

**Instituto Politécnico de Beja**  
**Escola Superior de Tecnologia e Gestão**  
**Licenciatura em Engenharia Informática**

**Trabalho de Grupo Nº 1**

**Sistemas Operativos**

*Francisco Manuel Baião Soudo Nº 14060*

*Miguel Pauzinho Nº 27131*

Orientado por:

*Luís Garcia, IPBeja*

**Relatório de Projeto**

Beja, 14 de Novembro de 2024



## Conteúdo

Capítulo 1 .....	4
Introdução.....	4
Capítulo 2 .....	5
2.1 Estrutura do trabalho.....	6
2.2 Processamento do corpus.....	7
2.3 Cálculo de Métricas.....	8
2.4 Criador de Palavras.....	9
2.5 Pares de palavras .....	10
2.6 Verificador de pares de palavras .....	11
2.7 Criador de frases.....	12
2.8 Pares de frases .....	13
2.9 Verificador de pares de frases.....	14
2.10 Limitador de palavras .....	15
2.11 Script Opcional (Mestre).....	16
2.12 Copiador de ficheiros .....	18
Capítulo 3 .....	20
Conclusão .....	20

# Capítulo 1

## Introdução

No âmbito da unidade curricular de Sistemas Operativos, este projeto visa o desenvolvimento de um conjunto de scripts em Shell Script (bash) para processamento de um corpus textual em italiano, com o objetivo de gerar dicionários para o software Eugénio V3.

O Eugénio V3 é um sistema de investigação em comunicação aumentativa e alternativa, capaz de realizar predição de palavras e frases. Para este efeito, necessita de dicionários com informação estatística sobre a utilização de palavras e frases.

O projeto envolve a criação, análise e transformação de um corpus textual de grande dimensão, neste caso o corpus italiano paisa.raw.utf8, composto por mais de 8 milhões de linhas, do qual serão utilizadas as primeiras 200.000 linhas. O objetivo é gerar quatro ficheiros fundamentais para o funcionamento do Eugénio:

- words.txt: lista de palavras e respetivas ocorrências;
- words\_pairs.txt: lista de pares de palavras e respetivas ocorrências;
- sentences.txt: lista de frases e respetivas ocorrências;
- sentences\_pairs.txt: pares de frases e respetivas ocorrências.

Adicionalmente, serão implementados scripts de verificação e controlo, garantindo que os dicionários respeitam as restrições do Eugénio (máximo de 250.000 entradas) e seguem os requisitos estabelecidos.

O desenvolvimento foi realizado utilizando a linguagem Bash, assegurando compatibilidade com os sistemas operativos Linux e Windows.

## Capítulo 2

### Objetivos do Trabalho

1. **Criar e organizar a estrutura de diretórios** necessária para o projeto.
2. **Obter o corpus italiano** (`paisa.raw.utf8`) e gerar um sub-corpus com 200.000 linhas.
3. **Analisar e extrair estatísticas** do corpus textual.
4. Criar dicionários de:
  - o palavras → `words.txt`
  - o pares de palavras → `words_pairs.txt`
  - o frases → `sentences.txt`
  - o pares de frases → `sentences_pairs.txt`
5. **Verificar a integridade dos dicionários**, assegurando que os pares existem nos ficheiros base.
6. **Limitar** o número de palavras/frases aos valores permitidos pelo Eugénio V3.
7. Criar um **script Windows** para instalar os dicionários no software Eugénio.
8. Criar um **Master Script** em bash que executa todos os requisitos pela ordem correta.
9. Produzir um **relatório académico** apresentando e explicando todas as soluções desenvolvidas.

## 2.1 Estrutura do trabalho

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1$ tree
.
└── tg1
    ├── corpus
    │   └── paisa.raw.utf8
    ├── corpus_info
    │   └── corpus_info.txt
    ├── corpus_txt
    │   └── paisa.raw.utf8
    ├── scripts
    │   ├── check_sentences_pairs_req9.sh
    │   ├── corpusInfo_req3.sh
    │   ├── corpus_sentences_req7.sh
    │   ├── corpustxt_req2.sh
    │   ├── create_word_pairs_req5.sh
    │   ├── limit_sentences_req10.sh
    │   ├── Master.sh
    │   ├── sentences_pairs_req8.sh
    │   ├── verifica_palavras_pares_req6.sh
    │   └── words_txt_req4.sh
    ├── sentences_dict
    │   ├── sentences_pairs.txt
    │   └── sentences.txt
    └── words_dict
        ├── words_pairs.txt
        └── words.txt
```

Cada diretoria contém:

- **corpus/** → corpus original (opcional).
- **corpus\_txt/** → subcorpus de 200.000 linhas.
- **corpus\_info/** → estatísticas sobre o corpus.
- **words\_dict/** → *words.txt* e *words\_pairs.txt*.
- **sentences\_dict/** → *sentences.txt* e *sentences\_pairs.txt*.
- **scripts/** → todos os scripts bash desenvolvidos para os requisitos.

## 2.2 Processamento do corpus

### Requisito 2 — Criação do Sub-Corpus (200.000 linhas)

O objetivo deste requisito é extrair apenas as primeiras 200.000 linhas do ficheiro original *paisa.raw.utf8*.

O script **corpustxt\_req2.sh** realiza:

1. Verificação da existência do ficheiro de origem.
2. Extração das 200.000 primeiras linhas usando head.
3. Validação automática da contagem de linhas.

Este ficheiro constituirá a base de trabalho para todos os requisitos seguintes.

```
$ create_sub_corpus_req2.sh
$ create_sub_corpus.req2.sh
1  #!/bin/bash
2
3  # Variáveis
4  Ficheiro_Original="/home/$USER/tg1/corpus/paisa.raw.utf8"
5  Ficheiro_Trabalho="/home/$USER/tg1/corpus_txt/paisa_200k.txt"
6
7  # Criar diretoria se não existir
8  mkdir -p "/home/$USER/tg1/corpus_txt"
9
10 if [ -f "$Ficheiro_Original" ]
11 then
12     echo "O ficheiro está na pasta"
13     echo "A extrair as primeiras 200000 linhas"
14     head -200000 $Ficheiro_Original > $Ficheiro_Trabalho
15
16     contar_linhas=$(wc -l < $Ficheiro_Trabalho)
17
18     if [ $contar_linhas -eq 200000 ]
19     then
20         echo "Linhas extraídas com sucesso"
21     else
22         echo "As linhas não foram extraídas com sucesso"
23     fi
24 else
25     echo "O ficheiro não está na pasta"
26 fi
27
```

#### 2.2.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ chmod +x corpustxt_req2.sh
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./corpustxt_req2.sh
O ficheiro está na pasta
A extrair as primeiras 200000 linhas
Linhas extraídas com sucesso
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

## 2.3 Cálculo de Métricas

Neste ponto o pretendido era que os alunos desenvolvessem um script com o objetivo de proceder à caracterização do corpus utilizado, calculando um conjunto de métricas essenciais para a sua análise. Este script foi projetado para extrair as seguintes informações do corpus: o número total de caracteres, a quantidade de linhas não vazias, o número total de palavras, o número de palavras distintas, o quociente entre o total de palavras distintas e o total de palavras, o número total de frases, o número de frases distintas e o quociente entre o total de frases distintas e o total de frases. Para alcançar estes objetivos, foi criado o script 2.2, que procede ao processamento do corpus e realiza os cálculos necessários para cada métrica indicada. O resultado desta operação é guardado num ficheiro de saída, localizado na subdiretoria `corpus_info`.

