

Polynomial algorithms for CFGs via semiring embeddings

Piotr Mikołajczyk

Theoretical Computer Science Department of Jagiellonian University

November 4, 2020

- **Parsing with Derivatives** (2011) – by M. Might, D. Darais, D. Spiewak
- **On the Complexity and Performance of Parsing with Derivatives** (2016) – by M. Adams, C. Hollenbeck, M. Might
- **A C++ implementation of Parsing With Derivatives** (2019) – ☺

Regular languages

- \emptyset and $\{\varepsilon\}$ are regular
- $\forall c \in \Sigma \{c\}$ is regular
- If A and B are regular, then $A \cup B$, $A \circ B$ and A^* are regular
- No other language is regular

Regular languages

- \emptyset and $\{\epsilon\}$ are regular
- $\forall c \in \Sigma \{c\}$ is regular
- If A and B are regular, then $A \cup B$, $A \circ B$ and A^* are regular
- No other language is regular

$$R = \{\epsilon\} \cup \{a\} \cdot (\{a\} \cup \{b\})^*$$

- They are just like recursive regular languages

Context-free languages

- They are just like recursive regular languages
- As $L^* \equiv \varepsilon \cup (L \circ L^*)$, we can give up using Kleene's star operator

Context-free languages

A context-free grammar is a tuple $\mathcal{G} = (\mathcal{N}, \Sigma, S, P)$. A language which it generates is denoted as $\mathcal{L}(\mathcal{G})$.

Context-free languages

A context-free grammar is a tuple $\mathcal{G} = (\mathcal{N}, \Sigma, S, P)$. A language which it generates is denoted as $\mathcal{L}(\mathcal{G})$.

Example:

- $\mathcal{N} = \{S\}$
- $S \rightarrow aSa|bSb|\varepsilon$
- $\mathcal{L}(\mathcal{G}) = \{\varepsilon, aa, bb, abba, aaaa, \dots\}$

Binarization

Each grammar \mathcal{G} can be transformed (in polynomial time) to an equivalent *binarized* grammar \mathcal{G}' ($\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$) – operators (concatenation and alternative) are considered as purely binary; operands must be from $\Sigma \cup \mathcal{N}$.

Binarization

Each grammar \mathcal{G} can be transformed (in polynomial time) to an equivalent *binarized* grammar \mathcal{G}' ($\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$) – operators (concatenation and alternative) are considered as purely binary; operands must be from $\Sigma \cup \mathcal{N}$.

- $C \rightarrow ccc$
- $B \rightarrow b|A$
- $A \rightarrow \varepsilon|Bb|C$

Binarization

Each grammar \mathcal{G} can be transformed (in polynomial time) to an equivalent *binarized* grammar \mathcal{G}' ($\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$) – operators (concatenation and alternative) are considered as purely binary; operands must be from $\Sigma \cup \mathcal{N}$.

- $C \rightarrow ccc$

- $C \rightarrow C'c$

- $C' \rightarrow cc$

- $B \rightarrow b|A$

- $B \rightarrow b|A$

- $A \rightarrow \varepsilon|Bb|C$

- $A \rightarrow \varepsilon|A'$

- $A' \rightarrow B'|C$

- $B' \rightarrow Bb$

Grammar graph

Given a binarized grammar \mathcal{G} , we can consider its **graph** $G(\mathcal{G}) = (\mathbb{V}, \mathbb{E})$, where: $\mathbb{V} = (\mathcal{N} \cup \Sigma \cup \{\emptyset, \varepsilon\})$ and $(u, v) \in \mathbb{E}$ if and only if there is a production in \mathcal{P} which has u on its left side and v on the right side.

The resulting graph is directed and has ordered edges, i.e. we distinguish the left child from the right one.

