

**Control-flow analysis**

Lecture 22

Thursday, April 17, 2014

**1 Dataflow analysis for functional programs**

Last lecture we looked at abstract interpretation of IMP, an imperative programming language. Abstract interpretation provides a way to statically approximate the runtime behavior of a program more precisely than the type systems that we have seen so far.

However, how would we perform abstract interpretation for a functional language? That is, in a functional language, functions are values, and to perform abstract interpretation on a functional language, we would need an abstract domain for functions.

In this lecture we will consider *dataflow analysis* for functional programs. Dataflow analysis is a technique for statically determining what “facts” are true at various points in the program. The dataflow facts often concern what the set of possible values that variables may have at certain program points. (Dataflow analysis and abstract interpretation are related to each other, and in some ways are similar to each other; we will not dig into the precise connection between them here.)

In order to perform dataflow analysis (i.e., to track possible values through a program), we need to know the control flow of a program. In IMP, this is straightforward: the lexical structure of the program tells us the control flow structure. But in a functional language, if we see an application  $f y$ , then control will jump to the function body of whatever function the variable  $f$  evaluates to. This makes determining the control flow of a functional language more complex than determining the control flow structure of a language like IMP, where functions are not first-class values. (Similar issues arise in object-oriented languages, where objects are first class and contain executable code.)

To determine which functions the variable  $f$  may evaluate to, we can perform dataflow analysis to approximate the values that  $f$  may evaluate to. So in a functional language, we need to perform dataflow analysis and control-flow analysis simultaneously.

We will consider the OCFA analysis. (CFA stands for Control-Flow Analysis; it is an instance of  $k$ -CFA analysis, where the parameter  $k$  determines the precision of the analysis.) Like our type inference analysis, we will generate a set of constraints from a program, and then find a solution to the set of constraints.

**1.1 Labeled lambda calculus**

Let’s consider a lambda calculus with integers. The syntax of the language is similar to what we’ve seen before, but every expression in our source program will have a unique label. Labels are used to uniquely identify program expressions. We use  $l$  to range over labels.

$$e ::= n^l \mid x^l \mid (\lambda x. e)^l \mid (e_1 e_2)^l \mid (e_1 + e_2)^l$$

Every expression has a label, and we write  $labelof(e)$  for the label of expression  $e$ . For convenience, we also write  $e^l$  to mean that  $l$  is a label such that  $labelof(e) = l$ .

Also, given a program  $e$ , we write  $expof(e, l)$  for the (unique) subexpression  $e'$  of  $e$  such that  $labelof(e') = l$ . That is  $expof(e, l)$  allows us to get the expression associated with label  $l$ .

We will use a large-step environment semantics for this analysis (see Lecture 17). Recall that environments  $\rho$  are maps from variables to values, and that judgment  $\langle e, \rho \rangle \Downarrow v$  means that expression  $e$  in environment  $\rho$  evaluates to value  $v$ . Note that in an environment semantics, a function  $\lambda x. e$  evaluates to a *closure*  $((\lambda x. e)^l, \rho_{lex})$ , where  $\rho_{lex}$  is the environment that was current when  $\lambda x. e$  was evaluated. That is,  $\rho_{lex}$  binds the free variables (except  $x$ ) in the function body  $e$ . We define the label of a closure to be the label of the function:  $labelof(((\lambda x. e)^l, \rho_{lex})) = l$ .

