

Linear Conjunctive Reachability as Tensor Completion

Anonymous Author(s)

Abstract

Brzozowski (1964) defines a regular expression derivative as the suffixes which complete a known prefix. In this work, we establish a Galois connection with Valiant’s (1975) fixpoint construction in the context-free setting, and further extend their work into the hierarchy of bounded context-sensitive languages realizable by finite CFL intersection, i.e. conjunctive languages. We illustrate how to lower conjunctive language recognition onto a system of multilinear equations over finite fields. In addition to its theoretical value, this connection has yielded surprisingly useful applications in incremental parsing, code completion and program repair.

1 Introduction

Recall that a CFG is a quadruple consisting of terminals (Σ), nonterminals (V), productions ($P: V \rightarrow (V \mid \Sigma)^*$), and a start symbol, (S). It is a well-known fact that every CFG is reducible to *Chomsky Normal Form*, $P': V \rightarrow (V^2 \mid \Sigma)$, in which every production takes one of two forms, either $w \rightarrow xz$, or $w \rightarrow t$, where $w, x, z: V$ and $t: \Sigma$. For example, the CFG, $P := \{S \rightarrow SS \mid (S) \mid ()\}$, corresponds to the CNF:

$$P' = \{ S \rightarrow QR \mid SS \mid LR, \quad L \rightarrow (, \quad R \rightarrow), \quad Q \rightarrow LS \}$$

Given a CFG, $\mathcal{G}' : \langle \Sigma, V, P, S \rangle$ in CNF, we can construct a recognizer $R : \mathcal{G}' \rightarrow \Sigma^n \rightarrow \mathbb{B}$ for strings $\sigma : \Sigma^n$ as follows. Let 2^V be our domain, 0 be \emptyset , \oplus be \cup , and \otimes be defined as:

$$X \otimes Z := \{ w \mid \langle x, z \rangle \in X \times Z, (w \rightarrow xz) \in P \} \quad (1)$$

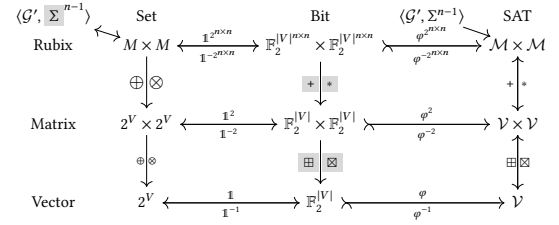
If we define $\sigma_r^\dagger := \{w \mid (w \rightarrow \sigma_r) \in P\}$, then initialize $M_{r+1=c}^0(\mathcal{G}', e) := \sigma_r^\dagger$ and solve for the fixpoint $M^* = M + M^2$,

$$M^0 := \begin{pmatrix} \emptyset & \sigma_1^{\rightarrow} & \emptyset & \cdots & \emptyset \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \emptyset & \cdots & \emptyset & \ddots & \sigma_n^{\uparrow} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \emptyset & \cdots & \emptyset & \cdots & \emptyset \end{pmatrix} \Rightarrow M^* = \begin{pmatrix} \emptyset & \sigma_1^{\rightarrow} & \Lambda & \cdots & \Lambda_{\sigma}^* \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \emptyset & \cdots & \emptyset & \ddots & \Lambda \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \emptyset & \cdots & \emptyset & \cdots & \emptyset \end{pmatrix}$$

we obtain the recognizer, $R(\mathcal{G}', \sigma) := S \in \Lambda_{\sigma}^*? \Leftrightarrow \sigma \in \mathcal{L}(\mathcal{G})?$ Full details of the bisimilarity between parsing and matrix multiplication can be found in Valiant [?] and Lee [?], who shows its time complexity to be $\mathcal{O}(n^{\omega})$ where ω is the least matrix multiplication upper bound (currently, $\omega < 2.77$).

