

Processamento de Linguagens (3<sup>o</sup> ano de MIEI)

**Trabalho Prático 1**

**Enunciado 1**

Relatório de Desenvolvimento

Sofia Santos  
(a89615)

Carolina Vila Chã  
(a89495)

5 de abril de 2021

## **Resumo**

Neste relatório encontra-se descrita a resolução do Enunciado 1 do Trabalho Prático 1 de 2021 para a UC Processamento de Linguagens. O foco deste trabalho é a utilização de expressões regulares (*ER*) para a descrição de padrões de frases, filtragem, e transformação de texto.

# Conteúdo

<b>1</b>	<b>Introdução</b>	<b>2</b>
<b>2</b>	<b>Análise do Problema</b>	<b>3</b>
2.1	Descrição . . . . .	3
2.2	Especificação . . . . .	3
2.2.1	Dados . . . . .	3
2.2.2	Alíneas . . . . .	3
<b>3</b>	<b>Resolução</b>	<b>4</b>
3.1	Estruturas de Dados . . . . .	4
3.2	Algoritmo e Decisões . . . . .	4
3.2.1	<i>main.py</i> . . . . .	4
3.2.2	<i>similar.py</i> . . . . .	12
<b>4</b>	<b>Conclusão</b>	<b>14</b>
<b>A</b>	<b>Exemplos de Utilização</b>	<b>15</b>
<b>B</b>	<b>Ficheiro HTML - Equipas</b>	<b>17</b>
<b>C</b>	<b>Ficheiro HTML - Equipa Minho Aventura</b>	<b>18</b>
<b>D</b>	<b>Website gerado</b>	<b>19</b>
<b>E</b>	<b>Código do Programa</b>	<b>21</b>

# Capítulo 1

## Introdução

O principal objetivo deste trabalho é consolidar os conhecimentos sobre ER e Python, adquiridos na UC de Processamento de Linguagens.

## Capítulo 2

# Análise do Problema

### 2.1 Descrição

Para o enunciado 1 pretende-se trabalhar com um arquivo desportivo criado por um Organizador de Provas de Orientação, realizadas em diferentes locais, e para diferentes graus de dificuldade, adaptadas a diferentes classes de participantes.

São apresentadas então várias alíneas que expõem os problemas a resolver pelo grupo.

### 2.2 Especificação

#### 2.2.1 Dados

Os dados a utilizar para este projeto encontram-se no ficheiro "inscritos-form.json", fornecido pelos docentes da UC.

#### 2.2.2 Alíneas

É-nos pedido para processar o ficheiro "inscritos-form.json" de modo a:

- a) imprimir o nome (convertido para maiúsculas) de todos os concorrentes que se inscrevem como 'Individuais' e são de 'Valongo'.
- b) imprimir o nome completo, o email<sup>1</sup> e a prova em que está inscrito cada concorrente cujo nome seja 'Paulo' ou 'Ricardo', desde que usem o Gmail.
- c) imprimir toda a informação dos atletas da equipe "TURBULENTOS".
- d) imprimir a lista dos escalões por ordem alfabética e para cada um indicar quantos atletas estão inscritos nesse escalão.
- e) gerar uma página HTML com a lista das equipas inscritas em qualquer prova, indicando o seu nome e o número dos atletas que a constituem e que se inscreveram pelo menos uma vez numa prova; essa lista deve estar ordenada por ordem decrescente do número de atletas; além disso, cada equipa deve ter um link para outra página HTML com a informação que achar interessante sobre cada atleta indicando as provas em que cada um participou.

---

<sup>1</sup>No enunciado original, era pedido o "telemóvel" – no entanto este valor não existe, pelo que se escolheu o email.

## Capítulo 3

# Resolução

### 3.1 Estruturas de Dados

Usamos várias estruturas de dados auxiliares ao longo do projeto para armazenar e manipular os dados obtidos pela leitura e *parsing* do ficheiro *.json*. A estrutura principal deste projeto é uma lista de dicionários, em que cada dicionário corresponde a um atleta e armazena a sua informação com uma organização semelhante à do ficheiro *.json*.

### 3.2 Algoritmo e Decisões

O código do nosso projeto está contido em dois ficheiros. O principal (*main.py*) contém o código relativo ao *parsing* do ficheiro de input e à resolução das alíneas do trabalho prático, para além de uma simples interface gráfica para poder visualizar os resultados das mesmas de forma eficaz. O outro ficheiro (*similar.py*) contém uma função auxiliar desenvolvida por nós para calcular a semelhança entre duas palavras. Devido à sua relativa complexidade e versatilidade, considerámos que seria mais lógico colocá-la num ficheiro separado.

#### 3.2.1 *main.py*

Antes de tudo, são importados os módulos necessários à resolução deste trabalho prático, utilizados pelo nosso grupo.

---

```
1 import re
2 import sys
3 from datetime import datetime, timedelta
4 import similar
```

---

O nosso grupo decidiu ler apenas uma vez o ficheiro de *input*, colocando a sua informação relevante numa estrutura de dados, que é lida posteriormente de modo a responder às alíneas do problema. Visto que as operações de leitura de ficheiros são mais lentas do que as de leitura a partir da memória, este método revela-se mais eficiente.

Primeiramente, apresentamos ao utilizador um pequeno menu, a partir do qual pode escolher qual o ficheiro que pretende utilizar.

---

```
1 def main():
2     fp = None
3     inscritos = list()
4
```

---

```

5     while not fp:
6         opt = input("1) Usar ficheiro predefinido (inscritos-form.json).\n\
7 2) Especificar ficheiro a usar.\n\n\
8 Escolha a opção pretendida: ").strip()
9         if opt == '1':
10            fp = "inscritos-form.json"
11        elif opt == '2':
12            fp = input("Introduza o caminho para o ficheiro que pretende utilizar: ").strip()

```

---

O ciclo `while` garante-nos que apenas pode ser introduzida uma opção válida.