Grammar graph

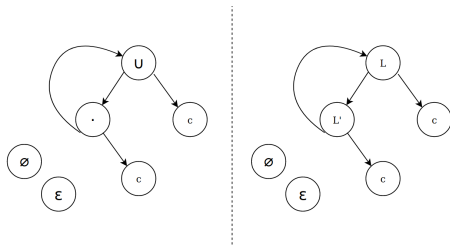
- $\mathcal{L} = (\mathcal{L} \cdot c) \cup c$
- $L = L'|c$
- $L' = Lc$

Grammar graph

- $\mathcal{L} = (\mathcal{L} \cdot c) \cup c$

- $L = L' | c$

- $L' = Lc$



Derivative of a language

Derivative by a symbol – b – from alphabet:

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$
- leave only these words, which start with the chosen symbol:
 $\{bar, bee\}$

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$
- leave only these words, which start with the chosen symbol:
 $\{bar, bee\}$
- remove first symbol: $\{ar, ee\}$

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$
- leave only these words, which start with the chosen symbol:
 $\{bar, bee\}$
- remove first symbol: $\{ar, ee\}$
- voilà

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$
- leave only these words, which start with the chosen symbol:
 $\{bar, bee\}$
- remove first symbol: $\{ar, ee\}$
- voilà ($D_b(\{foo, bar, bee\}) = \{ar, ee\}$)

Derivative of a language

Derivative by a symbol – b – from alphabet:

- take a language: $\{foo, bar, bee\}$
- leave only these words, which start with the chosen symbol: $\{bar, bee\}$
- remove first symbol: $\{ar, ee\}$
- voilà ($D_b(\{foo, bar, bee\}) = \{ar, ee\}$)

Derivative of a language

$$D_c(L) = \{w \mid cw \in L\}$$

Derivative of a language

Derivative of a language

$$D_c(L) = \{w \mid cw \in L\}$$

$$\text{So } cw \in L \iff w \in D_c(L)$$

$$w = a_1 \dots a_k \implies D_w(L) = D_{a_k}(\dots D_{a_1}(L) \dots)$$

$$\text{So } w \in L \iff \varepsilon \in D_w(L)$$

- $D_c(\emptyset) = \emptyset$
- $D_c(\varepsilon) = \emptyset$
- $D_c(a) = \begin{cases} \emptyset & a \neq c \\ \varepsilon & a = c \end{cases}$
- $D_c(L_1 \cup L_2) = D_c(L_1) \cup D_c(L_2)$
- $(D_c(L^*) = D_c(L) \circ L^*)$
- $D_c(L_1 \circ L_2) = \begin{cases} D_c(L_1) \circ L_2 & \varepsilon \notin L_1 \\ (D_c(L_1) \circ L_2) \cup D_c(L_2) & \varepsilon \in L_1 \end{cases}$

Nullability function

$$\delta : \mathcal{P}(\Sigma^*) \rightarrow \{\emptyset, \{\varepsilon\}\}$$

$$\delta(L) = \begin{cases} \emptyset & \varepsilon \notin L \\ \{\varepsilon\} & \varepsilon \in L \end{cases}$$

Nullability function

$$\delta : \mathcal{P}(\Sigma^*) \rightarrow \{\emptyset, \{\varepsilon\}\}$$

$$\delta(L) = \begin{cases} \emptyset & \varepsilon \notin L \\ \{\varepsilon\} & \varepsilon \in L \end{cases}$$

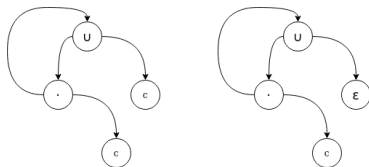
With this, we have: $D_c(L_1 \circ L_2) = (D_c(L_1) \circ L_2) \cup (\delta(L_1) \circ D_c(L_2))$

Nullability function

- $\delta(\emptyset) = \emptyset$
- $\delta(\varepsilon) = \varepsilon$
- $\delta(a) = \emptyset$
- $\delta(L_1 \cup L_2) = \delta(L_1) \cup \delta(L_2)$
- $\delta(L_1 \circ L_2) = \delta(L_1) \circ \delta(L_2)$
- $(\delta(L^*) = \varepsilon)$

Derivative on graphs

$$\mathcal{L} = (\mathcal{L} \cdot c) \cup c$$
$$D_c(\mathcal{L}) = (D_c(\mathcal{L}) \cdot c) \cup \varepsilon$$