Values in our language are thus the following.

$$v ::= n^l \mid ((\lambda x. e)^l, \rho)$$

$$\frac{}{\langle x^l, \rho \rangle \Downarrow \rho(x)} \quad \frac{}{\langle n^l, \rho \rangle \Downarrow n^l} \quad \frac{\langle e_1, \rho \rangle \Downarrow n_1^{l_1} \quad \langle e_2, \rho \rangle \Downarrow n_2^{l_2}}{\langle (e_1 + e_2)^l, \rho \rangle \Downarrow n^l} \quad n = n_1 + n_2$$

$$\frac{}{\langle (\lambda x. e)^l, \rho \rangle \Downarrow ((\lambda x. e)^l, \rho)} \quad \frac{\langle e_1, \rho \rangle \Downarrow ((\lambda x. e)^{l_1}, \rho_{lex}) \quad \langle e_2, \rho \rangle \Downarrow v_2 \quad \langle e, \rho_{lex}[x \mapsto v_2] \rangle \Downarrow v}{\langle (e_1 e_2)^l, \rho \rangle \Downarrow v}$$

Notice that every value is labeled with the label of the expression that created it. (For integer values, the label is either the label of the integer literal that appeared in the source program, or the label of an addition.)

## 1.2 Analysis

The aim of our analysis is to approximate the values that expressions can evaluate to. Since functions are values, this will also tell us about control flow, since given an application  $e_1 e_2$ , knowing which function values expression  $e_1$  may evaluate to tells us about the possible control flow of the function application.

Our analysis will find a function  $C : \mathbf{Label} \rightarrow \mathcal{P}(\mathbf{Label})$  that approximates for each label  $l$  the set of values that the expression labeled  $l$  may evaluate to. Specifically, if  $l' \in C(l)$ , then the expression labeled  $l$  may evaluate to a value labeled  $l'$ .

Our analysis will also find a function  $r : \mathbf{Var} \rightarrow \mathcal{P}(\mathbf{Label})$  that for each variable  $x$  will approximate the set of values that  $x$  may be bound to. That is,  $l \in r(x)$ , then the variable  $x$  may be bound to a value labeled  $l'$ . We assume without loss of generality that variable names in a program are unique. (If we don't have this assumption, the analysis will still work, it will just be less precise.)

Given a program  $e$ , we produce a set of constraints on functions  $C$  and  $r$  by examining the program. Intuitively, if we can find functions  $C$  and  $r$  that satisfy the constraints, then  $C$  and  $r$  will give us correct information about what values expressions and variables may evaluate to.

We generate the set of constraints using function  $\mathcal{C}[\cdot]_e : \mathbf{Expr} \rightarrow \mathcal{P}(\mathbf{Constraint})$ . Here,  $e$  is the program we are analyzing, and we use it in order to map labels to expressions.  $\mathcal{C}[\cdot]_e$  is defined as follows.

$$\begin{aligned} \mathcal{C}[n^l]_e &= \{l \in C(l)\} \\ \mathcal{C}[(e_1 + e_2)^l]_e &= \mathcal{C}[e_1]_e \cup \mathcal{C}[e_2]_e \cup \{l \in C(l)\} \\ \mathcal{C}[x^l]_e &= \{r(x) \subseteq C(l)\} \\ \mathcal{C}[(\lambda x. e_1)^l]_e &= \{l \in C(l)\} \cup \mathcal{C}[e_1]_e \\ \mathcal{C}[(e_1^{l_1} e_2^{l_2})^l]_e &= \mathcal{C}[e_1^{l_1}]_e \cup \mathcal{C}[e_2^{l_2}]_e \\ &\quad \cup \{l' \in C(l_1) \Rightarrow C(l_2) \subseteq r(x) \mid \text{expref}(e, l') = (\lambda x. e_0^{l_0})^{l'}\} \\ &\quad \cup \{l' \in C(l_1) \Rightarrow C(l_0) \subseteq C(l) \mid \text{expref}(e, l') = (\lambda x. e_0^{l_0})^{l'}\} \end{aligned}$$

Let's consider what each of these constraints mean. For values  $n^l$ ,  $(\lambda x. e_1)^l$  and addition  $(e_1 + e_2)^l$  we have constraint  $l \in C(l)$ . This means that the expression labeled  $l$  may evaluate to a value labeled  $l$ , and, if we look at the semantics, it is indeed the case.

