



Escola de Engenharia
Universidade do Minho

DEPARTAMENTO DE ENGENHARIA INFORMÁTICA
Mestrado Integrado em Engenharia Informática
Processamento de Linguagens

GAWK

Trabalho prático nº2

Enunciado 4



Marco Gonçalves
A75480



Ricardo Canela
A74568

Braga, 28 de Abril de 2019

Conteúdo

1	Introdução	2
2	Processador de CETEMPúblico	3
2.1	Enunciado 2.4	3
2.2	Descrição do problema	3
2.3	Resolução do problema	3
2.3.1	Alínea a)	6
2.3.2	Alínea b)	7
2.3.3	Alínea c)	8
2.3.4	Alínea d)	9
3	Estrutura HTML	10
4	Desfecho	11

1. Introdução

Foi proposto como trabalho prático da unidade curricular de Processamento de Linguagem, a utilização de *GAWK* para tratamento de informação. Este trabalho prático tem como principais objectivos: Com esta proposta o professor disponibilizou 5 enunciados distintos em que cada grupo escolheria um deles baseado na seguinte formula : $(\text{menor } n^{\circ} \text{ de aluno do grupo } \%5)+1$, que nos nosso caso $(74568\%5)+1 = 4$. No exercício 4 é nos pedido para recebendo o ficheiro *CETEMPúblicoAnotado2018-p1.txt*, (O *CETEMPúblico* é um corpus de aproximadamente 180 milhões de palavras em português europeu) e tratar a informação de forma a efetuar a análise do ficheiro e responder ao requisitos do enunciado:

- Aumentar a experiência de uso do ambiente Linux e de algumas ferramentas de apoio à programação;
- Aumentar a capacidade de escrever Expressões Regulares (ER) para descrição de padrões de frases;
- Desenvolver, a partir de ERs, sistemática e automaticamente *Processadores de Linguagens Regulares*, que filtrem ou transformem textos.
- Utilizar o sistema de produção para filtragem de texto *GWAK*.

Como complemento do problema, os resultados do projeto são apresentados numa estrutura HTML.

2. Processador de CETEMPúblico

2.1 Enunciado 2.4

O formato CETEMPúblico, é um ficheiro dividido por linhas e nessas linhas contém tags em XML que adicionam informação sobre anotação frásica para além dessas linhas contêm linhas com a informação detalhada representados palavras. Pretende-se, assim escrever um sistema de produção para filtragem de texto, GAWK.

- Conta o número de Extratos, Parágrafos e Frases;
- Extrai a lista das multi-word-expressions e respetivo número de ocorrências;
- Calcular a lista dos verbos PT: (Lema, para palavras com pos=V) e respetivo número de ocorrências;
- Determina o dicionário implícito no corpóra - calcule a lista das palavras associando-lhes os possíveis (lema, pos).

2.2 Descrição do problema

O nosso ficheiro está separado por linhas e existem 2 tipos de linhas as linhas para dar início ou fim a conjuntos (tipo xml). Neste trabalho temos que contar os Extratos que começam por `<ext informação desnecessária>` e acabam por `</ext>`, contar os parágrafos que começam por `<p informação desnecessária>` e acabam por `</p>` e contar frases que começam por `<s>` e terminam quando aparece uma linha com `</s>`. Temos ainda que extrair a lista das *multi-word-expression* e os respectivos números de ocorrências. Outra das finalidades deste trabalho é calcular a lista de verbos e respectivas ocorrências, bem como criar um dicionário de palavras onde lhes associamos o (lema, pos).

2.3 Resolução do problema

Para a resolução deste problema apenas precisamos de um programa. Neste programa comecemos no **BEGIN** por definir o *Field Separator (FS)* como `"\t"`, pois as linhas que contêm mais de 1 *Number of Fields (NF)* têm os campos separados por *tabs* e inicializamos a `mweF = 0` que serve para sabermos se estamos dentro de uma multi-word-expression, inicializamos ainda o `p` que representa o número de parágrafos, o `s` que representa o número de frases, e o `ext` que representa o número de extratos. Após o **BEGIN** decidimos dividir as linhas por `NF<5` e `NF>5` isto porque supostamente apenas existem linhas com 1 campo e linhas com 9 campos e como quando tem mais de um campo precisamos de utilizar como o maior campo o \$5.

```

NF<5{
    aux = first($0)
    if(aux == "<mwe" || aux == "<mwe>"){
        mweF = 1
        mweW = ""
    }
    if(aux == "<p" || aux == "<p>") p++
    if(aux == "<s>") s++
    if(aux == "<ext" || aux == "<ext>") ext++
    if(aux == "</mwe>"){
        mweF = 0;
        if(!mweList[mweW]){
            mweList[mweW] = 1;
        } else {
            mweList[mweW]++;
        }
    }
}
}

