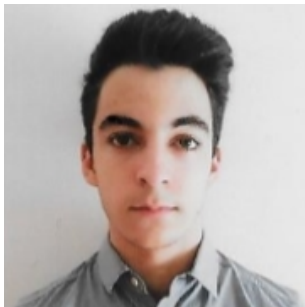


Processamento de Linguagens e Compiladores (3<sup>o</sup> Ano)

**Trabalho Prático 1**

Relatório de Desenvolvimento

Bruno Fernandes  
(a95972)



Nelson Almeida  
(a95652)



Nuno Costa  
(a97610)



01/11/2022

# Conteúdo

<b>1</b>	<b>Introdução</b>	<b>2</b>
<b>2</b>	<b>Análise e Especificações</b>	<b>4</b>
2.1	Processador de Pessoas listadas nos Róis de Confessados . . . . .	4
2.2	Processador de Registos de Exames Médicos Desportivos . . . . .	4
<b>3</b>	<b>Resolução do Exercício 1</b>	<b>5</b>
<b>4</b>	<b>Resolução do Exercício 2</b>	<b>10</b>
<b>5</b>	<b>Testes e Resultados obtidos</b>	<b>18</b>
5.1	Problema 1 . . . . .	18
5.1.1	Frequência de Processos Por Ano . . . . .	18
5.1.2	Frequência de nomes próprios e apelidos por séculos . . . . .	18
5.1.3	Frequência dos vários tipos de relação . . . . .	19
5.1.4	Imprimir os 20 primeiros registos em formatos Json . . . . .	19
5.2	Problema 2 . . . . .	22
5.2.1	index.html . . . . .	22
5.2.2	Datas extremas dos registos no dataset . . . . .	22
5.2.3	Distribuição por modalidade em cada ano e no total . . . . .	23
5.2.4	Distribuição por ano e género . . . . .	24
5.2.5	Distribuição por morada . . . . .	24
5.2.6	Porcentagem de aptos e não aptos por ano . . . . .	25
<b>6</b>	<b>Conclusão</b>	<b>26</b>
6.1	Bibliografia . . . . .	27

# Capítulo 1

## Introdução

No âmbito da disciplina de Processamentos de Linguagens e Compiladores foi-nos proposto pelo docente Pedro Rangel Henriques um trabalho de pesquisa cujo objetivo principal é consolidar a aprendizagem da utilização de linguagens regulares em programas em linguagem Python.

Decidimos então resolver os exercícios 1 e 2.

O exercício 1 tem como intuito a construção de vários programas em Python para processar o ficheiro de texto "processos.txt" com a intenção de calcular as frequências de alguns elementos pedidos ao longo do exercício.

O exercício 2 tem como objetivo a construção de vários programas para processar o dataset "emd.csv" e produzir o que é solicitado durante o exercício.

Neste documento apresentamos a nossa resolução para cada um dos problemas.

## Estrutura do Relatório

O relatório está organizado da seguinte forma:

Começamos por fazer uma pequena introdução, capítulo 1, onde referimos o objetivo deste trabalho e explicamos resumidamente o intuito de cada problema que decidimos resolver.

No capítulo 2 apresentamos o enunciado dos problemas escolhidos.

O capítulo 3 contém a resolução do primeiro problema e no 4 a resolução do segundo.

No capítulo 5 mostramos os resultados obtidos.

E, por fim, o ultimo capítulo contém a conclusão do trabalho realizado.

## Capítulo 2

# Análise e Especificações

### 2.1 Processador de Pessoas listadas nos Róis de Confessados

Construa agora um ou vários programas Python para processar o texto 'processos.txt' com o intuito de calcular frequências de alguns elementos (a ideia é utilizar arrays associativos para o efeito) conforme solicitado a seguir:

- a) Calcula a frequência de processos por ano (primeiro elemento da data);
- b) Calcula a frequência de nomes próprios (o primeiro em cada nome) e apelidos (o ultimo em cada nome) por séculos;
- c) Calcula a frequência dos vários tipos de relação: irmão, sobrinho, etc.
- d) Imprimir os 20 primeiros registos num novo cheiro de output mas em formato Json.

### 2.2 Processador de Registos de Exames Médicos Desportivos

Neste exercício pretende-se trabalhar com um dataset gerado no âmbito do registo de exames médicos desportivos. Construa, então, um ou vários programas Python para processar o dataset "emd.csv" e produzir o solicitado nas alíneas seguintes:

- 1. Página principal: de nome "index.html", contendo os seguintes indicadores estatísticos:
  - a) Datas extremas dos registos no dataset;
  - b) Distribuição por modalidade em cada ano e no total;
  - c) Distribuição por idade e género (para a idade, considera apenas 2 escalões:  $\leq 35$  anos e  $> 35$ );
  - d) Distribuição por morada;
  - e) Percentagem de aptos e não aptos por ano.
- 2. Página do indicador: clicando no indicador na página principal, devemos saltar para a página do indicador onde temos a informação que permitiu obter esse indicador. Por exemplo, para a distribuição por morada, a página deverá apresentar uma lista de moradas, ordenada alfabeticamente e para cada morada deverá apresentar uma sublista de registos, ordenada alfabeticamente por nome de atleta (com os dados: nome do atleta, modalidade).

