



UNIVERSITÄT
BAYREUTH

University of Bayreuth
Institute for Computer Science

Bachelor Thesis

in Applied Computer Science

Topic: Randomly Generated CFGs, The CYK-Algorithm
And A GUI

Author: Andreas Braun <andreasbraun5@aol.com>
Matrikel-Nr. 1200197

Version date: February 23, 2017

1. Supervisor: Prof. Dr. Wim Martens
2. Supervisor: Tina Trautner

To my parents.

Abstract

The abstract of this thesis will be found here.

Zusammenfassung

Hier steht die Zusammenfassung dieser Bachelorarbeit.

Contents

1	Introduction	9
2	Technology Overview	10
3	Algorithms	11
4	Course of Action	24
	Used software	26
	Listings	28
	Tables	30
	References	32

1 Introduction

2 Technology Overview

3 Algorithms

Analogue to the script:

grammar $G = (V, \Sigma, S, P)$.

V is a finite set of variables.

Σ is an alphabet.

S is the starting symbol and $S \in V$.

P is a finite set of rules: $P \subseteq V \times (V \cup \Sigma)^*$.

G is in CNF and therefore it holds, more specifically: $P \subseteq V \times (V^2 \cup \Sigma)$. Think about if it is ok.

V s are Variables like "A, B, ...".

$(V \cup \Sigma)^*$ are terminals like "a, b, ..." and compound variables like "AB, BS, AC, ...".

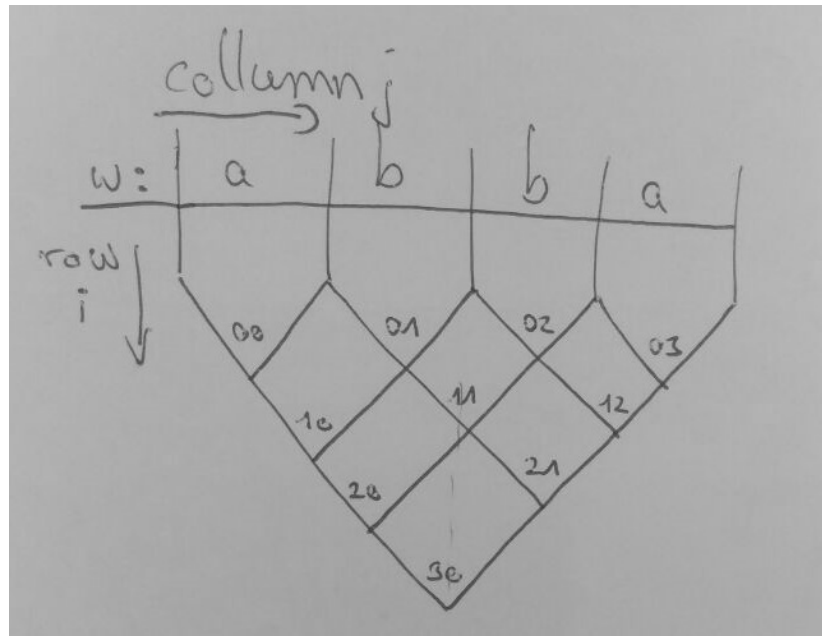
$i, j \in \mathbb{N}$; $[i, j] := \{i, i+1, \dots, j-1, j\} \subseteq \mathbb{N}_{\geq 0}$. $word \in \Sigma^*$

$Pyramid := \{cell_{i,j} \mid i \in \mathbb{N}_{\geq 0}, j \in [0, j_{max} - i], i_{max} = j_{max} = |word| - 1\}$.

$cell_{i,j} = \{Cell_{i,j} \mid Cell_{i,j} \subseteq V\}$. If each $cell_{i,j} = \emptyset$ then it is called *EmptyPyramid*.

