



Universidade do Minho
Escola de Ciências

Universidade do Minho
Licenciatura em Ciências da Computação

Processamento de Linguagens e Compiladores

Trabalho Prático 1

BibTeXPro, Um processador de BibTeX

Grupo 13

Pedro Faria
A72640

João Gouveia
A87995

João Goulart
A82643

Novembro 2021

Conteúdo

1	Introdução	2
2	Descrição da abordagem seguida	2
3	Código	8
4	Exemplos de utilização	10

1 Introdução

No âmbito da Unidade Curricular de Processamento de Linguagens e Compiladores, desenvolvemos o projeto **BibTeXPro, Um processador de BibTeX** que propõe a implementação de um processador capaz de filtrar todas as entradas de uma base de dados textual. Este documento encontra-se estruturado de forma a que seja feita a análise de cada uma das tarefas executadas, bem como os métodos e abordagens que o nosso grupo utilizou para as executar, seguida de exemplos de funcionamento (incluindo Inputs e respectivos Outputs)

Como forma de familiarização com o formato BibTeX foi sugerida a consulta do ficheiro exemplo-utf8.bib. É-nos ainda explicado que as categorias de referência são identificadas pela primeira palavra logo após o carácter "@". Assim sendo, começamos por efetuar a **Análise do documento BibTeX** acima referido.

2 Descrição da abordagem seguida

Após analisar o documento BibTeX, procedemos para a realização das tarefas que devíamos executar. Abordaremos, então, cada uma destas tarefas, explicando todos os pensamentos que tivemos e detalhes do seu desenvolvimento.

a) Contagem das categorias que ocorrem no documento

Tendo em mente que as categorias de referência são identificadas pela primeira palavra após o carácter "@", rapidamente estabelecemos um método para a identificação das categorias presentes no documento BibTeX. Para tal, formulámos uma expressão regular que nos permitisse encontrar todas as palavras do documento que se encontrassem logo a seguir a um carácter "@". Chegámos, então, à seguinte Expressão Regular:

```
1 pattern\_categorias = r'(@[\w]+{(.+\n)+}\n\n) '
```

Para nos auxiliar em algumas das alíneas seguintes à a), utilizámos um dicionário *dict_cat*, que tem o nome da categoria (*string*) como chave, como valor tem uma lista de dois elementos cujo primeiro elemento é a incidência (*integer*) e o segundo um dicionário com a chave da entrada como chave (*string*) e valor uma lista de dois elementos, cujo primeiro elemento é uma lista dos autores e o segundo é o título da entrada. Mas por enquanto, só nos iremos focar na incidência.

```

1 dict_cat = dict()
2 for categoria in categorias:
3     categoria = categoria[0]
4     categoria_nome = re.findall(pattern_nome_cat,categoria)[0].lower()
5     #(...)
6     incidencia = 1
7     #(...)
8     if categoria_nome in dict_cat:
9         incidencia += dict_cat[categoria_nome][0]
10    #(...)
11    dict_cat[categoria_nome] = [incidencia,entradas]

```

Como podemos observar no excerto de código mostrado, a variável de incidência é inicializada com valor 1. Caso esta categoria já exista no dicionário, é adicionado o valor já presente à variável de incidência, de forma a atualizá-la.

b) Filtrar a chave, os autores e o título associados a cada categoria

Utilizando o mesmo método da alínea a), definimos uma Expressão Regular que nos auxiliasse a filtrar apenas a chave, os autores e o título de cada categoria. Para esse efeito, utilizamos as Expressões Regulares

```

1 pattern_chave = r'(?<={} [\w:-]+(=?=\n) '
2 pattern_autor = r'author *= *({|"}{*(([\w.\s]|,|-|\\\\"|\\\\~)+)(?=}|",\n) '
3 pattern_titulos = r' title *= *({|"}((.*)|(.*\n.*)|(\n.*))({|"})'

```

para filtrar as chaves, autores e títulos, respectivamente.

