

Type analysis

Letterio Galletta

Introduction to type analysis

- Type analysis is one of the most famous and used static analysis
- Every static typed language performs type analysis when is compiled/interpret

The goal

The aim of type analysis is to ensure that no type errors occur at run-time. To do that a type system:

- classify program phrases according to values they compute
- ensure that each operation is applied to correct data

What is a type system?

What a type? (my answer)

A type of a term usually denotes some approximation of the term semantics:

- a set of values equipped with a set of operations
- something more complex describing the behavior of a term, e.g., effects

A type system is an instance of deduction systems where theorems to be proved are about program types

The recipe for defining a type system

Defining a new type system consists of 4 steps:

1. Define the syntax of types and the syntax of type environments
2. Define the format of typing judgments and the typing rules (the typing relation)
3. Prove the type system is sound with respect to the semantics of the language
4. Define the typing algorithm and prove it correct

Our running example: the fun language

$x, f \in Id$

$c \in Const$

$e \in Exp ::= c$

| x

| **if** e_1 **then** e_1 **else** e_2

| $e_1 \diamond e_2$

| **let** $x = e_1$ **in** e_2

| **fun** $x \rightarrow e$

| $e_1 e_2$

constant literal

identifier

conditionals

$\diamond \in \{+, -, *, \dots\}$ primitive operators

declarations

lambda abstraction

function application

Defining types

- Usually, programming languages provide us values that are scalar, e.g., integers, strings, and values that are compound, e.g., lists, arrays, maps
- A type can be considered as a summary/description of a collection of values that share properties, e.g., 1 and 2 share the property to be integer values and that can be added
- The syntax of types describes which the type for scalar values are and how we can define new types by composing those already existent one
- The syntax of types is specified using a grammar, whose language corresponds to all possible types

Type syntax of fun

$$\tau ::= \tau_c \mid \tau_1 \rightarrow \tau_2 \qquad \tau_c \in \{int, bool, string, \dots\}$$

- The type τ_c abstractly represents the type of a generic scalar value
- There is only one type constructor $\cdot \rightarrow \cdot$ that allows us to create functional types

Typing environments

A type environment records the association among variables and their types, formally it is a list of pairs

$$\Gamma ::= \emptyset \mid x : \tau, \Gamma$$

Usually, variables are assumed to be distinct and Γ is considered as a partial function, with the expected meaning for $\Gamma(x)$ and $\text{dom}(\Gamma)$

Type judgments

A judgment is a statement about the type of a program term typically it has the form

$$\Gamma \vdash M$$

where the free variables of M occur in Γ (we say Γ *closes* M)

A judgment usually express a relation between a term and a type (typing relation) or between two types (sub-typing relation)

Type judgments of fun

In our language we have only a kind of type judgment

$$\Gamma \vdash e : \tau$$

meaning that the expression e has type τ in Γ .

Typing rules

A typing rule specifies when a judgment is true:

$$\frac{\Gamma \vdash M_1 \cdots \Gamma \vdash M_n}{\Gamma \vdash M}$$

Rules without premises are axioms, thus, always true

A judgment is valid iff it can be derived from axioms, inductively applying judgments (derivation tree)

Saying that a term M has type τ in Γ means to find a derivation for the judgment $\Gamma \vdash M : \tau$

Typing rules of fun (1)

$$\begin{array}{c} \text{T}_{\text{CONST}} \\ \hline \Gamma \vdash c : \tau_c \end{array} \qquad \begin{array}{c} \text{T}_{\text{VAR}} \\ x \in \text{dom}(\Gamma) \quad \Gamma(x) = \tau \\ \hline \Gamma \vdash x : \tau \end{array}$$

$$\begin{array}{c} \text{T}_{\text{OP}} \\ \Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \diamond : \tau_1 \times \tau_2 \rightarrow \tau \\ \hline \Gamma \vdash e_1 \diamond e_2 : \tau \end{array}$$

$$\begin{array}{c} \text{T}_{\text{IF}} \\ \Gamma \vdash e_1 : \text{bool} \quad \Gamma \vdash e_2 : \tau \quad \Gamma \vdash e_3 : \tau \\ \hline \Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : \tau \end{array}$$

Typing rules of fun (2)

$$\frac{\text{T}_{\text{LET}} \quad \Gamma \vdash e_1 : \tau_1 \quad x : \tau_1, \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \mathbf{let} \, x = e_1 \mathbf{in} \, e_2 : \tau_2}$$

$$\frac{\text{T}_{\text{FUN}} \quad x : \tau_1, \Gamma \vdash e : \tau_2}{\Gamma \vdash \mathbf{fun} \, x \rightarrow e : \tau_1 \rightarrow \tau_2}$$

$$\frac{\text{T}_{\text{APP}} \quad \Gamma \vdash e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash e_2 : \tau_1}{\Gamma \vdash e_1 \, e_2 : \tau_2}$$

