

Soluzione dell'Esercitazione 1 del Corso di Elaborazione Linguaggio Naturale (seconda versione)

Luca Baronti*

25 marzo 2011

1 Espressioni Regolari

Le varie classi di stringhe riconosciute sono le seguenti:

[a-zA-Z]+ Tutte le parole all'interno di un testo separate da spazi. Escludendo i segni di punteggiatura, le lettere accentate e, in generale, qualsiasi simbolo che non sia le lettere upper-lower case

[A-Z][a-z]* Tutte le parole che cominciano con una lettera maiuscola, in quanto è necessario che sia presente innanzitutto un carattere nel range [A-Z] e successivamente zero o più caratteri in lowercase. Nel caso in cui ci sia una parola contenente due lettere uppercase, questa espressione effettuerà un doppio matching suddividendo la parola in due che hanno come iniziale il carattere uppercase

\d+(\.\d+)? Tutte le cifre intere o decimali (la cui parte decimale è separata da un punto)

(([bcdfghjklmnpqrstvwxyz][aeiou][bcdfghjklmnpqrstvwxyz])*) Tutte le triple costituite da una vocale in mezzo a due consonanti (tutto in lowercase)

\w+([\^\w\s]+) Tutte le parole (esclusi i segni di punteggiatura e gli spazi) oppure uno o più segni di punteggiatura.

Le espressioni regolari per far corrispondere le seguenti classi di stringhe sono:

1. Un singolo articolo determinato (“a”, “an”, “the”):

$\backslash b(a|n|the)\backslash b$

2. Un'espressione aritmetica che usa interi, somme e moltiplicazioni (2*3+8):
Nel caso non siano previsti spazi tra le cifre e gli operandi usiamo

*lbaronti@gmail.com

`[\d+|*|\+]`

Se invece vogliamo includere *anche* gli spazi nel riconoscimento possiamo usare

`\d+[\+| |*|\d+]*`

2 T9

Un esempio di codice che prende in ingresso un corpus ed una sequenza di caratteri digitati da una tastiera, e restituisce le **5** parole più probabili è il seguente

```
1 def T9(corpusPath, seq):
2     # definisco la tastiera
3     tastiera={
4         '1':(' '),
5         '2':('a','b','c'),
6         '3':('d','e','f'),
7         '4':('g','h','i'),
8         '5':('j','k','l'),
9         '6':('m','n','o'),
10        '7':('p','q','r','s'),
11        '8':('t','u','v'),
12        '9':('w','x','y','z')}
13    # recupero le informazioni del corpus
14    d={}
15
16    txtfile = open(corpusPath, "r")
17
18    for line in txtfile:
19        line=line.strip()
20        if line in d:
21            d[line]=d[line]+1
22        else:
23            d[line]=1.0
24    txtfile.close()
25
26    size=len(d)
27    # elimino le parole di lunghezza sbagliata
28    d=dict([(k,v) for k,v in d.items() if len(k)==len(seq)])
29    # scansiono la sequenza di numeri eseguendo i controlli con il dizionario
30    i=0
31    for c in seq:
32        if c=='1':
33            break
34
35        el=tastiera[c]
36        elcandidate=[]
37
38        d=dict([(k,v) for k,v in d.items() if k[i] in el])
39
40        i=i+1
41
42    # stampo il dizionario risultante
43    import operator
44    i=0
45    for e in sorted(d.items(), key=operator.itemgetter(1), reverse=True):
46        print e[0], e[1]/size
47        i=i+1
48        if i==5:
49            break
```

Naturalmente è possibile creare una versione equivalente che utilizza le espressioni regolari. Il seguente esempio costruisce una lista di pattern ordinata, ed esegue i match sui termini del dizionario, lettera per lettera.

```

1 def T9(corpusPath, seq):
2     import re
3     # creo la lista di pattern corrispondenti alla tastiera
4     tastiera=[
5         re.compile('\s'),
6         re.compile('[a-c]'),
7         re.compile('[d-f]'),
8         re.compile('[g-i]'),
9         re.compile('[j-l]'),
10        re.compile('[m-o]'),
11        re.compile('[p-s]'),
12        re.compile('[t-v]'),
13        re.compile('[w-z]')]
14
15    # recupero le informazioni del corpus
16    d={}
17
18    txtfile = open(corpusPath, "r")
19
20    for line in txtfile:
21        line=line.strip()
22        if line in d:
23            d[line]=d[line]+1
24        else:
25            d[line]=1.0
26    txtfile.close()
27
28    size=len(d)
29    i=0
30    # elimino le parole di lunghezza sbagliata
31    d=dict([(k,v) for k,v in d.items() if len(k)==len(seq)])
32
33    # seleziono il pattern corrispondente, ed eseguo il pruning sul dizionario
34    for num in seq:
35        currentPattern=tastiera[int(num)-1]
36        # pruning
37        d=dict([(k,v) for k,v in d.items() if currentPattern.match(k[i])])
38        i+=1
39
40    # stampo il dizionario risultante
41    import operator
42    i=0
43    for e in sorted(d.items(), key=operator.itemgetter(1), reverse=True):
44        print e[0], e[1]/size
45        i+=1
46        if i==5:
47            break

```

Per entrambe le versioni, la ricerca della parola inglese “*home*”, codificata in “4663” da luogo al seguente risultato

```

>>> from t9 import *
>>> T9("NPS_chat.txt", "4663")
good 0.0222606975019
home 0.00173138758348
gone 0.00148404650012
goof 0.000247341083354

```

dal quale si evince che la parola “*good*”, codificata allo stesso modo, è più frequente (di un ordine di grandezza) rispetto a quella voluta, all’interno del corpus scelto.

