

Valoració d'eficiència de diferents algorismes de cerca

Miquel Perello Nieto

7 de març de 2012

Resum

La finalitat d'aquesta pràctica es comprovar si avui en dia continua essent vàlida la gràfica dels algorismes mes eficients presentada a classe. Per tal motiu la feina pel proper dia serà avaluar, per a tres alfabet de mida $m-1, m$ i $m+1$ caràcters (en aquest cas $m = 3$), els algorismes mes eficients (en temps d'execució) per a patrons de longitud k ($1 \leq k \leq \infty$). Els algorismes que tindrem en compte són els de Força bruta, Horspool, BNDM, BOM (tota la informació dels diferents algorismes es pot veure a [1]).

A part es podrà observar en la pràctica els resultats de l'algoritme Karp-Rabin, el qual només tindrà caràcter informatiu, sense entrar en el calcul de resultats.

1 Introducció

Durant aquesta pràctica voldrem veure si els resultats esperats amb la gràfica 1 es correspon amb l'execució pràctica en un ordinador amb processador de 64 bits. Però primerament veurem els diferents algorismes que utilitzarem en la pràctica, per veure quins son els seus costos teòrics, quines característiques te i una implementació de cadascun en C.

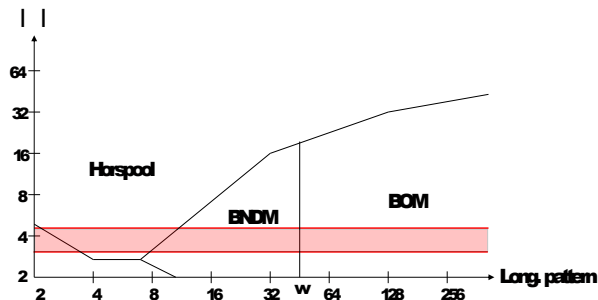


Figura 1: Gràfica de resultats que volen validar, amb marca en el rang dels nostres alfabet

1.1 Bruteforce

L'algoritme de força bruta consisteix en comprovar per totes les posicions del text entre 0 i $n-m$, si el principi del patró coincideix o no amb el text. Després de cada comprovació es shifta el patró una posició a la dreta.

Aquest algoritme no requereix cap fase de preprocés, i un espai extra constant adicionalment al patró i al text.

Durant la fase de cerca la comparació dels caràcters pot ser en qualsevol ordre. La complexitat temporal d'aquesta fase es $O(mn)$. I el nombre de comparacions esperat es de $2n$.

1.1.1 Característiques

- No hi ha fase de preprocés;
- Espai extra necessari constant;
- Sempre mou la finestra una posició a la dreta;
- La comparació es pot efectuar en qualsevol ordre;
- Fase de cerca en temps $O(mn)$;
- Comparacions de text esperat de $2n$;

1.1.2 Codi C

Codi C 1: codi Bruteforce

```
1 void BF(char *x, int m, char *y,  
2         int n) {  
3     int i, j;  
4     /* Searching */
```

```

5     for (j = 0; j <= n - m; ++j)
6     {
7         for (i = 0; i < m && x[i]
8             == y[i + j]; ++i);
9         if (i >= m)
10            OUTPUT(j);
11    }
12 }

```

```

9     /* Searching */
10    j = 0;
11    while (j <= n - m) {
12        c = y[j + m - 1];
13        if (x[m - 1] == c && memcmp(x,
14            y + j, m - 1) == 0)
15            OUTPUT(j);
16        j += bmBc[c];
17    }

```

1.2 Horspool

El shift per caràcter incorrecte de l'algoritme de Boyer-Moore no és molt eficient per alfabetos petits, però quan els alfabetos són grans comparats amb la longitud del patró és comporta molt bé.

Utilitzant-lo sol es produeix un algoritme molt eficient en la pràctica. El Horspool proposa utilitzar només el shift de caràcter incorrecte del caràcter més la la dreta de la finestra per calcular els shifts de l'algoritme de Boyer-Moore.

La fase de preprocés es en complexitat temporal $O(m+)$ i espacial $O()$

La fase de cerca té el pitjor cas quadràtic però es pot demostrar que el nombre mig de comparacions en un text es entre $1/$ i $2/(+1)$.