## Capítulo 3

# Resolução do Exercício 1

Inicialmente criamos um ficheiro principal "main.py", onde instalamos os packages necessários para a execução do programa pedindo ao utilizador o caminho (path) para o ficheiro que pretende utilizar, tendo como pré-definido o ficheiro "processos.txt". Neste mesmo ficheiro criamos a função "parser", que recebe como parâmetro o ficheiro pretendido, gerando um conjunto de arrays com o conteúdo do mesmo que será usado posteriormente para cada uma das 4 alíneas.

- a) Para resolver a primeira alínea, desenvolvemos o seguinte código para calcular a frequência de processos por ano.

```
1  def __init__(self, processos):
2      print(self.processesPerTimeSpan(processos))
3
4  def processesPerTimeSpan(self, processos):
5      processosPorData = {}
6      while (processos != []):
7          ano = re.match("(([0-9]{1,4})-*)", processos[0][1]).group(2)
8          processosNoAno, processosNoutrosAnos = self.filtrarPorAno(processos, ano)
9          processosPorData[ano] = processosNoAno
10         processos = processosNoutrosAnos
11     return processosPorData
12
13 def filtrarPorAno(self, processos, yearToFilter):
14     procPerYear = 0
15     processosNoutrosAnos = []
16     for pessoa in processos:
17         flag = 0
18         for elem in pessoa:
19             if re.search(yearToFilter+r"-[0-9]{2}-[0-9]{1,2}", elem):
20                 flag = 1
21             if flag == 1:
22                 procPerYear += 1
23             else:
24                 processosNoutrosAnos.append(pessoa)
25     return (procPerYear, processosNoutrosAnos)
26
```

Começamos por criar uma função "processesPerTimeSpan" para obter um dicionário com a frequência de processos para cada ano.

Para descobrir cada ano fizemos uso da função match com a expressão regular "(([0-9]{1,4})-\*)" e invocamos a função auxiliar "filtrarPorAno".

```

1 while (processos != []):
2     ano = re.match("([0-9]{1,4})-*)", processos[0][1]).group(2)
3     processosNoAno, processosNoutrosAnos = self.filtrarPorAno(processos, ano)
4

```

A função auxiliar "filtrarPorAno" tem como resultado o número de processos do ano passado como argumento e uma lista contendo todos os processos restantes existentes no ficheiro.

Assim sendo, para cada processo utilizamos a função search com a expressão regular "[0-9]{2}-[0-9]{1,2}" para encontrar todas as datas desse ano passado com argumento.

```

1 for pessoa in processos:
2     flag = 0
3     for elem in pessoa:
4         if re.search(yearToFilter+r"[0-9]{2}-[0-9]{1,2}", elem):
5             flag = 1
6     if flag == 1:
7         procPerYear += 1
8     else:
9         processosNoutrosAnos.append(pessoa)
10

```

- b) Para a resolução desta alínea, desenvolvemos o seguinte código de modo a obter a frequência de nomes próprios e apelidos por século.

```
1 def __init__(self, processos):
2     print(self.nomesPorSeculo(processos))
3
4 def nomesPorSeculo(self, processos):
5     dic = {"firstNames": {}, "lastNames": {}}
6     for processo in processos:
7         ano = re.match("([0-9]{1,4})-*)", processo[1]).group(2)
8         sec = int(ano)//100 + 1
9         if sec not in dic["firstNames"]:
10             dic["firstNames"][sec] = {}
11             dic["lastNames"][sec] = {}
12         i = 2
13         while (i <= 4):
14             names = re.split(" ", processo[i])
15             fstName = names[0]
16             lastName = names[-1]
17             if fstName not in dic["firstNames"][sec]:
18                 dic["firstNames"][sec][fstName] = 0
19             if lastName not in dic["lastNames"][sec]:
20                 dic["lastNames"][sec][lastName] = 0
21             dic["firstNames"][sec][fstName] += 1
22             dic["lastNames"][sec][lastName] += 1
23             i += 1
24         return dic
25
```

Começamos por criar uma função "nomesPorSeculo" com o objetivo de obter um dicionário com duas chaves, uma para os nomes próprios e a outra para os apelidos.