Estas expressões foram utilizadas de modo a conseguirmos extrair os dados para o dicionário criado na alínea a).

```

1 for categoria in categorias:
2     #(...)
3     chave = re.findall(pattern_chave,categoria)[0]
4     autor = re.findall(pattern_autor,categoria, re.IGNORECASE)[-1][1]
5     autores = re.split(r' *and +| +and\n *',autor,100,re.IGNORECASE)
6     autores = addAutor(dict_aut, autores)
7     #(...)
8     entradas = dict()
9
10    try:
11        titulo = re.findall(pattern_titulos,categoria, re.IGNORECASE)[0]

```

```

12     except:
13         continue
14
15     if categoria_nome in dict_cat:
16         #(...)
17         entradas = dict_cat[categoria_nome][1]
18
19     entradas[chave] = [autores,titulo]
20     dict_cat[categoria_nome] = [incidencia,entradas]

```

Usamos duas estratégias adicionais para realizarmos esta extração de dados:

- a função *addAutor*, que foi criada de forma a separar os vários autores extraídos com a expressão regular e guardá-los num dicionário (sendo a chave o nome de um autor e o valor os outros autores com quem ele trabalhou) e numa lista.

```

1     def addAutor(dic_aut,autores):
2         for x in range(len(autores)):
3             autores[x] = re.split(r', ',autores[x],re.IGNORECASE)
4             if len(autores[x]) == 2:
5                 autores[x] = autores[x][1] + ' ' + autores[x][0]
6             else:
7                 autores[x] = autores[x][0]
8
9             if '' in autores:
10                 autores.remove('')
11
12             for autor1 in autores:
13                 for autor2 in autores:
14                     autor1 = ' '.join(re.split(r'\n+ ',autor1,re.IGNORECASE))
15                     autor2 = ' '.join(re.split(r'\n+ ',autor2,re.IGNORECASE))
16                     if autor1 != autor2:
17                         if autor1 not in dic_aut:
18                             dic_aut[autor1] = list()
19                             dic_aut[autor1].append(autor2)
20
21         return autores

```

- lidar com uma exceção relacionada com a expressão regular dos títulos, que não funcionava para apenas uma entrada no ficheiro de exemplos.

```

1     try:
2         titulo = re.findall(pattern_titulos,categoria, re.IGNORECASE)[0]

```

```

3     except:
4         continue

```

Depois de obtermos o dicionário *dict_cat* completo, entregamo-lo como input na função *html_builder()* que estrutura o ficheiro HTML. A nossa decisão final foi criar uma única tabela que engloba as alíneas a) e b). A tabela contém duas filas com duas colunas para a informação das categorias pedida na alínea a) e imediatamente a baixo tem todos os livros relacionados com a categoria específica, organizados por chave, autores e título.

```

1     #(...)
2     for categoria in dict_cat:
3         html += '''
4             <tr>
5                 <th class="bortop">Categoria</th>
6                 <td class="tabcen bortop" colspan="2">{0}</td>
7             </tr>'''.format(categoria)
8         html += '''
9             <tr>
10                <th class="borbot">Incidencia</th>
11                <td class="tabcen borbot" colspan="2">{0}</td>
12            </tr>'''.format(str(dict_cat[categoria][0]))
13        html += '''
14            <tr>
15                <th>Chave</th>
16                <th>Autores</th>
17                <th>Titulos</th>
18            </tr>
19        #(...)

```

Escrevemos também uma pequena porção de código CSS para a tabela ser de mais fácil leitura.

```

1     .tabcen {
2         text-align: center;
3     }
4     td,th {
5         border: 1px solid grey;
6     }
7     th {
8         background-color: #80808054;
9     }

```

```

10     .bortop {
11         border-top: 5px solid black;
12     }
13     .borbot {
14         border-bottom: 2px solid black;
15     }

```

c) Criar um índice de autores, mapeando cada autor nos respectivos registros

Para realizarmos esta alínea, também utilizamos o mesmo ciclo utilizado nas alíneas anteriores, aproveitando também a lista de autores e as informações já retiradas (chave, título, nome da categoria).

```

1     for autor in autores:
2         if autor.lstrip() not in dict_aut_indice:
3             dict_aut_indice[autor.lstrip()] = list()
4             dict_aut_indice[autor.lstrip()].append([chave,titulo[1],categoria_nome])

```

Em seguida, utilizamos a informação adicionada ao dicionário para escrever uma *string* ordenada por autor alfabeticamente organizada da seguinte forma:

```

1     Nome de autor (x entradas):
2         "Título de livro 1", chave, nome da categoria
3         "Título de livro 2", chave, nome da categoria
4         ...
5         "Título de livro x", chave, nome da categoria

```

d) Construção de um grafo que mostre todos os autores que publicaram com um dado autor

Usando o dicionário *dict_aut* definido na alínea b), definimos a seguinte função que constroi um grafo em DOT com as relações entre autores.

