

Linguaggi e Traduttori

simi

February 2025

Chapter 1

Introduzione

Inizio 8.45 fine 10.15

Quello che vedremo noi sono gli interpreti e in particolare riusciremo a fare un eseguibile usando un LLVM (che compila in una sorta di linguaggio macchina IR). Esiste una libreria (che ha fatto lui) che si chiama LibLet che permette di visualizzare l'esecuzione degli algoritmi che vedremo. Tutta questa cosa sarà scritta in Python senza focalizzarci sull'orientazione ad oggetti per vedere che queste cose non si possono fare solo ad oggetti. Vedremo anche un ripassone di algoritmi per vedere alcuni aspetti cruciali nel corso.

1.1 Python

Ci sono alcune strutture dati super comode che si possono usare con una sintassi comoda che sono le liste (strutture non omogenee), gli insiemi e i dizionari. Una cosa comoda di questi oggetti è che sono iterabili, in particolare è possibile costruire questi oggetti attraverso un meccanismo di comprehension in cui racchiudiamo tra i simboli sintattici la costruzione degli oggetti attraverso l'iterazione su un altro oggetto:

```
1 s = [x * x for x in range(10)]
```

Se ci metto una clausola dopo questa viene valutata:

```
1 # set comprehension: i numeri pari tra gli interi in [0, 9]
2 s = {x * x for x in range(9) if x % 2 == 0}
```

Si può imporre l'iterazione usando iter, notiamo che è comodo passare una sentinella per cui quando è finita l'iterazione ci ritorna la sentinella (tipicamente None):

```
1 #iterazione tramite iter/next
2
3 it = iter('alcune parole divise da spazi'.split())
4 while True:
5     w = next(it, None)
6     if w is None: break
7     print(w)
```

Le funzioni sono cittadini del primo ordine, possiamo assegnarle a variabili e passarle ad altre funzioni. Le useremo per implementare in modo economico i visitor, algoritmi ricorsivi che navigano le strutture dati in modo ricorsivo per fare delle cose. In più vediamo le dispatch table (un modo comodo per fare l'object oriented) e la memorizzazione tramite i decorator.

```
1 def quadra(x):
2     return x * x
3
4 def applica(fun, lst):
5     return [fun(x) for x in lst]
6
7 applica(quadra, [1, 2, 3])
```

Noi useremo molto le liste di liste, perchè ci rappresentano gli alberi e su queste possiamo definire visite.

```
1 lol = [1, [2, 3], [4, [5, 6]]]
2
```

```

3 #applichiamo una funzione scalare f a tutti gli elementi
4
5 def visit(f, lol):
6     for elem in lol:
7         if instance(elem, list):
8             visit(f, elem)
9         else:
10            f(elem)

```

1.1.1 Dispatch table

Cominciamo a vedere un piccolo esempio di parsing di un'espressione. La prima parte del parsing è suddividere la struttura lineare del testo in chunk concettuali che non è un lavoro banalissimo:

```

1  expr = "3 + 12 * 4 + 1 * 2"
2  tokens = iter(expr.split())

```

dopo di che dobbiamo definire la semantica delle operazioni, in qualche modo dobbiamo riassociare quello che osserviamo nel flusso dei token con le nostre interpretazioni. Una dispatch table è quella che associa delle informazioni a delle funzioni:

```

1 #semantica delle operazioni, tramite la dispatch table
2 def somma(x, y):
3     return x + y
4 def prodotto(x, y):
5     return x*y
6
7 DT = {
8     '+': somma,
9     '*': prodotto
10 }

```

A questo punto se voglio valutare l'espressione usiamo in modo iterativo la dispatch table, occhio che non stiamo rispettando le regole aritmetiche, ma associamo sempre a sinistra:

```

1 result = int(next(tokens))
2
3 while True:
4     t = next(tokens, None)
5     if t is None: break
6     of = DT[t]
7     result = of(result, int(next(tokens)))
8
9 result

```

1.1.2 Memorizzazione

Spesso e volentieri ci capiterà di esplorare algoritmi ricorsivi per cui per risolvere un problema con un'istanza grande risolveremo il problema su sue istanze più piccole, può accadere che nel processo di spezzamento andiamo a risolvere un sottoproblema che è già

stato risolto da qualcun altro, quindi se non adopero accorgimenti particolari ricalcolo gli stessi risultati, l'esempio tipico è il calcolo di Fibonacci.