```

Figura 2.1: NF<5

Durante o processamento das linhas com NF<5, sabemos que estamos perante uma linha que contem informação sobre se começa ou acaba algum tipo de conjunto e como temos conjuntos que contém informação desnecessária separada por um espaço decidimos criar uma função *first* que utiliza a função split para separar uma string e ao receber o retornando assim o primeiro elemento da lista recebido pelo split.

```

function first(s){
    split(s,a," ",b)
    return a[1]
}

```

Figura 2.2: First

Tendo este primeiro elemento este pode ser igual a «mwe», «mwe> o que faz com que tenhamos que por a mweF=1, para quando processar-mos os NF>5 saibamos que estamos dentro de uma multi-word-expression, e ainda inicializamos a mweW=. Caso seja igual a «p» ou igual a «p>» incrementamos o p, caso seja «s» incrementamos o s, caso seja «ext» ou «ext>» incrementamos o ext e caso seja «/mwe>» que significa que acabou a multi-word-expression ponos a mweF=0 e adicionamos a mweW atual á nossa lista de expressões mweList, caso esta já exista na nossa lista incrementamos o numero de aparições.

```

NF>5 {
    if($5 == "V") {
        if(!verbos[$4]){
            verbos[$4] = 1;
        } else {
            verbos[$4]++;
        }
    }
    if($5 != "" && !dicionario[$1] && $5 != "PU" && $5 != "NUM") {
        dicionario[$1] = $4 " ", " $5 "."
    }
    if(mweF == 1){
        mweW = concat(mweW,$1)
    }
}

```

Figura 2.3: NF>5

Durante o tratamento do NF>5 estamos perante uma linha que contem uma palavra e informação á cerca da mesma. O 1º campo contem a palavra, o 4º contem o lema, e o 5º contem a pos(part of speech). Neste processamento fazemos 3 coisas, que são elas adicionar a palavra á lista de verbos caso o \$5 seja igual a "V" ou incrementá-la caso já lá esteja presente, adicionamos a palavra ao nosso dicionário caso o \$5 seja diferente de , "PU" ou "NUM_card", e caso a mweF esteja igual a 1 igualamos a mweW á nossa função concat definida por nós que recebe 2 strings e devolve apenas uma que é o resultados da sua junção e neste caso recebe a antiga mweW e o \$1.

```
function concat(a,b){
    return a " " b
}
```

Figura 2.4: concat

```
END {
    #Indice
    print indice() > "indice.html"
    #Estatísticas a)
    print begin() > "data/stats.html"
    print title("Estatísticas") > "data/stats.html"
    print "<li>Parágrafos: " p "</li>\n" > "data/stats.html"
    print "<li>Frases: " s "</li>\n" > "data/stats.html"
    print "<li>Extratos: " ext "</li>\n" > "data/stats.html"
    print end() > "data/stats.html"
    #MWE b)
    print begin() > "data/mwe.html"
    print title("Multi-word-expressions") > "data/mwe.html"
    for(k in mweList){
        print "<li>" k " : " mweList[k] "</li>" > "data/mwe.html"
    }
    print end() > "data/mwe.html"
    #Verbos c)
    print begin() > "data/verbos.html"
    print title("Verbos") > "data/verbos.html"
    for(k in verbos){
        print "<li>" k " : " verbos[k] "</li>" > "data/verbos.html"
    }
    print end() > "data/verbos.html"
    #Dicionário d)
    print title("Dicionário") > "data/dicionario.html"
    print begin() > "data/dicionario.html"
    for(k in dicionario){
        if(dicionario[k] != "") print "<li>" k " : " dicionario[k] "</li>" > "data/dicionario.html"
    }
}
```

Figura 2.5: Criação de html

```

function link(u,t){
    return "<li><a href='" u "'> t "</a></li>\n"
}
function title(u){
    return "<h1>" u "</h1>\n"
}
function begin(){
    return "<html>\n\t<head>\n\t\t<meta charset=\"UTF-8\"/></h
ead>\n\t<body>\n\t\t<ul>\n\t\t<link rel=\"stylesheet\" href=\"https://
stackpath.bootstrapcdn.com/bootstrap/4.3.1/css/bootstrap.min.css\"
    integrity=\"sha384-gg0yR0iXCbMQv3Xipma34MD+dH/1fQ784/j6cY/iJTQU0h
cWr7x9JvoRxT2MZw1T\" crossorigin=\"anonymous\">"
}
function end(){
    return "\n\t\t</ul>\n\t</body>\n</html>"
}
function indice(){
    return begin() "\n" title("Índice") link("data/stats.html"
,"Estatísticas") link("data/mwe.html","Multi-word-expressions") li
nk("data/verbos.html","Lista de Verbos") link("data/diccionario.htm
l","Diccionario") end()
}

```

Figura 2.6: Funções auxiliares aos prints

Por fim no **END** geramos todos o html recorrendo a prints em ficheiros da informação até aqui guardada.