(Consideramos como autores que colaboram frequentemente autores que tivessem escrito mais que dois livros juntos.)

```
1  def graph_builder(dict_aut,autor):
2      for autor1 in dict_aut:
3          dict_aut_temp = dict_aut[autor1]
4          for autor2 in dict_aut_temp:
5              if dict_aut[autor1].count(autor2) < 2:
6                  dict_aut[autor1].remove(autor2)
7          dict_aut[autor1] = set(dict_aut[autor1])
8
9      graph_dot = "digraph G {\n"
10
11     if autor:
12         for autor2 in dict_aut[autor]:
13             graph_dot += "{0}" -> "{1}"\n'.format(autor,autor2)
14     else:
15         for autor1 in dict_aut:
16             for autor2 in dict_aut[autor1]:
17                 graph_dot += "{0}" -> "{1}"\n'.format(autor1,autor2)
18
19     graph_dot += "}"
20
21     return graph_dot
```

Caso não tenha sido fornecido nenhum input na consola, irá ser gerado um grafo com todos os autores. Se for fornecido o input, o código irá só construir o grafo a partir do autor indicado.

3 Código

Abaixo encontra-se o código python do nosso trabalho na sua totalidade.

```
1  #!/usr/bin/python
2  # -*- coding: utf8 -*-
3
4  import re
5  from autor_builder import autor_builder
6  from html_builder import html_builder
7  from indice_builder import indice_builder
8  from graph_builder import graph_builder
9
10 f = open('exemplo-utf8.bib','r', encoding="utf8")
11
12 pattern_categorias = r'(@[\w]+{(.+\n)+}\n\n)'
13 pattern_nome_cat = r'(?<=@) [\w]+(?:={}) '
14 pattern_chave = r'(?<={}) [\w:-]+(?:=, \n)'
15 pattern_autor = r'author *= *({|"})*(((\w.\s|,|-|\\\'|\\\'~)+)(?="|",\n)'
16 pattern_titulos = r' title *= *({|"})((.*)|(.*\n.*)|(\n.*))({|"})'
17
18 categorias = re.findall(pattern_categorias,f.read())
19 f.close()
20 dict_aut = dict()
21 dict_cat = dict()
22 dict_aut_indice = dict()
23 #dict_cat {cat: [incidencia, {chave:[[autores],titulo]]}}
24
25 for categoria in categorias:
26     categoria = categoria[0]
27     categoria_nome = re.findall(pattern_nome_cat,categoria)[0].lower()
28
29     chave = re.findall(pattern_chave,categoria)[0]
30     autor = re.findall(pattern_autor,categoria, re.IGNORECASE)[-1][1]
31     autores = re.split(r' *and +| +and\n *',autor,100,re.IGNORECASE)
32     autores = autor_builder(dict_aut, autores)
33     incidencia = 1
34     entradas = dict()
35
36     try:
37         titulo = re.findall(pattern_titulos,categoria, re.IGNORECASE)[0]
38     except:
39         continue
40
41     if categoria_nome in dict_cat:
```

```

42     incidencia += dict_cat[categoria_nome][0]
43     entradas = dict_cat[categoria_nome][1]
44
45     entradas[chave] = [autores,titulo]
46     dict_cat[categoria_nome] = [incidencia,entradas]
47
48     for autor in autores:
49         if autor.lstrip() not in dict_aut_indice:
50             dict_aut_indice[autor.lstrip()] = list()
51             dict_aut_indice[autor.lstrip()].append([chave,titulo[1],categoria_nome])
52
53 html = html_builder(dict_cat)
54 open('index.html','w').write(html)
55 print("Ficheiro HTML criado com sucesso.")
56
57 indice_autores = indice_builder(dict_aut_indice)
58 open('indice_autores.txt','w').write(indice_autores)
59 print("Ficheiro TXT criado com sucesso.")
60
61 autor = input('Autor: ')
62 graph = graph_builder(dict_aut,autor)
63 open('graph.dot','w').write(graph)
64 print("Ficheiro DOT criado com sucesso.")

