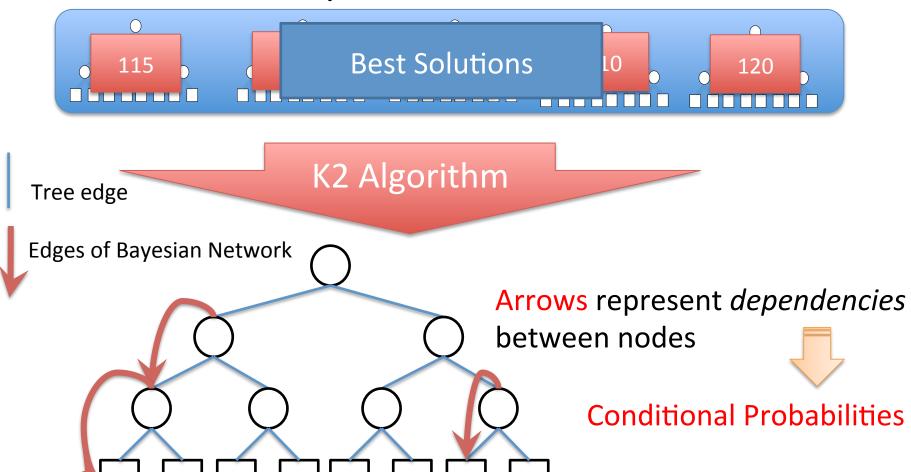


2. Evaluation: Evaluate solutions





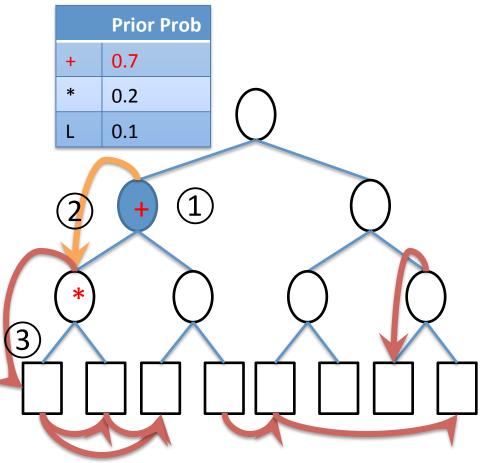
4. Construction of Bayesian Network on tree structure



5. Sampling M-1 individuals (M: population size):

Sample from ancestors

Parent	Child	Conditional Prob
+	+	0.6
+	*	0.2
+	L	0.2
*	+	0.05
*	*	0.15
*	L	0.8
L	+	0.25
L	*	0.2
L	L	0.55



6. Graph Transformation

Bayesian Network: Directed Graph

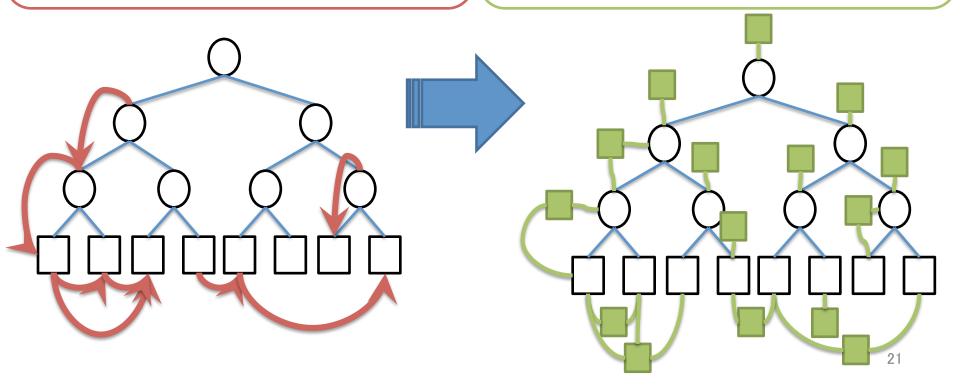
Directed Edge: conditional probability

Nodes without parents: prior probability

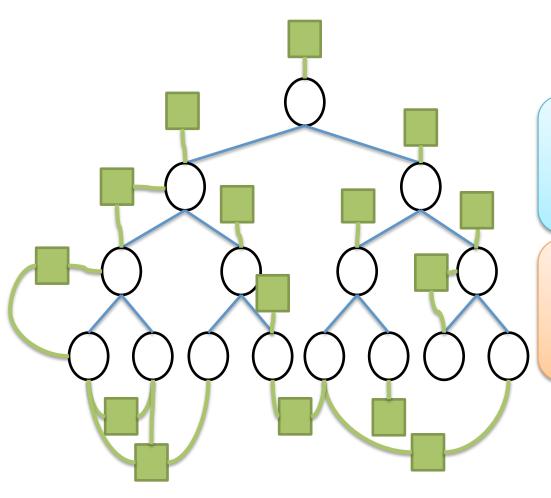
Factor Graph: Undirected Graph

Factor: prior and conditional probability

Undirected Edge: Only represent connection



7. Generate MPS by LBP on the factor graph



- 1. Initialize all messages to 0
- 2. Message Passing factor variable

Messages from factor to variable

$$\begin{aligned} & \text{leafs} \quad \mu_{f \to x}(x) = \ln f(x) \\ & \text{The others} \\ & \mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right] \end{aligned}$$