```
$ corpus_info_req3.sh
1 #!/bin/bash
2 # Variáveis
3 Ficheiro_Trabalho="/home/$USER/tg1/corpus_txt/paisa_200k.txt"
4 Resultado="/home/$USER/tg1/corpus_info/corpus_info.txt"
5 # Criar diretória se não existir
6 mkdir -p "/home/$USER/tg1/corpus_info"
7 # Conta corretamente o número de caracteres
8 numero_caracteres=$(wc -c < $Ficheiro_Trabalho)
9 # Conta apenas linhas não vazias (> "$") ignora linhas em branco
10 numero_linhas_nao_vazias=$(grep -v '^$' $Ficheiro_Trabalho | wc -l)
11 # Calcula o número total de palavras existentes no ficheiro do corpus
12 # cat: lê o conteúdo completo do ficheiro e envia-o para o "pipe"
13 # tr: Comando translate ou transform | -s opção de "squeeze", comprime repetições do mesmo caractere
14 numero_total_de_palavras=$(cat $Ficheiro_Trabalho | tr -s ' ' | wc -w)
15 # Converte para minúsculas, remove pontuação, quebra palavras em linhas, elimina vazias, ordena e conta únicas
16 numero_total_de_palavras_diferentes=$(tr 'A-Z' 'a-z' < $Ficheiro_Trabalho | sed 's/[[[:punct:]]]/ /' | tr -s ' ' | tr '\n' '' | grep -v '^$' | sort | uniq | wc -l)
17 # Calcula a proporção entre o número de palavras diferentes e o número total de palavras
18 total_de_palavras_diferentes_porc_total_de_palavras=$echo "scale=6; $numero_total_de_palavras_diferentes / $numero_total_de_palavras" | bc
19 # Define frases como linhas do corpus
20 numero_total_de_frases=$numero_linhas_nao_vazias
21 # Remove duplicadas e conta frases únicas
22 numero_total_de_frases_diferentes=$(grep -v '^$' $Ficheiro_Trabalho | sort | uniq | wc -l)
23 # Calcula a proporção (ou quociente) entre $numero_total_de_frases_diferentes / $numero_total_de_frases
24 total_de_frases_diferentes_porc_total_de_frases=$echo "scale=6; $numero_total_de_frases_diferentes / $numero_total_de_frases" | bc
25 # Bloco que organiza e limpa
26 {
27     echo "==== INFORMAÇÕES DO CORPUS ==="
28     echo ""
29     echo "Número de caracteres: $numero_caracteres"
30     echo "Quantidade de linhas não vazias: $numero_linhas_nao_vazias"
31     echo "Número total de palavras: $numero_total_de_palavras"
32     echo "Número total de palavras diferentes: $numero_total_de_palavras_diferentes"
33     echo "Quociente (palavras diferentes / total palavras): $total_de_palavras_diferentes_porc_total_de_palavras"
34     echo "Número total de frases: $numero_total_de_frases"
35     echo "Número total de frases diferentes: $numero_total_de_frases_diferentes"
36     echo "Quociente (frases diferentes / total frases): $total_de_frases_diferentes_porc_total_de_frases"
37 } > $Resultado
38 echo "Resultados guardados em $Resultado"
39 cat $Resultado
```

### 2.3.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ chmod +x corpusInfo_req3.sh
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./corpusInfo_req3.sh
Resultados guardados em /home/francisco/Sistemas-Operativos-TG1/tg1/corpus_info/corpus_info.txt
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

```
Open ▾ corpus.info.txt
=====
Data de criação do ficheiro : Fri Nov 14 06:54:30 PM WET 2025

Número de caracteres: 17973242
Quantidade de linhas não vazias: 200000
Número total de palavras: 2758049
Número total de palavras diferentes: 90324
Quociente (palavras diferentes / total palavras): .0327
Número total de frases: 200000
Número total de frases diferentes: 34558
Quociente (frases diferentes / total frases): .17275
```

## 2.4 Criador de Palavras

Neste ponto o pretendido era que os alunos desenvolvessem um script com o objetivo de criar um ficheiro de palavras que regista cada palavra do corpus juntamente com o seu número de ocorrências. Este ficheiro, denominado "words.txt", deve apresentar em cada linha uma palavra única seguida da sua frequência no corpus. A listagem foi estruturada de forma a ser apresentada em ordem alfabética. Para o desenvolvimento deste script foi feito uso de algumas comandos como por exemplo awk e o sort para fazer o processamento do corpus, todos os caracteres não alfabéticos em quebra de linha, transforma ainda todos os caracteres maiúsculos em minúsculos para desta forma garantir que a contagem de palavras não seja induzida em erro por ser key sensitive, de seguida é utilizado o uniq para contar as ocorrências, por fim ao organizar as palavras por ordem alfabética o ficheiro gerado inicialmente é ainda processado pelo script que limita o conteúdo .