Para tal começamos por encontrar o ano em cada processo usando a função match e a expressão regular "([0-9]{1,4})-\*)" e deste modo sabemos a que século pertence esse ano adicionando-o, caso não exista, ao dicionário.

```
1 for processo in processos:
2     ano = re.match("([0-9]{1,4})-*)", processo[1]).group(2)
3     sec = int(ano)//100 + 1
4     if sec not in dic["firstNames"]:
5         dic["firstNames"][sec] = {}
6         dic["lastNames"][sec] = {}
7
```

Para obter o nome próprio e o apelido recorreremos à separação dos nomes de cada pessoa por utilizando a função split, partindo-os nomes pelos espaços existentes entre cada um.

Após isso adicionamos os nomes próprios e apelidos aos respetivos séculos.

```
1 names = re.split(" ", processo[i])
2 fstName = names[0]
3 lastName = names[-1]
4 if fstName not in dic["firstNames"][sec]:
5     dic["firstNames"][sec][fstName] = 0
6 if lastName not in dic["lastNames"][sec]:
7     dic["lastNames"][sec][lastName] = 0
8 dic["firstNames"][sec][fstName] += 1
9 dic["lastNames"][sec][lastName] += 1
10
```



- c) Para a resolução deste problema, desenvolvemos o código que se segue e calcula a frequência dos vários tipos de relação.

```
1 def __init__(self, processos):
2     print(self.tiposParentesco(processos))
3
4 def tiposParentesco(self, pessoas):
5     tipos = {}
6     for pessoa in pessoas:
7         for elem in pessoa:
8             parent = re.findall(r'(?<=[a-z]\,)([PMTISANB][A-Za-z ]+|Filho|Filhos)(?=\. )',
9                                 elem)
10            for x in parent:
11                if x in tipos.keys():
12                    tipos[x] = tipos.get(x, 0) + 1
13                else:
14                    tipos[x] = 1
```

Começamos por criar uma função "tiposParentesco" que tem como retorno um dicionário que a cada tipo de parentesco associa o número de vezes que o mesmo aparece no ficheiro "processos.txt".

Primeiramente começamos por identificar cada grau de parentesco encontrado no ficheiro. E visto que podem existir mais do que um por cada pessoa, recorremos à função findall com a expressão regular "(?<=[a-z]\,)([PMTISANB][A-Za-z ]+|Filho|Filhos)(?=\. )".

```
1     tipos = {}
2     for pessoa in pessoas:
3         for elem in pessoa:
4             parent = re.findall(r'(?<=[a-z]\,)([PMTISANB][A-Za-z ]+|Filho|Filhos)(?=\. )',
5                                 elem)
```

- d) Para a resolução da última alínea, pretendíamos imprimir os 20 primeiros registos num novo ficheiro de output mas em formato Json e para tal desenvolvemos o código abaixo.

```
1  def __init__(self, processos, limit):
2      if (self.jsonFormatter(processos[0:limit])):
3          print("Ficheiro criado com sucesso")
4      else:
5          print("Ocorreu um erro a criar o ficheiro")
6
7  def jsonFormatter(self, pessoas):
8      file = open("./TP1/ex1/rois.json", "w")
9      file.write("{\"data\":[")
10     templateJson = '''{"folder_id":#,"date":#,"name":#,"father":#,"mother":#,"obs":#}
11     ,,,
12     pos = 0
13     for pessoa in pessoas:
14         pessoa.append("")
15         person = templateJson
16         for i in pessoa:
17             person = re.sub("#", f"\{i}\\"", person, count=1)
18             if (pos < len(pessoas)-1):
19                 file.write(person+",\n")
20             else:
21                 file.write(person)
22             pos += 1
23     file.write("]}\n")
24     file.close()
25     return True
```

Inicialmente criamos a função "jsonFormatter" que irá retornar o novo output em formato Json.

Esta função começa por escrever de acordo com o template Json no ficheiro criado por nós "rois.json" e retornar em caso de sucesso "Ficheiro criado com sucesso".

## Capítulo 4

# Resolução do Exercício 2

Começamos por criar um ficheiro principal "main.py" onde se começa instala os packages necessários para a execução do programa e pede ao utilizador o path para o ficheiro que pretende utilizar, tendo como pré-definido o ficheiro "emd.csv". Com a função "parser" recebendo como parâmetro o ficheiro pretendido, obtemos um conjunto de arrays com o conteúdo do mesmo que é usado posteriormente para cada uma das 5 alíneas.