Note that regarding one $cell_{i,j}$: $cell_{i,j} = cellDown$, $cell_{i-1,j} = cellUpperLeft$ and $cell_{i-1,j+1} = cellUpperRight$

Multisets allow duplicate elements. Notation of a multiset is done via the index b, that stands for bag: $multiset_b$



Algorithm 1: distributeRhse**Input:** $G, Rhse \subseteq (V^2 \cup \Sigma), i, j$ **Output:** Grammar G with uniform randomly distributed $Rhse$'s.

```

1 foreach  $rhse \in Rhse$  do
2   | choose  $n$  uniform randomly in  $[i, j]$ ;
3   | choose  $V_{add} :=$  uniform random subset of size  $n$  from  $V$ ;
4   |  $G.P = G.P \cup \{ "v \rightarrow rhse" \mid v \in V_{add} \}$ ;
5 end
6 return  $G$ ;

```

Algorithm 2: calculateSubsetForCell**Input:** $Pyramid, cell_{i,j}$ **Output:** $V_{i,j} \subseteq V^2$

```

1  $V_{i,j} = \emptyset$ ;
2 for  $m := i - 1 \rightarrow 0$  do
3   |  $V_{i,j} = V_{i,j} \cup \{ X \mid X \rightarrow YZ, Y \in V_{m,j}, Z \in V_{i-m-1,m+j+1} \}$ ;
4 end
5 return  $V_{i,j}$ ;

```

Algorithm 3: checkForceCombinationForCell**Input:** $cell_{i,j} \subseteq V^2, cell_{i-1,j} \subseteq V^2, cell_{i-1,j+1} \subseteq V^2, G.P$ **Output:** $varsForcing \subseteq V$

```

1  $varsForcing = \emptyset$ ;
2  $varComp = \{ XY \mid X \in cell_{i-1,j} \wedge Y \in cell_{i-1,j+1} \}$ ;
3 foreach  $v \in cellDown$  do
4   |  $prods = \{ p \mid v \in G.P.V \}$ ;
5   |  $rhse = \{ rhse \mid rhse \in prods.(V^2 \cup \Sigma) \}$ ;
6   | if  $varComp \not\subseteq rhse$  then
7     |  $varsForcing = varsForcing \cup v$ ;
8   | end
9 end
10 return  $varsForcing$ ;

```

Input: $cell_{i,j} = cellDown, cell_{i-1,j} = cellUpperLeft$ and $cell_{i-1,j+1} = cellUpperRight$

Algorithm 4: GeneratorGrammarDiceRollMartens

Input: $word \in \Sigma^*$, V , Σ , S , $P = \emptyset$, $minCount\Sigma$, $maxCount\Sigma$,
 $minCountVarComp$, $maxCountVarComp$

Output: G

```

1  $G = (V, \Sigma, S, P)$ ;
2  $G = distributeRhse(G, \Sigma, minCount\Sigma, maxCount\Sigma)$ ;
3  $Pyramid = CYK.calculatePyramid(G, word)$ ;
4 for  $i := 1 \rightarrow i_{max}$  do
5   for  $j := 0 \rightarrow j_{max} - i$  do
6      $sub = calculateSubsetForCell(Pyramid, i, j)$ ;
7     foreach  $vc \in sub$  do
8       while  $cell_{i_{max},0} = \emptyset$  do
9          $distributeRhse(G, vc, minCountVarComp,$ 
10            $maxCountVarComp)$ ;
11          $Pyramid = CYK.calculatePyramid(G, word)$ ;
12       end
13     end
14   end
15 return  $G$ ;

```

Line 3: Fills the $i=0$ row of the pyramid.

line 5: Instead of going from left to right, choose j uniform randomly with the restrictions that one cell is only visited one time.

Note: The algorithm tends to finish already within $i = 1$ loop.

Algorithm 5: GeneratorGrammarDiceRollMartens2

Input: $word \in \Sigma^*$, V , Σ , S , $P = \emptyset$, $minCount\Sigma$, $maxCount\Sigma$,
 $minCountVarComp$, $maxCountVarComp$

