

Compressor i descompressor d'altures de vol

Enric Carpintero Rico

Abril 2025

1 Introducció

Aquest document detalla la implementació del compressor i descompressor de fitxers d'altures de vol demanat per l'assignatura de Compresió de Dades i Imatges de la Facultat d'Informàtica de Barcelona, quadrimestre de primavera 2025.

A excepció de la classe auxiliar `BinTree.hh`, implementada pel professorat del departament de CS de la UPC, tot el codi de l'entrega ha estat desenvolupat per mi.

Per als fitxers de prova proporcionats inicialment a Atenea, `file11111.txt` i `file22222.txt` la ràtio de compressió obtinguda ha estat de 5,78 i 7,97, respectivament. La mida del programa compilat resultant és de 85 kB, mentre que el codi font sense compilar, excloent les llibreries de sistema i les llibreries estàndard, ocupa 31 kB. El temps de compressió pels dos fitxers de prova, als ordinadors de la facultat, és de 2,1 i 1,7 segons, respectivament. El temps de descompressió és 6,7 i 5,7 segons, respectivament.

2 Algorisme

Aquest algorisme està dissenyat únicament per al tipus de fitxers indicats per l'exercici: Cadenes de text pla en ASCII, representant seqüències de nombres naturals, o el nombre -1, separats entre ells per un únic espai o salt de línia. Per tant, aquest algorisme no contempla qualsevol altre tipus de fitxer o dades que no compleixi amb aquestes característiques, i per tant no es pot fer servir com a compressor i descompressor de fitxers de forma general.

L'algorisme de compressió està basat en:

- El fet que tots els elements dels fitxers a comprimir corresponen a nombres naturals.
- Increments/diferències entre elements consecutius.
- Codis de Huffman.

```

-1 -1 -1 -1 59359 59337 59316 59285 59242 59258 59257 59209 59185
-1 -1 -1 -1 59360 59346 59330 59309 59293 59298 59294 59255 59224
-1 -1 -1 -1 59368 59352 59339 59326 59321 59320 59310 59294 59255
-1 -1 -1 -1 59380 59368 59353 59341 59332 59324 59315 59303 59275
-1 -1 -1 -1 59399 59384 59360 59342 59332 59332 59314 59308 59286
-1 -1 -1 -1 59419 59382 59352 59339 59322 59314 59314 59303 59285
-1 -1 -1 -1 59401 59374 59350 59331 59315 59299 59289 59287 59283
-1 -1 -1 -1 59426 59387 59385 59356 59327 59315 59301 59284 59276
-1 -1 -1 -1 59438 59414 59420 59378 59347 59329 59313 59298 59277
-1 -1 -1 -1 59468 59461 59425 59387 59369 59348 59329 59305 59282
-1 -1 -1 -1 59539 59487 59422 59402 59394 59364 59337 59312 59287
-1 -1 -1 -1 59573 59501 59480 59442 59415 59380 59344 59309 59287
-1 -1 -1 -1 59617 59578 59525 59462 59421 59386 59354 59324 59291
-1 -1 -1 -1 59687 59618 59544 59474 59421 59373 59342 59319 59305
-1 -1 -1 -1 59709 59649 59582 59496 59433 59383 59357 59343 59335
-1 -1 -1 -1 59746 59680 59615 59549 59481 59431 59390 59368 59366
-1 -1 -1 -1 59806 59764 59711 59642 59587 59486 59431 59398 59396
-1 -1 -1 -1 59828 59795 59755 59704 59625 59514 59463 59453 59436
-1 -1 -1 -1 59864 59823 59778 59736 59698 59594 59528 59502 59463
-1 -1 -1 -1 59908 59841 59775 59742 59695 59643 59585 59541 59474
-1 -1 -1 -1 59939 59860 59780 59727 59693 59661 59634 59552 59475
-1 -1 -1 -1 59953 59849 59771 59720 59684 59655 59620 59550 59492
-1 -1 -1 -1 -1 59911 59828 59762 59710 59670 59638 59610 59523 595
-1 -1 -1 -1 -1 59872 59801 59739 59671 59639 59614 59558 59483 595

```

Figure 1: Exemple de seqüència de dades que comprimeix l'algorisme.