De seguida, o ficheiro especificado é aberto com *encoding* UTF-8 e utilizamos funções do módulo `re` para ler o conteúdo do ficheiro e armazenar a informação que pretendemos obter do mesmo. A partir da função `re.search()` somos capazes de guardar a informação relativa aos inscritos numa variável `string_inscritos`.

---

```

1     with open(fp, encoding="utf-8") as f:
2         m = re.search(r'\s*inscritos':[ \s*((?:.|\\n)*)\\]\s*', f.read())
3         if not m:
4             sys.exit("Erro - Ficheiro inválido")
5         string_inscritos = m.group(1)

```

---

Usamos ainda a função `sys.exit()` para sair do programa caso o ficheiro introduzido não seja válido. Para verificar a validade do ficheiro procuramos pela chave *"inscritos"*; se estiver presente consideramos que o ficheiro é válido, e o valor correspondente a esta chave diz respeito aos dados que vamos guardar.

A partir daqui somos capazes de percorrer a string obtida e passar os valores relativos a cada inscrito para uma lista de dicionários, neste caso a lista `inscritos`, em que cada dicionário diz respeito a uma pessoa inscrita.

---

```

1     lista_inscritos = re.split(r'}\s*,\s*{', string_inscritos.strip("\t\n{}"))
2     progress = 10
3     for i, inscrito in enumerate(lista_inscritos):
4         if i * 100 // len(lista_inscritos) > progress:
5             print(f"Lendo ficheiro ... {progress}% concluído")
6             progress += 10
7             insc_dict = dict()
8             for linha in inscrito.splitlines():
9                 linha_stripped = linha.strip()
10                if linha_stripped:
11                    m = re.search(r'"(\w+)"\s*:\s*"(.*)"\s*,?', linha_stripped)
12                    insc_dict[m.group(1)] = m.group(2)
13
14                not_duplicate = all(any(similar.similarity(inscrito[key].lower(), insc_dict[
15                    key].lower()) > 2 for key in inscrito) for inscrito in inscritos)
16
17                if not_duplicate: inscritos.append(insc_dict)

```

---

Primeiro dividimos a string `string_inscritos` usando a função `re.split()` com um separador correspondente ao separador usado no ficheiro *.json* para separar a informação dos inscritos. Ficamos assim com uma lista (`lista_inscritos`) em que cada elemento diz respeito à informação de um inscrito, numa string.

Depois apenas percorremos essa lista e para cada elemento, linha a linha, usamos a função `re.search()` para obter os dados relativos a cada inscrito, armazenando-os em dicionários, um por inscrito, dicionários estes que são depois agregados na lista mencionada anteriormente. Antes desta agregação, é feita a filtragem de entradas repetidas.

O grupo decidiu que apenas iria excluir a entrada de um atleta se já existisse outra entrada com os mesmos valores. Poderíamos filtrar entradas repetidas através do nome ou da data de nascimento, mas poderíamos estar a excluir pessoas com o mesmo nome e nascidas no mesmo dia, para além de que a mesma pessoa pode-se inscrever duas vezes com nomes ou datas de nascimento diferentes, algo que reparámos que acontece no ficheiro fornecido, por isso não seria um método eficaz.

O único dado que poderia ser usado para excluir repetidos seria o endereço de email, assumindo que todas as pessoas que se inscrevessem teriam que usar um endereço de email diferente, mas tal como referimos no parágrafo anterior, no ficheiro fornecido constam vários atletas com o mesmo endereço de email, e é-nos impossível saber qual deles se inscreveu primeiro. Logo, não somos capazes de usar um dado específico como filtro de repetidos, e decidimos assim que dois atletas são a mesma pessoa apenas se tiverem todos os seus dados iguais ou quase iguais.

Uma solução proposta por nós seria a inclusão de um elemento extra no perfil de cada atleta: o Número de Cartão de Cidadão (CC) ou o NIF. Uma vez que estes números são únicos, permitiriam excluir entradas repetidas. Outra opção era proibir a inscrição de mais do que um atleta com o mesmo endereço de email, e poderíamos assim usar o endereço de email para excluir repetidos.

Para verificar se dois dados são parecidos o suficiente para serem considerados iguais, por exemplo, no caso de haver erros ortográficos, desenvolvemos um simples algoritmo que calcula a semelhança de duas palavras ou expressões. Falamos deste algoritmo em detalhe na secção **3.2.2** deste relatório. Neste caso, como estamos a comparar os dados de cada atleta com os de todos os outros, este processo é um bocado lento. Pelos nossos testes demora cerca de 15 segundos com o ficheiro disponibilizado, o que é normal, tendo em conta a dimensão deste ficheiro e o funcionamento do algoritmo. Sem esta verificação o nosso programa consegue gerar a lista final (**inscritos**) em muito menos do que um segundo.

De qualquer das formas, ficamos assim com os dados do ficheiro *.json* numa estrutura de dados organizada, o que nos vai permitir percorrê-la de forma eficiente e obter os dados necessários para responder a cada alínea. Para obter esta estrutura recorreremos a Expressões Regulares, tal como nos foi pedido, e a funções e operadores *standard* do Python, em situações onde usar Expressões Regulares seria muito menos eficiente.