Recognizing algorithm

```
def recognize(G, w):  
    for c ∈ w:  
        G = Dc(G)  
    return δ(G)
```

Some problems:

Some problems:

- forgetting about parsed symbols (tokens) – $D_c(\{c\}) = \varepsilon$

Some problems:

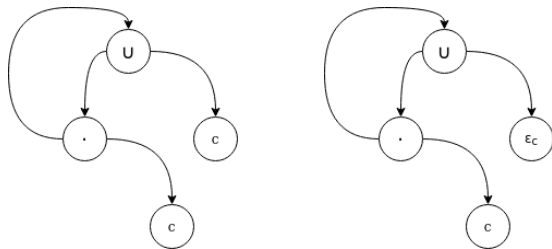
- forgetting about parsed symbols (tokens) – $D_c(\{c\}) = \varepsilon$
- loosing information by using δ –
$$D_c(L_1 \circ L_2) = (D_c(L_1) \circ L_2) \cup (\delta(L_1) \circ D_c(L_2))$$

Some problems:

- forgetting about parsed symbols (tokens) – $D_c(\{c\}) = \varepsilon$
- loosing information by using δ –
$$D_c(L_1 \circ L_2) = (D_c(L_1) \circ L_2) \cup (\delta(L_1) \circ D_c(L_2))$$
- skewing parse tree (by operators associativity)

$$\varepsilon_{\Sigma} = \{\varepsilon_a : a \in \Sigma\}$$

$$D_a(c) = \begin{cases} \varepsilon_a & a = c \\ \emptyset & a \neq c \end{cases}$$



$$D_c(L_1 \circ L_2) = (D_c(L_1) \circ L_2) \cup (\Delta(L_1) \circ D_c(L_2))$$

$$D_a(\Delta(P)) = \emptyset$$

$$\mathcal{M} = \{\vdash_i\}_{i \in \mathbb{N}}$$

- $C \rightarrow ccc \vdash_1$
- $B \rightarrow b \vdash_2 \mid A \vdash_3$
- $A \rightarrow \varepsilon \vdash_4 \mid Bb \vdash_5 \mid C \vdash_6$

$$\mathcal{M} = \{\vdash_i\}_{i \in \mathbb{N}}$$

- $C \rightarrow ccc \vdash_1$

- $B \rightarrow b \vdash_2 \mid A \vdash_3$

- $A \rightarrow \varepsilon \vdash_4 \mid Bb \vdash_5 \mid C \vdash_6$

- $C \rightarrow C' C''$

- $C'' \rightarrow c \vdash_1$

- $C' \rightarrow cc$

- $B \rightarrow B' \mid B''$

- $B'' \rightarrow A \vdash_3$

- $B' \rightarrow b \vdash_2$

- ...

A semiring \mathcal{R} is a triple $(R, +_{\mathcal{R}}, \cdot_{\mathcal{R}}, 0_{\mathcal{R}}, 1_{\mathcal{R}})$, where:

- R is a set of semiring's elements
- $(R, +_{\mathcal{R}})$ is a commutative monoid with $0_{\mathcal{R}}$ as an identity element
- $(R, \cdot_{\mathcal{R}})$ is a monoid with $1_{\mathcal{R}}$ as an identity element and 0 as an annihilator
- $\cdot_{\mathcal{R}}$ is distributive over $+_{\mathcal{R}}$ (on both sides)

A semiring \mathcal{R} is a triple $(R, +_{\mathcal{R}}, \cdot_{\mathcal{R}}, 0_{\mathcal{R}}, 1_{\mathcal{R}})$, where:

- R is a set of semiring's elements
- $(R, +_{\mathcal{R}})$ is a commutative monoid with $0_{\mathcal{R}}$ as an identity element
- $(R, \cdot_{\mathcal{R}})$ is a monoid with $1_{\mathcal{R}}$ as an identity element and 0 as an annihilator
- $\cdot_{\mathcal{R}}$ is distributive over $+_{\mathcal{R}}$ (on both sides)

We skip the \mathcal{R} subscript next to the operators and elements whenever possible.