```
$ create_words_req4.sh
1 #!/bin/bash
2
3 # Variáveis
4 Ficheiro_Trabalho="/home/$USER/tg1/corpus_txt/paisa_200k.txt"
5 Resultado="/home/$USER/tg1/words_dict/words.txt"
6
7 # Criar diretoria se não existir
8 mkdir -p "/home/$USER/words_dict"
9
10 # Processar palavras do corpus
11 # cat: lê o ficheiro completo
12 # tr: -s remove espaços duplicados, deixando apenas um espaço entre palavras
13 # sed: substitui toda a pontuação por espaços
14 # tr: converte maiúsculas para minúsculas
15 # tr: substitui espaços por quebras de linha
16 # grep: remove linhas vazias que podem surgir depois de substituir pontuação
17 # sort: ordena todas as palavras alfabeticamente - passo essencial para o uniq agrupar repetições
18 # uniq -c: conta as ocorrências de cada palavra (a opção -c adiciona a contagem à esquerda)
19 # sed: remove espaços iniciais que o uniq -c adiciona antes do número
20 # awk: reordena para formato 'palavra contagem'
21 # sort: ordena novamente alfabeticamente pela primeira coluna (a palavra)
22 cat $Ficheiro_Trabalho | \
23     tr -s '' | \
24     sed 's/[:punct:]/ /g' | \
25     tr 'A-Z' 'a-z' | \
26     tr ' ' '\n' | \
27     grep -v '^$' | \
28     sort | \
29     uniq -c | \
30     sed 's/^ *//g' | \
31     awk '{print $2, $1}' | \
32     sort -k1 > $Resultado
33
34 echo "Ficheiro de palavras criado: $Resultado"
35 echo "Total de palavras: $(wc -l < $Resultado)"
36
```

### 2.4.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./words_txt_req4.sh
Ficheiro de palavras criado: /home/francisco/Sistemas-Operativos-TG1/tg1/words_dict/words.txt
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

## 2.5 Pares de palavras

O objetivo deste ponto era que os alunos desenvolvessem um script que processa o corpus italiano de maneira que o output tivesse em cada linha um par de palavras e as suas ocorrências no corpus. O output deve ainda estar ordenado de forma alfabética e limitado a um máximo de 250 mil palavras e frases. Para a realização deste script foi feito uso de vários comandos à disposição do utilizador como se pode observar pela listagem 2.5. O script começa por caracteres que não sejam letras por newlines de maneira a isolar cada palavra em uma linha diferente, em seguida as palavras são convertidas para minúsculas de maneira a evitar possíveis erros na contagem. Após esses comandos serem executados as palavras são concatenadas na mesma linha com uma barra vertical como separador e é então feita a ordenação, contagem e estruturação do output. Por fim o ficheiro gerado inicialmente é processado pelo script de limitação de conteúdo.

```
$ create_word_pairs_req5.sh
1  #!/bin/bash
2
3  # Variáveis
4  Ficheiro_Trabalho="/home/$USER/tg1/corpus_txt/paisa_200k.txt"
5  Resultado="/home/$USER/tg1/words_dict/words_pairs.txt"
6
7  # Cria diretória se não existir
8  mkdir -p "/home/$USER/tg1/words_dict"
9  # Verificação do corpus
10 if [ ! -f "$Ficheiro_Trabalho" ]
11 then
12     echo "Erro: Ficheiro $Ficheiro_Trabalho não encontrado!"
13     exit 1
14 fi
15 # Processar pares de palavras
16 # cat: lê o ficheiro completo
17 # tr -s: remove espaços duplicados
18 # sed: substitui pontuação por espaços
19 # tr: converte para minúsculas
20 # tr -s: converte espaços para quebras de linha
21 # grep: remove linhas vazias
22 # awk: cria pares de palavras consecutivas (guarda palavra anterior e imprime par)
23 # sort: ordena os pares
24 # uniq <: conta ocorrências de cada par
25 # sed: remove espaços iniciais
26 # awk: reordena para formato "palavra1 palavra2 contagem"
27 # sort: ordena alfabeticamente pelos dois primeiros campos
28 cat $Ficheiro_Trabalho | \
29   tr -s '\n' | \
30   sed 's/[[:punct:]]//g' | \
31   tr 'À-Ã' 'à-ã' | \
32   tr -s ' ' '\n' | \
33   grep -v '^$' | \
34   awk '{NR>1 {print prev, $0} {prev=$0}}' | \
35   sort | \
36   uniq < | \
37   sed 's/ /|/' | \
38   awk '{print $2, $3, $1}' | \
39   sort -k1,2 > $resultado
40
41 echo "Ficheiro de pares de palavras criado: $Resultado"
42 echo "Total de pares: $(wc -l < $Resultado)"
```

### 2.5.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./create_word_pairs_req5.sh
Ficheiro de pares de palavras criado: /home/francisco/Sistemas-Operativos-TG1/tg1/words_dict/words_pairs.txt
Total de pares únicos: 746111
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

## 2.6 Verificador de pares de palavras

Foi pedido que após realizar o script para a formação de pares de palavras os alunos criassem outro script que verifica se cada palavra dos pares estava presente no ficheiro com a lista de palavras ("words.txt"). Para a criação deste script foi feito uso do comando "awk" e das suas funções para formatar as linhas como pretendido. O script começa por percorrer o ficheiro com pares de palavras e adiciona cada palavra a um array que então passa por uma condição if que verifica se cada palavra pertence ao array de palavras.

```
$ verify_word_pairs_req6.sh
1 #!/bin/bash
2
3 # Variáveis
4 Palavras="/home/$USER/tg1/words_dict/words.txt"
5 Pares="/home/$USER/tg1/words_dict/words_pairs.txt"
6
7 echo "A verificar pares de palavras..."
8
9 # Extraír palavras dos pares (primeira e segunda coluna)
10 awk '{print $1"\n"$2}' $Pares | sort -u > /tmp/palavras_dos_pares.txt
11
12 # Extraír palavras do dicionário (primeira coluna)
13 awk '{print $1}' $Palavras | sort -u > /tmp/palavras_dicionario.txt
14
15 # Verificar se todas as palavras dos pares existem no dicionário
16 # comm -23: mostra linhas que estão no primeiro ficheiro mas não no segundo
17 palavras_faltam=$(comm -23 /tmp/palavras_dos_pares.txt /tmp/palavras_dicionario.txt | wc -l)
18
19 if [ $palavras_faltam -eq 0 ]
20 then
21     echo "✓ Todas as palavras dos pares existem no dicionário"
22 else
23     echo "✗ Faltam $palavras_faltam palavras no dicionário"
24 fi
25
26 # Limpar ficheiros temporários
27 rm /tmp/palavras_dos_pares.txt /tmp/palavras_dicionario.txt
```

### 2.6.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./verifica_palavras_pares_req6.sh
A verificar se todas as palavras dos pares existem em words.txt...
```

Todas as palavras dos pares existem no ficheiro words.txt

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```



```
Open - F
corpus_info.txt | words.txt | corpus_info.txt
corpus_info.txt
words.txt
0000 2
000 522
000euro 1
000w 1
0001 1
001 1
0014 3
00187 2
002 1
004 2
005 1
00 523
007 4
0076108 1
008819 1
0103 1
018467 1
01209 1
01 342
014 1
01602 1
018 1
```