## Alínea a

**>> Imprimir o nome (convertido para maiúsculas) de todos os concorrentes que se inscrevem como 'Individuais' e são de 'Valongo'.**

---

```
1      print("\nNomes dos concorrentes individuais de Valongo:")
2      for inscrito in inscritos:
3          if similar.similarity(inscrito["equipa"].lower(), "individual") < 2
              and (n := re.search(r'(?i:valongo)', inscrito["morada"])):
4              print("-", inscrito["nome"].upper())
```

---

Graças à estrutura de dados que criámos anteriormente, responder a esta alínea é muito simples. Apenas percorremos a lista **inscritos** e para cada inscrito verificamos se a sua equipa é *"Individual"* e se a palavra *"Valongo"* aparece na sua morada.

Tal como na filtragem inicial de dados, usamos também aqui o nosso algoritmo de semelhança. Neste caso, usamo-lo para verificar se o nome da equipa de cada inscrito é próximo o suficiente de *"individual"* para poder ser considerado um erro ortográfico.

## Alínea b

**>> Imprimir o nome completo, o endereço de email e a prova em que está inscrito cada concorrente cujo nome seja 'Paulo' ou 'Ricardo', desde que usem o GMail.**



---

```

1      print("\nNome, email e prova dos concorrentes chamados \"Paulo\" ou \"Ricardo\" que usam o Gmail:")
2      for inscrito in inscritos:
3          if m := re.search(r'@gmail\.com$', inscrito["email"]):
4              if o := re.search(r'(?i:Paulo|Ricardo)', inscrito["nome"]):
5                  print(f'_{inscrito["nome"]} _{inscrito["email"]} _{inscrito["prova"]} ')

```

---

Tal como a alínea anterior, esta revela-se muito simples de resolver a partir da nossa estrutura de dados. Primeiro verificamos, para cada inscrito, se o seu endereço de email acaba em `"@gmail.com"`. Se tal for o caso, verificamos depois se o seu nome contém as palavras `"Paulo"` ou `"Ricardo"`.

Decidimos contabilizar atletas que tivessem estes nomes em qualquer parte do seu nome. Por exemplo, um atleta chamado `"André Ricardo"` será incluído nesta pesquisa. Se quiséssemos apenas contabilizar os atletas com estes nomes como primeiro nome, apenas teríamos que trocar a função `re.search()` para `re.match()`, ou colocar um `^` no início da Expressão Regular.

### Alínea c

**>> Imprimir toda a informação dos atletas da equipa "TURBULENTOS".**

---

```

1      turbulentos = list()
2      for inscrito in inscritos:
3          if m := re.fullmatch(r'(?i:turbulentos)', inscrito["equipa"]):
4              turbulentos.append(inscrito)

```

---

Para obter os atletas da equipa `"TURBULENTOS"`, apenas temos de, mais uma vez, percorrer a nossa estrutura de dados e filtrar os atletas cuja equipa tenha o nome `"TURBULENTOS"`. O *matching* é feito de forma *case insensitive*, tal como nas outras alíneas, de modo a encontrar todos os resultados que possam estar escritos de forma diferente.

Na nossa análise do ficheiro encontramos uma equipa chamada `"Os Turbulentos"`, que no contexto deste trabalho decidimos considerar como sendo uma equipa diferente, mas que poderia ser considerada a mesma equipa, com o nome escrito de forma diferente. Para tal, apenas teríamos que, de forma semelhante à alínea anterior, trocar a função `re.fullmatch()` para `re.search()`.

---

```

1      i = 0
2      while True:
3          print(f'\nNome: _{turbulentos[i]["nome"]} _\nData de nascimento: _{turbulentos[i]["dataNasc"]} _\nMorada: _{turbulentos[i]["morada"]} _\nEmail: _{turbulentos[i]["email"]} _\nProva: _{turbulentos[i]["prova"]} _\nEscalão: _{turbulentos[i]["escalao"]} ')
4          print(f'\nPágina {i+1} de {len(turbulentos)}\n\n[a] _ver anterior; [s] _ver seguinte; [e] _sair')
5          opt = input()
6          if opt == 'a' and i > 0: i -= 1
7          elif opt == 's' and i < len(turbulentos) - 1: i += 1
8          elif opt == 'e': break

```

---

Devido à grande quantidade e complexidade dos resultados desta alínea, criámos um ciclo que mostra a informação de cada atleta em `"páginas"`, como um livro, permitindo ao utilizador visualizar toda a informação de uma forma simples e eficaz.

## Alínea d

>> Imprimir a lista dos escalões por ordem alfabética e para cada um indicar quantos atletas estão inscritos nesse escalão.

---

```
1         esc_dict = dict()
2         for inscrito in inscritos:
3             if re.match(r'(\w+)', inscrito["escala"]):
4                 x = esc_dict.get(inscrito["escala"], 0)
5                 esc_dict[inscrito["escala"]] = x+1
6         for escala in sorted(esc_dict):
7             print("Escala:", escala, "; Número de atletas:", esc_dict[escala])
```

---

Nesta alínea a resolução torna-se ligeiramente mais complexa que as anteriores, uma vez que temos de contar o número de atletas pertencentes a cada escala.

Para tal, criamos uma nova estrutura - `esc_dict` - que, para cada escala, armazena o número de atletas pertencentes a esse escala.

Primeiro percorremos a estrutura de dados que temos usado nas outras alíneas, e para cada atleta incrementamos em uma unidade o valor no dicionário `esc_dict` correspondente ao seu escala. Caso ainda não exista uma entrada com esse escala, é criada uma com o valor inicial 0.

De seguida, o dicionário é ordenado alfabeticamente e são imprimidos para o terminal os valores obtidos.

## Alínea e

>> Gerar uma página HTML com a lista das equipas inscritas em qualquer prova, indicando o seu nome e o número dos atletas que a constituem e que se inscreveram pelo menos uma vez numa prova; essa lista deve estar ordenada por ordem decrescente do número de atletas; além disso, cada equipa deve ter um link para outra página HTML com a informação que achar interessante sobre cada atleta indicando as provas em que cada um participou.

