



Universidade do Minho

Processamento de Linguagens

Trabalho Prático 1

Março 2019

André Guilherme Nunes Viveiros, A80524

César Augusto da Costa Borges, A81644

Luís José Rodrigues da Silva Macedo, A80494

Resumo

Este documento consiste no relatório correspondente ao primeiro trabalho prático realizado no âmbito da Unidade Curricular de Processamento de Linguagens, do curso de Mestrado Integrado em Engenharia Informática da Universidade do Minho, no ano letivo de 2018/2019. Este primeiro trabalho consiste em desenvolver uma solução a um desafio proposto pelos docentes da Unidade Curricular.

Índice

1. Introdução	1
1.1. Objetivos	1
1.2. Contextualização e caso de estudo	1
1.3. Estrutura do relatório	2
2. Fundamentação	3
3. Normalização	4
3.1. Limpeza dos Dados	4
3.2. Estrutura de Dados	5
4. Expressões Regulares	7
5. HTML	11
5.1. Noticias	11
5.2. Tags	12
6. Conclusão	15
Lista de Siglas e Acrónimos	16

Índice de Figuras

1. Exemplo de uma noticia.	5
2. Exemplo de um ficheiro HTML com as noticias.	11
3. Exemplo de um ficheiro HTML com uma noticias.	12
4. Exemplo de um ficheiro HTML com as tags e as suas repetições. . . .	13
5. Exemplo de um ficheiro HTML com as noticias que contêm a tag "ignorância.	14

1. Introdução

Este relatório apresenta e documenta a totalidade do trabalho desenvolvido no âmbito do primeiro trabalho prático da Unidade Curricular de Processos de Linguagem, do curso de Mestrado Integrado em Engenharia Informática da Universidade do Minho, no ano letivo de 2018/2019.

Este capítulo indica os objetivos, contextualiza e apresenta o caso de estudo, descrevendo também as motivações e objetivos do projeto.

1.1. Objetivos

Este trabalho tem como objetivos:

- A resolução de um desafio usando *ER* e *Flex*;
- A familiarização dum ambiente *Unix*;
- A familiarização das *ER*;
- O desenvolvimento de *Processadores de Linguagens Regulares*;
- A utilização do *Flex*

1.2. Contextualização e caso de estudo

Dos 8 enunciados propostos pelos docentes da *UC*, este grupo escolheu o enunciado 8, de acordo com a formula proposta. Considera-se o caso de

estudo a limpeza e normalização de um ficheiro com milhares de notícias de um jornal angolano. Através dessa normalização, espera-se a criação de um ficheiro *HTML* para cada notícia, a criação de um ficheiro *HTML* com todas as *tags*, em que cada *tag* é um *link* para outro ficheiro com *links* para todas as notícias com essa *tag*. Por ultimo, a criação de uma lista com todas as *tags* encontradas e numero de ocorrências.

1.3. Estrutura do relatório

Este documento apresenta a seguinte estrutura:

- No Capítulo 2 (fundamentação) refere-se as ideais gerais para o desenvolvimento da solução;
- O relatório é concluído no Capítulo 6 com observações relevantes.

2. Fundamentação

Tendo-se apresentado o caso de estudo e identificado as motivações e objetivos para o desenvolvimento da solução ao desafio em questão, este capítulo fundamenta a sua construção.

De forma a limpar as notícias usou-se *ER* com o objetivo de encontrar todas as características pertinentes de cada notícia, tais como:

- *Tags*;
- *ID*;
- Categoria;
- Data;
- Autor;
- Artigo.

Em conjunto com a limpeza, ocorre a normalização, ou seja, guarda-se os dados numa estrutura de dados. Após ter todas as notícias normalizadas, cria-se todos os ficheiros *html* necessários.

3. Normalização

Tendo-se a fundamentação do sistema, este capítulo explica o processo de normalização do artigo.

3.1. Limpeza dos Dados

Após a observação dos artigos do Jornal Angolano reparou-se que todas as notícias têm uma estrutura semelhante á descrita na Figura 1, ou seja, as informações de cada noticia encontram-se por esta ordem:

1. *tags*;
2. *ID*;
3. categoria;
4. titulo;
5. autor e a data;
6. texto do artigo.