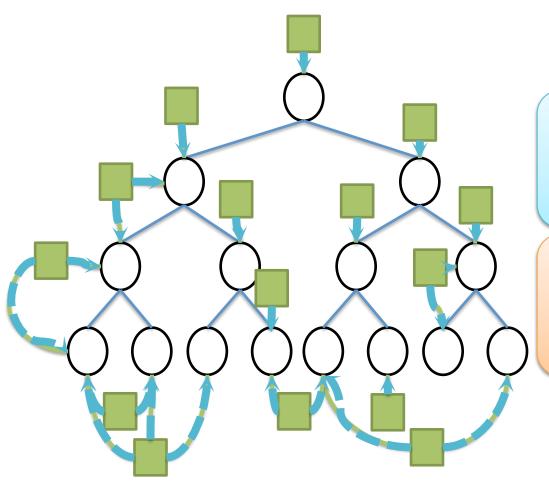
Messages from variable to factor

leafs
$$\mu_{x \to f}(x) = 0$$
 The others
$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

3. Get MPS

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)} \mu_{f_s \to x}(x) \right]$$

7. Generate MPS by LBP on the factor graph



- 1. Initialize all messages to 0
- 2. Message Passing factor variable

Messages from factor to variable

$$\begin{aligned} & \text{leafs} \quad \mu_{f \to x}(x) = \ln f(x) \\ & \text{The others} \\ & \mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right] \end{aligned}$$

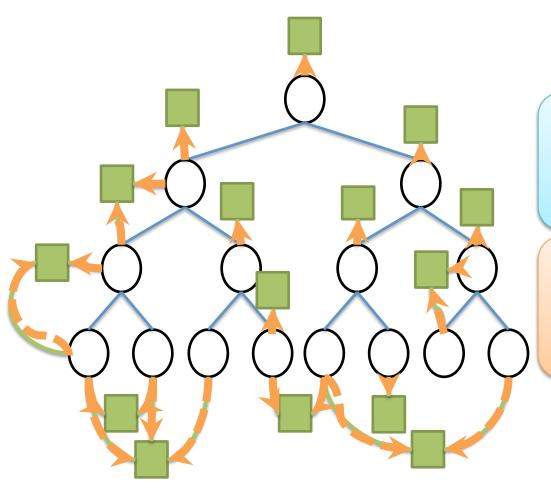
Messages from variable to factor

leafs
$$\mu_{x \to f}(x) = 0$$
 The others
$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

3. Get MPS

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)} \mu_{f_s \to x}(x) \right]$$

7. Generate MPS by LBP on the factor graph



- 1. Initialize all messages to 0
- 2. Message Passing factor variable

Messages from factor to variable

$$\begin{aligned} & \text{leafs} \quad \mu_{f \to x}(x) = \ln f(x) \\ & \text{The others} \\ & \mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right] \end{aligned}$$

Messages from variable to factor

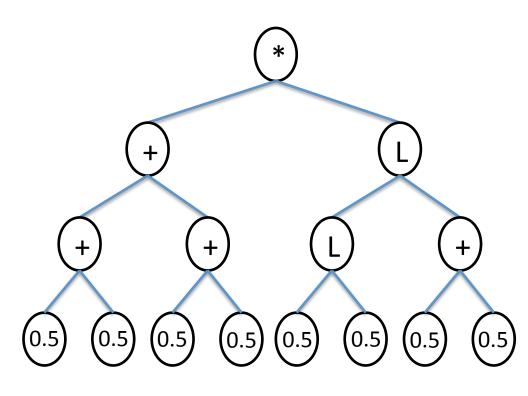
leafs
$$\mu_{x \to f}(x) = 0$$
 The others
$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

3. Get MPS

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)} \mu_{f_s \to x}(x) \right]$$

7. Generate MPS by LBP on the factor graph

Most Probable Solution



- 1. Initialize all messages to 0
- 2. Message Passing factor variable

Messages from factor to variable

$$\begin{aligned} & \text{leafs} \quad \mu_{f \to x}(x) = \ln f(x) \\ & \text{The others} \\ & \mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right] \end{aligned}$$

Messages from variable to factor

leafs
$$\mu_{x \to f}(x) = 0$$
 The others
$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

3. Get MPS

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)} \mu_{f_s \to x}(x) \right]$$

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

Experiments

- Benchmark on 3 problems with different characters:
 - MAX
 - Deceptive MAX (DMAX) (strong deceptiveness)
 - Royal Tree (many symbols)
- Compared the performance of POLE-BP(proposed),
 POLE and Simple GP(no mutation)
 - The relation between the average number of evaluations and tree size (the number of nodes)

POLE-BP(Proposed)