Note: in the rule T_{FUN} the type τ_1 of the formal parameter x is guessed

An example of derivation

$$\frac{\frac{\frac{\Gamma \vdash x : int}{\Gamma \vdash 3 : int} \quad < : int \times int \rightarrow bool}{\Gamma \vdash x < 3 : bool} \quad \frac{\frac{\Gamma \vdash x : int}{\Gamma \vdash 1 : int} \quad + : int \times int \rightarrow int}{\Gamma \vdash x + 1 : int} \quad \frac{\frac{\Gamma \vdash y : int}{\Gamma \vdash 1 : int} \quad + : int \times int \rightarrow int}{\Gamma \vdash y + 1 : int}}{\Gamma \vdash \text{if } x < 3 \text{ then } x + 1 \text{ else } y + 1 : int}$$

where $\Gamma = x : int, y : int, \emptyset$

An example of type error

$$\frac{\begin{array}{c} \Gamma \vdash x : int \\ \Gamma \vdash 3 : int \\ \hline < : int \times int \rightarrow bool \\ \hline \Gamma \vdash x < 3 : bool \end{array} \quad \begin{array}{c} \Gamma \vdash x : int \\ \Gamma \vdash 1 : int \\ \hline + : int \times int \rightarrow int \\ \hline \Gamma \vdash x + 1 : int \end{array} \quad \begin{array}{c} \Gamma \vdash y : real \\ \Gamma \vdash 1.0 : real \\ \hline \boxplus : real \times real \rightarrow real \\ \hline \Gamma \vdash y \boxplus 1.0 : real \end{array}}{\Gamma \vdash \mathbf{if} \ x < 3 \ \mathbf{then} \ x + 1 \ \mathbf{else} \ y \boxplus 1.0 : \mathbf{error}}$$

where $\Gamma = x : int, y : real, \emptyset$

The T_{IF} rule requires that both branches have exactly the same type: *int* and *real* are different types because there is no implicit conversion nor sub-typing

Soundness

- The soundness ensures that the type system does what it was supposed to do, i.e., it prevents type errors to happen at run-time (well-typed programs don't go wrong)
- A sound type system discards statically all programs that may cause a type error, but it may discard more due to approximations
- The precision of a type system can be improved by enriching the type language, e.g., by adding polymorphism or sub-typing

An example of approximation

$$\frac{\begin{array}{c} \emptyset \vdash 2 : int \\ \emptyset \vdash 1 : int \\ + : int \times int \rightarrow int \\ \hline \emptyset \vdash 2 + 1 : int \end{array} \quad \begin{array}{c} \emptyset \vdash "hello" : string \\ \emptyset \vdash "world" : string \\ \oplus : string \times string \rightarrow string \\ \hline \emptyset \vdash "hello" \oplus "world" : string \end{array}}{\emptyset \vdash \mathbf{if} \text{ true } \mathbf{then} 2 + 1 \mathbf{else} "hello" \oplus "world" : \mathbf{error}}$$

The expression does not type check, but it never raises a run-time type error, yet:

$$\mathbf{if} \text{ true } \mathbf{then} 2 + 1 \mathbf{else} "hello" \oplus "world" \rightarrow 2 + 1 \rightarrow 3$$

Statically, we cannot decide which branch of **if** to take, thus, we consider both of them

Typing algorithm

- So far, we have considered the logical presentation of a type system but the actually to run the analysis we need a type-checking algorithm
- There are different techniques for defining these algorithms: sometimes translating rules in code suffices, but often to deal with advanced features, e.g., polymorphism, we need other approaches
- However, the algorithm must be proved to compute what the type system prescribes
 - Soundness:** if the algorithm returns that a term M has type t , then the judgment $\vdash M : t$ must be valid
 - Completeness:** if the term M has type t , the algorithm must return t as type of M

Typing algorithm for fun

In the rule T_{FUN} we guess the type of the parameter x

$$\begin{array}{c} T_{\text{FUN}} \\ \frac{x : \tau_1, \Gamma \vdash e : \tau_2}{\Gamma \vdash \mathbf{fun} \, x \rightarrow e : \tau_1 \rightarrow \tau_2} \end{array}$$

but guessing is not feasible algorithmically

Two possible approaches:

1. Introducing type annotations, the programmer specifies the types
2. Inferring the type of values from their usage, no annotations

Adding type annotations to fun

We annotate the formal parameter of a functional value with its type

$$\begin{array}{l} e \in \text{Exp} ::= \dots \\ \quad | \text{fun } x : \tau_1 \rightarrow e \quad \text{lambda abstraction} \\ \quad | \dots \end{array}$$

The rule T_{FUN} becomes

$$\frac{\text{T}_{\text{FUN}} \quad x : \tau_1, \Gamma \vdash e : \tau_2}{\Gamma \vdash \text{fun } x : \tau_1 \rightarrow e : \tau_1 \rightarrow \tau_2}$$

Our type system becomes algorithmic and its rules can be directly implemented. The resulting algorithm is sound and complete.

Demo

See the file fun.ml

A more interesting and complete implementation is MiniML from PLZoo project:
<https://github.com/andrejbauer/plzoo/>

Conclusion

- Type systems are one of the most used technique for the analysis of programs
- We describe the case of simple type checking for a functional language: type annotations in functions
- Many other topics: type inference, polymorphism á la ML, parametric polymorphism, type system for security, type and effect systems, the different family of dependent types, objects and classes, flow-sensitive types.