For our purposes, we will also require the semirings to have an additional element $\infty_{\mathcal{R}}$ (or simply ∞) with the following properties:

$$\begin{aligned}\forall e \in R \quad e + \infty &= \infty, \\ 0 \cdot \infty &= \infty \cdot 0 = 0, \\ \forall e \in R - \{0\} \quad e \cdot \infty &= \infty \cdot e = \infty.\end{aligned}$$

Every nonterminal P can be associated with the language $\mathcal{L}(P) \subseteq \Sigma^*$. Thus we can perceive alternative and concatenation as respective operators in the semiring $\mathcal{R}_\Sigma = (\wp(\Sigma^*), \cup, \cdot, \emptyset, \{\varepsilon\})$.

Every nonterminal P can be associated with the language $\mathcal{L}(P) \subseteq \Sigma^*$. Thus we can perceive alternative and concatenation as respective operators in the semiring $\mathcal{R}_\Sigma = (\wp(\Sigma^*), \cup, \cdot, \emptyset, \{\varepsilon\})$.

Now we can generalize the function δ introduced previously. From now, $\delta_{\mathcal{R}}$ will represent any homomorphism between \mathcal{R}_Σ and an arbitrary semiring \mathcal{R} .

Generic algorithm

```
def recognize $\langle \mathcal{R} \rangle$ (G, w):  
    for c  $\in$  w:  
        G = Dc(G)  
    return  $\delta_{\mathcal{R}}$ (G)
```

Back to recognition

For recognition we can work with the Boolean semiring, i.e.

$\mathcal{R}_{\mathbb{B}} = (\mathbb{B}, \vee, \wedge, 0, 1)$ with $\mathbb{B} = \{0, 1\}$, $\infty = 1$.

Back to recognition

For recognition we can work with the Boolean semiring, i.e.

$\mathcal{R}_{\mathbb{B}} = (\mathbb{B}, \vee, \wedge, 0, 1)$ with $\mathbb{B} = \{0, 1\}$, $\infty = 1$.

$$\delta_{\mathbb{B}}(\emptyset) = 0,$$

$$\delta_{\mathbb{B}}(\epsilon) = 1,$$

$$\forall a \in \Sigma \quad \delta_{\mathbb{B}}(\epsilon_a) = 1,$$

$$\forall a \in \Sigma \quad \delta_{\mathbb{B}}(a) = 0.$$

Counting parse trees

For counting parse trees we can work with $\mathcal{R}_{\mathbb{N}} = (\mathbb{N} \cup \{\infty\}, +, \cdot, 0, 1)$, which is the standard semiring of non-negative integers, enriched with a special element ∞ behaving “naturally”, except that $\infty \cdot 0 = 0$.

Counting parse trees

For counting parse trees we can work with $\mathcal{R}_{\mathbb{N}} = (\mathbb{N} \cup \{\infty\}, +, \cdot, 0, 1)$, which is the standard semiring of non-negative integers, enriched with a special element ∞ behaving “naturally”, except that $\infty \cdot 0 = 0$.

$$\delta_{\mathbb{N}}(\emptyset) = 0$$

$$\delta_{\mathbb{N}}(\varepsilon) = 1$$

$$\forall a \in \Sigma \cup \mathcal{M} \quad \delta_{\mathbb{N}}(\varepsilon_a) = 1$$

$$\forall a \in \Sigma \quad \delta_{\mathbb{N}}(a) = 0$$

Parsing

For parsing we can work with $\mathcal{R}_{\mathbb{N}} = (\mathcal{Q} \times (\epsilon_{\Sigma} \cup \mathcal{M})^*, \oplus, \otimes)$, where the first coordinate, an element from $\mathcal{Q} = \{ \text{NONE}, \text{UNIQUE}, \text{FINITELY_MANY}, \text{INFINITELY_MANY} \}$ is one of the quantity categories.

For parsing we can work with $\mathcal{R}_{\mathbb{N}} = (\mathcal{Q} \times (\epsilon_{\Sigma} \cup \mathcal{M})^*, \oplus, \otimes)$, where the first coordinate, an element from $\mathcal{Q} = \{ \text{NONE}, \text{UNIQUE}, \text{FINITELY_MANY}, \text{INFINITELY_MANY} \}$ is one of the quantity categories.