1. Para exercício 1 criamos um ficheiro "index.html" que faz referência às restantes alíneas.

```
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4      <meta charset="UTF-8">
5      <meta http-equiv="X-UA-Compatible" content="IE=edge">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <title>Ex2 TP1 PLC</title>
8      <link rel="stylesheet" href="main.css">
9  </head>
10 <body>
11     <div className="mainBar" style="padding-top: 50vh;">
12         <a href="a.html">A</a>
13         <a href="b.html">B</a>
14         <a href="c.html">C</a>
15         <a href="d.html">D</a>
16         <a href="e.html">E</a>
17     </div>
18 </body>
19 </html>
20
```

Para as alíneas b), c), d), e) e f) utilizamos a biblioteca "pandas" para fazer a função "htmlGenerator" para através de um dataframe gerar a representação da tabela em html, esta é guardada numa variável para depois ser processada.

Utilizamos também a biblioteca "matplotlib" para podermos gerar gráficos.

Para adicionar o título da página, os gráficos gerados e o "inline styling" ao código HTML gerado, utilizamos o re.sub que faz a substituição da tag table gerada por o que pretendemos adicionar.

- a) Para resolver a primeira alínea, criamos um array com todas as datas do ficheiro seleccionado e determinamos o mínimo e o máximo deste.

```
1 def a(self, file):
2     datas = [i[2] for i in file]
3     file = open("./TP1/ex2/1/website/a.html", "w")
4     templateText =
5         '''<!DOCTYPE html >
6         <html lang = "en" >
7         <head >
8             <meta charset = "UTF-8" >
9             <meta http-equiv = "X-UA-Compatible" content = "IE=edge" >
10            <meta name = "viewport" content = "width=device-width, initial-scale
11            =1.0" >
12            <title >a</title >
13            <link rel="stylesheet" href="./main.css">
14        </head >
15        <body >
16        </body >
17        </html >'''
18     newText = re.sub(
19         "<body >", f'''\n<body >\n<p class="title">Datas extremas dos registos no
20         dataset</p>\n<p class="minData">Extremo inferior do dataset: {min(datas)}</p>\n
21         <p class="maxData">Extremo superior do dataset: {max(datas)}</p>''',
22         templateText)
23     file.write(newText)
24     file.close()
```

- b) Para a resolução da alínea b) e assim determinar a distribuição por modalidade em cada ano e no total, desenvolvemos o seguinte código.

```
1 def b(self, data):
2     sportsPerYear, totalSports = self.sportsPerYearandTotal(data)
3     indexes = sorted(sportsPerYear)
4     sportsPerYear = {i: sportsPerYear[i] for i in indexes}
5     totalSports = dict(sorted(totalSports.items(), key=lambda x: x[1]))
6     self.plotter(sportsPerYear, totalSports)
7     self.htmlGenerator(sportsPerYear, totalSports)
8
9 def sportsPerYearandTotal(self, data):
10    sportsPerYear = {}
11    totalSports = {}
12    while (data):
13        year = re.match("([0-9]{1,4})-*", data[0][2]).group(2)
14        sportsYear, updatedData = self.sportsPerYear(data, year)
15        sportsPerYear[year] = sportsYear
16        data = updatedData
17
18    for year in sportsPerYear:
19        for sport in sportsPerYear[year]:
20            if sport not in totalSports:
21                totalSports[sport] = 0
22                totalSports[sport] += sportsPerYear[year][sport]
23    return sportsPerYear, totalSports
24
25 def sportsPerYear(self, data, yearToFilter):
26    sportsPerY = {}
27    newData = []
28    for person in data:
29        if re.match(yearToFilter+r"(-[0-9]{1,2}){2}", person[2]):
30            sport = person[8]
31            if sport not in sportsPerY:
32                sportsPerY[sport] = 0
33                sportsPerY[sport] += 1
34            else:
35                newData.append(person)
36    return sportsPerY, newData
37
38 def htmlGenerator(self, data1, data2):
39    file = open(f"./TP1/ex2/1/website/b.html", "w")
40    templateText1 = pd.DataFrame(data1).to_html()
41    templateText1 = re.sub(
42        r'''<table border="1" class="dataframe">''', '''<h1 class="title">
43        Distribui o por modalidade em cada ano e no total</h1>\n<div class="images"
44        style="display: flex; width: 100%; padding-bottom: 2rem;">\n\n</div>\n<link rel="stylesheet" href="./main.css">\n<table border="1"
47        class="dataframe">''', templateText1
48    )
49    file.write(templateText1)
50    templateText2 = pd.DataFrame(data2, index=["Praticantes"]).to_html()
51    file.write(templateText2)
52    file.close()
53
54 def plotter(self, data1, data2):
55    df1 = pd.DataFrame(data1)
```

```

51     df1.plot(kind="bar")
52     plt.savefig("./TP1/ex2/1/website/src/b1.png")
53     df2 = pd.DataFrame(data2, index=["Praticantes"])
54     df2.plot(kind="bar")
55     plt.savefig("./TP1/ex2/1/website/src/b2.png")
56