```

4 Exemplos de utilização

Para demonstrar o funcionamento do nosso trabalho introduzimos o seguinte conteúdo como input para o nosso programa.

```
1 @techreport{Camila,  
2   author ={{projecto Camila}},  
3   editor ={L.S. Barbosa and J.J. Almeida and J.N. Oliveira and Luís Neves},  
4   title = "\textsc{Camila} - A Platform for Software Mathematical Development",  
5   url="http://camila.di.uminho.pt",  
6   type="(Páginas do projecto)",  
7   institution = "umdi",  
8   year=1998,  
9   keyword = "FS",  
10  }  
11  
12 @techreport{Barbosa95b,  
13   author = "L.S. Barbosa and J.J. Almeida",  
14   title = "Growing Up With \textsc{Camila}",  
15   institution = "umdi",  
16   year = 1995,  
17   number = "DI-CAM-95:7:1",  
18   url = "http://www.di.uminho.pt/~lsb/pub_camila/romantic.ps.gz",  
19   keyword ="Camila, formal specification, didactics",  
20  }  
21  
22 @inproceedings{Ramalho95,  
23   author = "J.C. Ramalho and J.J. Almeida and P.R. Henriques",  
24   title = "Algebraic Specification of Documents",  
25   booktitle = "TWT10 - Algebraic Methods in Language Processing",  
26   year = 1995,  
27   month = "6--8 Dec.",  
28   editor = "A. Nijholt and G. Scollo and R. Steetskamp",  
29   address = "Twente University, Netherlands",  
30   note = "AMiLP'95",  
31   series = "Twente Workshop on Language Technology",  
32   url="http://natura.di.uminho.pt/~jj/bib/amilp95.ps.gz",  
33   docpage="http://www.di.uminho.pt/~jcr/projectos/david/ARTIGOS/AMiLP95/amilp95.html",  
34   pages = "55--64",  
35   keyword ="PDavid, Camila, SGML",  
36  }  
37  
38 @article{sepln06,  
39   author = {Alberto Simões and J. João Almeida},  
40   title = {{NatServer:} A Client-Server Architecture for building Parallel
```

```

41     Corpora applications},
42     year = {2006},
43     journal = {Procesamiento del Lenguaje Natural},
44     address = {Zaragoza, Spain},
45     url = {http://alfarrabio.di.uminho.pt/~albie/publications/sepln06.pdf},
46     month = {September},
47     lang = {EN},
48     volume = {37},
49     pages = {91--97},
50     abstract = {Parallel corpora are important resources for most
51                 Natural Language processing tasks. From the common
52                 applications, like machine translation, to the
53                 usually mono-lingual tasks as paraphrase detection
54                 and word sense disambiguation, most researchers are
55                 using massive parallel corpora. Thus, the
56                 availability of an efficient way to manage them is
57                 very important. This paper presents a Client-Server
58                 architecture to query efficiently parallel corpora
59                 and probabilistic translation dictionaries.},
60 }
61
62 @Article{MSH05,
63     author = {Giovana Mendes and Nuno Alberto Silva and Pedro Rangel Henriques},
64     title = {Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de
65             erros num sistema de logs},
66     journal = {Sistemas de Informação},
67     editor = {},
68     publisher = {APSI: Associação Portuguesa de Sistemas de Informação},
69     year = {2005},
70     month = {},
71     volume = {17},
72     number = {},
73     pages = {91-100}
74 }
75
76 @Article{ALHF02,
77     author = {Gustavo Arnold and Giovani Librelotto and Pedro Rangel Henriques
78             and Jaime Fonseca},
79     title = {O Uso da Linguagem RS em Robótica},
80     journal = {Revista Electrónica e Telecomunicações},
81     editor = {},
82     publisher = {Departamento de Electrónica e Telecomunicações da Universidade de Aveiro},
83     year = {2002},
84     month = {Apr},

```

```
85     volume = {},
86     number = {},
87     pages = {501-508}
88 }
```