Dito isto, decidiu-se criar uma estrutura onde fosse possível guardar a informação mais pertinente da mesma.


```

<!--=====2014/140-milhoes-de-dolares-para-dois-petroleiros/index.html -->
<sub>
#TAG: tag:{Sonangol} tag:{coreia do sul} tag:{petroleiros}
#ID:{post-5702 post type-post status-publish format-standard has-post-thumbnail hentry category-economia category-destaque tag-sonangol tag-coreia-do-sul tag-petroleiros}
EconomiaDestaque

140 milhões de dólares para dois petroleiros

PARTILHE VIA:
#DATE: [116eb] Redacção F8 - 11 de Dezembro de 2014
140 milhões de dólares para dois petroleiros - Folha 8

A Sociedade Nacional de Combustíveis de Angola (SONANGOL) assinou hoje, em Seul, Coreia do Sul, um acordo para a construção de dois navios petroleiros, orçados em 140 milhões de dólares, que deverão ser entregues em 2017.

O acordo foi rubricado pelo Presidente do Conselho da Administração da SONANGOL, Francisco de Lemos Jose Maria e o seu homólogo sul-coreano, da DSME, Ko Jae-Ho.

Uma nota de imprensa da concessionária angolana refere que os navios designados Suezmax, terão 156,290 toneladas métricas cada um, 274 metros de comprimento, 48 metros de largura e uma altura de 23,7 metros.

O documento salienta que a empresa norueguesa DNV-GL está encarregue de acompanhar e classificar o processo de fabrico das referidas embarcações, garantindo que são respeitados os parâmetros convencionais.

A SONANGOL possui presentemente em construção nos mesmos estaleiros dois navios sonda de perfuração, que deverão ser entregues, respectivamente, no quarto trimestre de 2015 e no segundo trimestre de 2016.

O investimento, segundo a concessionária angolana, visa reforçar a exploração petrolífera em águas profundas, ultra-profundas e no pré-sal e o objectivo da construção desses dois navios, adianta a nota, sublinhando que o recurso à Coreia do Sul foi feito por se tratar de um país líder mundial na engenharia de construção naval.

Angola é o segundo maior produtor de petróleo da África subsaariana.

Etiquetas: Sonangolcoreia do sulpetroleiros
</sub>

```

Figura 1: Exemplo de uma noticia.

3.2. Estrutura de Dados

De forma a guardar todos os artigos criou-se uma estrutura de dados capaz de guardar toda a informação de um artigo, isto é, uma estrutura chamada *Noticia*, que guarda toda a informação da noticia e uma estrutura chamada *Tag*, que guarda uma *tag*, o numero de repetições e os *IDs* das noticias onde aparece.

Noticia:

- char *id, o *ID*;
- char *title, o titulo
- char *category, a categoria;
- char *date, o autor e a data;
- char **tags, o conjunto das *tags*;
- int lenght_tags, o numero de *tags*;

- char *text, o texto.

Tag:

- char *tag, a tag;
- int rept, o numero de repetições;
- char **id_noticia, o conjunto dos *IDs* das noticias em que aparece.

Para que seja possível guardar todas as noticias lidas do ficheiro, decidiu-se usar estruturas auxiliares, da biblioteca *GLib*. Essas estruturas são duas *GHashTable*, **tags** e **noticias**.

	Key	Value
tags	<i>tag(char *)</i>	<i>Tag</i> (estrutura de dados)
noticias	<i>ID da noticia(char *)</i>	<i>Noticia</i> (estrutura de dados)

De modo a inserir e percorrer estas estruturas externa usou-se as seguintes funções, da mesma biblioteca:

- *g_hash_table_insert*;
- *g_hash_table_lookup*;
- *g_hash_table_foreach*;

De forma a ter um conjunto de noticias por *tag*, o preenchimento da *GHashTable* tags é feito depois de ter percorrido todo o ficheiro e ter guardado todas as noticias, ou seja, percorre-se as noticias e por cada noticia percorre-se as *tags* dessa noticia, inserindo-as na *GHashTable* juntamente com o *ID* da noticia e aumentando o numero de repetições.

4. Expressões Regulares

Após as estruturas estarem bem definidas, construiu-se um autômato para modelar o sistema de acordo com a nossa estrutura.