Output: G

```

1  $G = (V, \Sigma, S, P)$ ;
2  $G = distributeRhse(G, \Sigma, minCount\Sigma, maxCount\Sigma)$ ;
3  $Pyramid = CYK.calculatePyramid(G, word)$ ;
4  $sub = \emptyset$ ;
5 for  $i := 1 \rightarrow i_{max}$  do
6   for  $j := 0 \rightarrow j_{max} - i$  do
7      $sub = sub \cup \{(A, i) \mid A \in calculateSubsetForCell(Pyramid, i, j)\}$ ;
8   end
9    $sub_b = \{B \mid B \in sub, sub_b \text{ models } i\text{-dependent priority mechanism}\}$ ;
10  while  $cell_{i_{max},0} = \emptyset \wedge threshold_i = false$  do
11    choose one vc uniform randomly  $\in sub_b$ ;
12     $distributeRhse(G, vc, minCountVarComp, maxCountVarComp)$ ;
13     $Pyramid = CYK.calculatePyramid(G, word)$ ;
14    evaluate and update  $threshold_i$ ;
15  end
16 end
17 return  $G$ ;

```

Line 3: Fills the $i=0$ row of the pyramid.

Line 7: $(AB, 1), (AB, 2), (BC, 3) \dots \in sub \rightarrow$ multiple occurrences of AB are allowed. This considers "more important" compound variables.

Line 11: One vc can be chosen several times.

Note: Threshold: Linear or log function $f(i)$?

Note: Priority mechanism: In line $i + 1$ the $k = \{(A, l) \mid (A, l) \in sub, l = i\}$ are preferred over the $m = \{(A, n) \mid (A, n) \in sub, n < i\}$. In what way are they preferred? Using some kind of factor to weight the i of (A, i) .

```
1 Algorithm: CYK.calculateSetVAdvanced
2 Input: grammar, word;
3 Output: Set<VariableK> [][] cYKMatrix;
4
5 Set<VariableK> [][] cYKMatrix = new Set<VariableK>[wordSize][wordSize];
6 cYKMatrix = calculateCYKMatrix;
7 return cYKMatrix;
```

Listing 1: CYK.calculateSetVAdvanced

```

1 Algorithm: GeneratorGrammarDiceRollOnly
2 Input: settings;
3 Output: grammar;
4 Note: A lot of productions are generated, that later on are not needed
5 for parsing the specific word.
6
7 Grammar grammar = new Grammar();
8 // Part1: Distribute the terminals.
9 grammar = distributeDiceRollRightHandSideElements(
10     grammar, settingsTerminals, minCountTerminals,
11     maxCountTerminals, settingsListVariables);
12 // Part2: Distribute the compound variables.
13 Set<Variables> vars = settings.getVariables();
14 Set<VariablesCompound> setVarComp;
15 setVarComp = calculate all the possible tupels of ({vars}, {vars});
16 grammar = distributeDiceRollRightHandSideElements(
17     grammar, settingsTerminals, minCountVariableCompound,
18     maxCountVariableCompound, setVarComp);
19 return grammar

```

Listing 2: GeneratorGrammarDiceRollOnly


```

1 Algorithm: GeneratorGrammarDiceRollOnlyBias
2 Input: settings;
3 Output: grammar;
4 Note: A lot of productions are generated, that later on are not needed
5 for parsing the specific word.
6
7 Grammar grammar = new Grammar();
8 // Distribute the terminals.
9 grammar = distributeDiceRollRightHandSideElementsBias(
10     grammar, settingsTerminals, settingsMinCountTerminals,
11     settingsMCountTerminals, settingsListVars, settingsFavouritism);
12
13 // Distribute the compound variables.
14 Set<VariablesCompound> setVarComp;
15 setVarComp = calculate all the possible tupels of ({vars}, {vars});
16 grammar = distributeDiceRollRightHandSideElementsBias(
17     grammar, varComp, settingsMinCountVars,
18     settingsMaxCountVars, settingsListVars, settingsFavouritism);
19 return grammar;

```

Listing 3: GeneratorGrammarDiceRollOnlyBias

```

1 Algorithm: distributeDiceRollRightHandSideElementsBias
2 Input: grammar, setRhse, minCount, maxCount, listVars, favouritismList;
3 Output: grammar;
4 Note: Because of dice rolling anyways and lots of grammars being
5 generated, no rhse is added if the production already exists.
6
7 // Calculate the bloated varSet.
8 List<Variable> varsBloated;
9 for(Variables varTemp : settings.getVariables()){
10     tempFavour = randomly pick favouritism[i];
11     varsBloated.add({tempFavour times varTemp});
12     favouritism.remove(tempFavour);
13 }
14 // Because of dice rolling anyways and lot of grammars being generated,
15 // just no rhse is added if the production already exists.
16 grammar = distributeDiceRollRightHandSideElements( grammar,
17     varsBloated, minCount, maxCount, listVars );
18 return grammar;