O código/ficheiro HTML para responder às alíneas a) e b) gerado foram os seguintes.

```
1      <!DOCTYPE html>
2      <html>
3      <head>
4      <style>
5          .tabcen {
6              text-align: center;
7          }
8          td,th {
9              border: 1px solid grey;
10         }
11         th {
12             background-color: #80808054;
13         }
14         .bortop {
15             border-top: 5px solid black;
16         }
17         .borbot {
18             border-bottom: 2px solid black;
19         }
20     </style>
21 </head>
22 <body>
23     <table>
24     <tbody>
25     <tr>
26         <th class="bortop">Categoria</th>
27         <td class="tabcen bortop" colspan="2">techreport</td>
28     </tr>
29     <tr>
30         <th class="borbot">Incidencia</th>
31         <td class="tabcen borbot" colspan="2">2</td>
32     </tr>
33     <tr>
34         <th>Chave</th>
35         <th>Autores</th>
36         <th>Titulos</th>
37     </tr>
38
39     <tr>
40         <td>Camila</td>
41         <td>
```

```

42     <ul>
43         <li>projecto Camila</li>
44     </ul></td>
45     <td>\textsc{Camila} - A Platform for Software Mathematical Development</td>
46 </tr>
47     <td>Barbosa95b</td>
48     <td>
49         <ul>
50             <li>L.S. Barbosa</li>
51             <li>J.J. Almeida</li>
52         </ul></td>
53     <td>Growing Up With \textsc{Camila}</td>
54 </tr>
55     <th class="bortop">Categoria</th>
56     <td class="tabcen bortop" colspan="2">inproceedings</td>
57 </tr>
58 <tr>
59     <th class="borbot">Incidencia</th>
60     <td class="tabcen borbot" colspan="2">1</td>
61 </tr>
62 <tr>
63     <th>Chave</th>
64     <th>Autores</th>
65     <th>Titulos</th>
66 </tr>
67
68 <tr>
69     <td>Ramalho95</td>
70     <td>
71         <ul>
72             <li>J.C. Ramalho</li>
73             <li>J.J. Almeida</li>
74             <li>P.R. Henriques</li>
75         </ul></td>
76     <td>Algebraic Specification of Documents</td>
77 </tr>
78     <th class="bortop">Categoria</th>
79     <td class="tabcen bortop" colspan="2">article</td>
80 </tr>
81 <tr>
82     <th class="borbot">Incidencia</th>
83     <td class="tabcen borbot" colspan="2">3</td>
84 </tr>
85 <tr>

```

```

86         <th>Chave</th>
87         <th>Autores</th>
88         <th>Titulos</th>
89     </tr>
90
91     <tr>
92         <td>sepln06</td>
93         <td>
94             <ul>
95                 <li>Alberto Simões</li>
96                 <li>J. João Almeida</li>
97             </ul></td>
98         <td>{NatServer:} A Client-Server Architecture for building Parallel Corpora application
99     </tr>
100    <tr>
101        <td>MSH05</td>
102        <td>
103            <ul>
104                <li>Giovana Mendes</li>
105                <li>Nuno Alberto Silva</li>
106                <li>Pedro Rangel Henriques</li>
107            </ul></td>
108        <td>Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de erros num sistema
109    </tr>
110    <tr>
111        <td>ALHF02</td>
112        <td>
113            <ul>
114                <li>Gustavo Arnold</li>
115                <li>Giovani Librelotto</li>
116                <li>Pedro Rangel Henriques</li>
117                <li>Jaime Fonseca</li>
118            </ul></td>
119        <td>O Uso da Linguagem RS em Robótica</td>
120    </tr>
121 </tbody>
</table>
</body>
</html>

```


Categoria	techreport	
Incidencia	2	
Chave	Autores	Titulos
Camila	<ul style="list-style-type: none"> projecto Camila 	\textsc{Camila} - A Platform for Software Mathematical Development
Barbosa95b	<ul style="list-style-type: none"> L.S. Barbosa J.J. Almeida 	Growing Up With \textsc{Camila}
Categoria	inproceedings	
Incidencia	1	
Chave	Autores	Titulos
Ramalho95	<ul style="list-style-type: none"> J.C. Ramalho J.J. Almeida P.R. Henriques 	Algebraic Specification of Documents
Categoria	article	
Incidencia	3	
Chave	Autores	Titulos
sepin06	<ul style="list-style-type: none"> Alberto Simões J. João Almeida 	(NatServer:) A Client-Server Architecture for building Parallel Corpora applications
MSH05	<ul style="list-style-type: none"> Giovana Mendes Nuno Alberto Silva Pedro Rangel Henriques 	Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de erros num sistema de logs
ALHF02	<ul style="list-style-type: none"> Gustavo Arnold Giovani Librelotto Pedro Rangel Henriques Jaime Fonseca 	O Uso da Linguagem RS em Robótica