```

Começamos por criar a função auxiliar "sportsPerYearandTotal" para obtermos dois dicionários, um com os desportos distribuídos por anos e outro com todos os desportos existentes no ficheiro. Para cada pessoa existente no ficheiro utilizamos o match com a expressão regular "([0-9]{1,4}-\*)" para encontrar o ano que vai ser utilizado como parâmetro em outra função auxiliar "sportsPerYear".

```

1     year = re.match("([0-9]{1,4}-*)", data[0][2]).group(2)
2     sportsYear, updatedData = self.sportsPerYear(data, year)
3

```

Na função "sportsPerYear" começamos por criar um dicionário com os desportos existentes nesse ano e a sua frequência.

```

1     if re.match(yearToFilter+r"(-[0-9]{1,2}){2}", person[2]):
2         sport = person[8]
3         if sport not in sportsPerY:
4             sportsPerY[sport] = 0
5             sportsPerY[sport] += 1
6

```

E um array onde vão ser adicionados apenas os dados relativos aos restantes anos (diferentes do ano passado como parâmetro), para assim, evitar repetições e garantir um melhor controlo.

```

1     else:
2         newData.append(person)
3

```

No fim vão ser retornados o dicionário "sportsPerY" e o array "newData" que são utilizados para atualizar o dicionário "sportsPerYear" e os dados do array "data" do método anterior.

```

1     sportsPerYear[year] = sportsYear
2     data = updatedData
3

```

Por último, com os dados colecionados no dicionário "sportsPerYear", soma o número de desportistas para cada ano e para cada desporto e coleciona no dicionário de "totalSports".

```

1     for year in sportsPerYear:
2         for sport in sportsPerYear[year]:
3             if sport not in totalSports:
4                 totalSports[sport] = 0
5                 totalSports[sport] += sportsPerYear[year][sport]
6

```

c) Para a resolução da alínea c) e assim determinar a distribuição por idade e género, desenvolvemos o seguinte código.

```
1  def c(self, processos):
2      data = processos[1:] # aqui ficam tudo menos os indicadores
3      data = self.distByAgeAndGender(data)
4      self.plotter(data)
5      self.htmlGenerator(data)
6
7  def distByAgeAndGender(self, data):
8      dist = {"<35": {}, ">=35": {}}
9      ageIndex = 5
10     genderIndex = 6
11     for person in data:
12         ageRange = ">=35"
13         if int(person[ageIndex]) < 35:
14             ageRange = "<35"
15         if person[genderIndex] not in dist[ageRange]:
16             dist[ageRange][person[genderIndex]] = 0
17             dist[ageRange][person[genderIndex]] += 1
18     return dist
19
20 def htmlGenerator(self, data1):
21     file = open(f"./TP1/ex2/1/website/c.html", "w")
22     templateText = pd.DataFrame(data1).to_html()
23     templateText = re.sub(
24         r'''<table border="1" class="dataframe">''', '''<p class="title">
Distribui o por idade e g nero </p>\n<div class="images" style="padding-
bottom: 2rem;">\n</div>\n<link rel="stylesheet"
href="./main.css">\n<table border="1" class="dataframe">''', templateText)
25     file.write(templateText)
26     file.close()
27
28 def plotter(self, data):
29     df1 = pd.DataFrame(data)
30     df1.plot(kind="bar")
31     plt.savefig("./TP1/ex2/1/website/src/c.png")
32
```

Criamos a função auxiliar "distByAgeAndGender" que começa por criar um dicionário onde vão ser guardadas as informações, vai percorrer cada pessoa do ficheiro, verificar se a sua idade é maior ou menor que 35 anos, verificar o género dessa pessoa e incrementar na contação.

- d) Para a resolução da alínea d) e assim determinar a distribuição por morada, desenvolvemos o seguinte código.

```
1  def d(self, processos):
2      data = processos[1:]
3      data = self.distByAddress(data)
4      self.plotter(data)
5      self.htmlGenerator(data)
6
7  def distByAddress(self, data):
8      distByAddress = {}
9      for person in data:
10         address = person[7]
11         if address not in distByAddress:
12             distByAddress[address] = 0
13             distByAddress[address] += 1
14     return distByAddress
15
16 def htmlGenerator(self, data1):
17     file = open(f"./TP1/ex2/1/website/d.html", "w")
18     templateText = pd.DataFrame(data1, index=[""]).to_html()
19     templateText = re.sub(
20         r'''<table border="1" class="dataframe">''', '''<p class="title">
Distribui o por morada</p>\n<div class="images" style="padding-bottom: 2rem
; ">\n</div>\n<link rel="stylesheet" href="./main.
css">\n<table border="1" class="dataframe">''', templateText
21     )
22     file.write(templateText)
23     file.close()
24
25 def plotter(self, data):
26     if len(data) >= 50:
27         df1 = pd.DataFrame(data, index=[""])
28         df1.plot(kind="bar")
29         plt.savefig("./TP1/ex2/1/website/src/d.png")
30     else:
31         df1 = pd.DataFrame(data, index=[""])
32         df1.plot(kind="bar")
33         plt.savefig("./TP1/ex2/1/website/src/d.png")
34
```