## 2.7 Criador de frases

Foi pedido que os alunos desenvolvessem um outro script que isolava cada frase do corpus italiano em uma linha e contasse a sua ocorrência no mesmo. Como é possível observar o script tem um funcionamento simples, fazendo uso da função "tr" para substituir espaços por barras verticais, o comando awk para isolar a frase e os comandos "sort" e "uniq -c" para fazer a ordenação e contagem da ocorrência das frases.

```

9 # create_sentences_retail
10 #!/bin/bash
11 # Variáveis
12 Ficheiro_Trabalho="/home/$USER/tgl/corpus_txt/paisa_200k.txt"
13 Resultado="/home/$USER/tgl/sentences_dict/sentences.txt"
14 MAX_WORDS=30
15 # Verifica se o diretório se não existir
16 mkdir -p "/home/$USER/tgl/sentences_dict"
17 # Verificação do corpus
18 if [ ! -f "$Ficheiro_Trabalho" ]
19 then
20     echo "Erro: Ficheiro $Ficheiro_Trabalho não encontrado!"
21     echo "Por favor, execute primeiro o script de requisito 2 para criar o corpus."
22     exit 1
23 fi
24
25 echo "A processar frases do corpus..." 
26 # Processar frases com limite de 30 palavras (evita crash de Eugénio)
27 # grep: remove linhas vazias | sed: divide texto por ponto final (\.) | tr: converte para minúsculas | sed: remove caracteres especiais, mantendo apenas letras e espaços
28 # sed: remove espaços em branco e adiciona espaço | sed: remove espaços no início e fim de linha | awk: limita cada frase a máximo 30 palavras e filtra frases muito curtas
29 # tr: substitui espaços por || grep: remove linhas vazias ou só com || | sort: ordena as frases | uniq -c: conta ocorrências | sed: remove espaços iniciais
30 # awk: reformata para "frase contagem"
31 # sorteia ordem alfabeticamente
32 grep -v "[[:space:]]*\t" "$Ficheiro_Trabalho" | \
33     sed 's/[[:space:]]*/\n/g' | \
34     tr '\r\n\t' '\n' | \
35     grep -v '^\[[:space:]\]*$' | \
36     sed 's/[[:alpha:]]/ \1/g' | \
37     sed 's/[[:space:]]//g' | \
38     sed 's/[[:space:]]//g' | \
39     awk '{if ($1 > 30) {next}}{if ($1 < 1) {next}}{if ($1 > max) {max=$1}}{if ($1 < min) {min=$1}}{result = result $1 "\n"}{if ($1 == max) {print result}}{if ($1 == min) {print result}}{if ($1 < max & $1 > min) {result = result $1 "\n"} }{if ($1 < min & $1 > max) {result = result $1 "\n"} }{if ($1 < min) {min=$1}}{if ($1 > max) {max=$1}}}' | \
40     awk -v max="$MAX_WORDS" -v min="1" | \
41     sort -n | \
42     grep -v '^$' | \
43     grep -v '^ *$' | \
44     sort | \
45     uniq -c | \
46     sed 's/^ *//g' | \
47     awk '{c=$1 $1=""; sub(/ /, ""); print $0, c}' | \
48     sort > "$Resultado"
49
50 echo "Ficheiro de frases criado: $Resultado"
51 echo "Total de Frases Únicas: $(wc -l < $Resultado)"

```

## 2.7.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./corpus_sentences_req7.sh  
Ficheiro de frases criado: /home/francisco/Sistemas-Operativos-TG1/tg1/sentences_dict/sentences.txt  
Total de frases únicas: 34558  
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

