

A formalisation of the Cocke-Younger-Kasami algorithm

Maksym Bortin

June 17, 2024

Abstract

The theory provides a formalisation of the Cocke-Younger-Kasami algorithm [1] (CYK for short), an approach to solving the word problem for context-free languages. CYK decides if a word is in the languages generated by a context-free grammar in Chomsky normal form. The formalized algorithm is executable.

Contents

1	Basic modelling	2
1.1	Grammars in Chomsky normal form	2
1.2	Derivation by grammars	2
1.3	The generated language semantics	3
2	Basic properties	3
2.1	Properties of generated languages	5
3	Abstract specification of CYK	5
3.1	Properties of <i>subword</i>	6
3.2	Properties of <i>CYK</i>	7
4	Implementation	7
4.1	Main cycle	7
4.2	Initialisation phase	9
4.3	The overall procedure	10

```
theory CYK
imports Main
begin
```

The theory is structured as follows. First section deals with modelling of grammars, derivations, and the language semantics. Then the basic properties are proved. Further, CYK is abstractly specified and its underlying recursive relationship proved. The final section contains a prototypical implementation accompanied by a proof of its correctness.

1 Basic modelling

1.1 Grammars in Chomsky normal form

A grammar in Chomsky normal form is here simply modelled by a list of production rules (the type CNG below), each having a non-terminal symbol on the lhs and either two non-terminals or one terminal symbol on the rhs.

datatype $('n, 't) \text{ RHS} = \text{Branch } 'n \ 'n$
 $\quad \quad \quad | \text{Leaf } 't$

type-synonym $('n, 't) \text{ CNG} = ('n \times ('n, 't) \text{ RHS}) \text{ list}$

Abbreviating the list append symbol for better readability

abbreviation $\text{list-append} :: 'a \text{ list} \Rightarrow 'a \text{ list} \Rightarrow 'a \text{ list}$ (**infixr** · 65)
where $xs \cdot ys \equiv xs @ ys$

1.2 Derivation by grammars

A *word form* (or sentential form) may be built of both non-terminal and terminal symbols, as opposed to a *word* that contains only terminals. By the usage of disjoint union, non-terminals are injected into a word form by *Inl* whereas terminals – by *Inr*.

type-synonym $('n, 't) \text{ word-form} = ('n + 't) \text{ list}$
type-synonym $'t \text{ word} = 't \text{ list}$

A single step derivation relation on word forms is induced by a grammar in the standard way, replacing a non-terminal within a word form in accordance to the production rules.

definition $\text{DSTEP} :: ('n, 't) \text{ CNG} \Rightarrow (('n, 't) \text{ word-form} \times ('n, 't) \text{ word-form})$
 set
where $\text{DSTEP } G = \{ (l \cdot [\text{Inl } N] \cdot r, x) \mid l \ N \ r \ \text{rhs } x. (N, \text{rhs}) \in \text{set } G \wedge$
 $\quad \quad \quad (\text{case rhs of}$
 $\quad \quad \quad \text{Branch } A \ B \Rightarrow x = l \cdot [\text{Inl } A, \text{Inl } B] \cdot r$
 $\quad \quad \quad | \text{Leaf } t \Rightarrow x = l \cdot [\text{Inr } t] \cdot r) \}$

abbreviation $\text{DSTEP}' :: ('n, 't) \text{ word-form} \Rightarrow ('n, 't) \text{ CNG} \Rightarrow ('n, 't) \text{ word-form}$
 $\Rightarrow \text{bool } (- \dashrightarrow - [60, 61, 60] 61)$
where $w \dashrightarrow w' \equiv (w, w') \in \text{DSTEP } G$

abbreviation $DSTEP\text{-reflc} :: ('n, 't) \text{ word-form} \Rightarrow ('n, 't) \text{ CNG} \Rightarrow ('n, 't) \text{ word-form} \Rightarrow \text{bool}$ $(- \dashrightarrow^= - [60, 61, 60] 61)$
where $w - G \dashrightarrow^= w' \equiv (w, w') \in (DSTEP\ G)^=$

