

# A Method to Evaluate CFG Comparison Algorithms

2014 14th International Conference on Quality Software  
Patrick P.F. Chan, Christian Collberg

김영철

2016. 6. 29.

# EVALUATION METHODOLOGY

- Graph Edit Distance 개념을 기초로 함
  - Add a zero-degree node
  - Delete a zero-degree node
  - Add an edge between two existing nodes
  - Delete an existing edge
- scores 대신에 rankings 을 사용

# EVALUATION METHODOLOGY

- Evaluation Algorithm
  - Algorithm 1은 evaluation process 정보를 줌
  - generate test case
  - compare seed CFG with test case
  - rank the results
  - compute a goodness score

```
procedure EVALUATE( $\langle \mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_k \rangle, G_0, d, n$ )  
  for  $a \leftarrow 1, k$  do  
    for  $t \leftarrow 1, n$  do  
      for  $g \leftarrow 1, d$  do  
         $scores_g \leftarrow \Delta_{\mathcal{A}_a}(G_0, \text{TestCases}_{t,g})$   
      end for  
       $rank \leftarrow \text{rank}(scores)$   
       $\text{corr}_{\mathcal{A}_a,t} \leftarrow \text{pearson}((G_{t,1}, G_{t,2}, \dots, G_{t,d}), rank)$   
    end for  
     $\text{avg}_{\text{corr}_{\mathcal{A}_a}} = \text{average}(\text{corr}_{\mathcal{A}_a,1}, \text{corr}_{\mathcal{A}_a,2}, \dots, \text{corr}_{\mathcal{A}_a,t})$   
  end for  
  return rank( $\text{avg}_{\text{corr}}$ )  
end procedure
```

Algorithm 1

# EVALUATION METHODOLOGY

- Benchmark Graph Generation
  - 하나의 EDG (Edit Distance Graph)는 seed CFG로부터 생성될 수 있는 가능한 모든 CFG를 나타냄
    - EDG는 무한하기 때문에 depth  $d$ 를 설정
  - EDG nodes의 level이 생성됨
    - ex)  $i$  edit operations로 만들어진 EDG node는 level  $i$ 에 위치한다
    - $ED(G_0, G_i) = i$

# CFG SIMILARITY ALGORITHMS

- A. Based on k-subgraphs Mining
- B. Based on Edit Distance
- C. Based on Neighbor Matching
- D. Based on Simulation

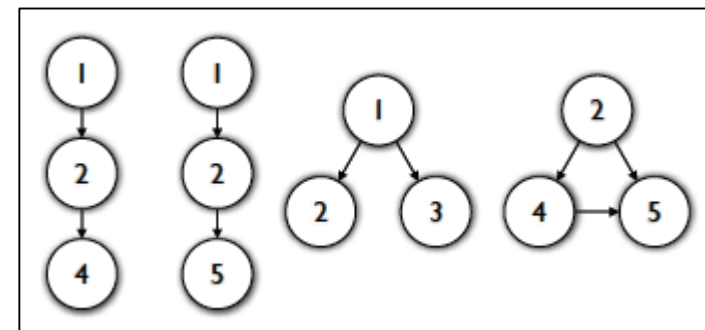
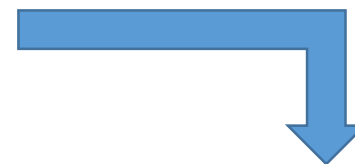
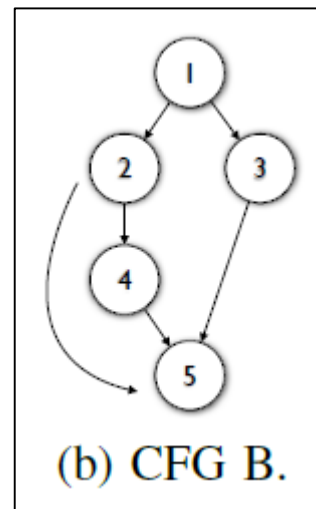
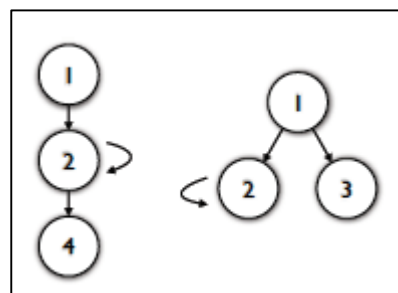
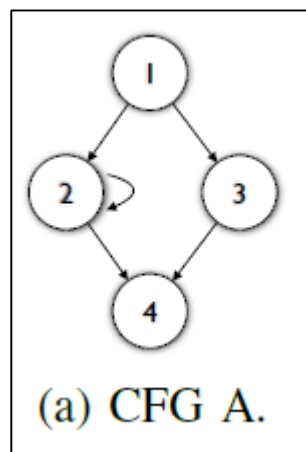
# Based on k-subgraphs Mining