```

Listing 4: distributeDiceRollRightHandSideElementsBias

Algorithm 6: distributeDiceRollRightHandSideElementsBias

Input: G , $rhse \subseteq RHSE$, $0 \leq minCount \leq maxCount \leq |G.V|$
 $, favouritism = \{x \mid x \in \mathbb{N} \wedge |favouritism| = |G.V|\}$

Output: G

- 1 $favour = \{(v, f) \mid v \in V \wedge f \in favouritism \wedge \text{tupel are created via dice roll}\};$
- 2 $varsBloated_b = \emptyset;$
- 3 **foreach** fav in $favour$ **do**
- 4 $varsBloated_b = varsBloated_b \uplus \{v^f\};$
- 5 **end**
- 6 **return**
 $distributeDiceRollRhse(G, MISTAKEHEREvarsBloated_b, minCount, maxCount);$
- 7 Still working on. One more Pparameter needed for
distributeDiceRollRighthandSideElement. Parameter V that defines the
variables the rhse are added to. OR make this algorithm independent.

Line 6: Note that $varsBloated_b$ is a multiset, but should actually be a set. Exceptions causing a duplicate production to the grammar are not relevant because G.P is a set.

Description of the checks here.

All test of the GrammarValidityChecker class are based on the simple setV matrix.

isValid = isWordProducible && isExamConstraints && isGrammarRestrictions

isWordProducible = CYK.algorithmAdvanced()

isExamConstraints = isRightCellCombinationsForced && isMaxSumOfProductionsCount
&& isMaxSumOfVarsInPyramidCount && countRightCellCombinationsForced

isGrammarRestrictions = isSizeOfWordCount && isMaxNumberOfVarsPerCellCount

```

1 Algorithm: checksumOfProductions
2 Input: grammar, maxSumOfProduction;
3 Output: isSumOfProductions;
4
5 return grammar.getProductionsAsList().size() <= maxSumOfProductions;

```

Listing 5: checksumOfProductions

```
1 Algorithm: checkMaxNumberOfVarsPerCell
2 Input: setVSimple, maxNumberOfVarsPerCell;
3 Output: isMaxNumberOfVarsPerCell;
4 Note: Checking for maxNumberOfVarsPerCell <= zero isn't allowed;
5
6 int tempMaxNumberOfVarsPerCell = 0;
7 int wordLength = tempSetV[0].length;
8 for ( int i = 0; i < wordLength; i++ ) {
9     for ( int j = 0; j < wordLength; j++ ) {
10         if ( tempSetV[i][j].size() > numberOfVarsPerCell ) {
11             numberOfVarsPerCell = tempSetV[i][j].size();
12         }
13     }
14 }
15 return tempMaxNumberOfVarsPerCell <= maxNumberOfVarsPerCell;
```

Listing 6: checkMaxNumberOfVarsPerCell

```
1 Algorithm: checkMaxSumOfVarsInPyramid
2 Input: setVSimple, maxSumOfVarsInPyramid;
3 Output: isMaxSumOfVarsInPyramid;
4
5 // put all vars of the matrix into one list and use its length.
6 List<Variable> allVarsList = new ArrayList<>();
7 for ( int i = 0; i < setVSimple.length; i++ ) {
8     for ( int j = 0; j < setVSimple.length; j++ ) {
9         tempVars.addAll( setVSimple[i][j] );
10    }
11 }
12 return allVarsList.size() <= maxSumOfVarsInPyramid;
```