L'idea è salvare in una cache i risultati parziali di una funzione,

```
1 def rendi_verbosa(f):
2     def f_verbosa(x):
3         result = f(x)
4         print(f'f({x}) = {result}')
5         return result
6     return f_verbosa
7
8 def quadrato(x):
9     x*x
10
11     quadrato_verboso = rendi_verbosa(quadrato)
12
13     q = quadrato_verboso(3)
14
15     #tenere da parte i risultati di una funzione
16     cache = {}
17
18     def memoize(f):
19         def f_memoized(x):
20             if x not in cache: cache[x] = f(x)
21             return cache[x]
22         return f_memoized
```

Esiste uno zucchero sintattico con cui possiamo annotare la funzione per memorizzare i risultati parziali @memoize

```
1 @memoize
2 def cubo(x):
3     return x ** 3
4
5     cache = {}
6
7     cubo(1), cubo(4), cubo(6)
8
9     cache
10
11 @memoize
12 def fib(n):
13     if n == 0 or n == 1: return 1
14     return fib(n - 1) + fib(n - 2)
```

1.2 Strutture dati

Le strutture dati che vedremo sono:

- alberi
- grafi
- pile/code

Gli alberi li useremo per rappresentare il testo che parseremo utilizzando una gerarchia del testo, e per visitare questa gerarchia utilizzeremo le visite. Un altro algoritmo che useremo è il backtracking che è molto utile per comprendere gli algoritmi di parsing, è una tecnica ricorsiva che tenta di risolvere i problemi con una sorta di forza bruta ma senza infilarsi nelle chiamate ricorsive anche inutili.

1.2.1 Alberi

Gli alberi sono la struttura dati più comune del corso perchè c'è una forte ricorrenza tra linguaggi/parsing e alberi, inoltre non è difficile immaginare che molti concetti che vedremo sono rappresentabili come alberi (l'espressione aritmetica dell'altra volta).

La forma più comune che useremo è la rappresentazione con liste di liste (lol), la prima è la radice e poi ci sono i sotto-alberi con tutti i loro figli.

```
1 # [radice]
2 # [radice alberi]
3
4 tree = [1, [11, [111]], [12, [121], [122]], [13]]
```

La cosa comoda di python è che c'è l'assegnamento ordinato, viene particolare comodo l'unpacking iterabile, se metto l'* prima della variabile mi prendo tutto quello che resta:

```
1 a, *b = 1, 2, 3
2 print(b) #[2, 3]
```

Questo è comodo perchè verrà comodo posso prendere i figli:

```
1 root, *children = tree
```

Come già detto useremo la libreria Liblet per visualizzare gli alberi partendo dalle lol (liste di liste). È possibile fare l'unpacking anche di questi alberi con la medesima sintassi.

```
1 from liblet import Tree, side_by_side
2
3 t = Tree.from_lol(tree)
4
5 root, *children = t
```

Visite

Quando facciamo una visita possiamo scegliere tante strategie, l'unica questione è decidere quando operare sul nodo, possiamo decidere di operare all'inizio (quando lo incontro) oppure occuparmi del nodo al termine della visita ai figli. Gli ordini vengono chiamati:

- pre-ordine: visito prima il nodo e poi i figli
- post-ordine: visito prima i figli e poi il nodo
- per livello

La visita in preordine si implementa in questo modo, posto che prende la funzione che fa i suoi conti:

```

1 def preorder(tree, root):
2     root, *children = tree
3     visitor(root)
4     for child in children: preorder(child, visitor)

```

Per fare il postordine scambio lo scarico ricorsivo con la funzione che fa i calcoli, ovviamente visualizzeremo nodi in ordine diverso.

```

1 def postorder(tree, visitor):
2     root, *children = tree
3     for child in children: preorder(child, visitor)
4     visitor(root)

```

Nel caso degli alberi è possibile visitarli anche per livelli, si chiama level order. È un algoritmo che non fa uso della ricorsione, generalmente si adopera una coda in cui accumulo i figli che vedo in modo tale che mi occupo per primo di quelli che metto dentro:

```

1 def levelorder(tree, visitor):
2     Q = queue()
3
4     Q.enqueue(tree)
5     while Q:
6         tree = Q.dequeue()
7         root, *children = tree
8         visitor(root)
9         for st in children: Q.enqueue(child)