abbreviation $DSTEP\text{-transc} :: ('n, 't) \text{ word-form} \Rightarrow ('n, 't) \text{ CNG} \Rightarrow ('n, 't) \text{ word-form} \Rightarrow \text{bool}$ $(- \dashrightarrow^+ - [60, 61, 60] 61)$
where $w - G \dashrightarrow^+ w' \equiv (w, w') \in (DSTEP\ G)^+$

abbreviation $DSTEP\text{-rtransc} :: ('n, 't) \text{ word-form} \Rightarrow ('n, 't) \text{ CNG} \Rightarrow ('n, 't) \text{ word-form} \Rightarrow \text{bool}$ $(- \dashrightarrow^* - [60, 61, 60] 61)$
where $w - G \dashrightarrow^* w' \equiv (w, w') \in (DSTEP\ G)^*$

1.3 The generated language semantics

The language generated by a grammar from a non-terminal symbol comprises all words that can be derived from the non-terminal in one or more steps. Notice that by the presented grammar modelling, languages containing the empty word cannot be generated. Hence in rare situations when such languages are required, the empty word case should be treated separately.

definition $Lang :: ('n, 't) \text{ CNG} \Rightarrow 'n \Rightarrow 't \text{ word set}$
where $Lang\ G\ S = \{w. [Inl\ S] - G \dashrightarrow^+ \text{ map } Inr\ w\}$

So, for instance, a grammar generating the language $a^n b^n$ from the non-terminal $"S"$ might look as follows.

definition $G\text{-anbn} =$
 $[(\text{"S"}, \text{Branch } \text{"A"}\ \text{"T"}),$
 $(\text{"S"}, \text{Branch } \text{"A"}\ \text{"B"}),$
 $(\text{"T"}, \text{Branch } \text{"S"}\ \text{"B"}),$
 $(\text{"A"}, \text{Leaf } \text{"a"}),$
 $(\text{"B"}, \text{Leaf } \text{"b"})]$

Now the term $Lang\ G\text{-anbn}\ \text{"S"}$ denotes the set of words of the form $a^n b^n$ with $n > 0$. This is intuitively clear, but not straight forward to show, and a lengthy proof for that is out of scope.

2 Basic properties

lemma $prod\text{-into-}DSTEP1 :$
 $(S, \text{Branch } A\ B) \in \text{set } G \implies$
 $L \cdot [Inl\ S] \cdot R - G \dashrightarrow L \cdot [Inl\ A, Inl\ B] \cdot R$
 $\langle \text{proof} \rangle$

lemma $prod\text{-into-}DSTEP2 :$
 $(S, \text{Leaf } a) \in \text{set } G \implies$
 $L \cdot [Inl\ S] \cdot R - G \dashrightarrow L \cdot [Inr\ a] \cdot R$

$\langle proof \rangle$

lemma *DSTEP-D* :

$s - G \rightarrow t \implies$
 $\exists L N R rhs. s = L \cdot [Inl N] \cdot R \wedge (N, rhs) \in set G \wedge$
 $(\forall A B. rhs = Branch A B \longrightarrow t = L \cdot [Inl A, Inl B] \cdot R) \wedge$
 $(\forall x. rhs = Leaf x \longrightarrow t = L \cdot [Inr x] \cdot R)$
 $\langle proof \rangle$

lemma *DSTEP-append* :

assumes $a: s - G \rightarrow t$
shows $L \cdot s \cdot R - G \rightarrow L \cdot t \cdot R$
 $\langle proof \rangle$

lemma *DSTEP-star-mono* :

$s - G \rightarrow^* t \implies length s \leq length t$
 $\langle proof \rangle$

lemma *DSTEP-comp* :

assumes $a: l \cdot r - G \rightarrow t$
shows $\exists l' r'. l - G \rightarrow^* l' \wedge r - G \rightarrow^* r' \wedge t = l' \cdot r'$
 $\langle proof \rangle$

theorem *DSTEP-star-comp1* :

assumes $A: l \cdot r - G \rightarrow^* t$
shows $\exists l' r'. l - G \rightarrow^* l' \wedge r - G \rightarrow^* r' \wedge t = l' \cdot r'$
 $\langle proof \rangle$

theorem *DSTEP-star-comp2* :

assumes $A: l - G \rightarrow^* l'$
 and $B: r - G \rightarrow^* r'$
shows $l \cdot r - G \rightarrow^* l' \cdot r'$
 $\langle proof \rangle$

lemma *DSTEP-trancl-term* :

assumes $A: [Inl\ S] -G \rightarrow^+ t$
and $B: Inr\ x \in set\ t$
shows $\exists N. (N, Leaf\ x) \in set\ G$
 $\langle proof \rangle$