## 2.8 Pares de frases

Similarmente ao script de criação de pares de palavras também foi pedido aos alunos que criassem um outro script que formava pares de frases, trocando os seus espaços por barras verticais e contando a sua ocorrência no corpus. Para além dos requisitos anteriormente mencionados o output deste script também deverá ficar limitado a 250 mil frases e palavras. Como é possível observar pela listagem 2.7 o script faz uso de uma variedade de comandos do terminal de Linux. O script começa por extrair cada frase e trocar os seus espaços por barras verticais, em seguida junta cada frase com a que antecede e prossegue a fazer a ordenação dos pares por ordem alfabética e contar as suas ocorrências. Por fim o script formata o par com a sua ocorrência de maneira a ficar de acordo com o formato aceite pelo Eugénio.

```
$ create_sentence_pairs_req8.sh
1 #!/bin/bash
2 # Autor: Francisco
3 Ficheiro_Trabalho="/home/SUSER/tg1/corpus_txt/paisa_200k.txt"
4 Resultado="/home/SUSER/tg1/sentences_dict/sentences_pairs.txt"
5 MAX_WORDS=30
6 # Cria diretória se não existir
7 if [ ! -d "$Ficheiro_Trabalho" ]
8 then
9     echo "Ficheiro $Ficheiro_Trabalho não encontrado!"
10    echo "Por favor, execute primeiro o script do requisito 2 para criar o corpus."
11    exit 1
12 fi
13 echo "A processar pares de frases do corpus..."
14 # Processar pares de frases com limite de 30 palavras (evita crash do Eugénio)
15 # grep: remove linhas vazias | sed: divide texto por pontuação (., , !, ?) | tr: converte para minúsculas | sed: remove caracteres especiais, mantendo apenas letras e espaços | sed: normaliza espaços múltiplos para um único espaço
16 # sed: remove espaços no inicio e fim da linha | awk: limita cada frase a máximo 30 palavras e filtra frases muito curtas | tr: substitui espaços por |
17 # grep: remove linhas vazias ou só com / | awk: cria pares de frases consecutivas | sort: ordena os pares
18 # grep -c conta ocorrências
19 # sed: remove espaços iniciais
20 # awk: reformata para "frase1 frase2 contagem"
21 # sort: ordena alfabeticamente
22 # sort: ordena alfabeticamente
23 grep -v '^$' "$Ficheiro_Trabalho" | \
24 sed 's/[[:upper:]]/`[[:lower:]]`/g' | \
25 tr '[:upper:]' '[:lower:]' | \
26 sed 's/[^[:alpha:]]*[[:space:]]*$/`/ | \
27 sed 's/`[^[:space:]]*/`/ | \
28 sed 's/`[^[:space:]]*/`/ | \
29 awk '{c=$1; prev=$0; words=split($0, words, " ")'
30         if (words > max) {
31             result = words[1]
32             for (i = 2; i <= max; i++) {
33                 result = result " " words[i]
34             }
35             print result
36         } else {
37             print $0
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    '
1001    '
1002    '
1003    '
1004    '
1005    '
1006    '
1007    '
1008    '
1009    '
1010    '
1011    '
1012    '
1013    '
1014    '
1015    '
1016    '
1017    '
1018    '
1019    '
1020    '
1021    '
1022    '
1023    '
1024    '
1025    '
1026    '
1027    '
1028    '
1029    '
1030    '
1031    '
1032    '
1033    '
1034    '
1035    '
1036    '
1037    '
1038    '
1039    '
1040    '
1041    '
1042    '
1043    '
1044    '
1045    '
1046    '
1047    '
1048    '
1049    '
1050    '
1051    '
1052    '
1053    '
1054    '
1055    '
1056    '
1057    '
1058    '
1059    '
1060    '
1061    '
1062    '
1063    '
1064    '
1065    '
1066    '
1067    '
1068    '
1069    '
1070    '
1071    '
1072    '
1073    '
1074    '
1075    '
1076    '
1077    '
1078    '
1079    '
1080    '
1081    '
1082    '
1083    '
1084    '
1085    '
1086    '
1087    '
1088    '
1089    '
1090    '
1091    '
1092    '
1093    '
1094    '
1095    '
1096    '
1097    '
1098    '
1099    '
1100    '
1101    '
1102    '
1103    '
1104    '
1105    '
1106    '
1107    '
1108    '
1109    '
1110    '
1111    '
1112    '
1113    '
1114    '
1115    '
1116    '
1117    '
1118    '
1119    '
1120    '
1121    '
1122    '
1123    '
1124    '
1125    '
1126    '
1127    '
1128    '
1129    '
1130    '
1131    '
1132    '
1133    '
1134    '
1135    '
1136    '
1137    '
1138    '
1139    '
1140    '
1141    '
1142    '
1143    '
1144    '
1145    '
1146    '
1147    '
1148    '
1149    '
1150    '
1151    '
1152    '
1153    '
1154    '
1155    '
1156    '
1157    '
1158    '
1159    '
1160    '
1161    '
1162    '
1163    '
1164    '
1165    '
1166    '
1167    '
1168    '
1169    '
1170    '
1171    '
1172    '
1173    '
1174    '
1175    '
1176    '
1177    '
1178    '
1179    '
1180    '
1181    '
1182    '
1183    '
1184    '
1185    '
1186    '
1187    '
1188    '
1189    '
1190    '
1191    '
1192    '
1193    '
1194    '
1195    '
1196    '
1197    '
1198    '
1199    '
1200    '
1201    '
1202    '
1203    '
1204    '
1205    '
1206    '
1207    '
1208    '
1209    '
1210    '
1211    '
1212    '
1213    '
1214    '
1215    '
1216    '
1217    '
1218    '
1219    '
1220    '
1221    '
1222    '
1223    '
1224    '
1225    '
1226    '
1227    '
1228    '
1229    '
1230    '
1231    '
1232    '
1233    '
1234    '
1235    '
1236    '
1237    '
1238    '
1239    '
1240    '
1241    '
1242    '
1243    '
1244    '
1245    '
1246    '
1247    '
1248    '
1249    '
1250    '
1251    '
1252    '
1253    '
1254    '
1255    '
1256    '
1257    '
1258    '
1259    '
1260    '
1261    '
1262    '
1263    '
1264    '
1265    '
1266    '
1267    '
1268    '
1269    '
1270    '
1271    '
1272    '
1273    '
1274    '
1275    '
1276    '
1277    '
1278    '
1279    '
1280    '
1281    '
1282    '
1283    '
1284    '
1285    '
1286    '
1287    '
1288    '
1289    '
1290    '
1291    '
1292    '
1293    '
1294    '
1295    '
1296    '
1297    '
1298    '
1299    '
1300    '
1301    '
1302    '
1303    '
1304    '
1305    '
1306    '
1307    '
1308    '
1309    '
1310    '
1311    '
1312    '
1313    '
1314    '
1315    '
1316    '
1317    '
1318    '
1319    '
1320    '
1321    '
1322    '
1323    '
1324    '
1325    '
1326    '
1327    '
1328    '
1329    '
1330    '
1331    '
1332    '
1333    '
1334    '
1335    '
1336    '
1337    '
1338    '
1339    '
1340    '
1341    '
1342    '
1343    '
1344    '
1345    '
1346    '
1347    '
1348    '
1349    '
1350    '
1351    '
1352    '
1353    '
1354    '
1355    '
1356    '
1357    '
1358    '
1359    '
1360    '
1361    '
1362    '
1363    '
1364    '
1365    '
1366    '
1367    '
1368    '
1369    '
1370    '
1371    '
1372    '
1373    '
1374    '
1375    '
1376    '
1377    '
1378    '
1379    '
1380    '
1381    '
1382    '
1383    '
1384    '
1385    '
1386    '
1387    '
1388    '
1389    '
1390    '
1391    '
1392    '
1393    '
1394    '
1395    '
1396    '
1397    '
1398    '
1399    '
1400    '
1401    '
1402    '
1403    '
1404    '
1405    '
1406    '
1407    '
1408    '
1409    '
1410    '
1411    '
1412    '
1413    '
1414    '
1415    '
1416    '
1417    '
1418    '
1419    '
1420    '
1421    '
1422    '
1423    '
1424    '
1425    '
1426    '
1427    '
1428    '
1429    '
1430    '
1431    '
1432    '
1433    '
1434    '
1435    '
1436    '
1437    '
1438    '
1439    '
1440    '
1441    '
1442    '
1443    '
1444    '
1445    '
1446    '
1447    '
1448    '
1449    '
1450    '
1451    '
1452    '
1453    '
1454    '
1455    '
1456    '
1457    '
1458    '
1459    '
1460    '
1461    '
1462    '
1463    '
1464    '
1465    '
1466    '
1467    '
1468    '
1469    '
1470    '
1471    '
1472    '
1473    '
1474    '
1475    '
1476    '
1477    '
1478    '
1479    '
1480    '
1481    '
1482    '
1483    '
1484    '
1485    '
1486    '
1487    '
1488    '
1489    '
1490    '
1491    '
1492    '
1493    '
1494    '
1495    '
1496    '
1497    '
1498    '
1499    '
1500    '
1501    '
1502    '
1503    '
1504    '
1505    '
1506    '
1507    '
1508    '
1509    '
1510    '
1511    '
1512    '
1513    '
1514    '
1515    '
1516    '
1517    '
1518    '
1519    '
1520    '
1521    '
1522    '
1523    '
1524    '
1525    '
1526    '
1527    '
1528    '
1529    '
1530    '
1531    '
1532    '
1533    '
1534    '
1535    '
1536    '
1537    '
1538    '
1539    '
1540    '
1541    '
1542    '
1543    '
1544    '
1545    '
1546    '
1547    '
1548    '
1549    '
1550    '
1551    '
1552    '
1553    '
1554    '
1555    '
1556    '
1557    '
1558    '
1559    '
1560    '
1561    '
1562    '
1563    '
1564    '
1565    '
1566    '
1567    '
1568    '
1569    '
1570    '
1571    '
1572    '
1573    '
1574    '
1575    '
1576    '
1577    '
1578    '
1579    '
1580    '
1581    '
1582    '
1583    '
1584    '
1585    '
1586    '
1587    '
1588    '
1589    '
1590    '
1591    '
1592    '
1593    '
1594    '
1595    '
1596    '
1597    '
1598    '
1599    '
1600    '
1601    '
1602    '
1603    '
1604    '
1605    '
1606    '
1607    '
1608    '
1609    '
1610    '
1611    '
1612    '
1613    '
1614    '
1615    '
1616    '
1617    '
1618    '
1619    '
1620    '
1621    '
1622    '
1623    '
1624    '
1625    '
1626    '
1627    '
1628    '
1629    '
1630    '
1631    '
1632    '
1633    '
1634    '
1635    '
1636    '
1637    '
1638    '
1639    '
1640    '
1641    '
1642    '
1643    '
1644    '
1645    '
1646    '
1647    '
1648    '
1649    '
1650    '
1651    '
1652    '
1653    '
1654    '
1655    '
1656    '
1657    '
1658    '
1659    '
1660    '
1661    '
1662    '
1663    '
1664    '
1665    '
1666    '
1667    '
1668    '
1669    '
1670    '
1671    '
1672    '
1673    '
1674    '
1675    '
1676    '
1677    '
1678    '
1679    '
1680    '
1681    '
1682    '
1683    '
1684    '
1685    '
1686    '
1687    '
1688    '
1689    '
1690    '
1691    '
1692    '
1693    '
1694    '
1695    '
1696    '
1697    '
1698    '
1699    '
1700    '
1701    '
1702    '
1703    '
1704    '
1705    '
1706    '
1707    '
1708    '
1709    '
1710    '
1711    '
1712    '
1713    '
1714    '
1715    '
1716    '
1717    '
1718    '
1719    '
1720    '
1721    '
1722    '
1723    '
1724    '
1725    '
1726    '
1727    '
1728    '
1729    '
1730    '
1731    '
1732    '
1733    '
1734    '
1735    '
1736    '
1737    '
1738    '
1739    '
1740    '
1741    '
1742    '
1743    '
1744    '
1745    '
1746    '
1747    '
1748    '
1749    '
1750    '
1751    '
1752    '
1753    '
1754    '
1755    '
1756    '
1757    '
1758    '
1759    '
1760    '
1761    '
1762    '
1763    '
1764    '
1765    '
1766    '
1767    '
1768    '
1769    '
1770    '
1771    '
1772    '
1773    '
1774    '
1775    '
1776    '
1777    '
1778    '
1779    '
1780    '
1781    '
1782    '
1783    '
1784    '
1785    '
1786    '
1787    '
1788    '
1789    '
1790    '
1791    '
1792    '
1793    '
1794    '
1795    '
1796    '
1797    '
1798    '
1799    '
1800    '
1801    '
1802    '
1803    '
1804    '
1805    '
1806    '
1807    '
1808    '
1809    '
1810    '
1811    '
1812    '
1813    '
1814    '
1815    '
1816    '
1817    '
1818    '
1819    '
1820    '
1821    '
1822    '
1823    '
1824    '
1825    '
1826    '
1827    '
1828    '
1829    '
1830    '
1831    '
1832    '
1833    '
1834    '
1835    '
1836    '
1837    '
1838    '
1839    '
1840    '
1841    '
1842    '
1843    '
1844    '
1845    '
1846    '
1847    '
1848    '
1849    '
1850    '
1851    '
1852    '
1853    '
1854    '
1855    '
1856    '
1857    '
1858    '
1859    '
1860    '
1861    '
1862    '
1863    '
1864    '
1865    '
1866    '
1867    '
1868    '
1869    '
1870    '
1871    '
1872    '
1873    '
1874    '
1875    '
1876    '
1877    '
1878    '
1879    '
1880    '
1881    '
1882    '
1883    '
1884    '
1885   
```