Criamos a função auxiliar "distByAddress" que começa por criar um dicionário onde vão ser guardadas as informações, vai percorrer cada pessoa do ficheiro, adicionar ao dicionário a sua morada caso esta não exista e incrementar o número de pessoas que vivem nessa morada.



- e) Para a resolução da alínea e) e assim determinar a percentagem de pessoas aptas e não aptas por ano, desenvolvemos o seguinte código.

```

1      def e(self, data):
2          data = data[1:]
3          readinessStats = {}
4          while (data):
5              year = re.match("([0-9]{1,4})-*)", data[0][2]).group(2)
6              (ready, total), newData = self.readinessPerYear(data, year)
7              readinessPerc = (ready/total)
8              readinessStats[year] = {
9                  "ready": float("{:.2f}".format(100*readinessPerc)),
10                 "not_ready": float("{:.2f}".format(100*(1-readinessPerc)))
11             }
12             data = newData
13             self.plotter(readinessStats)
14             self.htmlGenerator(readinessStats)
15
16     def readinessPerYear(self, data, year):
17         ready = 0
18         total = 0
19         newData = []
20         for person in data:
21             if re.match(year+r"([0-9]{1,2}){2}", person[2]):
22                 if person[-1] == "true":
23                     ready += 1
24                     total += 1
25             else:
26                 newData.append(person)
27         return (ready, total), newData
28
29     def htmlGenerator(self, data1):
30         file = open(f"./TP1/ex2/1/website/e.html", "w")
31         templateText = pd.DataFrame(data1).to_html()
32         templateText = re.sub(
33             r'''<table border="1" class="dataframe">''', '''<p class="title">
Percentagem de aptos e n o aptos por ano</p>\n<div class="images" style="
padding-bottom: 2rem;>\n</div>\n<link rel="
stylesheet" href="main.css">\n<table border="1" class="dataframe">''',
templateText
34         )
35         file.write(templateText)
36         file.close()
37
38     def plotter(self, data):
39         df1 = pd.DataFrame(data)
40         df1.plot(kind="bar")
41         plt.savefig("./TP1/ex2/1/website/src/e.png")
42

```

Começamos por percorrer cada pessoa do ficheiro e utilizando o match com a expressão regular "([0-9]{1,4})-\*)" encontrar o ano que vai ser utilizado como parâmetro na função auxiliar "readinessPerYear".

```

1      year = re.match("([0-9]{1,4})-*)", data[0][2]).group(2)
2      (ready, total), newData = self.readinessPerYear(data, year)
3

```

Esta função "readinessPerYear" começa por criar duas variáveis "ready" e "total" que vão servir para saber o número de pessoas aptas e o total de pessoas, respetivamente, e um array onde

vão ser adicionados apenas os dados relativos aos restantes anos (diferentes do ano passado como parâmetro), para assim, evitar repetições e garantir um melhor controlo.

Começamos por percorrer cada pessoa, verificar se o ano coincide com ano passado como parâmetro, incrementar na variável "ready" se tal se verificar e essa pessoa estiver apta.

```
1     if re.match(year+r"(-[0-9]{1,2}){2}", person[2]):
2         if person[-1] == "true":
3             ready += 1
4             total += 1
5
```

Na função principal, depois de receber o número de pessoas aptas e no total nesse ano, são calculadas as percentagens de atletas aptos e não aptos e colocadas no dicionário. Posteriormente os dados são atualizados.

```
1     readinessPerc = (ready/total)
2     readinessStats[year] = {
3         "ready": float("{:.2f}".format(100*readinessPerc)),
4         "not_ready": float("{:.2f}".format(100*(1-readinessPerc)))
5     }
6     data = newData
7
```