$\mathcal{C}[x^l]_e$  produces the constraint  $r(x) \subseteq C(l)$ , meaning that if  $x$  may be bound to a value labeled  $l'$  (i.e.,  $l' \in r(x)$ ), then the expression  $x^l$  may evaluate to that value (i.e.,  $l' \in C(l)$ ).

The constraints for application  $(e_1^{l_1} e_2^{l_2})^l$  are most interesting. The *conditional constraint*  $l' \in C(l_1) \Rightarrow C(l_2) \subseteq r(x)$  requires that if expression  $e_1$  may evaluate to function value  $(\lambda x. e_0^{l_0})^{l'}$  (i.e.,  $l' \in C(l_1)$ ), then  $x$ , the argument of that function, may be bound to anything that  $e_2$  can evaluate to (i.e.,  $C(l_2) \subseteq r(x)$ ).

Similarly, constraint  $l' \in C(l_1) \Rightarrow C(l_0) \subseteq C(l)$  requires that if expression  $e_1$  may evaluate to function value  $(\lambda x. e_0^{l_0})^{l'}$  (i.e.,  $l' \in C(l_1)$ ), then the application expression may evaluate to anything that function body  $e_0$  may evaluate to (i.e.,  $C(l_0) \subseteq C(l)$ ).

Solving these constraints is straightforward. We will start off with a very bad approximation to the functions:  $C_0(l) = \emptyset$  for all labels  $l$ , and  $r_0(x) = \emptyset$  for all variables  $x$ . We will then iteratively improve these approximations by adding labels to the sets to satisfy the constraints. We will keep iterating until we reach a fixed point. Since there are only a finite number of labels, and in each iteration we only add labels to the sets, we are guaranteed to reach a fixed point.

$$C_0 = \lambda l. \emptyset$$

$$r_0 = \lambda x. \emptyset$$

$$C_{i+1} = \lambda l. C_i(l)$$

$$\cup \{l \mid (l \in C(l)) \in \mathcal{C}[[e]]_e\}$$

$$\cup \bigcup \{r_i(x) \mid (r(x) \subseteq C(l)) \in \mathcal{C}[[e]]_e\}$$

$$\cup \bigcup \{C_i(l_0) \mid (l' \in C(l_1) \Rightarrow C(l_0) \subseteq C(l)) \in \mathcal{C}[[e]]_e \text{ and } l' \in C_i(l_1)\}$$

$$r_{i+1} = \lambda x. r_i(x)$$

$$\cup \bigcup \{C_i(l_2) \mid (l' \in C(l_1) \Rightarrow C(l_2) \subseteq r(x)) \in \mathcal{C}[[e]]_e \text{ and } l' \in C_i(l_1)\}$$

The least fixed point, which we will denote  $C^*$  and  $r^*$  is the bound of all of the approximations.

$$C^* = \lambda l. \bigcup_{i \in \mathbb{N}} C_i(l)$$

$$r^* = \lambda x. \bigcup_{i \in \mathbb{N}} r_i(x)$$

### 1.3 Example

Let's work through a simple example. Consider the following program.

$$e \equiv (((\lambda a. a^1)^2 (\lambda b. b^3)^4)^5 99^6)^7$$

The set of constraints for this simple program is as follows. (Exercise: make sure you understand how these constraints were derived.)

$$\begin{aligned} \mathcal{C}[[e]]_e = \{ & 2 \in C(2), 4 \in C(4), 6 \in C(6), \\ & r(a) \subseteq C(1), r(b) \subseteq C(3), \\ & 2 \in C(5) \Rightarrow C(6) \subseteq r(a), \\ & 4 \in C(5) \Rightarrow C(6) \subseteq r(b), \\ & 2 \in C(5) \Rightarrow C(1) \subseteq C(7), \\ & 4 \in C(5) \Rightarrow C(3) \subseteq C(7), \\ & 2 \in C(2) \Rightarrow C(4) \subseteq r(a), \\ & 4 \in C(2) \Rightarrow C(4) \subseteq r(b), \\ & 2 \in C(2) \Rightarrow C(1) \subseteq C(5), \\ & 4 \in C(2) \Rightarrow C(3) \subseteq C(5) \} \end{aligned}$$

Let's consider the result of solving it iteratively. In the following table, columns indicate the values of  $C_i(l)$  and  $r_i(x)$  over the various iterations.