Este autômato contém 9 autômatos mais pequenos:

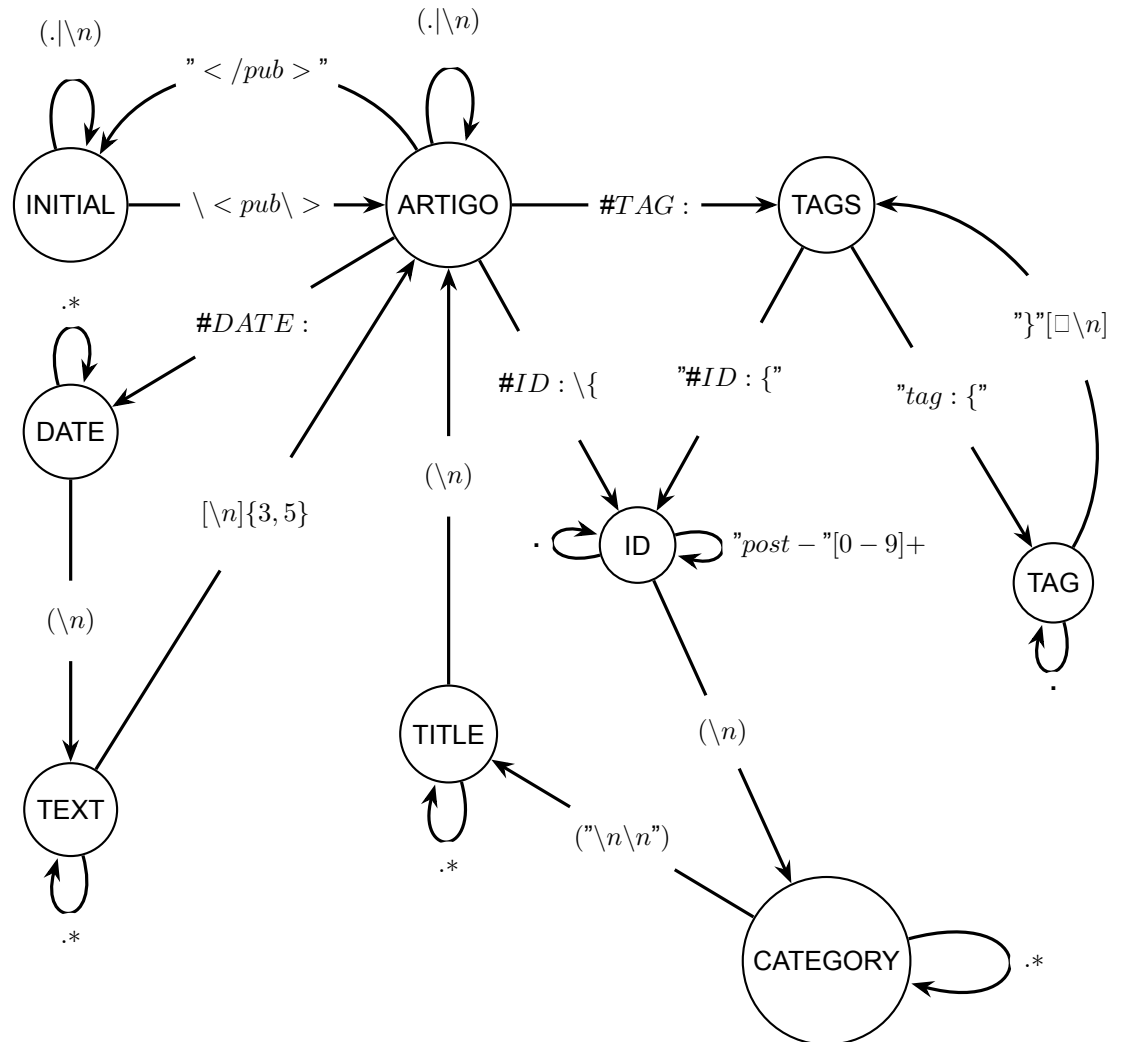
- INITIAL;
- ARTIGO;
- TAGS;
- TAG;
- ID;
- CATEGORY;
- TITLE;
- DATE;
- TEXT;

Quando o programa começa, tem à sua disposição as seguintes *ERs*:

1. $\backslash < pub \backslash >$
2. $< ARTIGO > " < /pub > "$
3. $< ARTIGO > \#TAG :$
4. $< ARTIGO > \#DATE :$
5. $< ARTIGO > \#ID : \backslash \{$
6. $< ARTIGO > (. \backslash n)$
7. $< TAGS > "tag : \{ "$
8. $< TAGS > "ID : \{ "$
9. $< TAG > " \} "[\square \backslash n]$, onde \square é igual a um espaço

10. $\langle TAG \rangle .$
11. $\langle ID \rangle "post - "[0 - 9] +$
12. $\langle ID \rangle (\backslash n)$
13. $\langle ID \rangle .$
14. $\langle CATEGORY \rangle .*$
15. $\langle CATEGORY \rangle (" \backslash n \backslash n ")$
16. $\langle TITLE \rangle .*$
17. $\langle TITLE \rangle (\backslash n)$
18. $\langle DATE \rangle (\backslash n)$
19. $\langle DATE \rangle .*$
20. $\langle TEXT \rangle .*$
21. $\langle TEXT \rangle [\backslash n]\{3, 5\}$
22. $(.|\backslash n)$