2.1 Properties of generated languages

lemma *Lang-no-Nil* :
 $w \in Lang\ G\ S \implies w \neq []$
 $\langle proof \rangle$

lemma *Lang-rtrancl-eq* :
 $(w \in Lang\ G\ S) = [Inl\ S] -G \rightarrow^* map\ Inr\ w$ **(is** $?L = (?p \in ?R^*)$ **)**
 $\langle proof \rangle$

lemma *Lang-term* :
 $w \in Lang\ G\ S \implies$
 $\forall x \in set\ w. \exists N. (N, Leaf\ x) \in set\ G$
 $\langle proof \rangle$

lemma *Lang-eq1* :
 $([x] \in Lang\ G\ S) = ((S, Leaf\ x) \in set\ G)$
 $\langle proof \rangle$

theorem *Lang-eq2* :
 $(w \in Lang\ G\ S \wedge 1 < length\ w) =$
 $(\exists A\ B. (S, Branch\ A\ B) \in set\ G \wedge (\exists l\ r. w = l \cdot r \wedge l \in Lang\ G\ A \wedge r \in Lang\ G\ B))$
(is $?L = ?R$ **)**
 $\langle proof \rangle$

3 Abstract specification of CYK

A subword of a word w , starting at the position i (first element is at the position 0) and having the length j , is defined as follows.

definition *subword* $w\ i\ j = take\ j\ (drop\ i\ w)$

Thus, to any subword of the given word w CYK assigns all non-terminals from which this subword is derivable by the grammar G .

definition *CYK* $G w i j = \{S. \text{subword } w i j \in \text{Lang } G S\}$

3.1 Properties of *subword*

lemma *subword-length* :

$i + j \leq \text{length } w \implies \text{length}(\text{subword } w i j) = j$
<proof>

lemma *subword-nth1* :

$i + j \leq \text{length } w \implies k < j \implies$
 $(\text{subword } w i j)!k = w!(i + k)$
<proof>

lemma *subword-nth2* :

assumes *A*: $i + 1 \leq \text{length } w$
shows $\text{subword } w i 1 = [w!i]$
<proof>

lemma *subword-self* :

$\text{subword } w 0 (\text{length } w) = w$
<proof>

lemma *take-split*[*rule-format*] :

$\forall n m. n \leq \text{length } xs \longrightarrow n \leq m \longrightarrow$
 $\text{take } n xs \cdot \text{take } (m - n) (\text{drop } n xs) = \text{take } m xs$
<proof>

lemma *subword-split* :

$i + j \leq \text{length } w \implies 0 < k \implies k < j \implies$
 $\text{subword } w i j = \text{subword } w i k \cdot \text{subword } w (i + k) (j - k)$
<proof>

lemma *subword-split2* :

assumes *A*: $\text{subword } w i j = l \cdot r$
 and *B*: $i + j \leq \text{length } w$
 and *C*: $0 < \text{length } l$
 and *D*: $0 < \text{length } r$

shows $l = \text{subword } w i (\text{length } l) \wedge r = \text{subword } w (i + \text{length } l) (j - \text{length } l)$
<proof>

3.2 Properties of CYK

lemma *CYK-Lang* :
 $(S \in \text{CYK } G \text{ } w \text{ } 0 \text{ } (\text{length } w)) = (w \in \text{Lang } G \text{ } S)$
(*proof*)

lemma *CYK-eq1* :
 $i + 1 \leq \text{length } w \implies$
 $\text{CYK } G \text{ } w \text{ } i \text{ } 1 = \{S. (S, \text{Leaf } (w!i)) \in \text{set } G\}$
(*proof*)

theorem *CYK-eq2* :
assumes $A: i + j \leq \text{length } w$
and $B: 1 < j$
shows $\text{CYK } G \text{ } w \text{ } i \text{ } j = \{X \mid X \text{ } A \text{ } B \text{ } k. (X, \text{Branch } A \text{ } B) \in \text{set } G \wedge A \in \text{CYK } G \text{ } w \text{ } i \text{ } k \wedge B \in \text{CYK } G \text{ } w \text{ } (i + k) \text{ } (j - k) \wedge 1 \leq k \wedge k < j\}$
(*proof*)