En primer pas, es calculen les diferències entre els nombres naturals consecutius, les seqüències de nombres -1 seguits, o la presència dels salts de línia.

A cada ocurrència, se li assigna un codi de longitud de bits variable indicant quin element o elements hem trobat, seguit d'un paràmetre, també de longitud variable, complementant el primer codi.

Finalment, aquestes ocurrències es divideixen en blocs de mida variable, i es codifiquen amb codis de Huffman.

2.1 Codis per diferència sobre el fitxer base

Al text pla original, podem trobar:

Nombres naturals: Per a cada natural que trobem, calculem la diferència de valor respecte a l'anterior nombre natural trobat. Segons el valor absolut d'aquesta diferència, el codifiquem de diferents maneres:

- **Els nombres son iguals:** Codifiquem amb 0b10.
- **La diferència és 1..16:** Codifiquem amb 0b00 seguit d'un bit pel signe (0: nou nombre és més gran. 1: nou nombre és més petit), seguit per la codificació en binari sense signe de restar 1 al valor absolut de la diferència, amb 4 bits.
- **La diferència és 17..272:** Codifiquem amb 0b010 seguit d'un bit pel signe (0: nou nombre és més gran. 1: nou nombre és més petit), seguit per la codificació en binari sense signe de restar 17 al valor absolut de la diferència, amb 8 bits.

- El nombre -1:** Es codifica amb 0b111 seguit de seguit per la codificació en binari sense signe de restar 1 al nombre de -1 seguits trobats de forma no interrompuda per cap natural ni salt de línia, amb 13 bits.

Els espais entre nombres no apareixen codificats a l'algorisme de compressió. Per cada dos nombres enters seguits, assumim que hi ha un sol caràcter d'espai entre els dos, tret que hi hagi un salt de línia.

La seqüència:

$$1031 \ 1031 \ 1125 \ -1 \ -1 \ -1 \ 1117 \setminus n$$

La codificaríem com a:

- [illegible]

11000000 00000000 00000000 01000000 01111001 00010011 01111000 00000000 10001000 1(0000000)

Que, expressat en hexadecimal, veuríem com a

0x C0 00 00 40 79 13 78 00 88 80

2.2 Codis de Huffman

Aquesta primera codificació es divideix en blocs de fins a 65536 bytes. Aquests blocs es tornen a codificar fent servir codis de Huffman. Un bloc no inclou necessàriament un nombre sencer de símbols codificats: un mateix símbol pot quedar dividit amb x bits al bloc i , i amb y bits al bloc $i+1$. Sí que és necessari que un bloc inclogui un nombre sencer de bytes, a excepció de l'últim bloc, pel qual el contingut dels bits no utilitzats de l'últim byte és irrellevant.

Tanmateix, els blocs ja codificats amb codis de Huffman no estan confinats a ocupar un nombre sencer de bytes. Un bloc codificat pot començar i acabar a qualsevol bit dins un byte.

Aquests codis de Huffman compleixen les mateixes característiques establertes per l'algorisme DEFLATE[1]:

1. Tots els codis d'una mateixa longitud tenen valors lexicogràficament consecutius, en el mateix ordre que els valors que representen.
2. Els codis més curts precedeixen lexicogràficament als codis més llargs.

Per a cada bloc es genera un arbre de Huffman independent. El diccionari sobre el qual es crea cada arbre de Huffman és el format pels símbols 0..255, amb el pes de cadascun depenent de la seva freqüència al bloc; i el símbol especial 256, de pes 1, que representa el codi de final de bloc.

Cada bloc es codifica amb un llistat de les longituds dels codis de Huffman de cada símbol, per ordre del símbol, seguit del codi de cada símbol del bloc, seguit del codi per al símbol de fi de bloc 256.