---

```
1         today = datetime.today()
2         equipas = dict()
3         provas = dict()
4         for i, inscrito in enumerate(inscritos):
5             prova = inscrito["prova"]
6             equipa = inscrito["equipa"].upper()
7             for eqp in equipas:
8                 if similar.similarity(eqp.upper(), equipa) < 2:
9                     equipa = eqp.upper()
10                    break
11             if not (m := re.match(r'|N/D|S/_CLUBE', equipa)):
12                 equipas.setdefault(equipa, list()).append(inscrito)
13                 provas.setdefault(prova, dict()).setdefault(equipa, list()).
14                     append(inscrito)
15             else:
16                 provas.setdefault(prova, dict()).setdefault("INDIVIDUAL", list())
17                     .append(inscrito)
```

---

Para esta alínea, devido à sua complexidade, criámos duas estruturas de dados auxiliares. A primeira - `equipas` - armazena, para cada equipa, os atletas que dela fazem parte, e a segunda - `provas` - armazena as equipas que participam em cada prova. Por sua vez, cada equipa contém os atletas das mesmas que participaram nas respetivas provas.

Usamos mais uma vez o algoritmo de semelhança criado por nós para verificar se dois nomes de equipa se referem à mesma equipa. Para além disso, excluimos determinados nomes considerados "inválidos", como ",", "n/d" ou "s/ clube". Como num dicionário em Python as chaves devem ter um único valor imutável, passamos os nomes das equipas para maiúsculas, de modo a tornar a capitalização igual para todos os nomes e não haver duplicados.

Com estas estruturas criadas, podemos criar os ficheiros HTML.

---

```

1         with open("equipas.html", "w", encoding="utf-8") as f:
2             f.write(f"""<!DOCTYPE html>\n
3 <html>
4     <head>
5         <title>Equipas</title>
6         <style>{style_provas}</style>
7     </head>
8     <body>
9         <h1>PROVAS</h1>
10        <div class="provas">
11            """
12                for prova in provas:
13                    f.write(f"""        <div class="prova">
14            <h2>{prova}</h2>
15            """
16                for equipa in sorted(provas[prova].keys(), key=lambda x: len(
17                    provas[prova][x]), reverse=True):
18                    f.write(f"""                <a href="/equipas/{''.join(x for x in
19                        equipa if x.isalnum())}.html">
20                    <span class="equipa">
21                        <h2>{equipa if equipa != "INDIVIDUAL" else "Sem equipa"}</h2>
22                        <div class="num_a"><p>{len(provas[prova][equipa])} atleta{'s' if len(
23                            provas[prova][equipa]) > 1 else ''}</p></div>
24                        @ -128,23 +133,23 @@ Escolha a opção pretendida: """).strip()
25                    </a>
26            """
27                f.write("        </div>\n")
28                f.write("    </div>
29    </body>
30 </html>
31 """

```

---

Este primeiro ficheiro diz respeito às provas e, para cada prova, apresenta as equipas e o número de atletas das mesmas que nela participaram.

Tal como pedido, as equipas são ordenadas por ordem decrescente do número de atletas.

Para além disso, é possível clicar em qualquer equipa para ir para outro ficheiro HTML, que contém a informação sobre os atletas que constituem essa equipa, incluindo as provas em que cada atleta participou.

Estes ficheiros HTML das equipas são criados da seguinte forma:

---

```

1         for equipa in equipas:
2             with open(f"equipas/{''.join(x for x in equipa if x.isalnum())}.html",
3                 "w", encoding="utf-8") as ff:
4                 ff.write(f"""<!DOCTYPE html>
5 <html>
6     <head>
7         <title>{equipa}</title>

```

---

```

7         <style>{style_equipas}</style>
8     </head>
9     <body>
10         <h1>{"Equipa: " + equipa if equipa != "INDIVIDUAL" else "Atletas sem
            equipa"}</h1>
11 @ -152,14 +157,14 @@ Escolha a opção pretendida: "").strip()
12         <div class="atletas">
13 """)
14             for atleta in equipas[equipa]:
15                 try:
16                     if "ou" not in atleta["dataNasc"]:
17                         birth = datetime.strptime(atleta["dataNasc"], "%d/%m/%
                            y")
18                         if birth > datetime.today(): birth = birth.replace(
                            year = birth.year - 100)
19                     except ValueError:
20                         birth = None
21                     ff.write(f"""" <div class="atleta">
22 <h3>{atleta["nome"]}</h3>
23 <ul>
24 <li>Idade: {str(today.year - birth.year - ((today.month, today
                    .day) < (birth.month, birth.day))) + " anos" if birth else
                    "-"}</li>
25 @ -168,14 +173,104 @@ Escolha a opção pretendida: "").strip()
26 </ul>
27 </div>
28 """)
29             ff.write(""" </div>
30 </body>
31 </html>
32 """)

```

---

Usámos o módulo *datetime* do Python para calcular a idade dos atletas com base na sua data de nascimento. Utilizou-se, também, *css* para tornar a consulta de informação sobre as provas e os seus concorrentes mais fácil, utilizando variáveis. Deste modo, conseguimos incluir todo o código relacionado com as páginas web no menor número possível de documentos.

