

Compressão de dados

AULA 17



Fonte: Best VPN for Data Compression

Referências: Entrada e saída binárias (PF), Compressão de dados (PF), Data Compression (SW), slides (SW), video (SW)

Introdução

Problema: representar um arquivo **GRANDE** por outro **menor**.

Exemplo:
arquivo **GRANDE**: ababababababababababababab
arquivo **menor**: 12ab

Por que comprimir?
Menor espaço de armazenamento.
Menor tempo de transmissão.

Software: gzip, bzip, 7z, etc.

Codificador e decodificador

Fluxo **B** é **original** e o fluxo **C(B)** é **codificado**.
Fluxo produzido pelo **expansor** é **decodificado**.
Fluxo **decodificado** é **idêntico** ao **original**, a compressão não perde informação (**lossless compression**).

Notação: $|B|$ é o número de bits de **B**.

Taxa de compressão: $|C(B)| / |B|$.

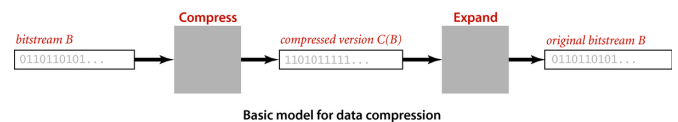
Desafio: obter a **menor** taxa de compressão possível.

Esquema

representar um dado **fluxo de bits** (*bitstream*) por outro mais curto.

Esquema básico de compressão de dados:

- ▶ um **compressor** transforma um fluxo de bits **B** em um fluxo **C(B)** e
- ▶ um **expansor** transforma **C(B)** de volta em **B**.



Considerações teóricas

Fato. Nenhum algoritmo pode garantir taxa de compressão estritamente **menor** que 1 para **todo** e qualquer fluxo de bits.

Fluxos de bits **aleatórios** são **pouco** compressíveis.

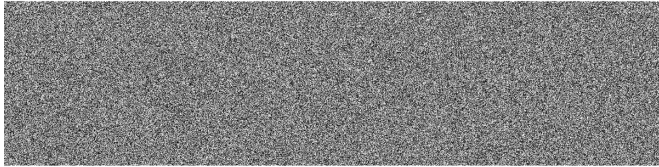
Mesmo fluxos **pseudo-aleatórios** podem ser difíceis de comprimir.

Considerações teóricas

Muitos fluxos de bits **parecem aleatórios** mas não são.

Exemplo: o fluxo de bits parece aleatório...

```
% java RandomBits | java PictureDump 2000 500
```



1000000 bits

A difficult file to compress: 1 million (pseudo-) random bits

Navigation icons: back, forward, search, etc.

Considerações teóricas

Outro exemplo:

$$4 \left(\sum_{i=0}^{\infty} (-1)^i (1/(2i + 1)) \right)$$

é uma representação **muito comprimida** da expansão decimal do número π .

Teoria da compressão de dados tem ligações fascinantes com a **Teoria da Informação** e os conceitos de **Aleatoriedade** e **Entropia**.

Navigation icons: back, forward, search, etc.

Cadeia de bits e fluxo de bits

Cadeia de bits (= *bitstring*) é uma sequência de bits:

```
110001000001110101010011111011001110001010000  
1100000000001111110000011101000111
```

fluxo de bits (= *bitstream*) é uma cadeia de bits na **entrada** ou na **saída** de um programa.

A classe **BinaryStdIn** lê um fluxo de bits a partir da **entrada padrão**.

A classe **BinaryStdOut** escreve um fluxo de bits na **saída padrão**.

Navigation icons: back, forward, search, etc.

Considerações teóricas

mas foi produzido pelo código

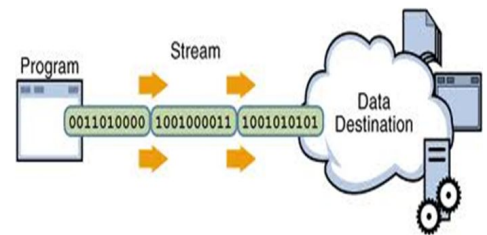
```
public class RandomBits {  
    public static void main(String[] args){  
        int x = 11111;  
        for(int i = 0; i < 1000000; i++) {  
            x = x * 314159 + 218281;  
            BinaryStdOut.write(x>0);  
        }  
        BinaryStdOut.close();  
    }  
}
```

Taxa de compressão: 0.002

Navigation icons: back, forward, search, etc.