Resumindo as listas anteriores, tem-se o seguinte automato com todas as *ERs* e as respectivas ações:



Este autômato procura o início da notícia, entrando no estado ARTIGO e ignora tudo o resto. Dentro do ARTIGO ele procura o início das *tags*, entrando em TAGS, dos *IDs*, entrando em ID e da data e autor, entrando em DATE, ignorando o resto.

O autómato TAGS procura o início de uma *tag*, que se for encontrada entra em TAG, e o início do *IDs*. O início do *ID* é procurado duas vezes pois, durante a análise do ficheiro com as notícias foram encontradas algumas exceções, em que uma delas não continha *tags*, e desta forma se não houver *tags* os *IDs* são encontrados no autómato ARTIGO.

Dentro do autómato TAG é procurado uma *tag* carácter a carácter devido ao problema intitula "*Black Hole*", ou seja, é encontrado o texto com mais caracteres e o pretendido é o de menor caracteres. Quando encontra o fim da *tag* volta a TAGS.

Se forem encontrados os *IDs* é procurado o *ID* que siga a forma "*post-X*", onde X é um número positivo. Os restantes *IDs* são ignorados. Se for encontrado uma mudança de linha entra-se em CATEGORY.

A categoria entra-se numa só linha, logo guarda-se essa linha toda. Quando for encontrado duas mudanças de linha seguidas entra-se em TITLE.

O processo de procura pelo título é o mesmo que o da categoria com uma diferença, quando encontra apenas uma mudança de linha volta ao autómato ARTIGO.

Como a data e o autor está contida numa só linha, é apanhada a linha toda removendo os primeiro nove caracteres (não são pertinentes). Após ser encontrado uma mudança de linha, entra-se em TEXT para procurar o texto da notícia.

Para guardar o texto da notícia são apanhadas todas as linhas até encontrar três ou quatro ou cinco mudanças de linha seguidas, voltando ao autómato ARTIGO.

5. HTML

5.1. Notícias

Um dos casos de estudo era a criação de ficheiros *HTML* para apresentar as notícias e as *tags*. Para tal foram criados um pagina *HTML* com uma lista de links apresentados pelos títulos das notícias, para os ficheiros individuais *HTML* das notícias. Pode ser observado um exemplo na Figura 2.

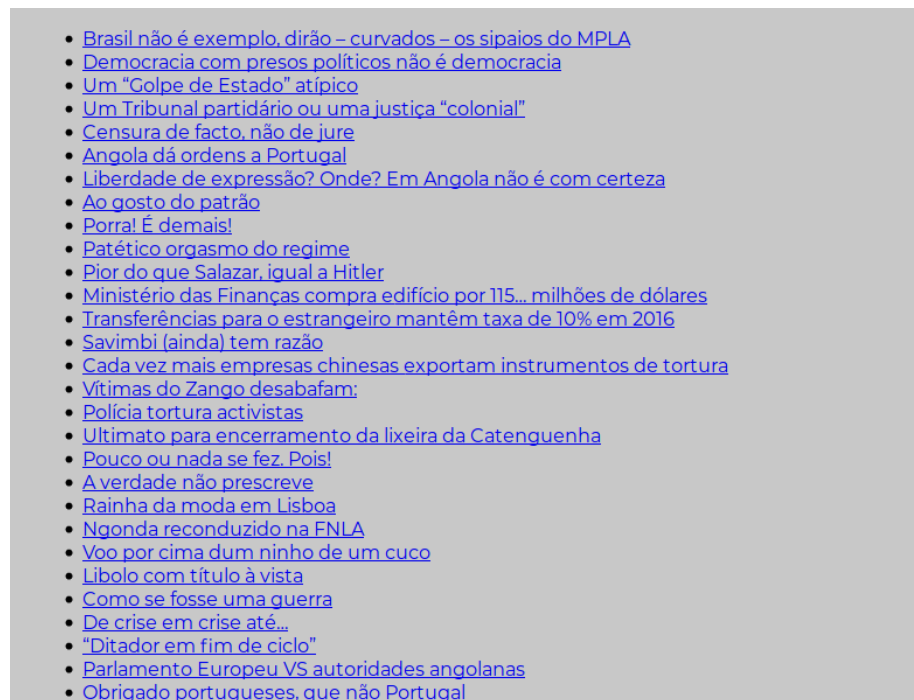


Figura 2: Exemplo de um ficheiro HTML com as notícias.