- Step1) Generating k-subgraphs
  - maximum outdegree는 2, maximum indegree가 1인 graphs의 set 생성
- Step2) Graph Fingerprinting
  - subgraph들은 숫자로 맵핑
- Step3) Computing Similarity
  - Jaccard index를 사용하여 비교

$$\frac{|\text{fingerprint}(A) \cap \text{fingerprint}(B)|}{|\text{fingerprint}(A) \cup \text{fingerprint}(B)|}$$

Jaccard index

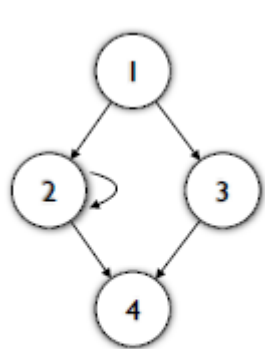
# Based on k-subgraphs Mining



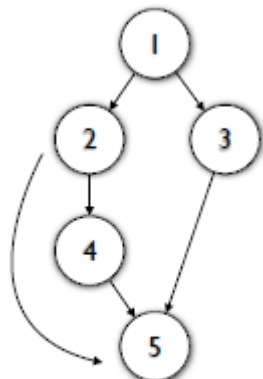
두 그래프의 유사도는 0

# Based on Edit Distance

- Step1) Building The Cost Matrix
  - $G_1, G_2$ 의 Vertex 수  $V_1, V_2$ 를 가지고  $(|V_1|+|V_2|)*(|V_1|+|V_2|)$  matrix 생성



(a) CFG A.



(b) CFG B.

|   | Cost of matching nodes          |          |          |          |          | Cost of deleting nodes in CFG A |          |          |          |
|---|---------------------------------|----------|----------|----------|----------|---------------------------------|----------|----------|----------|
| Cost of matching node 1 of CFG A to node 1 of CFG B | 0                               | 1        | 2        | 2        | 5        | 3                               | $\infty$ | $\infty$ | $\infty$ |
|   | 2                               | 1        | 2        | 2        | 3        | $\infty$                        | 5        | $\infty$ | $\infty$ |
|   | 2                               | 1        | 0        | 0        | 3        | $\infty$                        | $\infty$ | 3        | $\infty$ |
|   | 4                               | 3        | 2        | 2        | 1        | $\infty$                        | $\infty$ | $\infty$ | 3        |
| Cost of deleting node 1 of CFG B                    | 3                               | $\infty$ | $\infty$ | $\infty$ | $\infty$ | 0                               | 0        | 0        | 0        |
|   | $\infty$                        | 4        | $\infty$ | $\infty$ | $\infty$ | 0                               | 0        | 0        | 0        |
|   | $\infty$                        | $\infty$ | 3        | $\infty$ | $\infty$ | 0                               | 0        | 0        | 0        |
|   | $\infty$                        | $\infty$ | $\infty$ | 3        | $\infty$ | 0                               | 0        | 0        | 0        |
|   | $\infty$                        | $\infty$ | $\infty$ | $\infty$ | 4        | 0                               | 0        | 0        | 0        |
|   | Cost of deleting nodes in CFG B |          |          |          |          | Cost of matching dummy nodes    |          |          |          |
|   |                                 |          |          |          |          |                                 |          |          |          |

Cost of deleting node 4 of CFG B



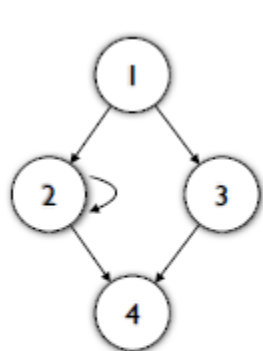
# Based on Edit Distance

- Step2) Finding Best Match of the Nodes
  - 가장 낮은 코스트를 갖는 매칭 비용을 찾는 것이 목적
  - Hungarian algorithm이 assignment problem의 optimal solution을 찾음
- Step3) Computing Similarity

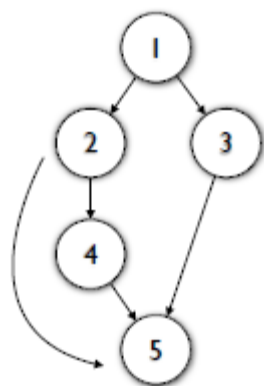
$$\text{sim}(G_1, G_2) = 1 - \frac{\text{edit distance}}{|V_1| + |E_1| + |V_2| + |E_2|}$$

# Based on Edit Distance

- Example
  - Hungarian algorithm을 이용하여 best matching 비용을 구함



(a) CFG A.



