

Progetto di Linguaggi di programmazione 2014-15
Prima parte
Analizzatore lessicale del Linguaggio LispKit

G.Filè

14 ottobre 2014

Il Linguaggio ListKit

Il ListKit è un linguaggio funzionale molto semplice, ma sufficientemente complesso da illustrare molte cose interessanti dei linguaggi di programmazione che vengono trattate durante il corso. La principale semplificazione del ListKit è che non ha tipi, ma può comunque manipolare valori diversi che sono i seguenti:

1. valori interi 12, 0, ~3 ...
2. valori booleani true e false
3. stringhe costanti `"ecco una stringa"`
4. liste come `cons(0 cons(true cons("abc" nil)))`.

Si deve subito osservare che le liste sono costruite con l'uso dell'operatore `cons` che corrisponde all'operatore `::` ML. Quindi in ListKit non si possono scrivere liste con le parentesi quadrate come in ML, ma solo con espressioni che usano l'operatore `cons` che costruisce le liste stesse. Questa scelta è volta a mantenere il linguaggio semplice senza perdere la possibilità di manipolare liste. Inoltre l'esempio mostra che in ListKit le liste non sono omogenee, ma possono contenere oggetti di tipo completamente diverso (come nell'esempio). D'altra parte i tipi non esistono in ListKit. Oltre a questo, è importante osservare che gli elementi di una lista sono separati semplicemente da uno o più spazi.

Ogni programma ListKit è costituito da una serie di dichiarazioni locali seguite da una espressione che usa le dichiarazioni. Questa espressione può essere costituita da altre dichiarazioni locali seguite da un'espressione che le usa e così via con annidamento di profondità arbitraria. Vediamo qualche esempio di programma ListKit.

```
val D= "let x=cons(\"ab\" cons(\"cd\" nil))
in if true then cons(\"01\" x) else nil end $";
```

La variabile `D` ha come valore una stringa che è un programma ListKit che dichiara in `x` una lista che consiste di 2 stringhe costanti (`"ab"` `"cd"`) che è seguita da un'espressione che aggiunge come primo elemento della lista la stringa costante `"01"`. In ListKIT le stringhe costanti sono iniziate e terminate da `"` per facilitarne il riconoscimento. Infatti in questo modo sono chiaramente diverse dalle keyword e dalle variabili. Inoltre l'esempio mostra nuovamente che il ListKit, per costruire liste, usa l'operatore `cons` che corrisponde al `::` dell'ML ed usa anche il costrutto condizionale `if_then_else`. Si osservi inoltre che, come per le liste, così i parametri formali e attuali delle funzioni sono separati da spazi. ListKit non ha tipi. Sta al programmatore scrivere cose che mescolano valori solo in modo sensato. Ogni programma termina sempre con il simbolo speciale `$`. Il programma `D` verrà esploso in una lista di caratteri (`explode(D)`) e dato in pasto all'analizzatore lessicale che chiameremo `lexi`. La seguente è la ben

nota funzione `reverse` che abbiamo già studiato in varie forme per introdurre il linguaggio ML.

```
val R="letrec rev = lambda(x y) if eq(x nil)
then nil else rev(cdr(x) cons(car(x) y))
in rev(cons(0 cons(1 cons(2 nil))) nil) end $";
```

Per prima cosa noteremo che la dichiarazione locale inizia ora con la keyword `letrec` (anziché `let`) ad indicare che si sta per definire una funzione ricorsiva. L'esempio mostra che la keyword che in `LispKit` definisce le espressioni funzionali è `lambda` (anziché `fn` come in ML). Inoltre la sintassi prevede che i parametri formali seguano la keyword `lambda` e siano racchiusi tra parentesi tonde e siano separati da spazi, mentre il corpo della funzione semplicemente segue la lista dei parametri formali. Le operazioni `head` e `tail` in `LispKit` sono indicate col nome che queste operazioni hanno nel linguaggio Lisp e cioè, `car` e `cdr`, rispettivamente. Il test di uguaglianza del `LispKit` è indicato con la funzione binaria `eq`.

Esaminiamo ora un esempio più complesso con dichiarazioni di funzioni ricorsive e di ordine superiore.

```
val C="letrec FACT = lambda (X) if eq(X 0)
then 1 else X*FACT(X-1)and
G = lambda (H L) if eq (L nil)
then L else cons(H(car( L ) ) G (H cdr ( L )))
in G ( FACT cons(1 cons(2 cons(3 nil))) ) end $";
```

Si noti nuovamente l'uso di `letrec` ad indicare che si stanno per dichiarare funzioni ricorsive. Nel `letrec` vengono dichiarate due funzioni ricorsive (separate dalla keyword `and`). La prima è la funzione fattoriale, mentre la seconda `G` è una funzione di ordine superiore in quanto aspetta una funzione `H` come parametro. Il fatto che il parametro `H` debba essere una funzione, lo si deduce dal suo uso nel corpo di `G`. Questo fatto viene anche confermato dalla successiva espressione che invoca `G` passandole proprio `FACT` come primo parametro.