1.2.1 Característiques

- Simplificació de l'algoritme Boyer-Moore;
- Fàcil de implementar;
- Temps de preprocés $O(m+)$ i espai $O()$;
- Temps de fase de cerca en $O(mn)$;
- El nombre de comparacions esperat per caràcter és de $1/$ i $2/(+1)$.

1.2.2 Codi C

Codi C 2: codi Horspool

```

1 void HORSPOOL(char *x, int m,
2     char *y, int n)
3 {
4     int j, bmBc[ASIZE];
5     char c;
6     /* Preprocessing */
7     preBmBc(x, m, bmBc);
8

```

1.3 BNDM

L'algoritme BNDM (Backward Nondeterministic Dawg Matching) utilitza una taula B, la qual per cada caràcter guarda una mascara de bits. La màscara del caràcter es crea si i només si $xi = c$.

La fase de cerca es manté en una paraula $d = dm - 1 \dots d0$, on la mida del patró m es inferior o igual a la mida de paraula de l'ordinador.

El bit di a la iteració k es posa a 1 si i només si $x[m - i..m - 1 - i + k] = y[j + m - k..j + m - 1]$. A la primera iteració "d" s'inicialitza a $1m - 1$. La formula per actualitzar "d" segueix $d' = (d \& B[yj]) \ll 1$.

Existeix una coincidència si i només si després de la iteració m el bit $dm - 1 = 1$.

Quan apareix una coincidència el prefix més llarg vist fins la coincidència determina el salt de finestra per la nova posició inicial.

1.3.1 Característiques

- Variant de l'algoritme Reverse Factor;
- Utilitza la simulació paral·lela de bit;
- Eficient si la mida del patró es inferior a la mida de paraula de l'ordinador.

1.3.2 Codi C

Codi C 3: codi BNDM

```

1 void BNDM(char *x, int m, char *
2     y, int n) {
3     int B[ASIZE];
4     int i, j, s, d, last;
5     if (m > WORD_SIZE)
6         error("BNDM");
7     /* Pre processing */

```

```

8  memset(B,0,ASIZE*sizeof(int));
9  s=1;
10 for (i=m-1; i>=0; i--){
11 B[x[i]] |= s;
12 s <<= 1;
13 }
14
15 /* Searching phase */
16 j=0;
17 while (j <= n-m){
18 i=m-1; last=m;
19 d = ~0;
20 while (i>=0 && d!=0) {
21     d &= B[y[j+i]];
22     i--;
23     if (d != 0){
24 if (i >= 0)
25     last = i+1;
26 else
27     OUTPUT(j);
28     }
29     d <<= 1;
30 }
31 j += last;
32 }
33 }

```

1.4 BOM

Els algorismes de tipus Boyer-Moore troben alguns sufixes del patró, però és possible trobar alguns prefixes del patró comprovant els caràcters de la finestra de dreta a esquerra i determinar la llargada del salt. Això és possible utilitzant el sufixe d'Oracle en el patró invertit. Aquesta estructura de dades es un Automat molt compacte que es capaç de reconèixer al menys tots els sufixes del patró i algunes altres paraules. La cerca de patrons amb l'algoritme utilitzant el patró invertit amb oracle es diu Backward Oracle Matching algorithm (les inicials del qual donen nom a aquest algoritme).

El sufix oracle d'una paraula "w" es un Automat Finit Determinista $O(w) = (Q, q_0, T, E)$.

El llenguatge acceptat per $O(w)$ es tal que $\{u \text{ en } * : \text{ existeix } v \text{ en } * \text{ tal que } w = vu\}$ en $L(O(w))$.

La fase de preprocés del Backward Oracle Matching algorithm consisteix en calcular el sufix d'oracle per l'invers del patró xR. A pesar del fet que aquest es capaç de reconèixer paraules que no son del patró, el sufix d'oracle pot ser utilitzat per fer

una comparació de paraules tenint en compte que la única paraula de mida major o igual a m és reconeguda per el patró invers de oracle.

Durant la fase de cerca l'algoritme Backward Oracle Matching busca els caràcters de la finestra de dreta a esquerra amb l'autòmat $O(xR)$ començant a l'estat q_0 . Aquest segueix fins que no hi ha una transició definida per el caràcter actual. En aquest moment la mida del prefix més llarg del patró, el qual és un sufix de la part comprovada es inferior que la mida del camí agafat de $O(xR)$ des de l'estat q_0 i l'últim estat trobat. Sabent aquesta mida, es trivial calcular la mida del salt a executar.