Entrada e saídas binárias

Input / Output Streams



Fonte: [Input / Output Streams Byte streams](#)
Referências: [BinaryStdIn](#), [BinaryStdOut](#), [BinaryIn](#), [BinaryOut](#)

Navigation icons: back, forward, search, etc.

BinaryStdIn

```
public class BinaryStdIn  
  
boolean    readBoolean()    lê 1 bit e devolve  
                                     o booleana correspondente  
char       readChar()       lê 8 bits  
char       readChar(int r)  lê r (entre 1 e 16)  
                                     bits  
String     readString()     lê fluxo em blocos  
                                     de 8 bits  
int        readInt()        lê 32 bits  
int        readInt(int r)   lê r (entre 1 e 32)  
                                     bits  
boolean    isEmpty()        fluxo está vazio?  
void       close()          feche o fluxo
```

Navigation icons: back, forward, search, etc.

Genomas

Virus de verdade, taxa de compressão de $12536/50000 = 0.25$:

```
% java PictureDump 512 100 < genomeVirus.txt
```



50000 bits

```
% java Genome - < genomeVirus.txt | java  
PictureDump 512 25
```



12536 bits

Navigation icons: back, forward, search, etc.

Codificação de comprimento de carreira

Em inglês: *run-length encoding* (=RLE).

Exemplo: cadeia de bits abaixo tem uma carreira de 15 0s, uma carreira de 7 1s, uma de 7 0s, e uma de 11 1s:

```
000000000000000011111111000000111111111111
```

Pode ser representada pela sequência 15 7 7 11. Usando 8 bits para cada um desses números, teremos uma cadeia de apenas 32 bits (ignore os espaços):

```
00001111000011110000111100001011
```

Navigation icons: back, forward, search, etc.

Codificação de comprimento de carreira

Para texto ASCII, a compressão é ruim

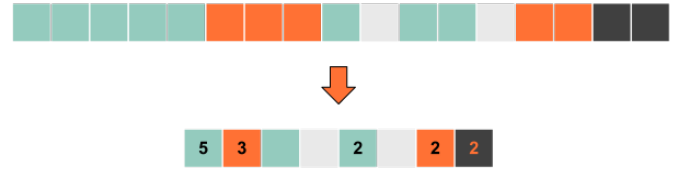
Exemplo: "ABRACADABRA!" tem carreiras curtas

```
% java BinaryDump 32 < abra.txt  
01000001010000100101001001000001  
01000011010000010100010001000001  
01000010010100100100000100100001  
96 bits
```

Navigation icons: back, forward, search, etc.

Codificação de comprimento de carreira

Lossless pixel compression



Fonte: Lossy Image Compression with Run-Length Encoding

Navigation icons: back, forward, search, etc.

Codificação de comprimento de carreira

Decisões de projeto:

- ▶ Use 8 bits (valores de 0 a 255) para cada comprimento de carreira.
- ▶ Use carreiras de comprimento 0 para dividir carreiras muito longas em blocos de comprimento menor que 256.
- ▶ A primeira carreira é sempre de 0s (e pode ser vazia).

Navigation icons: back, forward, search, etc.

Codificação de comprimento de carreira

Para texto ASCII, a compressão é ruim

Exemplo: "ABRACADABRA!" tem carreiras curtas

```
% java RunLength - < abra.txt | java  
HexDump 13  
01 01 05 01 01 01 04 01 02 01 01 01 02  
01 02 01 05 01 01 01 04 02 01 01 05 01  
01 01 03 01 03 01 05 01 01 01 04 01 02  
01 01 01 02 01 02 01 05 01 02 01 04 01  
416 bits
```

Taxa de compressão: $416/96 = 4.33$

Navigation icons: back, forward, search, etc.

