

A Data-Oriented Approach to Program Synthesis

Fritz Obermeyer

2016-04-06

Overview

Motivation:

Given a curated knowledge base,
fill holes in programs subject to tests.

Machine Learning Methods:

Tune the knowledge base's language.
Find common phrases to add to language.
Search for plausible conjectures.

Synthesize programs by sketching

Solar-Lezama (2008)

Fill holes in a program
with terms of a specified language
subject to constraints
modulo simplification.

Synthesize programs by sketching

Solar-Lezama (2008)

Fill holes in a program
with terms of a specified language
subject to constraints
modulo simplification.

How?

Best-first search
modulo simplification
filtered by constraint propagation.

Example program with one hole

conj = $\lambda s. \lambda r. \lambda f. \text{COMP } r \text{ COMP } f \ s$

conj = $\lambda s. \lambda r. \lambda f. \text{COMP COMP } r \ f \ s$

raise = $\lambda x. \lambda y. x$

lower = $\lambda x. \text{APP } x \ \top$

pull = $\lambda x. \lambda y. \text{JOIN } x \ \text{APP DIV } y$

push = $\lambda x. \text{APP } x \ \perp$

move1 = $\langle l, l \rangle$

move2 = $\langle \text{push}, \text{pull} \rangle$

move3 = $\langle \text{lower}, \text{raise} \rangle$

move4 = $\lambda r1, s1. \lambda r2, s2. \langle \text{COMP } r2 \ r1, \text{COMP } s1 \ s2 \rangle$

move5 = $\lambda r1, s1. \lambda r2, s2. \langle \text{APP APP conj } s1 \ r2, \text{APP APP conj } r1 \ s2 \rangle$

hole \sqsubseteq A

APP hole B \sqsubseteq I

APP hole B = I

APP hole CB $\not\sqsubseteq$ I

hole $\not\sqsubseteq$ move1 move1 \sqsubseteq A

hole $\not\sqsubseteq$ move2 move2 \sqsubseteq A

hole $\not\sqsubseteq$ move3 move3 \sqsubseteq A

hole $\not\sqsubseteq$ move4 move4 \sqsubseteq A

hole $\not\sqsubseteq$ move5 move5 \sqsubseteq A

Example language

```
language = {  
    'APP': 1.0,  
    // 'COMP': 1.6,  
    'JOIN': 3.0,  
    'B': 1.0,  
    'C': 1.3,  
    'A': 2.0,  
    '⊥': 2.0,  
    'T': 2.0,  
    'I': 2.2,  
    // 'Y': 2.3,  
    'K': 2.6,  
    'S': 2.7,  
    'J': 2.8,  
    'DIV': 3.0,  
}
```

Approximate with an abstract interpretation

Given a set P of $\sim 30K$ probe terms, define a 4-tuple of sets

$$\llbracket x \rrbracket = \left\{ \begin{array}{l} \text{above: } \{p \in P \mid x \sqsubseteq p\}, \\ \text{below: } \{p \in P \mid p \sqsubseteq x\}, \\ \text{nabove: } \{p \in P \mid x \not\sqsubseteq p\}, \\ \text{nbelow: } \{p \in P \mid p \not\sqsubseteq x\} \end{array} \right\}$$

Approximate with an abstract interpretation

Given a set P of $\sim 30K$ probe terms, define a 4-tuple of sets

$$\llbracket x \rrbracket = \left\{ \begin{array}{l} \text{above: } \{p \in P \mid x \sqsubseteq p\}, \\ \text{below: } \{p \in P \mid p \sqsubseteq x\}, \\ \text{nabove: } \{p \in P \mid x \not\sqsubseteq p\}, \\ \text{nbelow: } \{p \in P \mid p \not\sqsubseteq x\} \end{array} \right\}$$

Growing P leads to more accurate approximations $\llbracket x \rrbracket$.

Verify by constraint propagation

Constraints are simple predicates: $=$, \sqsubseteq , $\not\sqsubseteq$.
This language is Π_0^2 -complete.

Verify by constraint propagation

Constraints are simple predicates: $=$, \sqsubseteq , $\not\sqsubseteq$.
This language is Π_0^2 -complete.

$\llbracket x \sqsubseteq y \rrbracket$ is satisfied iff
 $\text{below}(x) \cap \text{nbelow}(y) = \emptyset = \text{nabove}(x) \cap \text{above}(y)$.

Verify by constraint propagation

Constraints are simple predicates: $=$, \sqsubseteq , $\not\sqsubseteq$.
This language is Π_0^2 -complete.