2 Method

Note that $\bigoplus_{c=1}^n M_{r,c} \otimes M_{c,r}$ has cardinality bounded by $|V|$ and is thus representable as a fixed-length vector using the characteristic function, $\mathbb{1}$. In particular, \oplus, \otimes are redefined as \boxplus, \boxtimes over bitvectors so the following diagram commutes,¹



where \mathcal{V} is a function $\mathbb{F}_2^{|V|} \rightarrow \mathbb{F}_2$. Note that while always possible to encode $\mathbb{F}_2^{|V|} \rightarrow \mathcal{V}$ using the identity function, φ^{-1} may not exist, as an arbitrary \mathcal{V} might have zero, one, or in general, multiple solutions in $\mathbb{F}_2^{|V|}$. Although holes may occur anywhere, let us consider two cases in which Σ^+ is strictly left- or right-constrained, i.e., $\mathbf{x} \mathbf{z} \mathbf{x}, \mathbf{z} \mathbf{z} : \Sigma^{|x|+|z|}$.

Valiant's \otimes operator, which yields the set of productions unifying known factors in a binary CFG, naturally implies the existence of a left- and right-quotient, which yield the set of nonterminals that may appear the right or left side of a known factor and its corresponding root. In other words, a known factor not only implicates subsequent expressions that can be derived from it, but also adjacent factors that may be composed with it to form a given derivation.

Left Quotient

$$\frac{\partial}{\partial \vec{x}} = \{ z \mid (w \rightarrow xz) \in P \}$$



Right Quotient

$$\frac{\partial}{\partial \vec{z}} = \{ x \mid (w \rightarrow xz) \in P \}$$



The left quotient coincides with the derivative operator first proposed by Brzozowski [?] and Antimirov [?] over regular languages, lifted into the context-free setting (our work). When the root and LHS are fixed, e.g., $\frac{\partial S}{\partial x} : (\vec{V} \rightarrow S) \rightarrow \vec{V}$ returns the set of admissible nonterminals to the RHS. One may also consider a gradient operator, $\vec{\nabla} S : (\vec{V} \rightarrow S) \rightarrow \vec{V}$, which simultaneously tracks the partials with respect to a set of multiple LHS nonterminals produced by a fixed root.

¹Hereinafter, we use gray highlighting to distinguish between expressions containing only **constants** from those which may contain free variables.

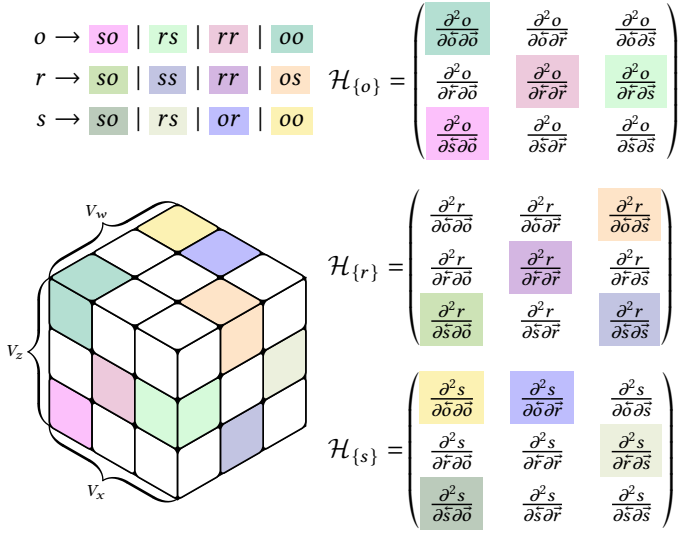


Figure 1. CFGs are witnessed by a rank-3 tensor, whose nonempty inhabitants indicate CNF productions. Gradients in this setting effectively condition the parse tensor M by constraining the superposition of admissible parse forests.

3 Context-sensitive reachability

It is well-known that the family of CFLs is not closed under intersection. For example, consider $\mathcal{L}_\cap := \mathcal{L}(\mathcal{G}_1) \cap \mathcal{L}(\mathcal{G}_2)$:

$$P_1 := \{ S \rightarrow LR, \quad L \rightarrow ab \mid aLb, \quad R \rightarrow c \mid cR \}$$

$$P_2 := \{ S \rightarrow LR, \quad R \rightarrow bc \mid bRc, \quad L \rightarrow a \mid aL \}$$