POLE(Conventional)

Simple GP

MAX

The number 1.0E+05 of evaluations 1.0E+04 Versus Tree size(the 1.0E+03 number of nodes)

Statistical

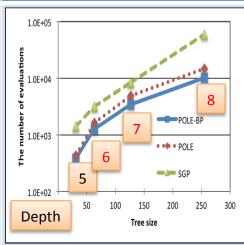
significance

&

Decreasing

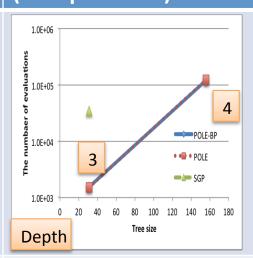
rate

(Normal)



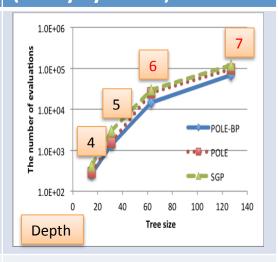
- POLE-BP is better for depth = 6, 7, 8at 1% significance level.
- Average 23.7% decreased.

DMAX (Deceptiveness)



No difference.

Royal Tree (Many symbols)



- POLF-BP is better for depth = 6, 7 at 1%significance level.
- Average 23.6% decreased.

POLE-BP(Proposed)

PC E(Conventional)

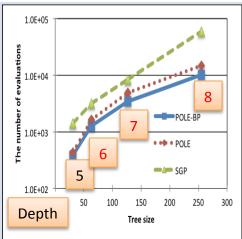
Simple GP

MAX

The number of evaluations Versus Tree size(the number of nodes)

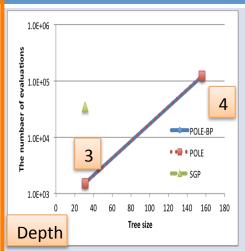
Statistical significance & Decreasing rate

(Normal)



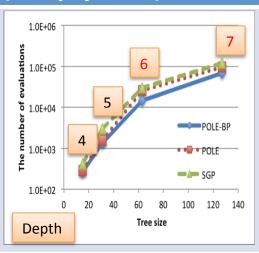
- POLF-BP is better for depth = 6, 7, 8at 1% significance level.
- Average 23.7% decreased.

DMAX (Deceptiveness)



No difference.

Royal Tree (Many symbols)



- POLF-BP is better for depth = 6, 7 at 1%significance level.
- Average 23.6% decreased.

POLE-BP(Proposed)

POLE(Conventiona'

Simple GP

MAX

(Normal)

The number of evaluations

Versus

Tree size(the number of nodes)

1.0E+03

1.0E+04

1.0E+03

1.0E+03

7

POLE-BP

1.0E+02

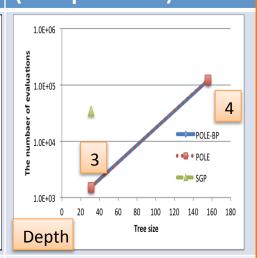
50
100
150
200
250
300

Tree size

Statistical significance & Decreasing rate

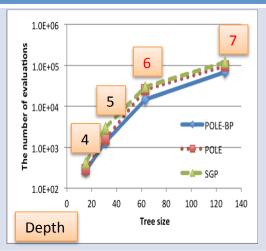
- POLE-BP is better for depth = 6, 7, 8 at 1% significance level.
- Average 23.7% decreased.

DMAX (Deceptiveness)



No difference.

Royal Tree (Many symbols)



- POLE-BP is better for depth = 6, 7 at 1% significance level.
- Average 23.6% decreased.

POLE-BP(Proposed)

POLE(Conventice al)

Simple GP

MAY

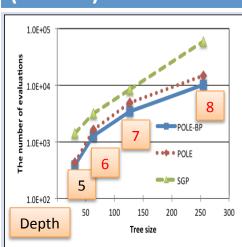
The number of evaluations

Versus

Tree size(the number of

nodes)

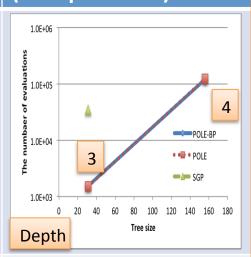
MAX (Normal)



Statistical significance & Decreasing rate

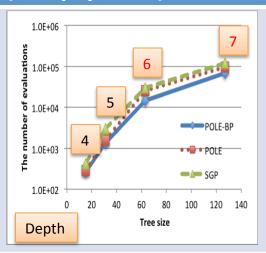
- POLE-BP is better for depth = 6, 7, 8 at 1% significance level.
- Average 23.7% decreased.

DMAX (Deceptiveness)



No difference.

Royal Tree (Many symbols)



- POLE-BP is better for depth = 6, 7 at 1% significance level.
- Average 23.6% decreased.