Listing 7: checkMaxSumOfVarsInPyramid

```

1 Algorithm: rightCellCombinationsForced
2 Input: setVSimple, minCountForced, grammar;
3 Output: isForced, countForced, setVSimpleVarsThatForce;
4 Note: Keep in mind that the setV matrix is a upper right matrix. But the
5 description of how the algorithm works is done, as if the setV pyramid
6 points downwards (reflection on the diagonal + rotation to the left).
7 Regarding one cell, its upper left cell and its upper right cell
8 are looked at. setV[i][j] = down cell. setV[i + 1][j] = upper right cell
9 setV[i][j - 1] = upper left cell.
10
11 int countForced = 0;
12 Set<Variable> [][] setVMarkedVarsThatForce;
13 for(Cell cell : setVSimple){
14     // Trivial cases that would fulfil the restrictions each time.
15     Ignore the upper two rows of the pyramid;
16     isRightCellCombinationForced = true;
17     if(!upperLeftCell.isEmpty() && !upperRightCell.isEmpty()) {
18         break;
19     }
20     setVariableCompound = calculate all the possible tupels of
21         ({varLeft}, {varRight});
22     for(Variable var : cellToBeVisited) {
23         varDownProdList = grammar.getProdList(varDown);
24         for(VariableCompound varComp : setVariableCompound) {
25             for(Production prod : varDownProdList){
26                 if(prod.getRhse() == varComp) {
27                     isForced = false;
28                 }
29             }
30         }
31         if(isRightCellCombinationForced) {
32             rightCellCombinationsForced++;
33             // Cell has index i and j.
34             setVMarkedVarsThatForce[i][j].add(var)
35         }
36     }
37 }
38 boolean isForced = countForced >= minCountRightCellCombinationsForced;
39 return isForced, countForced, setVMarkedVarsThatForce;

```

Listing 8: rightCellCombinationsForced

```
1 Algorithm: Util.removeUselessProductions
2 Input: grammar, setVSimple, word
3 Output: grammar
4 Note: Very similar to the calculateSetVAdvanced algorithm. Additionally
5 to storing the k, it is also saved, which production have been used.
6 All productions that haven't been need are removed, from the grammar.
7
8 Set<Production> allProductions = grammar.getProductions();
9 Set<Production> onlyUsefulProductions;
10 onlyUsefulProductions = calculate useful productions with the input of
11     grammar, setVSimple and word ;
12 grammar.remove(allProductions);
13 return grammar.add(onlyUsefulProductions);
```

Listing 9: Util.removeUselessProductions

4 Course of Action

The informal goal is to find a suitable combination of a grammar and a word that meets the demands of an exam exercise. Also the CYK pyramid and one derivation tree of the word must be generated automatically as a "solution picture".

Firstly the exam exercise must have an upper limit of variables per cell while computing the CYK-pyramid.

Secondly the exam exercise must have one or more "special properties" so that it can be checked if the students have clearly understood the algorithm, e.g. "Excluding the possibility of luck."

The more formal goal is identify and determine parameters that in general can be used to define the properties of a grammar, so that the demanded restrictions are met. Also parameters could be identified for words, but "which is less likely to contribute, than the parameters of the grammar." [Second appointment with Martens]

Some introductory stuff:

Possible basic approaches for getting these parameters are the Rejection Sampling method and the "Tina+Wim" method.

Also backtracking plays some role, but right now I don't know where to put it. Backtracking is underapproach to Rejection Sampling.

Note: Starting with one half of a word and one half of a grammar.

Identify restrictions (=parameters) regarding the grammar.

Maybe find restrictions regarding the words, too.

Procedures for automated generation. Each generation procedure considers different restrictions and restriction combinations. The restrictions within one generation procedure can be optimized on its own. Up till now:

Generating grammars: DiceRolling, ...

Generating words: DiceRolling, ...

Parameter optimisation via theoretical and/or benchmarking approach.

Benchmarking = generate N grammars and test them, (N=100000).

Define a success rate and try to increase it.

The overall strategy is as following:

- 1.) Identify theoretically a restriction/parameter for the grammar. Think about the influence it will have. Think also about correlations between the restrictions.
- 2.) Validate the theoretical conclusion with the benchmark. Test out the influence of this parameter upon the success rate. Try different parameter settings.

The ordering of step 1 and step 2 can be changed.

Used software

Listings

Following are some interesting classes referenced in the thesis that were too long to fit into the text.

Tables

This section contains all tables referenced in this thesis.

References

- [1] JSR 220: Enterprise Java Beans 3.0 <https://jcp.org/en/jsr/detail?id=220>, 09/09/2015