O ficheiro de texto com o índice dos autores gerado foi o seguinte.

```
1 A
2 Alberto Simões (1 entrada):
3     "{NatServer:} A Client-Server Architecture for building Parallel Corpora
4     applications", sepln06, article
5
6
7 G
8 Giovana Mendes (1 entrada):
9     "Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de erros num
10     sistema de logs", MSH05, article
11
12
13 Giovanni Librelotto (1 entrada):
14     "O Uso da Linguagem RS em Robótica", ALHF02, article
15
16
17 Gustavo Arnold (1 entrada):
18     "O Uso da Linguagem RS em Robótica", ALHF02, article
19
20
21 J
22 J. João Almeida (1 entrada):
23     "{NatServer:} A Client-Server Architecture for building Parallel Corpora
24     applications", sepln06, article
25
26
27 J.C. Ramalho (1 entrada):
28     "Algebraic Specification of Documents", Ramalho95, inproceedings
29
30
31 J.J. Almeida (2 entradas):
32     "Growing Up With \textsc{Camila}", Barbosa95b, techreport
33     "Algebraic Specification of Documents", Ramalho95, inproceedings
34
35
36 Jaime Fonseca (1 entrada):
37     "O Uso da Linguagem RS em Robótica", ALHF02, article
38
39
40 L
41 L.S. Barbosa (1 entrada):
42     "Growing Up With \textsc{Camila}", Barbosa95b, techreport
```

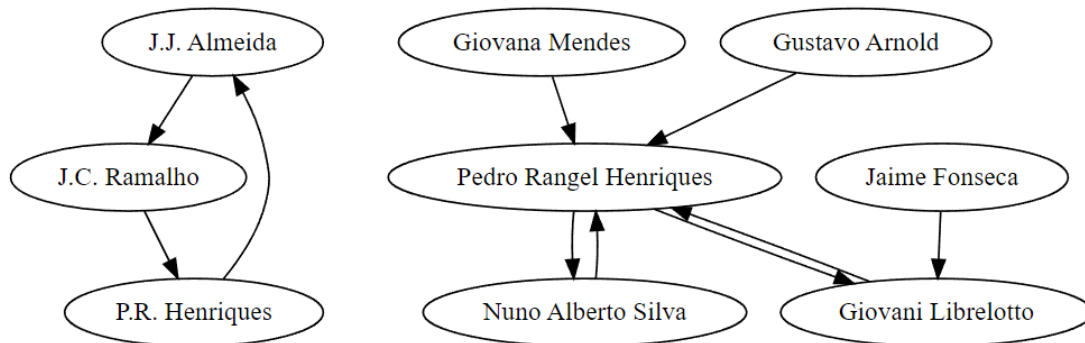
43
44
45 N
46 Nuno Alberto Silva (1 entrada):
47 "Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de erros num
48 sistema de logs", MSH05, article
49
50
51 P
52 P.R. Henriques (1 entrada):
53 "Algebraic Specification of Documents", Ramalho95, inproceedings
54
55
56 Pedro Rangel Henriques (2 entradas):
57 "Utilizando uma Base de Dados XML Nativa aplicada ao tratamento de erros num
58 sistema de logs", MSH05, article
59 "O Uso da Linguagem RS em Robótica", ALHF02, article
60
61
62 projecto Camila (1 entrada):
63 "\textsc{Camila} - A Platform for Software Mathematical Development", Camila, techreport

Em seguida, mostramos o código do grafo gerado sem introduzir o nome de nenhum autor e a sua representação gráfica; e o código do grafo gerado se for introduzido o nome "Pedro Rangel Henriques" e a sua representação gráfica.

```

1 digraph G {
2   "J.J. Almeida" -> "J.C. Ramalho"
3   "J.C. Ramalho" -> "P.R. Henriques"
4   "P.R. Henriques" -> "J.J. Almeida"
5   "Giovana Mendes" -> "Pedro Rangel Henriques"
6   "Nuno Alberto Silva" -> "Pedro Rangel Henriques"
7   "Pedro Rangel Henriques" -> "Giovani Librelotto"
8   "Pedro Rangel Henriques" -> "Nuno Alberto Silva"
9   "Gustavo Arnold" -> "Pedro Rangel Henriques"
10  "Giovani Librelotto" -> "Pedro Rangel Henriques"
11  "Jaime Fonseca" -> "Giovani Librelotto"
12 }

```



```

1 digraph G {
2   "Pedro Rangel Henriques" -> "Giovani Librelotto"
3   "Pedro Rangel Henriques" -> "Nuno Alberto Silva"
4 }

```