---

Listing 3.1: Variável `style_equipas` utilizada para a exibição da página das provas

---

```

1 style_provas = ""
2 html {
3     background: #5bcefa;
4     color: #166888;
5     font-family: fantasy;
6     min-height: 100%;
7 }
8
9 h1, h2 {
10     text-align: center;
11     width: 100%;
12 }
13
14 .num_a {
15     position: absolute;
16     bottom: 0;

```

```

17     left: 0px;
18     right: 0px;
19     text-align: center;
20 }
21
22 .provas {
23     display: flex;
24     flex-direction: column;
25 }
26
27 .prova {
28     display: flex;
29     flex-wrap: wrap;
30     justify-content: center;
31     background: #f5a9b8;
32     color: #97273e;
33     border-radius: 20px;
34     margin: 50px;
35 }
36
37 .equipa {
38     position: relative;
39     flex-grow: 1;
40     display: block;
41     /*background-color:#ffa4d2;
42     color: #ff339a;*/
43     background: white;
44     color: rgb(82, 82, 82);
45     padding: 10px;
46     width: 250px;
47     border-radius: 8px;
48     min-height: 200px;
49 }
50
51 a {
52     margin: 15px;
53     text-decoration: none;
54 }
55 """

```

---

Listing 3.2: Variável `style_equipas` utilizada para a exibição da página de cada equipe

---

```

1 style_equipas = """
2 html {
3     background:          #a3a3a3;
4     color: black;
5     font-family: fantasy;
6     min-height: 100%;
7 }
8
9 h1,h2,h3 {
10     text-align: center;
11 }
12
13 .atletas {

```

```

14     display: flex;
15     flex-wrap: wrap;
16     justify-content: center;
17 }
18
19 .atleta {
20     /*background-color: #ffccf5;
21     color: #7e3a70;*/
22     background: #800080;
23     color: white;
24     padding: 10px;
25     margin: 15px;
26     width: 20%;
27     border-radius: 15px;
28 }
29 """

```

---

Todos estes ficheiros HTML contém secções (delimitadas com a tag `<div>`) de modo a sermos capazes de usar ficheiros CSS para tornar os *websites* gerados mais bonitos e apelativos.

### 3.2.2 *similar.py*

Este algoritmo, apesar de parecer relativamente complexo, é bastante simples. Recebe duas palavras ou expressões como input e tenta calcular a semelhança entre elas. Neste caso, definimos "semelhança" como a quantidade de letras que têm em comum, e o valor devolvido pela função é decrescente, por outras palavras, quanto mais próximo de zero for, mais próximas as duas palavras são. Por exemplo, duas palavras iguais terão semelhança igual a 0, enquanto que duas palavras completamente diferentes terão semelhança igual ao comprimento da menor palavra. Pode parecer contra-intuitivo, visto que duas palavras semelhantes deveriam ter um valor alto de semelhança, mas trabalhando com uma escala decrescente somos capazes de mais facilmente detetar palavras "parecidas". Dando um exemplo mais concreto, as palavras "correr" e "acorrer" têm semelhança 1, pois apenas uma letra é diferente.

Um exemplo da utilização deste algoritmo no nosso trabalho seria a palavra "indivudual", que aparece no ficheiro disponibilizado e que consideramos ser um erro ortográfico, apesar de também poder haver uma equipa com este nome. Usando o nosso algoritmo, verificamos que a "distância" entre esta palavra e a palavra "individual" é de 1, ou seja, apenas uma letra é diferente. Assim, nas alíneas *a)* e *e)* contabilizamos atletas com o valor "indivudual" no campo "equipa" como sendo atletas sem equipa.

---

```

1 cache = dict()
2
3 def similarity(word1 : str, word2 : str):
4     if (word1, word2) in cache:
5         return cache[(word1, word2)]
6     if word1 == word2: return 0
7     if not word1:
8         return len(word2)
9     elif not word2:
10        return len(word1)
11    check0 = similarity(word1[1:], word2[1:])
12    if check0 == 0:
13        cache[(word1, word2)] = 1
14        return 1
15    if word1[0] == word2[0]:

```

```
16         cache[(word1,word2)] = check0
17         return check0
18     check1 = similarity(word1[1:], word2)
19     if check1 < check0:
20         cache[(word1,word2)] = 1 + check1
21         return 1 + check1
22     check2 = similarity(word1, word2[1:])
23     if check2 < check0:
24         cache[(word1,word2)] = 1 + check2
25         return 1 + check2
26     cache[(word1,word2)] = 1 + check0
27     return 1 + check0
```

---

Este algoritmo conta com uma cache pois, como funciona de forma recursiva, acaba por se tornar bastante lento para palavras ou expressões muito grandes. A cache ajuda-nos a reduzir de forma exponencial o seu tempo de execução.

## Capítulo 4

# Conclusão

A realização deste projeto permitiu consolidar a matéria lecionada não só sobre ER, mas também sobre Python. Possibilitou, especificamente, a prática e consolidação de: dicionários em python, filtragem de texto utilizando *Expressões Regulares*, o desenvolvimento de *Processadores de Linguagens Regulares*, e a utilização do módulo *'re'* e as suas funções, para além de, em geral, aumentar a nossa capacidade de escrever ER para padrões de texto e trabalhar com Python.



## Apêndice A

### Exemplos de Utilização

```
1) Usar ficheiro predefinido (inscritos-form.json).
2) Especificar ficheiro a usar.

Escolha a opção pretendida: █
```

Figura A.1: Menu inicial.

```
1) Nomes dos concorrentes individuais de Valongo.
2) Nome, email e prova dos concorrentes chamados "Paulo" ou "Ricardo" que usam o GMail.
3) Informação dos atletas da equipa "TURBULENTOS"
4) Lista ordenada de escalões e atletas por escalão.
5) Gerar página HTML com informação sobre provas e equipas.

0) Sair.

Escolha a opção pretendida: █
```

Figura A.2: Menu principal.

```
Nomes dos concorrentes individuais de Valongo:
- VERA CRISTINA MOREIRA DELGADO
- PAULO DOMINGUES
- DULCE MOREDA
```

Figura A.3: Resposta à alínea *a*).