```

Alberi con attributi

Una cosa che potrebbe essere molto comoda è che potremmo avere bisogno di usare nodi con informazioni più strutturate per poter visualizzare alberi arricchiti di attributi. Per fare questo useremo un albero che abbia dict come valori e che conservi il valore numerico come valore della chiave val:

```

1 def add_attr(tree):
2     root, *children = tree
3     return [{'val': root}] + [add_attr(child) for child in children]
4
5 tree = [1, [11, [111]], [1200, [121], [122]], [13]]
6
7 add_attr(tree)

```

Gli attributi generalmente vengono calcolati, ad esempio la profondità a cui siamo. Il modo in cui si calcolano segue due direzioni:

- Attributi ereditati: calcolati per il papà e ereditati dai figli. Si usa una visita in pre-ordine.
- Attributi sintetizzati: in qualche modo otteniamo prima le informazioni sui figli e poi sintetizziamo verso l'alto (tipicamente la valutazione delle espressioni è un attributo sintetizzato). Si usa una visita in post-ordine, perché prima devo visitare i figli.

Ad esempio vediamo come spingere in giù la profondità del nodo, (nella visita ricevo che sono figlio di un nodo profondo 4), quello che fa il visitor è aggiornare quella profondità di 1 per vedere quanto è profondo lui:

```

1 def preorder_with_value(tree, visitor, value = None):
2     root, *children = tree
3     visitor(root, value)
4     for child in children: preorder_with_value(child, visitor,
5         root['depth'])
6
7 # visitor che aggiunge l'attributo depth (pari a 1 + il valore
8 # ereditato, il caso None riguarda la radice)
9
10 def add_depth(root, value):
11     root['depth'] = value + 1 if value is not None else 0
12
13 attr_tree = add_attr(tree)
14
15 # la radice riceverà None perché è il valore di default di value
16
17 preorder_with_value(attr_tree, add_depth)
18
19 Tree.from_lol(attr_tree)

```

Nella sintesi il codice è più semplice da scrivere perché possiamo usare il return per ritornare il valore al padre:

```

1 def postorder_with_return(tree, visitor):
2     root, *children = tree
3     values = [postorder_with_return(child, visitor) for child in
4         children] # sarà la lista vuota se non ci sono figli
5     return visitor(root, values)
6
7 # visitor che aggiunge l'attributo max (pari al massimo tra il valore
8 # del nodo e quelli sintetizzati dai figli)
9
10 def add_max(root, values):
11     root['max'] = max([root['val']] + values)
12     return root['max']
13
14 attr_tree = add_attr(tree)
15
16 postorder_with_return(attr_tree, add_max)
17
18 Tree.from_lol(attr_tree)

```

1.2.2 Grafi

Li avremo i grafi come strumenti, meno concretamente ma avremo delle visite per passare da uno stato all'altro ma non scriveremo esplicitamente grafi. Esistono grafi non orientati e orientati. Idealmente noi useremo i grafi orientati. Tipicamente i grafi si rappresentano con:

- Matrice di adiacenza, $N \times N$ dove metto true e false dove c'è un collegamento. È $O(1)$ per il costo computazionale ma $O(n^2)$ nello spazio.

- Tengo tutte le coppie di nodi collegati se la matrice è sparsa (lista dei nodi)
- Tengo le liste di adiacenza di ogni nodo, per ogni nodo tengo in una lista tutti i nodi che ha collegati

Le rappresentazioni più comode che usiamo sono la rappresentazione con lista di archi (una tupla di tuple) oppure le liste di adiacenza rappresentata da un dict di set:

```

1  # dagli archi alla mappa delle adiacenze
2
3
4  # per ogni nodo n (sia s o t), adjacency[n] = set()
5
6  adjacency = dict()
7  for s, t in arcs:
8      adjacency[s] = set()
9      adjacency[t] = set()
10
11 # aggiungo gli outlink
12
13 for s, t in arcs: adjacency[s] |= {t}
14
15 adjacency

```

Algoritmi

La visita di un grafo è un procedimento che ci porta da un nodo ad esplorare tutti gli altri, ho in qualche modo bisogno di salvarmi da qualche parte i nodi che ho visitato e quelli che devo visitare. Immaginiamo di avere bisogno di una sacca in cui mettere le cose che abbiamo bisogno di fare (ci mettiamo dentro i figli dei nodi), dobbiamo stare attenti a non rimettere dentro quello che ci è già entrato e ha senso tirare fuori gli elementi in un certo ordine:

- FIFO: il primo nodo che metto è il primo che tolgo, i nodi sono messi in una coda. È la storia del level order. Così effettuiamo una visita in ampiezza.
- LIFO: last in first out, è una pila. Effettuiamo così una visita in profondità.