4 Implementation

One of the particularly interesting features of CYK implementation is that it follows the principles of dynamic programming, constructing a table of solutions for sub-problems in the bottom-up style reusing already stored results.

4.1 Main cycle

This is an auxiliary implementation of the membership test on lists.

fun *mem* :: 'a \Rightarrow 'a list \Rightarrow bool
where
 $\text{mem } a \ [] = \text{False} \mid$
 $\text{mem } a \ (x\#\text{xs}) = (x = a \vee \text{mem } a \ \text{xs})$

lemma *mem[simp]* :
 $\text{mem } x \ \text{xs} = (x \in \text{set } \text{xs})$
(*proof*)

The purpose of the following is to collect non-terminals that appear on the lhs of a production such that the first non-terminal on its rhs appears in the first of two given lists and the second non-terminal – in the second list.

fun *match-prods* :: ('n, 't) CNG \Rightarrow 'n list \Rightarrow 'n list \Rightarrow 'n list
where *match-prods* [] *ls rs* = [] |
 $\text{match-prods } ((X, \text{Branch } A \text{ } B)\#\text{ps}) \ \text{ls } \ \text{rs} =$
 $(\text{if } \text{mem } A \ \text{ls} \wedge \text{mem } B \ \text{rs} \text{ then } X \ \#\ \text{match-prods } \text{ps } \ \text{ls } \ \text{rs}$
 $\text{else } \text{match-prods } \ \text{ps } \ \text{ls } \ \text{rs}) \mid$

$match-prods ((X, Leaf a)\#ps) ls rs = match-prods ps ls rs$

lemma *match-prods* :

$(X \in set(match-prods G ls rs)) =$
 $(\exists A \in set ls. \exists B \in set rs. (X, Branch A B) \in set G)$
 <proof>

The following function is the inner cycle of the algorithm. The parameters i and j identify a subword starting at i with the length j , whereas k is used to iterate through its splits (which are of course subwords as well) all having the length greater 0 but less than j . The parameter T represents a table containing CYK solutions for those splits.

function *inner* :: ('n, 't) CNG \Rightarrow (nat \times nat \Rightarrow 'n list) \Rightarrow nat \Rightarrow nat \Rightarrow nat \Rightarrow 'n list

where *inner* G T i k j =
 (if $k < j$ then $match-prods G (T(i, k)) (T(i + k, j - k)) @ inner G T i (k + 1) j$
 else [])
 <proof>

termination

<proof>

declare *inner.simps*[simp del]

lemma *inner* :

$(X \in set(inner G T i k j)) =$
 $(\exists l. k \leq l \wedge l < j \wedge X \in set(match-prods G (T(i, l)) (T(i + l, j - l))))$
 (is ?L G T i k j = ?R G T i k j)
 <proof>

Now the main part of the algorithm just iterates through all subwords up to the given length len , calls *inner* on these, and stores the results in the table T . The length j is supposed to be greater than 1 – the subwords of length 1 will be handled in the initialisation phase below.

function *main* :: ('n, 't) CNG \Rightarrow (nat \times nat \Rightarrow 'n list) \Rightarrow nat \Rightarrow nat \Rightarrow nat \Rightarrow (nat \times nat \Rightarrow 'n list)

where *main* G T len i j = (let $T' = T((i, j) := inner G T i 1 j)$ in
 if $i + j < len$ then $main G T' len (i + 1) j$
 else if $j < len$ then $main G T' len 0 (j + 1)$
 else T')

