Conditional Program Generation for Bimodal Program Synthesis

Swarat Chaudhuri Rice University www.cs.rice.edu/~swarat

(Joint work with Chris Jermaine, Vijay Murali, and Letao Qi)





Program synthesis

[Simon 1963, Summers 1977, Manna-Waldinger 1977, Pnueli-Rosner 1989]

Specification



Specification: Logical constraint that must be satisfied exactly

Algorithm: Search for a program that satisfies the specification.





"Bimodal" program synthesis



Neural Sketch Learning for Conditional Program Generation. Murali, Qi, Chaudhuri, and Jermaine. Arxiv 2017.



"Bimodal" program synthesis



- API calls or types that the program uses
- "Soft" I/O examples or constraints
- Natural language description of what the program does
- ...



The Bayou synthesizer: A demo

http://bit.ly/2zgP5fj



Conditional program generation

Assume random variables X and Prog, over **labels** and **programs** respectively, following a joint distribution Q(X, Prog).

Offline:

- You are given a set {(X_i, Prog_i)} of samples from Q(X, Prog). From this, learn a function g that maps evidence to programs.
- Learning goal: maximize $E_{(X,Prog)\sim Q}[I]$, where

$$I = \begin{cases} 1 & \text{if } g(X) \equiv Prog \\ 0 & \text{otherwise.} \end{cases}$$

Online: Given *X*, produce g(X).



In what we actually do

The map g is probabilistic.

Learning is maximum conditional likelihood estimation:

• Given { $(X_i, Prog_i)$ }, solve $\arg \max_{\theta} \sum_i \log P(Prog_i | X_i, \theta)$.



Programs

Language capturing the essence of API usage in Java.

Prog ::= skip | Prog₁; Prog₂ | call Call |
let
$$x = Call$$
 |
if Exp then Prog₁ else Prog₂ |
while Exp do Prog₁ | try Prog₁ Catch
Exp ::= Sexp | Call | let $x = Call : Exp_1$
Sexp ::= $c \mid x$
Call ::= Sexp₀. $a(Sexp_1, ..., Sexp_k)$
Catch ::= catch(x_1) Prog₁ ... catch(x_k) Prog_k
API method name

Labels

Set of API calls

• readline, write,...

Set of API datatypes

• BufferedReader, FileReader,...

Set of keywords that may appear while describing program actions in English

- read, file, write,...
- Obtained from API calls and datatypes through a camel case split



Challenges

Directly learning over source code simply doesn't work

- Source code is full of low-level, program-specific names and operations.
- Programs need to satisfy structural and semantic constraints such as type safety. Learning to satisfy these constraints is hard.



Language abstractions to the rescue!

Learn not over programs, but typed, syntactic models of programs.





The sketch of a program is obtained by applying an **abstraction function** α .

From sketch Y to program Prog: a fixed **concretization distribution** P(Prog | Y).

Learning goal changes to

• Given { (X_i, Y_i) }, solve arg max $\sum_i \log P(Y_i | X_i, \theta)$.













CAPER







P(Z) = Normal(0, I) $P(f(X) | Z) = Normal(Z, \sigma^{2}I)$

During learning, use Jensen's inequality to get smooth loss function





P(Z) = Normal(0, I) $P(f(X) | Z) = Normal(Z, \sigma^{2}I)$

During inference, get P(Z | X) using normal-normal conjugacy



Neural decoder





Concretization





Results

- Trained method on 100 million lines of Java/Android code. ~2500 API methods, ~1500 types.
- Synthesis of method bodies from scratch, given 2-3 API calls and types.
- Sketch learning critical to accuracy.
- Good performance compared to GSNNs (state of the art conditional generative model).
- Good results on label-sketch pairs not encountered in training set.



Thank you!

Questions?

swarat@rice.edu
http://www.cs.rice.edu/~swarat

(Research funded by the DARPA MUSE award #FA8750-14-2-0270)