(b) CFG B.

|          |          |          |          |          |          |          |          |          |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 0        | 1        | 2        | 2        | 5        | 3        | $\infty$ | $\infty$ | $\infty$ |
| 2        | 1        | 2        | 2        | 3        | $\infty$ | 5        | $\infty$ | $\infty$ |
| 2        | 1        | 0        | 0        | 3        | $\infty$ | $\infty$ | 3        | $\infty$ |
| 4        | 3        | 2        | 2        | 1        | $\infty$ | $\infty$ | $\infty$ | 3        |
| 3        | $\infty$ | $\infty$ | $\infty$ | $\infty$ | 0        | 0        | 0        | 0        |
| $\infty$ | 4        | $\infty$ | $\infty$ | $\infty$ | 0        | 0        | 0        | 0        |
| $\infty$ | $\infty$ | 3        | $\infty$ | $\infty$ | 0        | 0        | 0        | 0        |
| $\infty$ | $\infty$ | $\infty$ | 3        | $\infty$ | 0        | 0        | 0        | 0        |
| $\infty$ | $\infty$ | $\infty$ | $\infty$ | 4        | 0        | 0        | 0        | 0        |

# Based on Neighbor Matching

- neighbor의 similarity에 근거하여 similarity matrix를 빌드함

- Step1) Topological Similarity

- $k=0, a_{ij}^0=1$

- Step2) Node Similarity

$$a_{ij}^{k+1} = \frac{s_{\text{in}}^{k+1}(i, j) + s_{\text{out}}^{k+1}(i, j)}{2}$$

where

$$s_{\text{in}}^{k+1}(i, j) = \frac{1}{m_{\text{in}}} \sum_{l=1}^{n_{\text{in}}} a_{f_{ij}^{\text{in}}(l)}^k g_{ij}^{\text{in}}(l)$$

$$s_{\text{out}}^{k+1}(i, j) = \frac{1}{m_{\text{out}}} \sum_{l=1}^{n_{\text{out}}} a_{f_{ij}^{\text{out}}(l)}^k g_{ij}^{\text{out}}(l)$$

$$m_{\text{in}} = \max(|\text{IE}_i|, |\text{IE}_j|)$$

$$n_{\text{in}} = \min(|\text{IE}_i|, |\text{IE}_j|)$$

$$m_{\text{out}} = \max(|\text{OE}_i|, |\text{OE}_j|)$$

$$n_{\text{out}} = \min(|\text{OE}_i|, |\text{OE}_j|)$$

$$a_{ij}^{k+1} = \sqrt{y_{ij} \cdot \frac{s_{\text{in}}^{k+1}(i, j) + s_{\text{out}}^{k+1}(i, j)}{2}}$$

# Based on Simulation

- Labeled Transition Systems (LTS)를 사용하여 CFG를 만듦
- Recursively match the most similar outgoing nodes starting from the entry nodes
- Sum up the similarity of the matched nodes and edges

# Based on Simulation

- entry  $n_1, n_2$ 를 갖는 두 CFG의 유사도

$$S(n_1, n_2) = \begin{cases} N(n_1, n_2) & , \text{ if } OD_{n_1} = 0 \\ (1 - p) \cdot N(n_1, n_2) + p \cdot \max(W_1, W_2) & , \text{ otherwise,} \end{cases}$$

where

$$p \in (0, 1)$$

$$W_1 = \max_{n_2 \xrightarrow{b} n'_2} L(b, \epsilon) \cdot S(n_1, n'_2)$$

$$W_2 = \frac{1}{OD_{n_1}} \cdot \sum_{n_1 \xrightarrow{a} n'_1} \max(\max_{n_2 \xrightarrow{b} n'_2} (L(a, b) \cdot S(n'_1, n'_2)), L(a, \epsilon) \cdot S(n'_1, n_2))$$