$i$	$C_i(1)$	$C_i(2)$	$C_i(3)$	$C_i(4)$	$C_i(5)$	$C_i(6)$	$C_i(7)$	$r_i(a)$	$r_i(b)$
0									
1		2		4		6			
2		2		4		6		4	
3	4	2		4		6		4	
4	4	2		4	4	6		4	
5	4	2		4	4	6		4	6
6	4	2	6	4	4	6		4	6
7	4	2	6	4	4	6	6	4	6
8	4	2	6	4	4	6	6	4	6

At the 8th iteration, we have  $C_7 = C_8$  and  $r_7 = r_8$ , and we have reached a fixed point, and so  $C^* = C_8$  and  $r^* = r_8$ .

Let's double check that this analysis returned reasonable results. For example,  $C^*(5) = \{4\}$ , meaning that the expression  $((\lambda a. a^1)^2 (\lambda b. b^3)^4)^5$  may evaluate to a value labeled 4, i.e., to the value  $(\lambda b. b^3)^4$ . That is indeed consistent with the actual execution of the program. Another example:  $C^*(7) = \{6\}$ , meaning that the whole program may evaluate to a value labeled 6, i.e., to the labeled integer 99<sup>6</sup>.

Note that  $2 \notin C^*(5)$ . That is, the analysis correctly says that expression  $((\lambda a. a^1)^2 (\lambda b. b^3)^4)^5$  can not evaluate to  $(\lambda a. a^1)^2$ .

#### 1.4 Soundness

The analysis is sound, meaning that, given a program  $e_0$ , if  $r^*$  and  $C^*$  satisfy the set of constraints  $\mathcal{C}[[e_0]]_{e_0}$ , then  $C^*$  conservatively describes what expressions may evaluate to, and  $r^*$  conservatively describes what variables may be bound to.

**Theorem (Soundness).** *Let  $e_0$  be a program. Let  $r^*$  and  $C^*$  satisfy the set of constraints  $\mathcal{C}[[e_0]]_{e_0}$ . If  $\langle e_0, \emptyset \rangle \Downarrow v_0$  and  $\langle e, \rho \rangle \Downarrow v$  appears in the derivation of  $\langle e_0, \emptyset \rangle \Downarrow v_0$ , then  $\text{labelof}(v) \in C^*(\text{labelof}(e))$ .*

In order to prove the soundness theorem, we need a stronger lemma. To state the lemma, we will extend the set-constraint generation function to environments and closures.

$$\mathcal{C}[[\rho]]_e = \bigcup_{x \in \text{dom}(\rho)} \{C(\text{labelof}(\rho(x))) \subseteq r(x)\} \cup \mathcal{C}[[\rho(x)]]_e$$

$$\mathcal{C}[[\langle \lambda x. e_1 \rangle^l, \rho]]_e = \mathcal{C}[[\langle \lambda x. e_1 \rangle^l]]_e \cup \mathcal{C}[[\rho]]_e$$

With this extended definition in hand, we can state the lemma, which can then be proved by induction on derivations  $\langle e, \rho \rangle \Downarrow v$ .

**Lemma.** *Let  $e_0$  be a program,  $e$  an expression and  $\rho$  an environment. Let  $r^*$  and  $C^*$  satisfy the constraints  $\mathcal{C}[[\langle e, \rho \rangle]]_{e_0}$ . If  $\langle e, \rho \rangle \Downarrow v$  then  $\text{labelof}(v) \in C^*(\text{labelof}(e))$  and  $r^*$  and  $C^*$  satisfy the constraints  $\mathcal{C}[[v]]_{e_0}$ .*