Outline

- Introduction
- POLE: Program Optimization with Linkage Estimation
- LBP: Loopy Belief Propagation
- POLE-BP: The proposed method
- Experiments & Discussion
- Conclusion

Conclusion

- We proposed the hybrid PMBGP combining POLE and Loopy Belief Propagation
- LBP reduces the average number of evaluations to get the optimum solution in PMBGP as well as GA-type EDAs
- LBP enhances the search ability in small population

Thank you for your attention

Hiroyuki SATO sato@iba.t.u-tokyo.ac.jp

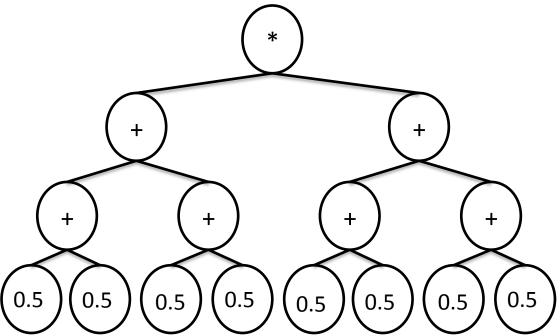
Future work

- Analyze how LBP works in PMBGP
- Real world application

MAX problem

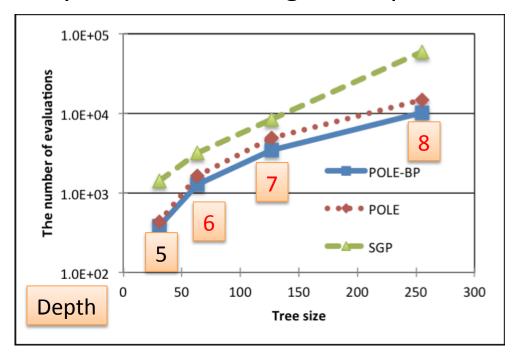
• Search the function which returns the biggest value $F = \{+, *\}, T = \{0.5\}$

The optimum solution: depth 4



MAX: Results and discussion

- Difference between POLE-BP and POLE for depth = 6, 7, 8 is statically significant at 1% significant level
- MAX problem has no deceptiveness
 - MPS always accelerates to get an optimum solution



DMAX: Deceptive MAX problem

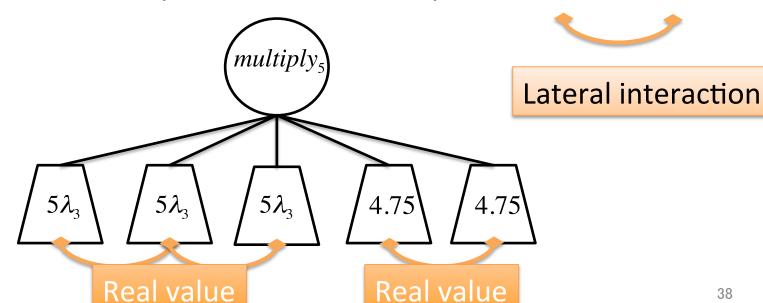
 Search the function which returns the biggest real value

$$F = \{add_5, multiply_5\}$$

$$T = \{\lambda_3, 0.95\}$$

$$\lambda_3 = \left(-\frac{1}{2} + \frac{i\sqrt{3}}{2}\right)$$

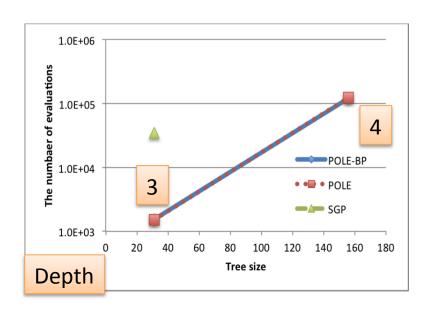
An optimum solution: depth 3

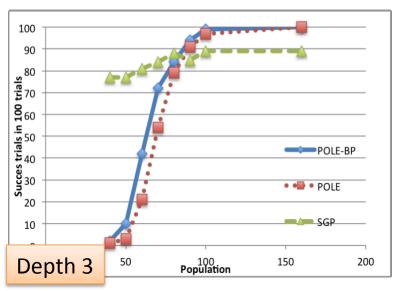


38

DMAX: Results and discussion

- No difference between POLE-BP and POLE in the average number of evaluations at 1% significant level
 - Because of deceptiveness, MPS does not work well to get an optimum
- In small population size, POLE-BP tends to get an optimal solution more than POLE
 - Because of sampling bias, sampling does not work well in small population. Therefore, LBP is effective in small population.