## Capítulo 5

# Testes e Resultados obtidos

### 5.1 Problema 1

#### 5.1.1 Frequência de Processos Por Ano

```
{'1894': 74, '1909': 39, '1867': 54, '1896': 75, '1904': 52, '1901': 59, '1883': 34, '1900': 50, '1902': 83, '1880': 62, '1889': 71, '1908': 51, '1869': 39, '1862': 40, '1906': 63, '1856': 72, '1875': 15, '1892': 57, '1733': 1188, '1778': 1066, '1691': 1084, '1730': 1200, '1899': 80, '1898': 91, '1877': 47, '1910': 27, '1881': 64, '1907': 47, '1884': 44, '1879': 55, '1895': 78, '1897': 71, '1707': 116, '1689': 628, '1713': 266, '1824': 251, '1703': 143, '1720': 172, '1890': 44, '1732': 2217, '1683': 142, '1863': 30, '1729': 47, '1694': 49, '1765': 9, '1754': 303, '1690': 278, '1755': 303, '1823': 213, '1708': 150, '1757': 32, '1699': 97, '1759': 90, '1712': 70, '1687': 97, '1738': 213, '1717': 267, '1684': 238, '1704': 358, '1688': 110, '1888': 70, '1734': 888, '1786': 352, '1798': 56, '1773': 482, '1821': 312, '1822': 303, '1809': 292, '1722': 457, '1680': 158, '1695': 44, '1728': 440, '1716': 238, '1849': 120, '1777': 1142, '1851': 66, '1785': 797, '1857': 80, '1686': 189, '1784': 361, '1780': 218, '1727': 169, '1788': 416, '1719': 387, '1847': 130, '1799': 103, '1829': 162, '1787': 743, '1805': 109, '1819': 242, '1844': 196, '1891': 73, '1731': 900, '1760': 280, '1741': 32, '1725': 272, '1802': 90, '1827': 97, '1885': 8, '1807': 602, '1710': 199, '1692': 246, '1706': 117, '1858': 46, '1739': 74, '1826': 171, '1714': 317, '1762': 342, '1743': 53, '1724': 110, '1697': 28, '1852': 74, '1740': 13, '1855': 86, '1723': 129, '1859': 87, '1811': 237, '1817': 207, '1685': 235, '1905': 42, '1893': 78, '1865': 40, '1848': 125, '1911': 16, '1882': 50, '1735': 153, '1871': 35, '1903': 20, '1850': 115, '1825': 155, '1843': 296, '1860': 77, '1812': 259, '1846': 86, '1845': 127, '1701': 108, '1746': 49, '1868': 22, '1715': 90, '1803': 167, '1830': 174, '1761': 381, '1766': 9, '1672': 17, '1876': 40, '1698': 153, '1726': 85, '1679': 83, '1750': 53, '1711': 75, '1810': 47, '1756': 75, '1721': 135, '1873': 46, '1831': 193, '1781': 208, '1818': 229, '1816': 352, '1783': 309, '1752': 80, '1866': 33, '1747': 58, '1771': 26, '1700': 82, '1682': 106, '1794': 57, '1742': 22, '1887': 42, '1836': 3, '1878': 76, '1886': 46, '1702': 139, '1767': 15, '1709': 148, '1779': 198, '1751': 107, '1681': 101, '1748': 46, '1678': 12, '1839': 85, '1770': 27, '1673': 19, '1792': 46, '1796': 86, '1828': 172, '1736': 47, '1737': 65, '1696': 42, '1853': 90, '1772': 122, '1718': 92, '1776': 36, '1804': 123, '1774': 38, '1832': 301, '1806': 106, '1820': 253, '1808': 311, '1854': 112, '1861': 62, '1833': 194, '1841': 47, '1705': 76, '1872': 42, '1874': 46, '1782': 155, '1693': 28, '1838': 10, '1797': 46, '1660': 1, '1864': 56, '1870': 12, '1744': 41, '1800': 107, '1655': 2, '1795': 55, '1840': 38, '1764': 24, '1763': 29, '1676': 9, '1768': 27, '1753': 37, '1842': 36, '1671': 14, '1801': 58, '1814': 11, '1815': 45, '1758': 22, '1749': 27, '1775': 30, '1745': 47, '1834': 13, '1769': 16, '1793': 45, '1665': 2, '1813': 12, '1635': 1, '1616': 1, '1791': 6, '1631': 2, '1658': 2, '1664': 5, '1789': 12, '1623': 3, '1675': 9, '1677': 6, '1666': 6, '1663': 3, '1674': 3, '1668': 5, '1661': 1, '1630': 4, '1628': 7, '1670': 3, '1633': 2, '1632': 3, '1669': 3, '1636': 1, '1667': 1, '1634': 1, '1625': 2, '1656': 2, '1620': 1, '1629': 1, '1622': 1, '1627': 2, '1790': 1, '1650': 1, '1662': 1}
```