2.3.1 Alínea a)

Estatísticas

- Parágrafos: 273177
- Frases: 616634
- Extratos: 115504

Figura 2.7: Estatísticas

2.3.2 Alínea b)

Multi-word-expressions

- pecados mortais : 8
- em favor : 139
- a longo prazo : 146
- receita bruta : 2
- contavam-se pelos dedos : 1
- De mistura : 1
- às três pancadas : 1
- Por volta de : 7
- Maré baixa : 1
- câmara escura : 1
- caixa torácica : 1
- tudo isto : 481
- previdência social : 3
- ao cabo : 59
- alma penada : 1
- prata da casa : 15
- trabalhos forçados : 24
- um tal de : 2
- a céu aberto : 54
- automóveis ligeiros : 28
- além disso : 331
- frente a : 206
- em vias de : 268
- largou mão : 2

Figura 2.8: Pequena amostra de Multi-word-expressions

2.3.3 Alínea c)

Verbos

- auto-exilar : 1
- vastovasto : 1
- de=Xir : 4
- aedificander : 2
- pderer : 1
- delapidar : 24
- auto-representar : 1
- catrapiscar : 1
- confluir : 21
- alcandorar : 10
- coincir : 1
- exagerar : 162
- confudar : 1
- desobrigar : 8
- esfiapar : 1
- accionar : 179
- trancar : 9
- descriar : 1
- malograr : 22
- s/ : 10
- semear : 99
- homolgar : 1
- tele-surfar : 1
- mal-regressar : 1

Figura 2.9: Pequena amostra dos Verbos

2.3.4 Alínea d)

Dicionário

- Sureno: Sureno, PROP.
- angélicos: angélico, ADJ.
- indignasse: indignar, V.
- Delauney: Robert=Delauney, PROP.
- Desligou-se: desligar+se, V+PERS.
- mongolismo: mongolismo, N.
- Cadáveres: Cadáveres=Esquitos, PROP.
- valerá: valer, V.
- estatelar: estatelar, V.
- ouriço: ouriço, N.
- neutro: neutro, ADJ.
- Tenciona: tencionar, V.
- Laffite176: Laffite176, PROP.
- Waterford: Waterford, PROP.
- entalaria: entalar, V.
- quadrigas: quadriga, N.
- Gaucher: Gaucher, PROP.
- paracetamol: paracetamol, N.
- Inger: Inger, PROP.
- Conforlimpa: Conforlimpa, PROP.
- sustentou: sustentar, V-REGENTE.
- Doelp: Michael=Doelp, PROP.
- proibimos: proibir, V.
- submeterem: submeter, V.

Figura 2.10: Pequena amostra do Dicionário

3. Estrutura HTML

A escrita dos ficheiros HTML é começada pelo programa *esquema1* que coloca no ficheiro *index.html* quatro ancoras. A primeira reporta para o ficheiro *stats.html*, que contem a informação, para responder á alinea a), respetivamente:

- Parágrafos: numero de paragrafos;
- Frases: numero de frases;
- Extratos: numero de extratos;

Da mesma maneira é necessário fazer a referencia para o ficheiro *mwe.html* que contem todas as multi-world-expressions e o respetivo numero de ocorrências. As próximas duas ancoras dizem respeito aos ficheiros que contêm a lista de todos os verbos e o respetivo numero de ocorrências. Por ultimo uma ligação ao ficheiro *dicionario.html* que contém a lista de palavras e os possíveis (lema,pos).

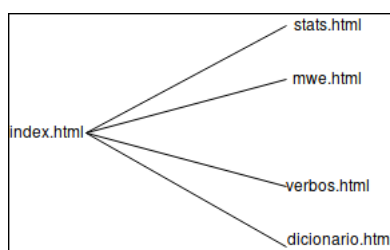


Figura 3.1: Árvore html

4. Desfecho

A realização de um projeto prático numa unidade curricular é sempre benéfica, neste projeto desenvolveram-se expressões regulares para identificar e alterar o texto, ajudando a consolidar toda a matéria lecionada, tão ou mais importante, dá-nos uma melhor perceção do impacto que o que estudamos tem em sistemas mais completos, construindo-se assim um sistema de produção para filtragem de texto GAWK, para além da experiência adquirida no uso do ambiente Linux.

hapterConclusão