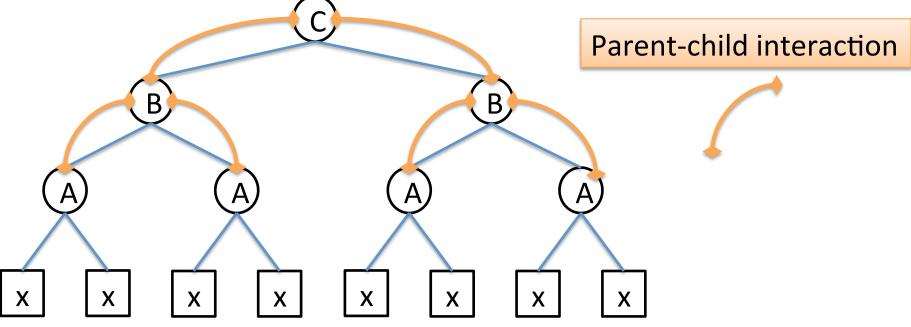




Royal tree problem

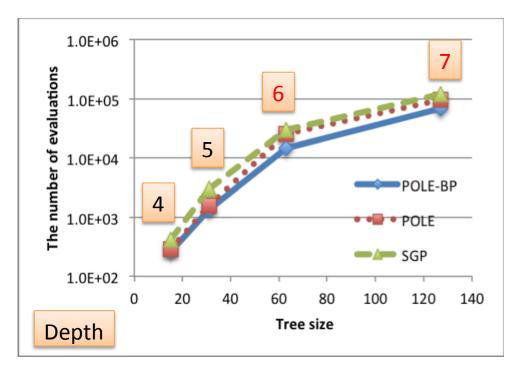
• Search *Perfect Tree* $F = \{A, B, C, D, E, F\}$ $T = \{x\}, x = 1$

The optimum solution(Perfect Tree): depth 4

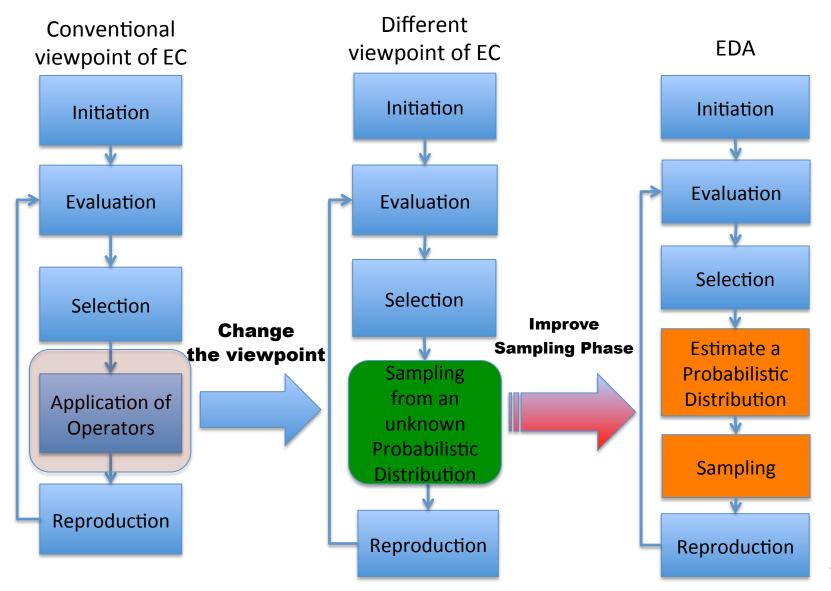


Royal tree: Results and discussion

- Difference between POLE-BP and POLE for depth 6, 7 is statically significant at 1% significant level.
- Royal tree uses more symbols (=depth) than MAX (3) and DMAX (4)
 - MPS generation by sampling is difficult, so MPS generation by LBP at each generation is effective



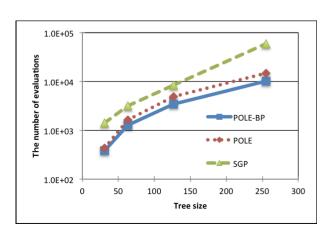
EDAs: Estimation Distribution Algorithms