- if the first coordinate is `UNIQUE`, the second one is the postorder of this unique parse tree; otherwise there could be anything - we do not care
- both operators \oplus, \otimes when given two elements from $\mathcal{R}_{\mathbb{N}}$ firstly look at the quantity coordinates and depending on them determine the resulting quantity. Then, if the resulting quantity is `UNIQUE`, they combine the second coordinates.
- 0-element is (NONE, ϵ) (in fact the second coordinate can contain anything)
- 1-element is $(\text{UNIQUE}, \epsilon)$
- ∞ -element is $(\text{INFINITELY_MANY}, \epsilon)$ (in fact the second coordinate can contain anything)

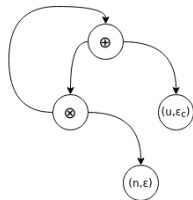
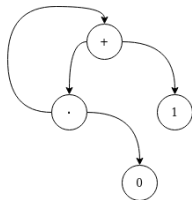
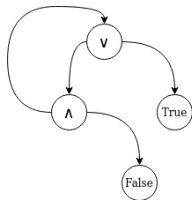
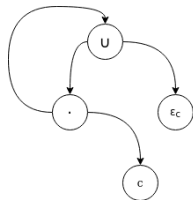
$$\delta_N(\emptyset) = (\text{NONE}, \varepsilon)$$

$$\delta_N(\varepsilon) = (\text{UNIQUE}, \varepsilon)$$

$$\forall a \in \Sigma \cup \mathcal{M} \quad \delta_N(\varepsilon_a) = (\text{UNIQUE}, \varepsilon_a)$$

$$\forall a \in \Sigma \quad \delta_N(a) = (\text{NONE}, \varepsilon)$$

All of them



$$\delta : \mathbb{V}(\mathcal{G}) \rightarrow R$$

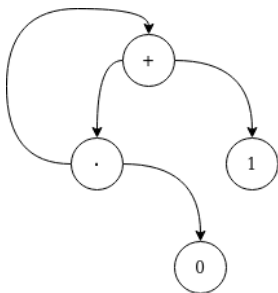
$$\delta(X) = \begin{cases} 0 & \text{if } X : \textit{Empty} \\ 0 & \text{if } X : \textit{Token} \\ 1 & \text{if } X : \textit{Epsilon} \\ e_X \in R & \text{if } X : \textit{ParsedToken} \\ \delta(X.ref) & \text{if } X : \textit{Delta} \\ \delta(X.left) +_{\mathcal{R}} \delta(X.right) & \text{if } X : \textit{Alternative} \\ \delta(X.left) \cdot_{\mathcal{R}} \delta(X.right) & \text{if } X : \textit{Concatenation} \end{cases}$$

Any ambiguities (that may appear due to cycles) can be dealt with by the following two rules:

$$\begin{aligned} \delta(X) = \alpha_1 \cdot \delta(X) \cdot \alpha_2 &\implies \delta(X) = 0 \\ \delta(X) = \alpha_1 \cdot \delta(X) \cdot \alpha_2 + \beta \wedge \alpha_i, \beta \neq 0 \wedge \delta(X) \notin \beta &\implies \delta(X) = \infty \end{aligned}$$

Algorithm

- 1 find all nodes X for which $\delta_{\mathcal{R}}(X) = 0_{\mathcal{R}}$,
- 2 propagate all “finite” values from \mathcal{R} as far as it is possible
- 3 what remains, should be equal to $\infty_{\mathcal{R}}$.



See in the full version.

- finding all zeros (induction on the number of connected components)

- finding all zeros (induction on the number of connected components)
- after marking all 0s and propagating finite values as far as possible, any value not yet calculated must be ∞ (by *recursive alternatives*)

- finding all zeros (induction on the number of connected components)
- after marking all 0s and propagating finite values as far as possible, any value not yet calculated must be ∞ (by *recursive alternatives*)
- sum disjointness (by *(semi-)parse-words*)