## 2.9 Verificador de pares de frases

Este script tem como objetivo verificar se os ficheiros **sentences.txt** e **sentences\_pairs.txt** foram criados corretamente e quantas entradas contêm. Ele começa por definir os caminhos destes dois ficheiros e, em seguida, informa o utilizador de que vai iniciar a verificação. Depois disso, o script conta quantas linhas existem em **sentences.txt**, que correspondem ao número total de frases únicas, e quantas linhas existem em **sentences\_pairs.txt**, que representam o total de pares de frases gerados. Por fim, apresenta estes valores no ecrã e mostra uma mensagem indicando que os ficheiros foram criados com sucesso. Trata-se de uma verificação simples que permite confirmar que o processamento das frases e dos seus pares foi concluído sem problemas.

```

1  #!/bin/bash
2
3  # Caminhos fixos
4  SentencesFile="/home/$USER/Sistemas-Operativos-TG1/tg1/sentences_dict/sentences.txt"
5  PairsFile="/home/$USER/Sistemas-Operativos-TG1/tg1/sentences_dict/sentences_pairs.txt"
6
7  # Verificação dos ficheiros necessários
8  if [ ! -f "$SentencesFile" ]; then
9      echo "Erro: Ficheiro $SentencesFile não encontrado!"
10     echo "Por favor, executa primeiro o script que cria o ficheiro de frases (Req7)."
11     exit 1
12 fi
13
14 if [ ! -f "$PairsFile" ]; then
15     echo "Erro: Ficheiro $PairsFile não encontrado!"
16     echo "Por favor, executa primeiro o script que cria o ficheiro de pares de frases (Req8)."
17     exit 1
18 fi
19
20 echo "A verificar se todas as frases dos pares existem no ficheiro sentences.txt..."
21 echo ""
22
23 # Extrai a lista de frases únicas (coluna 1)
24 awk '{$NF=""; sub(/[:space:]]+$/, ""); print}' "$SentencesFile" > /tmp/frases_tmp.txt
25
26 # Percorre cada par e verifica se ambas as frases existem
27 awk '
28 BEGIN {
29     while ((getline < "/tmp/frases_tmp.txt") > 0) {
30         frases[$0] = 1
31     }
32 }
33 {
34     n = NF
35     freq = $n
36     $n = ""
37     sub(/[:space:]]+$/, "")
38     split($0, parts, " ")
39     frase1 = parts[1]
40     frase2 = parts[2]
41     if (!(frase1 in frases)) {
42         print "⚠ Aviso: A frase 1 do par \"${frase1}\" não existe em sentences.txt"
43         erro = 1
44     }
45     if (!(frase2 in frases)) {
46         print "⚠ Aviso: A frase 2 do par \"${frase2}\" não existe em sentences.txt"
47         erro = 1
48     }
49 }
50 END {
51     if (erro != 1)
52         print "✅ Todas as frases dos pares existem no ficheiro sentences.txt"
53     else
54         print "⚠ Existem frases nos pares que não estão em sentences.txt"
55 }' "$PairsFile"
56
57 # Limpeza do ficheiro temporário
58 rm -f /tmp/frases_tmp.txt
59

```

## 2.9.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./check_sentences_pairs_req9.sh
A verificar se todas as frases dos pares existem no ficheiro sentences.txt...
```

Todas as frases dos pares existem no ficheiro sentences.txt  
 francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts\$

## 2.10 Limitador de palavras

O sistema Eugénio V3 apenas suporta ficheiros de texto até 250 mil frases e palavras, assim sendo o docente pediu que se desenvolvesse um ou mais scripts que limitassem o output gerado pelos scripts de criação de lista de palavras, lista de frases e dos seus pares. Como é possível observar pela listagem 2.9 o script percorre o ficheiro de texto passado como argumento da execução do script e adiciona cada linha do mesmo se não forem excedidos os limites de palavras ou frases.