<proof>

termination

<proof>

declare *main.simps*[simp del]

lemma *main* :
assumes $1 < j$
and $i + j \leq \text{length } w$
and $\bigwedge i' j'. j' < j \implies 1 \leq j' \implies i' + j' \leq \text{length } w \implies \text{set}(T(i', j')) = \text{CYK } G w i' j'$
and $\bigwedge i'. i' < i \implies i' + j \leq \text{length } w \implies \text{set}(T(i', j)) = \text{CYK } G w i' j$
and $1 \leq j'$
and $i' + j' \leq \text{length } w$
shows $\text{set}((\text{main } G T (\text{length } w) i j)(i', j')) = \text{CYK } G w i' j'$
 $\langle \text{proof} \rangle$

4.2 Initialisation phase

Similarly to *match-prods* above, here we collect non-terminals from which the given terminal symbol can be derived.

fun *init-match* :: $('n, 't) \text{CNG} \Rightarrow 't \Rightarrow 'n \text{ list}$
where *init-match* [] $t = []$ |
init-match $((X, \text{Branch } A B)\#ps) t = \text{init-match } ps t$ |
init-match $((X, \text{Leaf } a)\#ps) t = (\text{if } a = t \text{ then } X \# \text{init-match } ps t$
else init-match ps t)

lemma *init-match* :
 $(X \in \text{set}(\text{init-match } G a)) =$
 $((X, \text{Leaf } a) \in \text{set } G)$
 $\langle \text{proof} \rangle$

Representing the empty table.

definition *emptyT* = $(\lambda(i, j). [])$

The following function initialises the empty table for subwords of length 1, i.e. each symbol occurring in the given word.

fun *init'* :: $('n, 't) \text{CNG} \Rightarrow 't \text{ list} \Rightarrow \text{nat} \Rightarrow \text{nat} \times \text{nat} \Rightarrow 'n \text{ list}$
where *init'* G [] $k = \text{emptyT}$ |
init' $G (t\#ts) k = (\text{init}' G ts (k + 1))((k, 1) := \text{init-match } G t)$

lemma *init'* :
assumes $i + 1 \leq \text{length } w$
shows $\text{set}(\text{init}' G w 0 (i, 1)) = \text{CYK } G w i 1$
 $\langle \text{proof} \rangle$

The next version of initialization refines *init'* in that it takes additional account of the cases when the given word is empty or contains a terminal symbol that does not have any matching production (that is, *init-match* is an empty list). No initial table is then needed as such words can immediately be rejected.

```

fun init :: ('n, 't) CNG ⇒ 't list ⇒ nat ⇒ (nat × nat ⇒ 'n list) option
where init G [] k = None |
      init G [t] k = (case (init-match G t) of
        [] ⇒ None
        | xs ⇒ Some(emptyT((k, 1) := xs))) |
      init G (t#ts) k = (case (init-match G t) of
        [] ⇒ None
        | xs ⇒ (case (init G ts (k + 1)) of
          None ⇒ None
          | Some T ⇒ Some(T((k, 1) := xs))))

```

```

lemma init1:
  ⟨init' G w k = T⟩ if ⟨init G w k = Some T⟩
  ⟨proof⟩

```

```

lemma init2 :
  (init G w k = None) =
  (w = [] ∨ (∃ a ∈ set w. init-match G a = []))
  ⟨proof⟩

```

4.3 The overall procedure

```

definition cyk G S w = (case init G w 0 of
  None ⇒ False
  | Some T ⇒ let len = length w in
    if len = 1 then mem S (T(0, 1))
    else let T' = main G T len 0 2 in
      mem S (T'(0, len)))

```

```

theorem cyk :
  cyk G S w = (w ∈ Lang G S)
  ⟨proof⟩

```

```

value [code]
  let G = [(0::int, Branch 1 2), (0, Branch 2 3),
    (1, Branch 2 1), (1, Leaf "a"),
    (2, Branch 3 3), (2, Leaf "b"),
    (3, Branch 1 2), (3, Leaf "a")]
  in map (cyk G 0)
    [[ "b", "a", "a", "b", "a"],
     [ "b", "a", "b", "a"]

```

end

References

- [1] D. H. Younger. Recognition and parsing of context-free languages in time n^3 . *Information and Control*, 10(2):189 – 208, 1967.