Note that \mathcal{L}_\cap generates the language $\{ a^d b^d c^d \mid d > 0 \}$, which according to the pumping lemma is not context-free. We can encode $\bigcap_{i=1}^c \mathcal{L}(\mathcal{G}_i)$ as a polygonal prism with upper-triangular matrices adjoined to each rectangular face. More precisely, we intersect all terminals $\Sigma_\cap := \bigcap_{i=1}^c \Sigma_i$, then for each $t_\cap \in \Sigma_\cap$ and CFG, construct an equivalence class $E(t_\cap, \mathcal{G}_i) = \{ w_i \mid (w_i \rightarrow t_\cap) \in P_i \}$ and bind them together:

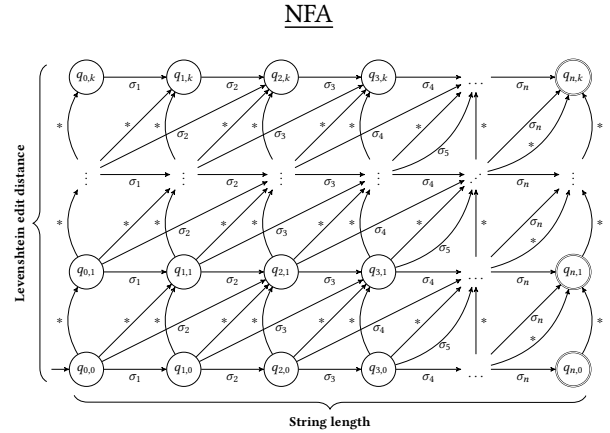
$$\bigwedge_{t \in \Sigma_\cap} \bigwedge_{j=1}^{c-1} \bigwedge_{i=1}^{|\sigma|} E(t_\cap, \mathcal{G}_j) \equiv_{\sigma_i} E(t_\cap, \mathcal{G}_{j+1}) \quad (2)$$



Figure 2. Orientations of a $\bigcap_{i=1}^4 \mathcal{L}(\mathcal{G}_i) \cap \Sigma^6$ configuration. As $c \rightarrow \infty$, this shape approximates a circular cone whose symmetric axis joins σ_i with orthonormal unit productions $w_i \rightarrow t_\cap$, and $S_i \in \Lambda_\sigma^*$ represented by the outermost bitvector inhabitants. Equations of this form are equiexpressive with the family of CSLs realizable by finite CFL intersection.

4 Levenshtein Reachability

Levenshtein reachability is recognized by the nondeterministic infinite automaton (NIA) whose topology $\mathcal{L} = \mathbb{Z} \times \mathbb{Z}$ can be factored into a product of (a) the monotone Chebyshev topology \mathbb{Z} , equipped with horizontal transitions accepting σ_i and vertical transitions accepting Kleene stars, and (b) the monotone knight's topology $\mathbb{Z} \times \mathbb{Z}$, equipped with transitions accepting σ_{i+2} . The structure of this space is representable as an acyclic NFA [?], populated by accept states within radius k of $q_{n,0}$, or equivalently, a left-linear CFG whose productions finitely instantiate the transition dynamics:



By intersection with a conjunctive language, we obtain a language $\mathcal{L}(\mathcal{G})$ that is both a subset of Σ^n and accepted by \mathcal{G} . We can then define the Levenshtein reachability problem as follows: given a string $\sigma \in \Sigma^n$, find the smallest k such that $\sigma \in \mathcal{L}(\mathcal{G})$.

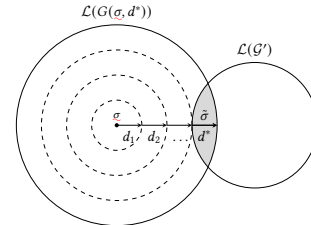


Figure 3. LED is computed gradually by incrementing d until $\mathcal{L}_d^\cap \neq \emptyset$.

5 Conclusion

Not only is linear algebra over finite fields an expressive language for inference, but also an efficient framework for inference on languages themselves. We illustrate a few of its applications for parsing incomplete strings and repairing syntax errors in context-free and sensitive languages. In contrast with LL and LR-style parsers, our technique can recover partial forests from invalid strings by examining the structure of M^* . In future work, we hope to extend our method to more natural grammars like PCFG and LCFRS.