L'analizzatore lessicale

L'analizzatore lessicale riceve come input un programma in `LispKit`, cioè una lista di caratteri e deve riconoscere le componenti elementari del linguaggio e deve metterle in una forma che sia semplice da manipolare nella successiva fase di analisi sintattica. Per esempio, deve riconoscere le costanti (per esempio i numeri interi oppure il valore `true`, eccetera), le parole chiave, gli identificatori, gli operatori ed i simboli di separazione.

Ognuna di queste componenti elementari viene rappresentata da un valore del seguente tipo ML `token`:

```

data Keyword_T = LET | IN | END | LETREC | AND | IF | THEN |
                ELSE | LAMBDA
                deriving (Show,Eq)

data Operator_T = EQ | LEQ | CAR | CDR | CONS | ATOM
                deriving (Show,Eq)

data Symbol_T = LPAREN | RPAREN | EQUALS | PLUS | MINUS |
               TIMES | DIVISION | VIRGOLA | DOLLAR
               deriving (Show,Eq)

data Token = Keyword Keyword_T | Operator Operator_T |
            Id String | Symbol Symbol_T | Number Integer |
            String String | Bool Bool | Nil | PROG_OK
            deriving (Show,Eq)

```

Per semplicità nel seguito chiameremo **token** i valori del tipo Token. Vediamo alcuni esempi di token corrispondenti a componenti elementari di programmi Lisp Kit:

- numeri interi: `Number(n)`, dove `n` è un intero;
- stringhe: `String(s)`, dove `s` è una stringa costante (senza `\`");
- identificatori: `Id(s)`, dove `s` è la stringa corrispondente all'identificatore;
- corrispondentemente a ciascuna keyword `K`, verrà prodotto il token `Keyword(K')`, dove `K'` è il valore `Keyword_T` corrispondente a `K`, per esempio, se `K=let`, allora il token corrispondente è `Keyword(LET)` e `Keyword(THEN)` rappresenta la keyword `then` e così via.
- nil: `Nil`;
- true e false : `Bool(true)` e `Bool(false)`;
- i simboli come `+`, `=`, `(` eccetera sono rappresentati da `Symbol(PLUS)`, `Symbol(EQUALS)`, `Symbol(LPAREN)` eccetera, dove `PLUS`, `EQUALS`, `LPAREN` sono costruttori del tipo `Symbol_T`.
- per gli operatori come `eq`, `car`, `cdr` eccetera, essi sono rappresentati da `Operator(EQ)`, `Operator(CAR)`, `Operator(CDR)`, eccetera, dove `EQ`, `CAR`, `CDR` sono costruttori del tipo `Operator_T`.

La segnatura della funzione ML `lexi` che realizza l'analizzatore è dunque la seguente:

`lexi :: [char] -> [token]`

Questa funzione deve implementare un automa a stati finiti che funziona secondo i seguenti principi. Partendo da uno stato iniziale I considera un carattere alla volta della sua lista in input. Chiamiamo questo carattere *s*:

- se *s* è uno spazio, lo si salta;
- se *s* è il dollaro, l'analisi termina;
- se *s* è numerico o è il carattere tilde che indica il meno unario, allora si deve passare ad un nuovo stato N atto a riconoscere i numeri; in N si continueranno a leggere caratteri numerici fino a che non si legga un non numerico e alla fine si produrrà il token `Number (x)` dove *x* è l'intero corrispondente alla stringa di caratteri letta;
- se invece *s* è un simbolo di parentesi o di = allora si deve produrre la corrispondente coppia `Symbol (LPAREN)` o `Symbol (RPAREN)` o `Symbol (EQUALS)`;
- infine se *s* è un simbolo alfabetico allora si deve passare ad un nuovo stato S che legge tutti i caratteri che seguono, fino ad un carattere che non può stare in una stringa (cioè né alfabetico né numerico) ed a quel punto dovrà “capire” se ha letto una keyword o una costante del Lisp Kit (`true`, `false` e `nil`) oppure un identificatore del programma. Ovviamente le possibilità sono mutuamente esclusive: non ci possono essere variabili che sono anche keyword o costanti del Lisp Kit. A seconda che si sia letta una keyword, o una costante o un identificatore, si produce il token corrispondente. Si osservi che una stringa come “falsetto” deve venire interpretata come la stringa “falsetto” e non la costante booleana “false” seguita da “tto”. Per avere questo secondo risultato sarebbe necessario scrivere “false tto”.
- se *s* è il carattere `'` che indica l'inizio di una stringa costante, allora si deve passare ad uno stato SC nel quale leggere tutti i caratteri che seguono fino al prossimo `'` che chiude la stringa costante. Sia *s* la stringa racchiusa dai 2 simboli `'`, allora il token da produrre è `String (s)`;
- Una volta che si è prodotto un token si ritorna allo stato I e si ricomincia.

Il fatto che si debba implementare un automa a stati finiti non deve essere preso troppo alla lettera. Gli stati possono venire realizzati con funzioni (generalmente ricorsive per ripetere le operazioni di scansione della lista di caratteri) e la transizione da I verso un altro stato X è realizzata semplicemente da un'invocazione della funzione che realizza lo stato X. Il ritorno da X a I è realizzato semplicemente con il return delle funzioni invocate. È importante osservare che diverse funzioni (S, SC ed N) *consumano* la lista di char e quindi quando ritornano ad I devono riportargli quello che è rimasto della lista di char, cioè la parte della lista che deve ancora venire esaminata.