L'algoritme Backward Oracle Matching té un cost quadràtic respecte el temps en el pitjor dels casos, però és optim en la mitjana. En aquesta el seu cost esdevé $O(n \cdot (\log m) / m)$.

1.4.1 Característiques

- Versió de l'algoritme *Reverse Factor* utilitzant el sufix d'oracle de xR en comptes de l'automat de xR.
- Ràpid en entorns amb mides de patró molt llargs i petits alfabetes;
- Fase de preprocés en complexitat temporal i espacial $O(m)$;
- Fase de cerca en complexitat temporal $O(mn)$;
- Optim en la mitja;

1.4.2 Codi C

Només la part externa de les transicions d'oracle s'emmagatzema en una llista de enllaços (un per estat). Els noms de les transicions no s'emmagatzemen però són computades per la paraula x.

Codi C 4: codi BOM

```

1 #define FALSE      0
2 #define TRUE       1
3
4 int getTransition(char *x, int p
5     , List L[], char c) {
6     List cell;
7
8     if (p > 0 && x[p - 1] == c)
9         return(p - 1);
10    else {

```

```

10     cell = L[p];
11     while (cell != NULL)
12         if (x[cell->element] ==
13             c)
14             return(cell->element
15                 );
16         else
17             cell = cell->next;
18     return(UNDEFINED);
19 }
20
21 void setTransition(int p, int q,
22     List L[]) {
23     List cell;
24     cell = (List)malloc(sizeof(
25         struct _cell));
26     if (cell == NULL)
27         error("BOM/setTransition")
28         ;
29     cell->element = q;
30     cell->next = L[p];
31     L[p] = cell;
32 }
33 void oracle(char *x, int m, char
34     T[], List L[]) {
35     int i, p, q;
36     int S[XSIZE + 1];
37     char c;
38     S[m] = m + 1;
39     for (i = m; i > 0; --i) {
40         c = x[i - 1];
41         p = S[i];
42         while (p <= m &&
43             (q = getTransition(
44                 x, p, L, c)) ==
45                 UNDEFINED) {
46             setTransition(p, i - 1,
47                 L);
48             p = S[p];
49         }
50         S[i - 1] = (p == m + 1 ? m
51             : q);
52     }
53     p = 0;
54 }
55
56 while (p <= m) {
57     T[p] = TRUE;
58     p = S[p];
59 }
60
61 void BOM(char *x, int m, char *y
62     , int n) {
63     char T[XSIZE + 1];
64     List L[XSIZE + 1];
65     int i, j, p, period, q, shift
66     ;
67
68     /* Preprocessing */
69     memset(L, NULL, (m + 1)*
70         sizeof(List));
71     memset(T, FALSE, (m + 1)*
72         sizeof(char));
73     oracle(x, m, T, L);
74
75     /* Searching */
76     j = 0;
77     while (j <= n - m) {
78         i = m - 1;
79         p = m;
80         shift = m;
81         while (i + j >= 0 &&
82             (q = getTransition(
83                 x, p, L, y[i + j
84                 ])) !=
85                 UNDEFINED) {
86             p = q;
87             if (T[p] == TRUE) {
88                 period = shift;
89                 shift = i;
90             }
91             --i;
92         }
93         if (i < 0) {
94             OUTPUT(j);
95             shift = period;
96         }
97         j += shift;
98     }
99 }

```

1.5 Karp-Rabin

Hashejar ens permet de manera simple evitar un nombre de comparacions de caràcter quadràtic en la majoria de situacions. En comptes de comprovar en cada posició del text si el patró coincideix, sembla més eficient comprovar només el contingut de les finestres que més *s'assemblen* al patró. Per comprovar la semblança de les dues paraules s'utilitza una funció de hash.

La fase de preprocés de l'algoritme Karp-Rabin consisteix en calcular el $hash(x)$. Aquest es pot resoldre amb un espai constant i un temps $O(m)$.

Durant la fase de cerca es suficient amb comparar el $hash(x)$ amb el $hash(y[j..j+m-1])$ for $0 \leq j < n-m$. Si es troba alguna semblança es necessari comprovar la igualtat de $x = y[j..j+m-1]$ caràcter per caràcter.

La complexitat temporal de la fase de cerca es $O(mn)$. El nombre esperat de comparacions per caracter es de $O(n+m)$.