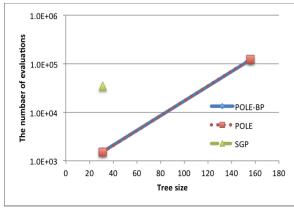


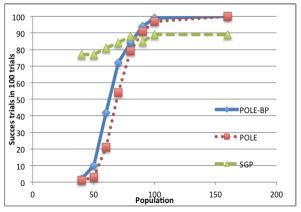
Experiments on 3 benchmarks

DMAX

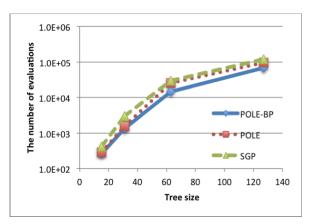
MAX





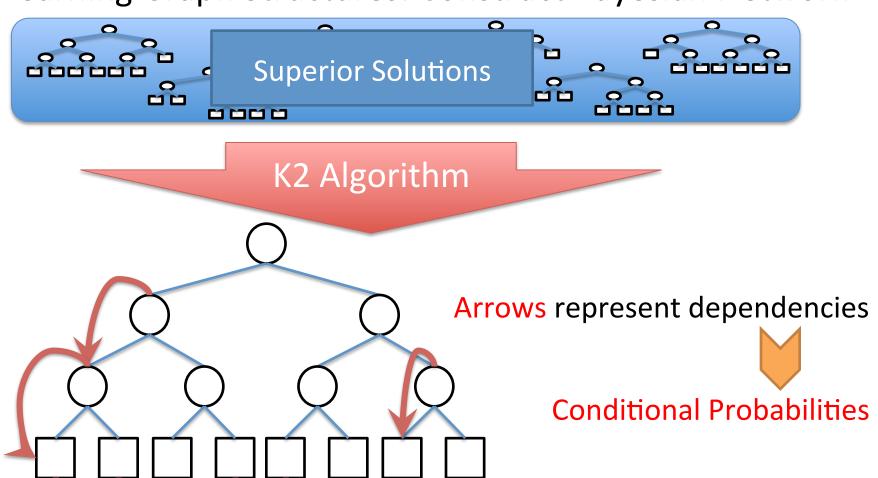


Royal tree



Estimation of PMBGP(Probabilistic Model Building GP) (GP version EDA, Bayesian Network)

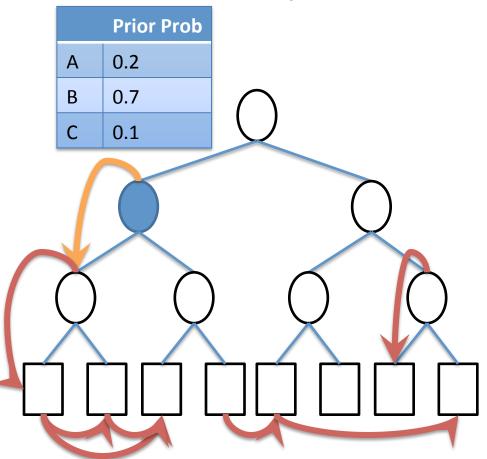
1. Learning Graph Structures: Construct Bayesian Network



Estimation of PMBGP

2. Parameter Estimation: Estimate conditional probabilities

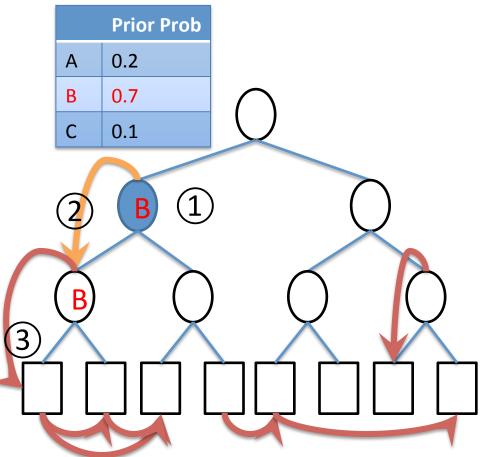
Parent	Child	Conditional Prob
А	А	0.1
А	В	0.7
А	С	0.2
В	А	0.05
В	В	0.15
В	С	0.8
С	А	0.25
С	В	0.2
С	С	0.55



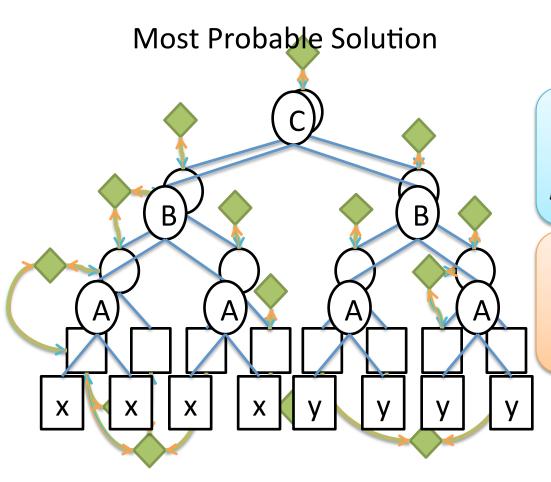
Sampling of PMBGP

Probabilistic Logic Sampling (PLS): Sample from ancestors

Parent	Child	Conditional Prob
А	А	0.1
А	В	0.7
А	С	0.2
В	А	0.05
В	В	0.15
В	С	0.8
С	А	0.25
С	В	0.2
С	С	0.55



Get MPS by Loopy Belief Propagation



- 1. Initialize all messages to 0
- 2. Message Passing

Messages from factor to variable

leafs
$$\mu_{f \to x}(x) = \ln f(x)$$

The others
$$\mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right]$$

Messages from variable to factor

leafs
$$\mu_{x \to f}(x) = 0$$
 The others
$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

3. Get MPS

Joint probability

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)} \mu_{f_s \to x}(x) \right]$$

Graph transformation

Bayesian Network: Directed Graph

Factor Graph: Undirected Graph



Directed Edge: conditional probability

Nodes without parents: prior probability

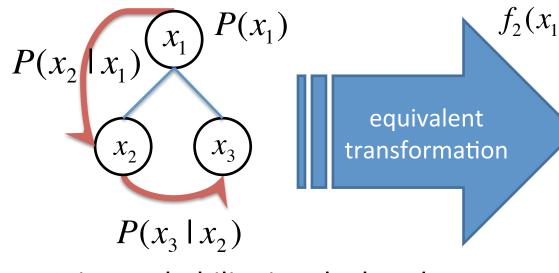
Blue edge: Tree structure

Factor: prior and conditional probability

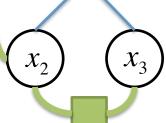
Undirected Edge: Only represent connection

$$f_1(x_1) = P(x_1)$$

Bayesian Network constructed on tree (depth=2)