3 Legge di Zipf

Come testo di riferimento si è usata la fonte di wikipedia su un argomento molto comune (l’acqua), dalla quale sono stati rimosse eventuali parole *outlier* come i

riferimenti o i pulsanti di modifica.

Si è utilizzata la seguente funzione per calcolarci le frequenze:

```
1 def p4(text):
2     txtfile = open(text, "r")
3
4     d={}
5     # estrapolo le parole con la loro frequenza
6     for line in txtfile.readlines():
7         line=line.strip()
8         for word in line.split():
9             if word in d:
10                d[word]=d[word]+1
11            else:
12                d[word]=1.0
13        size=len(d)
14        # normalizzo le frequenze
15        d=dict([(k,v/size) for k,v in d.items()])
16
17    import operator
18    ordered=sorted(d.items(), key=operator.itemgetter(1), reverse=True)
19    i=0
20    outfile = open("acquaFreq.dat", "w")
21    outfile2=open("costante.dat", "w")
22    outfile.write("id_parola_frequenza\n")
23    outfile2.write("id_parola_k\n")
24    for e in ordered:
25        i=i+1
26        outfile.write(str(i)+"_"+str(e[0])+"_"+str(e[1])+"\n")
27        outfile2.write(str(i)+"_"+str(e[0])+"_"+str(e[1]*i)+"\n")
28        if i==500: break
29
30    outfile.close()
31    outfile2.close()
```

Come era facile aspettarsi le parole più comuni sono state gli articoli e le preposizioni:

id	Parola	Frequenza (normalizzata)
1	di	0.126698262243
2	e	0.0581358609795
3	la	0.0429699842022
4	in	0.0382306477093
5	è	0.0379146919431
6	che	0.0366508688784
7	il	0.0309636650869
8	a	0.0296998420221
9	un	0.0284360189573
10	della	0.0249605055292
11	le	0.0243285939968
12	per	0.0233807266983

Una delle parole più significative del corpus analizzato, ovvero “acqua” si posiziona invece 22-esima, con una frequenza pari a ≈ 0.013902 .

Dall’andamento generale delle 500 parole più frequenti, osservabile in Fig.1, emerge come il valore della frequenza delle parole tenda a diminuire molto nelle prime posizioni, mentre per quanto riguarda l’eventuale riscontro della legge di Zipf, come si evince in Fig.2 le parole non sembrano seguire un andamento perfettamente costante rispetto al loro rank, tuttavia è possibile notare che il

loro prodotto oscilla attorno a 0.25 e a 0.5, e questo intervallo potrebbe rientrare entro i parametri di incertezza previsti dalla suddetta legge.

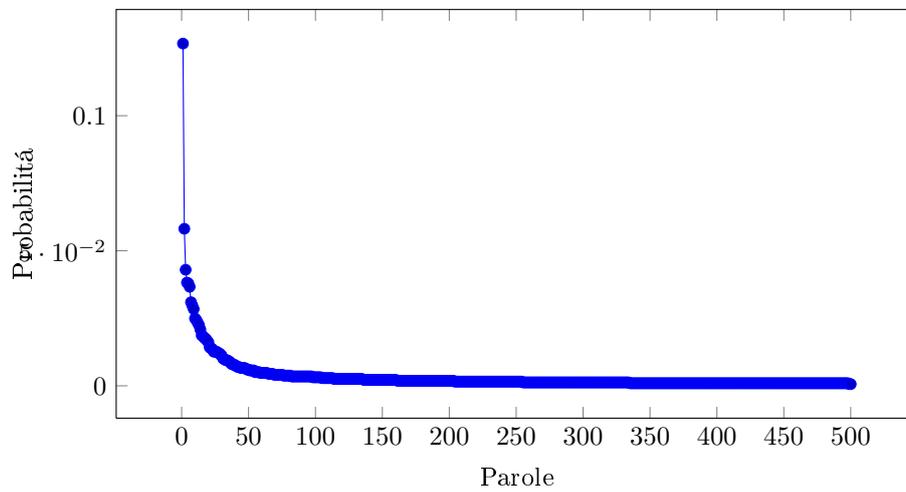


Figura 1: Frequenza delle parole in relazione al loro rank

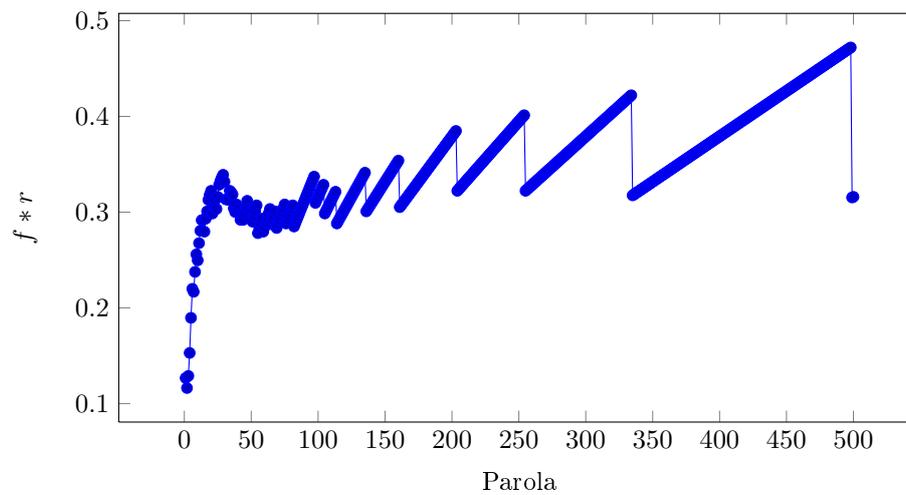


Figura 2: Prodotto delle frequenze delle parole con il loro rank