1.5.1 Característiques

- Utilitza una funció de hash.
- Fase de preprocés en temps $O(m)$ i espai constant.
- Fase de cerca en temps $O(mn)$.
- Temps d'execució esperat $O(n+m)$.

1.5.2 Codi C

```
Codi C 5: codi Karp-Rabin
1 #define REHASH(a, b, h) (((h) -
   (a)*d) << 1) + (b))
2
3 void KR(char *x, int m, char *y,
   int n) {
4     int d, hx, hy, i, j;
5
6     /* Preprocessing */
7     /* computes d = 2^(m-1) with
8        the left-shift operator */
9     for (d = i = 1; i < m; ++i)
10         d = (d<<1);
11
12     for (hy = hx = i = 0; i < m;
        ++i) {
```

```
13         hx = ((hx<<1) + x[i]);
14         hy = ((hy<<1) + y[i]);
15     }
16
17     /* Searching */
18     j = 0;
19     while (j <= n-m) {
20         if (hx == hy && memcmp(x,
21             y + j, m) == 0)
22             OUTPUT(j);
23         hy = REHASH(y[j], y[j + m
24             ], hy);
25         ++j;
26     }
```

2 Resultats

En aquesta secció veurem els resultats dels temps de cerca de cada un dels diferents algoritmes a analitzar. Es podrà veure de forma visual mitjançant plots generals. I en els casos en els que no es pugui apreciar la diferència de millora entre un algoritme i l'altre, es podrà veure un cas d'execució particular per fer desempatar els dos algoritmes, i una gràfica en el que sí es podrà observar diferència.

Finalment per cada una de les mides d'alfabet es presenta una petita gràfica lineal amb els rangs de mida de patró ideals per a cada algoritme.

2.1 Alfabet de mida 3

Per comprovar els resultats en l'execució dels diferents algoritmes en la cerca de patrons de l'alfabet de mida 3 s'ha realitzat primerament unes proves d'execució amb les següents característiques:

mida del text	1,4GB
mida patrons	$[2^0 \dots 2^9]$
repeticions per patró	4

Els resultats d'aquestes execucions donen una visió molt clara en la distribució dels diferents algoritmes (figura 2), però existeix un punt de conflicte entre el Horspool i el BNDM en el qual no es pot apreciar amb les repeticions efectuades el punt a partir del qual un supera a l'altre.

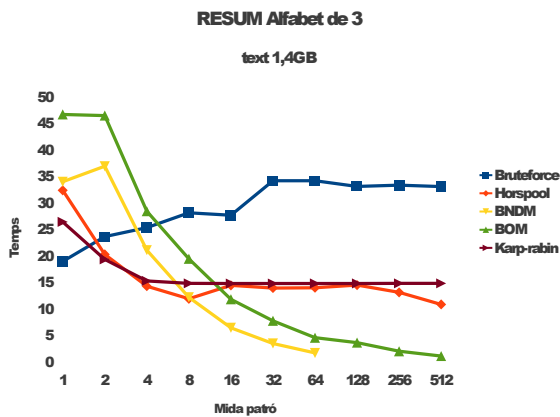


Figura 2: Grafica del temps per mida de patró per diferents algorismes amb un alfabet de mida 3

Per aquesta raó s'ha hagut d'efectuar una prova més concreta i focalitzada en el punt de conflicte, amb més repeticions per mida de patró, així que les següents proves han seguit el següent patró:

mida del text	2GB
mida patrons	[5...9]
repeticions per patró	10

Per tant després de l'execució d'aquestes proves es pot veure clarament el punt en el que el BNDM supera en temps d'execució a l'algoritme de Horspool (figura 2).

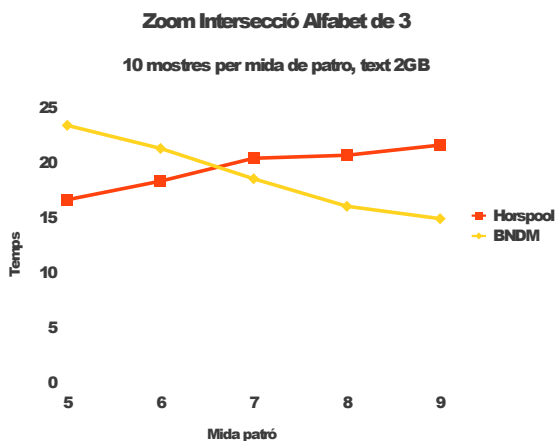


Figura 3: Zoom de conflicte entre Horspool i BNDM

Gracies a aquesta aproximació podem apreciar que amb patrons de mida sis l'algoritme Horspool segueix sent millor, en canvi en patrons de mida set el millor algoritme es millor el BNDM.