 $f_2(x_1, x_2) = P(x_2 | x_1)$



$$f_3(x_2, x_3) = P(x_3 \mid x_2)$$

Joint Probability is calculated as

$$P(x_1, x_2, x_3) = P(x_1)P(x_2 \mid x_1)P(x_3 \mid x_2) = f_1(x_1)f_2(x_1, x_2)f_3(x_2, x_3)$$

Effective calculation of the highest joint probability

• Let the number of cardinalities of each node be $N_{x_1} = |F|, N_{x_2} = N_{x_3} = |T|$

$$\max_{x_1, x_2, x_3} (P(x_1, x_2, x_3)) = \max_{x_1, x_2, x_3} (P(x_1)P(x_2 \mid x_1)P(x_3 \mid x_2))$$
 The number of multiplication
$$= \max_{x_1, x_2, x_3} (f_1(x_1)f_2(x_1, x_2)f_3(x_2, x_3)) \qquad N_{x_1}N_{x_2} + N_{x_1}N_{x_2}N_{x_3}$$

$$= \max_{x_3} \left[\max_{x_2} \left[\max_{x_1} \left[f_1 f_2 \right] f_3 \right] \right] \qquad N_{x_1}N_{x_2} + N_{x_2}N_{x_3}$$

Permuting max and multiplication reduces the number of multiplication More intuitional representation of this effective procedure?

→ Message Passing

Message Passing on a *chain graph* (Loopy Belief Propagation)

 Send messages in order, finally get the highest joint probability

Message definition

$$\mu_{f_1 \to x_1}(x_1) = f_1(x_1),$$

$$\mu_{x_1 \to f_2}(x_1) = \mu_{f_1 \to x_1}(x_1),$$

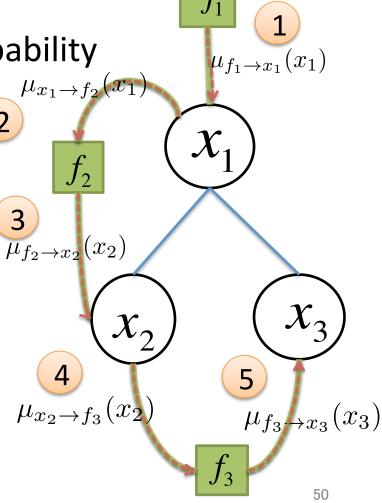
$$\mu_{f_2 \to x_2}(x_2) = \max_{x_1} \left[f_2(x_1, x_2) \mu_{x_1 \to f_2}(x_1) \right],$$

$$\mu_{x_2 \to f_3}(x_2) = \mu_{f_2 \to x_2}(x_2),$$

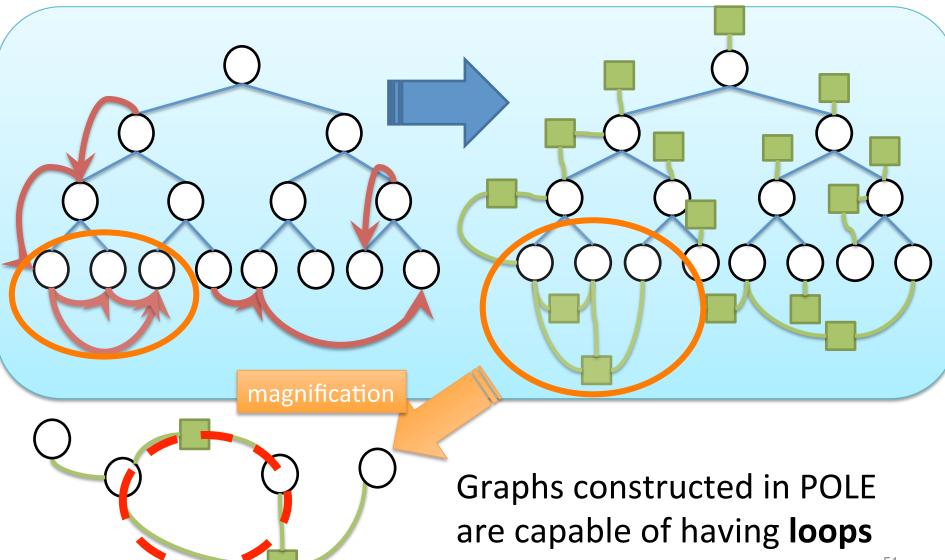
$$\mu_{f_3 \to x_3}(x_3) = \max_{x_2} \left[f_3(x_2, x_3) \mu_{x_2 \to f_3}(x_2) \right],$$