La codificació de les longituds dels codis depèn del valor de la longitud:

- **Longituds 1..31:** Es codifica amb la representació en binari sense signe de la longitud amb 5 bits.
- **Longituds 32..:** Es codifica amb 0b00000 seguit de la representació en binari sense signe de restar 31 a la longitud del codi, amb 8 bits.

Per a codificar la fi del fitxer, un cop codificats tots els blocs, s'afegeixen els bits 0b00000000000000.

3 Anàlisi

La motivació darrera l'algorisme de compressió base gira entorn poder aprofitar els valors numèrics, i no només els bytes purs en format ASCII.

La decisió de codificar els increments o decrements entre valors consecutius prové de la natura de la font de les dades: En tractar-se de les altures sobre nivell del mar d'un avió en vol, és natural que els canvis entre naturals consecutius siguin petits, i que tendeixin a un creixement o decreixement constant, cosa que dificultaria trobar elements repetits per a algorismes de diccionari.

Com podem observar a la Figura 2, el 99% d'enters consecutius es troben a una distància inferior a 128. Per aquest motiu, escollim codificar els increments



Figure 2: Distribució dels valors absoluts de les diferències entre naturals consecutius.



Figure 3: Freqüències absolutes pels bytes 0x00 a 0xFF per les dades codificades només pel primer pas del compressor, ordenades per freqüència.

o decrements petits en la quantitat menor de bits possible, per no malbaratar espai amb 0 redundants.

D'aquesta manera, reduïm la mida de cada nombre natural (*suposant un nombre de 5 dígits seguit d'un espai, 6 bytes = 48 bits*) en...

- 2400% per a naturals iguals (*1~10% dels casos*).
- 685% per a diferències menors o iguals a 16 (*55~85% dels casos*).
- 400% per a diferències menors o iguals a 272. (*5~25% dels casos*).
- 300% per a diferències menors o iguals a 4368. (*<1% dels casos*).

La compressió de les seqüències de -1 està destinada a cadenes llargues. Amb només 15 bits, es codifiquen seqüències de fins a 8192 elements. Tot i semblar un nombre excessiu per a seqüències curtes, aquesta codificació aconsegueix reduir el nombre de bits necessaris per a codificar qualsevol nombre de -1 seguits, fins i tot per a un de sol (*assumint l'espai seguint el nombre com a part de la codificació en text pla del mateix*). Per a cadenes llargues, com per exemple les que trobem als fitxers de proves, del voltant de 5500 elements seguits, i que signifiquen un 15% del total de nombres enters, aconseguim una reducció del 880000%.

Tal com es pot intuir a l'exemple del punt 2.1.1, o com veiem a la Figura 3, els valors com a bytes de la codificació inicial són molt desproporcionats. Concretament, hi ha una gran abundància dels bytes 0x00, i bytes fets per set

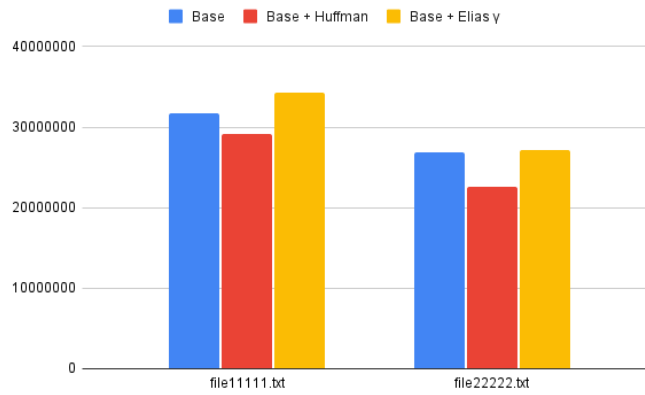


Figure 4: Comparació de les mides de fitxer per a les codificacions base, base + Huffman i base + Elias Gamma.



Figure 5: Freqüències absolutes pels bytes 0x00 a 0xFF per les dades codificades completament pel compressor fent servir codis de Huffman, ordenades per freqüència.

bits 0 i un d1. Això ens indica que hi ha marge de millora per la ràtio de compressió si tractem de nou aquestes dades.