Tot seguit es pot veure un resum visual en el que es poden apreciar els intervals de mida de patró amb l'algoritme que hauríem d'utilitzar per trigar el mínim temps en analitzar el text (figura 4).

Alfabet de cardinalitat 3

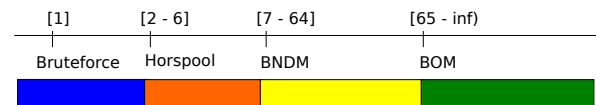


Figura 4: Valoració per mida del patró per un alfabet de cardinalitat 3

2.2 Alfabet de mida 4

Per comprovar els resultats en l'execució dels diferents algorismes en la cerca de patrons de l'alfabet de mida 4 s'ha realitzat primerament unes proves d'execució amb les següents característiques:

mida del text	1,4GB
mida patrons	[2 ⁰ ...2 ⁹]
repeticions per patró	4

Els resultats d'aquestes execucions donen una visió molt clara en la distribució dels diferents algorismes (figura 5), però igual que en el cas de l'alfabet de mida tres, en l'alfabet de mida quatre també existeix un punt de conflicte entre el Horspool i el BNDM en el qual no es pot apreciar amb les repeticions efectuades el punt a partir del qual un supera a l'altre.

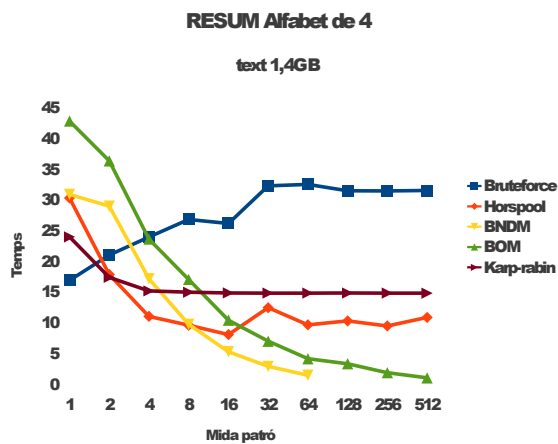


Figura 5: Grafica del temps per mida de patró per diferents algorismes amb un alfabet de mida 4

Per aquesta raó s'ha hagut d'efectuar una prova més concreta i focalitzada en el punt de conflicte, amb més repeticions per mida de patró, així que les següents proves han seguit el següent patró:

mida del text	2GB
mida patrons	[5...9]
repeticions per patró	10

Per tant després de l'execució d'aquestes proves es pot veure clarament el punt en el que el BNDM supera en temps d'execució a l'algoritme de Horspool (figura 5).

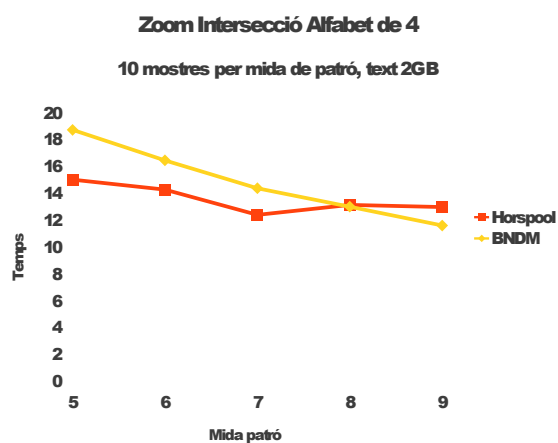


Figura 6: Zoom de conflicte entre Horspool i BNDM

Gracies a aquesta aproximació podem apreciar que amb patrons de mida sis l'algoritme Horspool segueix sent millor, en canvi en patrons de mida set el millor algoritme es millor el BNDM.

Tot seguit es pot veure un resum visual en el que es poden apreciar els intervals de mida de patró amb l'algoritme que hauríem d'utilitzar per trigar el mínim temps en analitzar el text (figura 7).