Tal como referido anteriormente no relatório, o tamanho do resultado desta alínea levou-nos a criar um sistema de páginas, sendo possível ao utilizador ver a informação relativa a um atleta de cada vez, mas com a capacidade de percorrer os resultados de forma simples.

```

Nome, email e prova dos concorrentes chamados "Paulo" ou "Ricardo" que usam o Gmail:
- paulo de castro rocha; pcastrorocha@gmail.com; Ultra Trail
- J Paulo Marques; pmarques269@gmail.com; Ultra Trail
- Paulo Serra; paulo.serra@gmail.com; Ultra Trail
- paulo Vilaça; pmv777@gmail.com; Ultra Trail
- Paulo Domingues; p.j.p.domingues@gmail.com; Ultra Trail
- Ricardo Jorge Dias Oliveira; Ricardo.transportesabranco@gmail.com; Ultra Trail
- Ricardo Reis; ricardoreiis@gmail.com; Ultra Trail
- paulo félix; pauloalexteixeirafelix@gmail.com; Ultra Trail
- João Ricardo Tavares Neves; rwichers4@gmail.com; Ultra Trail
- Ricardo Sousa; jose.ricardo.sousa@gmail.com; Corrida da Geira
- Ricardo Jorge Dias Oliveira; helderfva@gmail.com; Ultra Trail
- Ricardo Couto; rjcouto@gmail.com; Ultra Trail
- Ricardo Ernesto dos Santos Geraldes Domingues; santosgeraldes@gmail.com; Corrida da Geira

```

Figura A.4: Resposta à alínea *b*).

```

Nome: João Costa;
Data de nascimento: 04/04/70;
Morada: Rua Carolina Rosa Alves Nº27 Braga;
Email: jfscosta@gmail.com;
Prova: Ultra Trail;
Escalão: M40

Página 1 de 39

[a] ver anterior; [s] ver seguinte; [e] sair

```

Figura A.5: Resposta à alínea *c*).

```

Escalão: F40 ; Número de atletas: 3
Escalão: M40 ; Número de atletas: 36
Escalão: M50 ; Número de atletas: 10
Escalão: SENIOR Fem ; Número de atletas: 42
Escalão: SENIOR Masc ; Número de atletas: 164

```

Figura A.6: Resposta à alínea *d*).

## Apêndice B

# Ficheiro HTML - Equipas

Devido à grande extensão do ficheiro HTML, decidimos colocar neste relatório apenas um pequeno excerto, mas que representa a sua globalidade. Foi excluído o código relacionado com **css** uma vez que este já foi exposto anteriormente.

Listing B.1: Excerto do ficheiro HTML gerado

---

```
1 <!DOCTYPE html>
2
3 <html>
4   <head>
5     <title>Equipas</title>
6     <style>(...)</style>
7   </head>
8   <body>
9     <h1>PROVAS</h1>
10    <div class="provas">
11      <div class="prova">
12        <h2>Ultra Trail</h2>
13        <a href="/equipas/INDIVIDUAL.html">
14          <span class="equipa">
15            <h2>Sem equipa</h2>
16            <div class="num_a"><p>35 atletas</p></div>
17          </span>
18        </a>
19        <a href="/equipas/NAST.html">
20          <span class="equipa">
21            <h2>NAST</h2>
22            <div class="num_a"><p>9 atletas</p></div>
23          </span>
24        </a>
25        [...]
26      </div>
27    </div>
28  </body>
29 </html>
```

---

## Apêndice C

# Ficheiro HTML - Equipa Minho Aventura

Foi excluído o código relacionado com `css` uma vez que este já foi exposto anteriormente.

Listing C.1: Ficheiro HTML gerado

---

```
1 <!DOCTYPE html>
2   <html>
3     <head>
4       <title>MINHO AVENTURA</title>
5       <style> (...)</style>
6     </head>
7     <body>
8       <h1>Equipa: MINHO AVENTURA</h1>
9       <h2>Constituição: 2 atletas</h2>
10      <div class="atletas">
11        <div class="atleta">
12          <h3>Pedro Nóvoa</h3>
13          <ul>
14            <li>Idade: 58 anos</li>
15            <li>Escalão: M50</li>
16            <li>Prova: Ultra Trail</li>
17          </ul>
18        </div>
19        <div class="atleta">
20          <h3>Pedro Nóvoa</h3>
21          <ul>
22            <li>Idade: 58 anos</li>
23            <li>Escalão: M50</li>
24            <li>Prova: Ultra Trail</li>
25          </ul>
26        </div>
27      </div>
28    </body>
29  </html>
```

---

## Apêndice D

### Website gerado



PROVAS			
Ultra Trail			
<b>Sem equipa</b> 35 atletas	<b>NAST</b> 9 atletas	<b>EDV-VIANA TRAIL</b> 7 atletas	<b>PORTO RUNNERS</b> 6 atletas
<b>TURBULENTOS</b> 5 atletas	<b>TUGAS NA ESTRADA</b> 5 atletas	<b>COMPANHIA DO BAZÓFIAS</b> 5 atletas	<b>CLUBE ATLETISMO DE LAMAS</b> 4 atletas