Cada noticia foi estruturada com o titulo no topo da pagina, seguindo o autor e data, a categoria, a lista das *tags* e por ultimo o texto da noticia. Na Figura 3 é apresentado um exemplo de uma noticia.

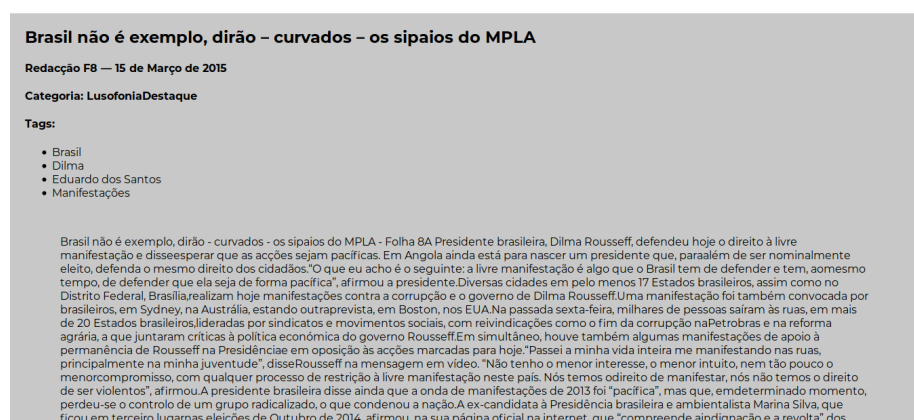


Figura 3: Exemplo de um ficheiro HTML com uma noticia.

5.2. Tags

A pagina com todas as *tags* contem uma lista com todas as *tags* e o numero de repetições, Figura 4. Para alem do nome da *tag* e do numero de repetições, o nome da *tag* está ancorado a uma pagina com uma lista das noticias onde essa *tag* aparece, Figura 5.

Tag	Number of repetitions
Processo de Paz	1
confiança	2
casa-ce. cabinda	1
gás	6
construtoras	1
slazar	1
um morto	1
genocida	1
puto	1
nativos	1
Chico Pobre	1
resultados	15
Mpindi André	1
duplicar	1
convocação	1
morro do moco	1
relógio	1
tony neves	1
Belas Shopping	1
José Saramago	1
Narciso Benedito	1
mobilidade	1
prejuízos	4
Matchedje	1

Figura 4: Exemplo de um ficheiro HTML com as tags e as suas repetições.

- [Nada pior do que ser africano e... albino](#)
- [A ecolália escatológica do Bolha](#)
- [É crime ser albino?](#)
- [Se a ignorância pagasse impostos, A Bola seria o maior pagador](#)
- [O professor Goebbels](#)
- [O papel de embrulho da loja do chinês](#)
- [Os fundamentos macroeconómicos](#)
- [Os "intelequetais"](#)
- [Ignorância voluntária afasta chilenos da realidade africana](#)
- [Pequenas e médias empresas têm medo de África](#)
- [Ser africano não é crime](#)
- [O pelotão do disparate](#)
- [A imbecilidade e o processo](#)
- [Burrice judicial ordena a prisão do morto e de Lídia Amões](#)

Figura 5: Exemplo de um ficheiro HTML com as notícias que contêm a tag "ignorância".

6. Conclusão

Em conclusão, o programa desenvolvido é capaz de processar um ficheiro com notícias, que seguem uma estrutura específica, guardando as informações pertinentes de cada notícia numa estrutura de dados, através de *ERs* de da ferramenta *Flex* e criando ficheiros *HTML* com:

- Uma lista com o título das notícias, ancorado aos ficheiros individuais das notícias;
- As informações guardadas de uma notícia, para todas as notícias;
- Uma lista com o nome das *tags*, ancorado a um ficheiro com a lista de notícias dessa *tag*;
- Uma lista com o título das notícias com uma certa *tag*, ancorado ao ficheiro da notícia.

Para além disso, o grupo consolidou os conhecimentos aprendidos sobre *ERs*, uso da ferramenta *Flex*.

Lista de Siglas e Acrónimos

ER	Expressões Regulares (<i>Regular Expressions</i>)
UC	Unidade Curricular
ID	Identificador (<i>Identity</i>)
XML	<i>Extensible Markup Language</i>