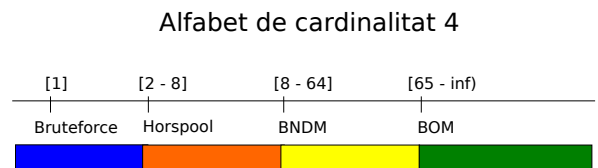


Figura 7: Valoració per mida del patró per un alfabet de cardinalitat 4

2.3 Alfabet de mida 5

Per comprovar els resultats en l'execució dels diferents algorismes en la cerca de patrons de l'alfabet de mida 5 s'ha realitzat primerament unes proves d'execució amb les següents característiques:

mida del text	1,4GB
mida patrons	[2 ⁰ ...2 ⁹]
repeticions per patró	4

Els resultats d'aquestes execucions donen una visió molt clara en la distribució dels diferents algorismes (figura 8), però igual que en el cas de l'alfabet de mida tres i quatre, en l'alfabet de mida cinc també existeix un punt de conflicte entre el Horspool i el BNDM en el qual no es pot apreciar amb les repeticions efectuades el punt a partir del qual un supera a l'altre.

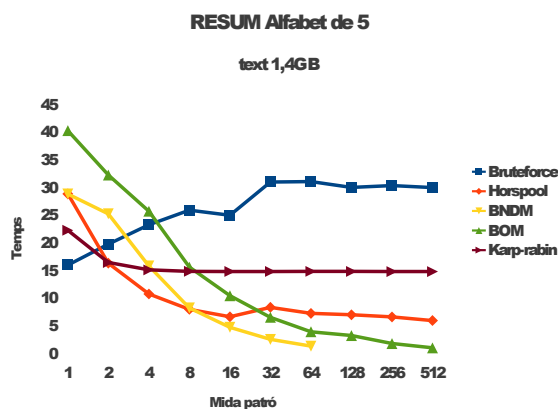


Figura 8: Grafica del temps per mida de patró per diferents algoritmes amb un alfabet de mida 5

Per aquesta raó s'ha hagut d'efectuar una prova més concreta i focalitzada en el punt de conflicte, amb més repeticions per mida de patró, així que les següents proves han seguit el següent patró:

mida del text	2GB
mida patrons	[5...13]
repeticions per patró	20

Per tant després de l'execució d'aquestes proves es pot veure clarament el punt en el que el BNDM supera en temps d'execució a l'algoritme de Horspool (figura 8).

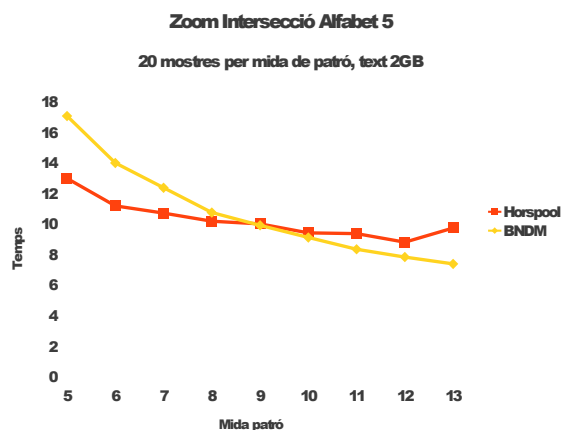


Figura 9: Zoom de conflicte entre Horspool i BNDM

Gracies a aquesta aproximació podem apreciar que amb patrons de mida sis l'algoritme Horspool segueix sent millor, en canvi en patrons de mida set el millor algoritme es millor el BNDM.

Tot seguit es pot veure un resum visual en el que es poden apreciar els intervals de mida de patró amb l'algoritme que hauríem d'utilitzar per trigar el mínim temps en analitzar el text (figura 10).

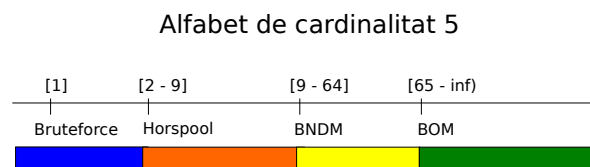


Figura 10: Valoració per mida del patró per un alfabet de cardinalitat 5

3 Conclusions

Un cop finalitzades totes les proves ja hem pogut contrastar els resultats, i gràcies a les gràfiques generades es pot observar de manera senzilla si els resultats han estat propers als esperats, així que detallem el que podem comprovar amb cada punt d'intersecció de la gràfica original. Ens centrarem exclusivament en la línia vermella que creua la gràfica 1 amb mida d'alfabet compresa entre tres i cinc.