$\llbracket x \sqsubseteq y \rrbracket$ is satisfied iff
 $\text{below}(x) \cap \text{nbelow}(y) = \emptyset = \text{nabove}(x) \cap \text{above}(y)$.

$\llbracket \text{APP } x \ y = z \rrbracket$ is a 3-way propagation node

$$\frac{x \sqsubseteq x' \quad y \sqsubseteq y'}{z \sqsubseteq \text{APP } x' \ y'}$$

$$\frac{x' \sqsubseteq x \quad y' \sqsubseteq y}{\text{APP } x' \ y' \sqsubseteq z}$$

etc.

And focus on data

1. Invest in a curated knowledge base.
2. Leverage the knowledge base in synthesis.

Invest in a Knowledge Base

- ~ 1 engineer year, ~ 10 CPU year,
- ~ 1GB compressed, ~ 10GB in memory.

Invest in a Knowledge Base

~ 1 engineer year, ~ 10 CPU year,
~ 1GB compressed, ~ 10GB in memory.

Automatically deduce many small facts.

Todd-Coxeter forward-chaining enumeration of
~ 30K equivalence classes.

VM-based with an inference rule compiler,
with ~ 100 inference rules.

ML Method 1: Tune the language

Fit a probabilistic grammar to a corpus.

Why?

- ▶ to learn facts about *relevant* terms
- ▶ to synthesize from *popular* components

How?

EM or Empirical Bayes

ML Method 2: Find common phrases

Suggest new atoms for the grammar.

Why?

“SK is a bad basis”

How?

Greedily suggest APP pairs from among $\sim 1M$.

ML Method 3: Find conjectures [1/3]

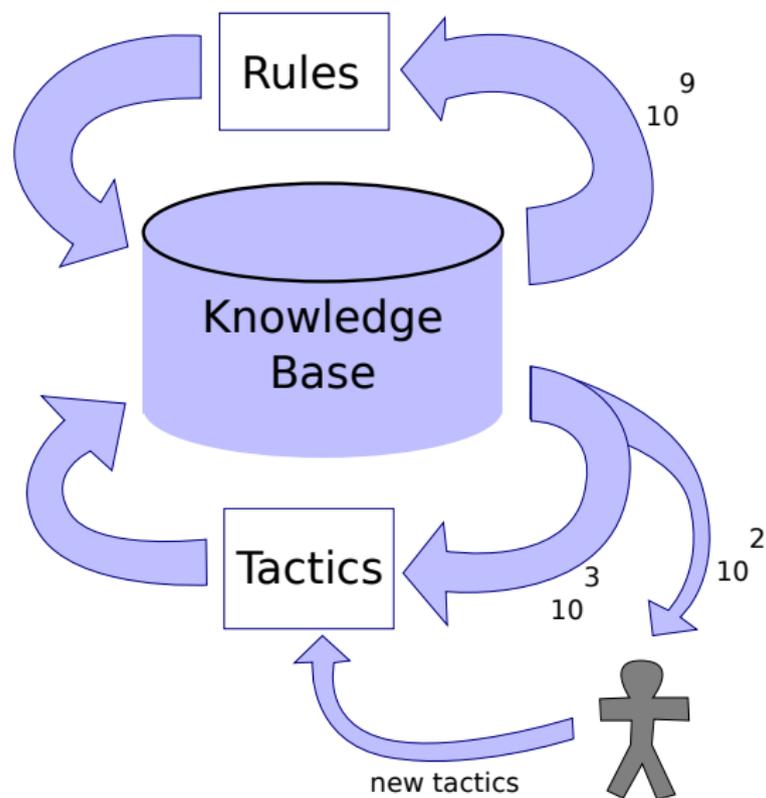
Search for plausible conjectures.

Among millions of statements

$$x = y, \quad x \sqsubseteq y, \quad x \not\sqsubseteq y$$

- ▶ Find 1000 to apply expensive tactics to
- ▶ Find 100 to send to a human

ML Method 3: Find conjectures [2/3]



ML Method 3: Find conjectures [3/3]

How?

Quantify evidence based on extensionality

$$\frac{}{\top \not\equiv \perp}$$

$$\frac{f\ x \not\equiv f\ y}{x \not\equiv y}$$

Questions?

Implemented at

<http://github.com/fritzo/pomagma>

Language

Combinatory algebra is a simple language:
 $S, K, I, AP(-, -)$.

There are no types, no variables, no binding.

Extensionally collapse to reduce search space.

Add nondeterminism to quotient out
implementation.

Implement types-as-closures via
nondeterminism.