Figura D.1: Página das Equipas

Atletas sem equipa			
Constituição: 66 atletas			
<b>MARIO PIRES</b> <ul style="list-style-type: none"> <li>Idade: 61 anos</li> <li>Escalação: M50</li> <li>Prova: Ultra Trail</li> </ul>	<b>Francisco Neto Silva</b> <ul style="list-style-type: none"> <li>Idade: 31 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Artur Bernardo</b> <ul style="list-style-type: none"> <li>Idade: 48 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Vera Cristina Moreira Delgado</b> <ul style="list-style-type: none"> <li>Idade: 40 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>
<b>Jorge Yong</b> <ul style="list-style-type: none"> <li>Idade: 50 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Paulo Serra</b> <ul style="list-style-type: none"> <li>Idade: 39 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>	<b>Tiago José Cadima Borges</b> <ul style="list-style-type: none"> <li>Idade: 38 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>	<b>jorge manuel martins silva</b> <ul style="list-style-type: none"> <li>Idade: 52 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>
<b>Paulo Domingues</b> <ul style="list-style-type: none"> <li>Idade: 50 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Dulce Moreda</b> <ul style="list-style-type: none"> <li>Idade: 51 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>	<b>António Fernandes</b> <ul style="list-style-type: none"> <li>Idade: 43 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Vera Cristina Moreira Delgado</b> <ul style="list-style-type: none"> <li>Idade: 40 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>
<b>Angelo Santa</b> <ul style="list-style-type: none"> <li>Idade: 48 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Paulo Jorge</b> <ul style="list-style-type: none"> <li>Idade: 48 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>António Fernandes</b> <ul style="list-style-type: none"> <li>Idade: 43 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Vitor Marques</b> <ul style="list-style-type: none"> <li>Idade: 48 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>

Figura D.2: Página dos atletas que não pertencem a nenhuma equipa

Equipa: NAST			
Constituição: 12 atletas			
<b>FERNANDO JORGE MENDES CARNEIRO</b> <ul style="list-style-type: none"> <li>Idade: 38 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>	<b>Gonçalo Mota</b> <ul style="list-style-type: none"> <li>Idade: 48 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Julio Antonio Cunha Peixoto Braga</b> <ul style="list-style-type: none"> <li>Idade: 50 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Helder Fernando Vieira Azevedo</b> <ul style="list-style-type: none"> <li>Idade: 43 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>
<b>Sérgio Paulo dos Santos Sá</b> <ul style="list-style-type: none"> <li>Idade: 44 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>	<b>Analberto Barros dos Santos</b> <ul style="list-style-type: none"> <li>Idade: 50 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>	<b>Ricardo Jorge Dias Oliveira</b> <ul style="list-style-type: none"> <li>Idade: 38 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>	<b>Carlos Manuel Oliveira Bras</b> <ul style="list-style-type: none"> <li>Idade: 41 anos</li> <li>Escalação: SENIOR Masc</li> <li>Prova: Ultra Trail</li> </ul>
<b>Ana Carolina Lopes Ferreira</b> <ul style="list-style-type: none"> <li>Idade: 47 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Sónia Dolores Ferreira Alves</b> <ul style="list-style-type: none"> <li>Idade: 41 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Deborah Alexandra Lopes Moreira</b> <ul style="list-style-type: none"> <li>Idade: 49 anos</li> <li>Escalação: SENIOR Fem</li> <li>Prova: Corrida da Geira</li> </ul>	<b>Julio Antonio Cunha Peixoto Braga</b> <ul style="list-style-type: none"> <li>Idade: 50 anos</li> <li>Escalação: M40</li> <li>Prova: Ultra Trail</li> </ul>

Figura D.3: Página dos atletas que pertencem à equipa NAST

## Apêndice E

# Código do Programa

---

```
1 import re
2 import sys
3 from datetime import datetime, timedelta
4 import similar
5
6 fp = None
7 inscritos = list()
8
9 while not fp:
10     opt = input("1) Usar ficheiro predefinido (inscritos-form.json).\n\
11 2) Especificar ficheiro a usar.\n\n\
12 Escolha a opção pretendida: ").strip()
13     if opt == '1':
14         fp = "inscritos-form.json"
15     elif opt == '2':
16         fp = input("Introduza o caminho para o ficheiro que pretende utilizar: ").strip()
17
18 with open(fp, encoding="utf-8") as f:
19     m = re.search(r'\s*inscritos: \[\s*((?:.\n)*)\]\s*', f.read())
20     if not m:
21         sys.exit("Erro - Ficheiro inválido")
22     string_inscritos = m.group(1)
23
24 lista_inscritos = re.split(r'\s*,\s*{', string_inscritos.strip("\t\n{}"))
25 for i, inscrito in enumerate(lista_inscritos):
26     print(f"Lendo ficheiro ... {round(i*100/len(lista_inscritos),2)}% concluído")
27     insc_dict = dict()
28     for linha in inscrito.splitlines():
29         linha_stripped = linha.strip()
30         if linha_stripped:
31             m = re.search(r'"(\w+)"\s*:\s*"(.*)"\s*,?', linha_stripped)
32             insc_dict[m.group(1)] = m.group(2)
33
34     not_duplicate = all(any(similar.similarity(inscrito[key].lower(), insc_dict[key].lower()) > 2 for key in inscrito) for inscrito in inscritos)
35
36     if not_duplicate: inscritos.append(insc_dict)
37
```

```

38 print("\nFicheiro carregado com sucesso!")
39
40 while True:
41     opt = input("""\n1) Nomes dos concorrentes individuais de Valongo.
42 2) Nome, email e prova dos concorrentes chamados "Paulo" ou "Ricardo" que usam o
    GMail.
43 3) Informação dos atletas da equipa "TURBULENTOS"
44 4) Lista ordenada de escalões e atletas por escalão.
45 5) Gerar página HTML com informação sobre provas e equipas.\n
46 0) Sair.\n
47 Escolha a opção pretendida: """).strip()
48     if opt == '0': break
49
50     elif opt == '1': # alinea a
51         print("\nNomes dos concorrentes individuais de Valongo:")
52         for inscrito in inscritos:
53             if similar.similarity(inscrito["equipa"].lower(), "individual") < 2 and (n
                    := re.search(r'(?i:valongo)', inscrito["morada"])):
54                 print("-", inscrito["nome"].upper())
55
56     elif opt == '2': # alinea b
57         print("\nNome, email e prova dos concorrentes chamados "Paulo" ou "Ricardo"
            "\nque usam o GMail:")
58         for inscrito in inscritos:
59             if m := re.search(r'@gmail\.com$', inscrito["email"]):
60                 if o := re.search(r'(?i:Paulo|Ricardo)', inscrito["nome"]):
61                     print(f'-{inscrito["nome"]};{inscrito["email"]};{inscrito["
                        prova"]}')
62
63     elif opt == '3': # alinea c
64         turbulentos = list()
65         for inscrito in inscritos:
66             if m := re.fullmatch(r'(?i:turbulentos)', inscrito["equipa"]):
67                 turbulentos.append(inscrito)
68         i = 0
69         while True:
70             print(f'\nNome:{inscrito[i]["nome"]};\nData de nascimento:{inscrito[i]
                    ["dataNasc"]};\nMorada:{inscrito[i]["morada"]};\n
                    Email:{inscrito[i]["email"]};\nProva:{inscrito[i]["prova"]};\n
                    Escalão:{inscrito[i]["escalao"]}')
71             print(f"\nPágina {i+1} de {len(turbulentos)}\n\n[a] ver anterior; [s] ver
                    seguinte; [e] sair")
72             opt = input()
73             if opt == 'a' and i > 0: i -= 1
74             elif opt == 's' and i < len(turbulentos) - 1: i += 1
75             elif opt == 'e': break
76
77     elif opt == '4': # alinea d
78         esc_dict = dict()
79         for inscrito in inscritos:
80             if re.match(r'(\w+)', inscrito["escalao"]):
81                 x = esc_dict.get(inscrito["escalao"], 0)
82                 esc_dict[inscrito["escalao"]] = x+1
83         for escalao in sorted(esc_dict):