```
# limit_entries_req10.sh
1 #!/bin/bash
2
3 # Verificações
4 WORDS_FILE="/home/SUSER/tg1/words.dict.words.txt"
5 WORDS_PAIRS_FILE="/home/SUSER/tg1/words.dict.words_pairs.txt"
6 SENTENCES_FILE="/home/SUSER/tg1/sentences.dict.sentences.txt"
7 SENTENCES_PAIRS_FILE="/home/SUSER/tg1/sentences.dict.sentences_pairs.txt"
8 MAX_ENTRIES=250000
9
10 echo "A aplicar limite de $MAX_ENTRIES entradas..."
11
12 # Limitar palavras se necessário
13 total=$(wc -l < "$WORDS_FILE")
14 if [ $total -gt $MAX_ENTRIES ]
15 then
16   head -n $MAX_ENTRIES "$WORDS_FILE" > "$WORDS_FILE.tmp"
17   mv "$WORDS_FILE.tmp" "$WORDS_FILE"
18   echo "Palavras limitadas a $MAX_ENTRIES"
19 fi
20
21 # Limitar pares de palavras se necessário
22 total=$(wc -l < "$WORDS_PAIRS_FILE")
23 if [ $total -gt $MAX_ENTRIES ]
24 then
25   head -n $MAX_ENTRIES "$WORDS_PAIRS_FILE" > "$WORDS_PAIRS_FILE.tmp"
26   mv "$WORDS_PAIRS_FILE.tmp" "$WORDS_PAIRS_FILE"
27   echo "Pares de palavras limitados a $MAX_ENTRIES"
28 fi
29
30 # Limitar frases se necessário
31 total=$(wc -l < "$SENTENCES_FILE")
32 if [ $total -gt $MAX_ENTRIES ]
33 then
34   head -n $MAX_ENTRIES "$SENTENCES_FILE" > "$SENTENCES_FILE.tmp"
35   mv "$SENTENCES_FILE.tmp" "$SENTENCES_FILE"
36   echo "Frases limitadas a $MAX_ENTRIES"
37 fi
38
39 # Limitar pares de frases se necessário
40 total=$(wc -l < "$SENTENCES_PAIRS_FILE")
41 if [ $total -gt $MAX_ENTRIES ]
42 then
43   head -n $MAX_ENTRIES "$SENTENCES_PAIRS_FILE" > "$SENTENCES_PAIRS_FILE.tmp"
44   mv "$SENTENCES_PAIRS_FILE.tmp" "$SENTENCES_PAIRS_FILE"
45   echo "Pares de frases limitados a $MAX_ENTRIES"
46 fi
47
48 echo ""
49 echo "Resumo final:"
50 echo " Palavras: $(wc -l < "$WORDS_FILE")"
51 echo " Pares de palavras: $(wc -l < "$WORDS_PAIRS_FILE")"
52 echo " Frases: $(wc -l < "$SENTENCES_FILE")"
53 echo " Pares de frases: $(wc -l < "$SENTENCES_PAIRS_FILE")"
```

## 2.10.1 Output esperado

```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$ ./limit_sentences_req10.sh
Entradas atuais em sentences.txt: 34558
Nothing to do: já tem <= 250000 entradas.
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts$
```

## 2.11 Script Opcional (Mestre)

O script **master.sh** funciona como um menu que permite ao utilizador executar facilmente todos os scripts do trabalho, sem ter de os correr um a um manualmente. Quando o script é iniciado, ele mostra uma lista numerada com todas as opções, cada uma correspondente a um requisito do TG1.

O utilizador escolhe um número, e o *master.sh* executa automaticamente o script correto. Depois de cada execução, o programa pausa e espera que o utilizador pressione ENTER, para que consiga ler as mensagens antes de continuar.

O script também inclui uma opção especial que prepara os ficheiros finais para instalar no software Eugénio. Esta opção copia os dicionários criados, muda-lhes o nome para os ficheiros que o Eugénio usa (geral.pal, geral.par, etc.) e coloca tudo numa pasta própria.

Existe ainda uma opção que executa **todos os scripts pela ordem certa**, desde o início até ao fim, garantindo que tudo é feito corretamente.

Em resumo, o **master.sh** serve para:

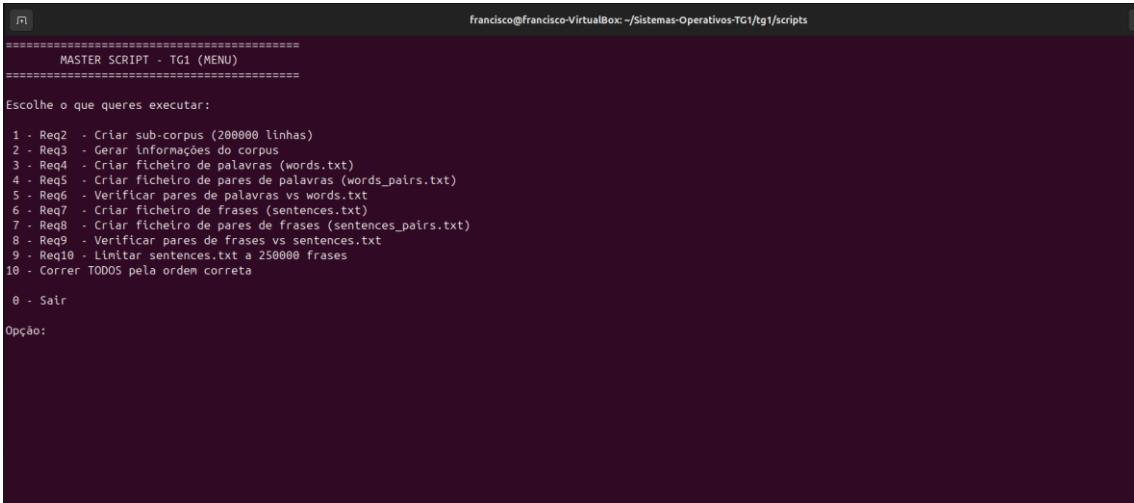
- facilitar a execução dos requisitos;
- evitar erros ao correr scripts manualmente;
- preparar automaticamente os ficheiros finais para o Eugénio;
- organizar todo o processo de forma simples e rápida.