Després de considerar una codificació amb codis de Huffman o amb *Run Length Encoding* fent servir la codificació Elias Gamma. Després de provar de codificar per segon cop els fitxers de prova, trobem que Huffman aconsegueix una ràtio de compressió al voltant d'1,1 sobre la codificació base, mentre que Elias Gamma empitjora els resultats. Per aquest motiu, decidim que la segona passada es realitzarà amb Huffman.

4 Codi

El codi de la pràctica està escrit en C++. Està dividit en 4 fitxers:

- **main.cpp** - Fitxer principal del codi. Únicament s'encarrega de fer una crida al compressor o descompressor.

- **CompressorPractica.cpp** - Implementa el compressor de fitxers. Consta d'una funció pública, que es crida amb dos paràmetres, strings indicant el camí al fitxer a codificar i el fitxer on desar-lo, que implementa el bucle principal; i una dotzena de funcions privades dedicades a codificar els blocs en codis de Huffman i escriure a disc.
 - **properNumero()**; A partir d'un string d'entrada i un iterador, retorna el pròxim nombre enter.
 - **diferenciaAlBloc()**; A partir de dos nombres naturals, afegeix al bloc la codificació de la diferència entre els dos.
 - **literalAlBloc()**; Afegeix al bloc la codificació d'un nombre natural.
 - **eolAlBloc()**; Afegeix al bloc la codificació d'un salt de línia.
 - **negatiuAlBloc()**; Afegeix al bloc la codificació d'una seqüència de -1 .
 - **bitsAlBloc()**; Afegeix al bloc una seqüència de bits de la llargària demanada.
 - **bitAlBloc()**; Afegeix al bloc un bit.
 - **byteAlBloc()**; Funció de suport a l'anterior que afegeix grups de bits al bloc.
 - **endOfFile()**; Buida el bloc actual a disc i escriu el símbol de fi de fitxer.
 - **blocAHuffman()**; Codifica el bloc actual amb codis de Huffman i els escriu a disc.
 - **longitudArbre()**; Donat un arbre binari, genera un vector amb les longituds dels codis de Huffman per a cada símbol.
 - **bitsAlDisc()**; Escriu al fitxer de disc una seqüència de bits de la llargària demanada.
 - **bitAlDisc()**; Escriu al fitxer de disc un bit.
 - **byteAlDisc()**; Escriu al fitxer de disc un byte.
 - **flushDisc()**; Escriu a disc la memòria cau que no s'ha escrit encara.
- **DescompressorPractica.cpp** - Implementa el descompressor de fitxers. Consta d'una funció pública, que es crida amb dos paràmetres, strings indicant el camí al fitxer a descodificar i el fitxer on desar-lo, que implementa el bucle principal; i set funcions privades dedicades a descodificar els blocs en codis de Huffman i escriure a disc.
 - **llegeixBitDeDisc()**; Retorna el pròxim bit del fitxer de disc.
 - **llegeixBitsDeDisc()**; Retorna els pròxims n bits demanats del fitxer de disc.
 - **llegeixByteDeDisc()**; Retorna el pròxim byte del fitxer de disc.

- `llegeixBlocDeDisc()`; Llegeix i descodifica un bloc codificat en Huffman del fitxer de disc.
- `llegeixBitDelBloc()`; Retorna el pròxim bit del bloc.
- `llegeixBitsDelBloc()`; Retorna els pròxims n bits demanats del bloc.
- `llegeixByteDelBloc()`; Retorna el pròxim byte del bloc.
- `BinTree.hpp` - Classe que implementa els arbres binaris. Aquesta classe ha estat distribuïda pel professorat del Departament de Ciències de la Computació de la UPC[2].

References

- [1] L. Peter Deutsch. *DEFLATE Compressed Data Format Specification version 1.3*. URL: <https://datatracker.ietf.org/doc/html/rfc1951>.
- [2] J. Petit. *BinTree.hh*. URL: <https://pro2.cs.upc.edu/>.