```



```

84         print("Escalão:", escalao, ";_Número_de_atletas:", esc_dict[escalao])
85
86     elif opt == '5': # alinea e
87         today = datetime.today()
88         equipas = dict()
89         provas = dict()
90         for i, inscrito in enumerate(inscritos):
91             prova = inscrito["prova"]
92             equipa = inscrito["equipa"].upper()
93             for eqp in equipas:
94                 if similar.similarity(eqp.upper(), equipa) < 2:
95                     equipa = eqp.upper()
96                 break
97             if not (m := re.match(r'|N/D|S/_CLUBE', equipa)):
98                 equipas.setdefault(equipa, list()).append(inscrito)
99                 provas.setdefault(prova, dict()).setdefault(equipa, list()).append(
100                     inscrito)
101             else:
102                 provas.setdefault(prova, dict()).setdefault("INDIVIDUAL", list()).
103                     append(inscrito)
104
105 # -----
106 #                               FICHEIRO PROVAS
107 # -----
108
109         with open("equipas.html", "w", encoding="utf-8") as f:
110             f.write("""<!DOCTYPE html>\n
111 <html>
112     <head>
113         <title>Equipas</title>
114         <link rel="stylesheet" href="style.css">
115     </head>
116     <body>
117         <h1>PROVAS</h1>
118         <div class="provas">
119             """
120             for prova in provas:
121                 f.write(f"""         <div class="prova">
122 <h2>{prova}</h2>
123             """
124                 for equipa in sorted(provas[prova].keys(), key=lambda x: len(provas[
125                     prova][x]), reverse=True):
126                     f.write(f"""                 <a href=". /equipas/{''.join(x for x in equipa
127                     if x.isalnum())}.html">
128 <span class="equipa">
129                     <h2>{equipa if equipa != "INDIVIDUAL" else "Sem equipa"}</h2>
130                     <div class="num_a"><p>{len(provas[prova][equipa])} atleta{'s' if len(
131                     provas[prova][equipa]) > 1 else ''}</p></div>
132 </span>
133 </a>
134             """
135                 f.write("        </div>\n")

```

```

133         f.write( """          </div>
134     </body>
135 </html>
136 """ )
137
138 # -----
139 #                               FICHEIROS EQUIPAS
140 # -----
141
142     for equipa in equipas:
143         with open(f"equipas/{''.join(x_ for x_ in equipa_ if x_.isalnum())}.html", "w",
144                 , encoding="utf-8") as ff:
145             ff.write(f """<!DOCTYPE html>
146
147 <html>
148     <head>
149         <title>{equipa}</title>
150         <link rel="stylesheet" href="style.css">
151     </head>
152     <body>
153         <h1>{"Equipa: " + equipa if equipa != "INDIVIDUAL" else "Atletas sem
154             equipa"}</h1>
155         <h2>Constituição: {len(equipas[equipa])} atleta{'s' if len(equipas[equipa
156             ]) != 1 else ''}</h2>
157         <div class="atletas">
158
159             for atleta in equipas[equipa]:
160                 try:
161                     if "ou" not in atleta["dataNasc"]:
162                         birth = datetime.strptime(atleta["dataNasc"], "%d/%m/%y")
163                         if birth > datetime.today(): birth = birth.replace(year =
164                             birth.year - 100)
165                     except ValueError:
166                         birth = None
167                     ff.write(f """          <div class="atleta">
168 <h3>{atleta["nome"]}</h3>
169 <ul>
170     <li>Idade: {str(today.year - birth.year - ((today.month, today
171         .day) < (birth.month, birth.day))) + " anos" if birth else
172         "-"}</li>
173     <li>Escalão: {atleta["escalao"] or "-"}</li>
174     <li>Prova: {atleta["prova"]}</li>
175 </ul>
176 </div>
177 """ )
178             ff.write( """          </div>
179
180     </body>
181 </html>
182 """ )
183
184     print("\nFicheiro HTML gerado com sucesso!")
185
186 else:

```

```
181         continue
182     input("\\nPrima ENTER para continuar.")
```

---