Calculation of the highest joint probability

$$\max(P(x_1, x_2, x_3)) = \max_{x_3} \left[\mu_{f_3 \to x_3} \right]$$
$$= \max_{x_3} \left[\max_{x_2} \left[\max_{x_1} \left[f_1 f_2 \right] f_3 \right] \right]$$

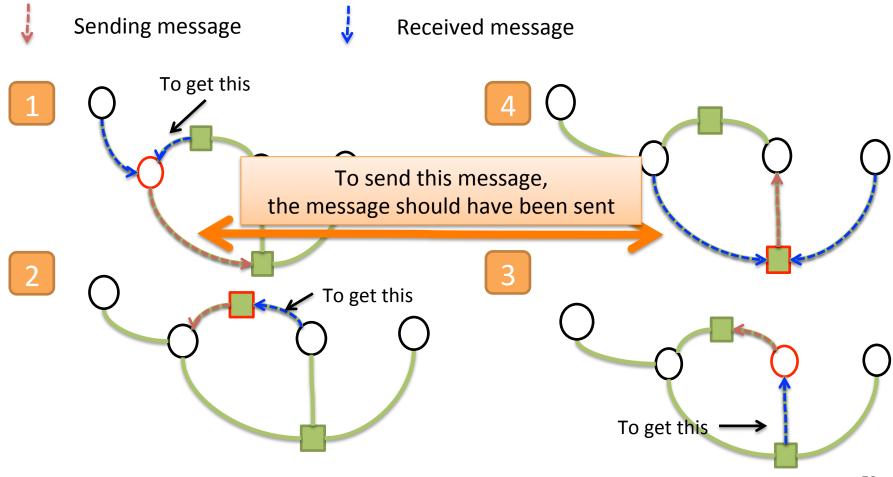


Problems of Message Passing in practice



Contradiction of Message Passing on graphs with loops

To send the message, the node should receive the other messages



Generalization of Message Passing

- To adapt loopy graphs,
 - 1. Set all messages with proper initial values
 - 2. Repeatedly send messages in arbitrary order
 - 3. Generally, messages converge to the right value
- To prevent underflow, substitute log and summation for multiplication in message
 - This substitution doesn't effect max

Graph Transformation

Bayesian Network : Directed Graph

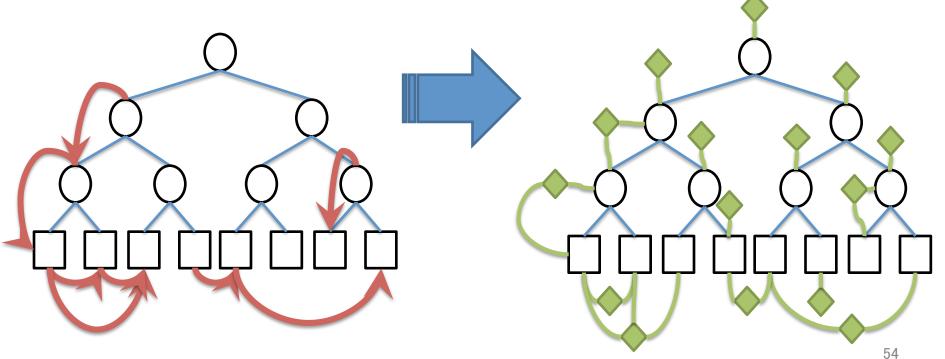
Directed Edge: conditional probability

Nodes without parents: prior probability

Factor Graph: Undirected Graph

Factor: prior and conditional probability

Undirected Edge: Only represent connection



Message Passing on loopy graphs LBP: Loopy Belief Propagation

1. Initialize all messages to 0

factor

variable





2. Message Passing repeatedly

Messages from factor to variable

Leafs

The others

$$\mu_{f \to x}(x) = \ln f(x)$$

$$\mu_{f \to x}(x) = \ln f(x) \qquad \mu_{f \to x}(x) = \max_{x, x_1, \dots, x_M} \left[\ln f(x, x_1, \dots x_M) + \sum_{m \in ne(f_s) \setminus x} \mu_{x_m \to f}(x_m) \right]$$

Messages from variable to factor

Leafs

The others

$$\mu_{x \to f}(x) = 0$$

$$\mu_{x \to f}(x) = \alpha_{xf} + \sum_{l \in ne(x) \setminus f} \mu_{f_l \to x}(x)$$

$$x^{\max} = \arg\max_{x} \left[\sum_{s \in ne(x)}^{\text{Joint probability}} \mu_{f_s \to x}(x) \right]$$

The weak point of prototype tree based PMBGP

- The solution which has the highest joint probability (MPS: Most Probable Solution) reflects the constructed Bayesian Network, which is usually used as probabilistic model.
- 2. However, traditional sampling does not always generate MPS
- 3. This problem wastes a part of learning and increases the number of evaluations to get an optimum solution.

The same problem is pointed out in GA-type EDAs, but sampling using Loopy Belief Propagation is known as a way to overcome it