```

# master.sh
1 #!/bin/bash
2
3 # Script Master - Executa scripts do TG1 à escolha do utilizador
4 # Trabalho de Grupo 1 - Sistemas Operativos
5
6 clear
7 echo "*****"
8 echo " MASTER SCRIPT - TG1 Eugénio Italiano"
9 echo "*****"
10 echo
11 echo "Scripts disponíveis:"
12 echo ""
13 echo " 1) Requisito 2 - Criar sub-corpus (200k linhas)"
14 echo " 2) Requisito 3 - Calcular métricas do corpus"
15 echo " 3) Requisito 4 - Criar ficheiro de palavras"
16 echo " 4) Requisito 5 - Criar ficheiro de pares de palavras"
17 echo " 5) Requisito 6 - Verificar palavras nos pares"
18 echo " 6) Requisito 7 - Criar ficheiro de frases"
19 echo " 7) Requisito 8 - Criar ficheiro de pares de frases"
20 echo " 8) Requisito 9 - Verificar frases nos pares"
21 echo " 9) Requisito 10 - Limitar dicionários a 250.000 entradas"
22 echo ""
23 echo " 10) Preparar ficheiros para instalação no Eugénio"
24 echo ""
25 echo " 11) Executar TODOS os scripts (ordem correta)"
26 echo ""
27 echo " 0) Sair"
28 echo ""
29 echo "*****"
30 echo ""
31
32 # Função para executar um script
33 execute_script()
34 {
35     local script_name=$1
36     local description=$2
37
38     echo ""
39     echo "*****"
40     echo ">>> Executando: $description"
41     echo "*****"
42     echo ""
43
44     if [ -f "$script_name" ]; then
45         ./"$script_name"
46         local exit_code=$?
47         echo ""
48         if [ $exit_code -eq 0 ]; then
49             echo "/> Script executado com sucesso!"
50         else
51             echo "X Erro ao executar script (código: $exit_code)"
52         fi
53     else
54         echo "X Erro: Script $script_name não encontrado!"
55     fi
56
57 }
58

```

## 2.11.1 Output esperado



```
francisco@francisco-VirtualBox:~/Sistemas-Operativos-TG1/tg1/scripts
=====
MASTER SCRIPT - TG1 (MENU)
=====

Escolhe o que queres executar:

1 - Req2 - Criar sub-corpus (200000 linhas)
2 - Req3 - Gerar informações do corpus
3 - Req4 - Criar ficheiro de palavras (words.txt)
4 - Req5 - Criar ficheiro de pares de palavras (words_pairs.txt)
5 - Req6 - Verificar pares de palavras vs words.txt
6 - Req7 - Criar ficheiro de frases (sentences.txt)
7 - Req8 - Criar ficheiro de pares de frases (sentences_pairs.txt)
8 - Req9 - Verificar pares de frases vs sentences.txt
9 - Req10 - Limitar sentences.txt a 250000 frases
10 - Correr TODOS pela ordem correta

0 - Sair

Opção:
```

## 2.12 Copiador de ficheiros

Por fim foi pedido que os alunos criassem um script Windows que copiasse os principais ficheiros gerados pelos scripts mencionados anteriormente para os seus respetivos ficheiros no software Eugénio. O objetivo deste requisito é que o Eugénio possa entender e sugerir palavras e frases em italiano ao utilizador.

Para tal foi criado um script no formato ".bat" visto ser o suportado pelo sistema operativo "Windows". O script começa por armazenar o caminho dos ficheiros a copiar em variável, de seguida verifica se o utilizador instalou a versão 64 bits ou 32 e então copia os ficheiros anteriormente guardados em variáveis para os que lhes correspondem na pasta do Eugénio. Há que ter em conta que de maneira ao script funcionar corretamente o utilizador deve substituir o caminho da variável "SOURCE\_DIR" pelo caminho para a pasta onde os outputs gerados anteriormente estão.

```

1  @echo off
2  REM =====
3  REM Instala dicionários Italianos no Eugénio 3.0
4  REM Copia os ficheiros já prontos da pasta C:\Eugenio_IT
5  REM para:
6  REM   - pasta de instalação (Program Files)
7  REM   - pasta do utilizador (AppData\Roaming)
8  REM =====
9
10 chcp 65001 >NUL
11
12 REM Pasta onde estão os ficheiros já prontos
13 set "SOURCE_DIR=C:\Eugenio_IT"
14
15 REM -----
16 REM Detectar diretoria do Eugénio
17 REM -----
18 if exist "C:\Program Files\Eugénio" (
19   set "TARGET_DIR=C:\Program Files\Eugénio"
20 ) else if exist "C:\Program Files (x86)\Eugénio" (
21   set "TARGET_DIR=C:\Program Files (x86)\Eugénio"
22 ) else (
23   echo ERRO: O Eugénio não foi encontrado.
24   pause
25   exit /b
26 )
27
28 REM Diretoria AppData do utilizador
29 set "USER_DIR=%AppData%\LabST2-INESC-ID\Eugénio 3.0"
30
31 echo.
32 echo A copiar ficheiros...
33 echo.
34
35 REM Copiar para Program Files
36 copy /Y "%SOURCE_DIR%\*" "%TARGET_DIR%"
37
38 REM Criar pasta AppData caso não exista
39 if not exist "%USER_DIR%" mkdir "%USER_DIR%"
40
41 REM Copiar para AppData
42 copy /Y "%SOURCE_DIR%\*" "%USER_DIR%"
43
44 echo.
45 echo Ficheiros copiados:
46 echo    -> %TARGET_DIR%
47 echo    -> %USER_DIR%
48 echo.
49 pause

```

## 2.12.1 Output esperado



## Capítulo 3

### Conclusão

A realização deste trabalho permitiu desenvolver um conjunto completo de scripts em bash capazes de transformar um corpus de texto italiano nos dicionários necessários para o funcionamento do software Eugénio V3. Ao longo do projeto foram aplicados conceitos essenciais de sistemas operativos, como automação, manipulação de ficheiros, uso de comandos Unix e criação de pipelines para processamento eficiente de dados.

Todos os requisitos definidos no enunciado foram cumpridos: desde a criação do sub-corpus, à análise estatística do texto, à geração dos ficheiros de palavras, frases e respetivos pares, até à verificação de consistência e limitação de entradas para garantir compatibilidade com o Eugénio. A implementação do **Master Script** permitiu ainda integrar todo o processo num sistema simples e automático, facilitando a execução sequencial de todos os passos.

O trabalho contribuiu para uma compreensão prática sobre como construir dicionários linguísticos a partir de grandes volumes de texto, reforçando competências em automação, análise de dados e programação em shell. No final, foi possível gerar corretamente os ficheiros finais em italiano, prontos para instalação e utilização no sistema Eugénio V3.