En la gràfica es podia observar que amb alfabet d'aquest rang existeix una petita àrea en els patrons de longitud inferior a tres en la qual el Horspool es comporta pitjor que altres algoritmes, en el nostre cas es compleix amb l'algoritme de força bruta, justament amb patrons de un caràcter.

Després es pot observar un àrea que pertany al Horspool que arriba fins als patrons de mida nou. Aquí es torna a veure que la similitud amb els nostres resultats es molt alta.

Tot seguit el BNDM es converteix en el millor algoritme fins arribar a la mida de patró igual a la mida de la paraula de l'ordinador; en el nostre cas utilitzàvem un ordinador amb processador de 64 bits, i per aquest motiu a partir d'aquesta mida de patró hem descartat els següents resultats.

Per ultim amb patrons de mida superior l'algoritme BOM es el més eficient tant en la gràfica anterior, com en els resultats d'aquesta pràctica.

Per tant amb la realització d'aquesta pràctica finalment hem pogut observar que els resultats de l'execució dels diferents algoritmes de cerca segueixen la relació que existia en la gràfica que volíem comprovar en els casos d'alfabets de mida tres, quatre i cinc.

A Codis

Codi C 6: Generació de text aleatori

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include <fcntl.h>
5 #include <sys/types.h>
6
7 #define EOS '\0'
8

```



```

9 char alfabet[7] = {'T', 'G', 'A',
10   , 'C', 'B', 'D', 'E'};
11 int main(int argc, char *argv[])
12 {
13     char filename[256];
14     int size;
15     int alf_size;
16     int fd;
17     unsigned long i;
18
19     srand ( time(NULL) );
20
21     size = atoi(argv[1]);
22     alf_size = atoi(argv[2]);
23
24     for (i = 0; i < size*1000000;
25         i++)
26     {
27         printf("%c",alfabet[(rand()
28             % alf_size)]);
29         if ((time(NULL)%(rand()+1))
30             == 0)
31         {
32             srand ( time(NULL) );
33         }
34     }
35     printf(EOS);
36 }
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

Codi C 7: Main de cerca en text

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include "list.h"
5 #define EOS '\0'
6 #define ASIZE 256
7 #define FALSE 0
8 #define TRUE 1
9
10 char alfabet[7] = {'T', 'G', 'A',
11   , 'C', 'B', 'D', 'E'};
12 unsigned long int readfile(char
13     *filename, char *text[])
14 {
15     FILE *fin;
16     unsigned long int len;

```

```

52
53     if (argc != 4) exit(0);
54
55     alf_size = atoi(argv[2]);
56     sprintf(filename, "%s", argv[1])
57         ;
58     n = readfile(filename, &text);
59     printf("n=%d\n", n);
60     printf("size \t pattern \t
61         trobats \t segons \n");
62     srand(0);
63     for (k = 1; k <= 512; k = k*2)
64     {
65         pattern_size = k;
66         for (i = 0; i < 4; i++)
67         {
68             for (j = 0; j <
69                 pattern_size; j++)
70             {
71                 pattern[j] = alfabet[
72                     rand()%alf_size];
73             }
74             pattern[pattern_size] =
75                 EOS;
76             //printf("%s\t", pattern);
77
78             initemps=clock();
79             //trobats = BF(pattern,
80                 strlen(pattern), text, n)
81                 ;
82             trobats = HORSPOOL(pattern
83                 , strlen(pattern), text, n
84                 );
85             //trobats = BNDM(pattern,
86                 strlen(pattern), text, n)
87                 ;
88             //XSIZE = pattern_size;
89             //trobats = BOM(pattern,
90                 strlen(pattern), text, n)
91                 ;
92             //trobats = KR(pattern,
93                 strlen(pattern), text, n)
94                 ;
95             endtemps=clock();
96             printf( "%d\t%s\t%d\t%f\n"
97                 , pattern_size, pattern
98                 , trobats, (endtemps -
99                     initemps)/(double)
100                     CLOCKS_PER_SEC );

```

Referències

- [1] "EXACT STRING MATCHING ALGORITHMS", Christian Charras - Thierry Lecroq, <http://www-igm.univ-mlv.fr/~lecroq/string/index.html>, Faculté des Sciences et des Techniques