Nella visita in profondità si può dare una implementazione ricorsiva (perchè la pila può essere la pila delle chiamate), quello che dobbiamo avere è una struttura dati che ci tiene quello che abbiamo già visto, per evitare di visitare di nuovo gli stessi nodi:

```

1  def depthfirst(adjacency, start, visit):
2      def walk(src):
3          visit(src)
4          seen.add(src)
5          for dst in adjacency[src]:
6              if dst not in seen:
7                  walk(dst)
8      seen = set()
9      walk(start)

```

La visita in ampiezza usa la coda, toglia dalla coda e lo visita e quando lo visita lo aggiunge:

```

1 from liblet import Queue
2
3 def breadthfirst(adjacency, start, visit):
4
5     Q = Queue()
6
7     seen = set()
8     Q.enqueue(start)
9     while Q:
10         src = Q.dequeue()
11         visit(src)
12         seen.add(src)
13         for dst in adjacency[src]:
14             if dst not in seen:
15                 Q.enqueue(dst)

```

Quello da notare nelle visite è che quello che succede nelle visite in profondità è che la pila non aumenta tanto di dimensione, l'unico problema è che se devo vedere nodi vicini è possibile che per visitarlo potrei dover fare un sacco di conti prima di arrivare a visitare quello. La visita in ampiezza d'altro canto è vero che procede per livelli, ma se il grafo è denso la coda può diventare molto grande. Quindi quello che dovremo vedere negli algoritmi di parsing è questa scelta, scendo velocemente nel grafo sperando nel successo oppure è meglio scendere per passi? quello che vedremo è che in generale se non ci sono grandi strategie l'ampiezza è l'unica strada mentre se ho qualcuno che mi dice qualcosa è meglio andare giù veloce nella strada giusta.

1.2.3 backtracking

L'idea è che possiamo decidere che alcuni rami dell'albero sono buoni e altri sono cattivi. Quindi l'idea è che se mi trovo in una strada in cui posso prevedere che se sono in un certo punto in cui non c'è speranza che andando avanti nella ricorsione non troverò mai il risultato allora posso tagliare quel ramo dell'albero. Ho una specie di soluzione parziale che posso piazzare.

Ho una soluzione candidata e se sono già capace di sapere se quella soluzione non va bene allora ritorno, se ho trovato il risultato lo stampo. Altrimenti vado avanti e faccio le chiamate ricorsive.

```

1 def backtrack(candidate):
2     if reject(candidate): return
3     if accept(candidate): output(candidate)
4     s = first(candidate)
5     while s:
6         backtrack(s)
7         s = next(candidate)

```

Come esempio vediamo come si segmentano le parole, che è un problema fondamentale nei motori di ricerca. Inanzitutto mi procuro un elenco di parole, se segmenti non è vuota e l'ultima parola del dizionario non è in WORDS butto via, altrimenti se non mi rimane niente ho trovato tutti, altrimenti provo a spaccare in due la seguente parole in tutti i modi possibili:

```

1 from urllib.request import urlopen
2
3 # WORDS sono le parole di almeno 2 caratteri (3 conta anche l'a-capo)
4
5 with
6     urlopen('https://raw.githubusercontent.com/napolux/paroleitaliane/master/parolei
7         as url:
8     WORDS = {word.decode().strip().upper() for word in url if len(word)
9         >= 3}
10
11 print(len(WORDS))
12
13 def segmenta(segmenti, resto):
14     if segmenti and not segmenti[-1] in WORDS: return
15     if not resto:
16         print(segmenti)
17         return
18     for i in range(1, 1 + len(resto)):
19         segmenta(segmenti + [resto[:i]], resto[i:])
20
21 segmenta([], 'ILCORRIEREDELLASERAEDIZIONENOTTURNA')

```