#### 5.1.2 Frequência de nomes próprios e apelidos por séculos

```
{'firstNames': {19: {'Aarao': 1, 'Antonio': 3618, 'Francisca': 311, 'Abel': 3, 'Maria': 3695, 'Francisco': 1779, 'Antonia': 428, 'Joao': 2626, 'Cecilia': 5, 'Abelardo': 2, 'Jose': 3810, 'Leopoldina': 9, 'Abilio': 8, 'Sebastiao': 91, 'Ieresa': 509, 'Ana': 1270, 'Flora': 4, 'Acacio': 3, 'Albina': 18, 'Adelino': 13, 'Custodia': 333, 'Bernardo': 203, 'Rosa': 660, 'Manuel': 3469, 'Rosalia': 25, 'Joana': 328, 'Adolfo': 3, 'Carlota': 18, 'Custodio': 241, 'Emilia': 29, 'Adriano': 9, 'Miguel': 157, 'Adriao': 1, 'Bernardina': 29, 'Afonso': 6, '': 171, 'Agostinho': 71, 'Luisa': 345, 'Domingos': 918, 'Helena': 40, 'Benta': 53, 'Margarida': 65, 'Mateus': 19, 'Gregorio': 17, 'Jeronima': 40, 'Joaquina': 249, 'Cleto': 6, 'Bernardino': 76, 'Aires': 2, 'Luis': 478, 'Albano': 6, 'Albertino': 1, 'Alberto': 14, 'Florinda': 23, 'Gracinda': 5, 'Benedita': 2, 'Henrique': 38, 'Albino': 26, 'Engracia': 29, 'Mariana': 233, 'Violante': 22, 'Cipriana': 7, 'Inacia': 38, 'Josefa': 318, 'Alexandre': 99, 'Marinha': 8, 'Perpetua': 20, 'Caetana': 39, 'Teodora': 8, 'Feliciano': 10, 'Bento': 353, 'Isabel': 202, 'Senhorinha': 49, 'Catarina': 75, 'Martina': 1, 'Lucia': 1, 'Alexandrino': 3, 'Vicente': 38, 'Clarina': 1, 'Caetano': 114, 'Alfredo': 15, 'Apolonia': 2, 'Lidorio': 1, 'Alvaro': 21, 'Vitorino': 25, 'Julio': 26, 'Rita': 82, 'Lourenco': 56, 'Amador': 1, 'Amaro': 6, 'Gaspar': 68, 'Ambrosio': 7, 'Americo': 1, 'Lino': 16, 'Anacleto': 10, 'Estevas': 20, 'Ana,Exposta': 2, 'Anastacio': 10, 'Narciso': 36, 'Andre': 27, 'Joaquim': 768, 'Inacio': 53, 'Angelo': 4, 'Anibal': 3, 'Ludovina': 9, 'Aniceto': 3, 'Anselmo': 6, 'Eugenia': 15, 'Escolastica': 3, 'Antao': 2, 'Andresa': 5, 'Antonino': 2, 'Rosendo': 3, 'Angelica': 35, 'Simao': 42, 'Balbina': 9, 'Pedro': 156, 'Jeronimo': 95, 'Genoveva': 25, 'Domingas': 46, 'Ventura': 7, 'Felizarda': 7, 'Vitoria': 33, 'Claudina': 10, 'Brigida': 6,...}}, {'lastNames': {19: {'Silva': 1318, 'Barroso': 138, 'oliveira': 456, 'Rebelo': 121, 'Freitas': 194, 'Pereira': 1197, 'Araujo': 719, 'Guerreiro': 47, 'Monteiro': 164, 'Alves': 245, 'Rocha': 235, 'Arantes': 28, 'Sousa': 804, 'Santos': 244, 'Jesus': 156, 'Leite': 199, 'Fernandes': 626, 'Guimaraes': 161, 'Borges': 94, 'Teixeira': 443, 'Carneiro': 142, 'Barbosa': 357, 'Matos': 160, 'Coelho': 201, 'Dias': 356, 'Ferreira': 592, 'Peixoto': 155, 'Costa': 995, 'Cardoso': 143, 'Eiras': 10, 'Goncalves': 789, 'Afonso': 179, 'Almeida': 256, 'Campos': 123, 'Pimenta': 65, 'Miranda': 135, 'Pedrosa': 16, 'Balazeiro': 1, 'Barros': 309, 'Abreu': 253, 'Vieira': 266, 'Meireles': 65, 'Moura': 161, 'Segura': 4, 'Vaz': 128, 'Casimira': 2, 'Pinheiro': 126, 'Freire': 20, 'catalao': 2, '': 171, 'Sanfins': 1, 'Soares': 218, 'Penteado': 5, 'Brito': 173, 'Azevedo': 320, 'Reis': 119, 'Mendes': 113, 'Antunes': 158, 'Nogueira': 76, 'Botelho': 53, 'Vale': 99, 'Ribeiro': 456, 'Sotomaior': 48, 'Felizarda': 3, 'Rodrigues': 741, 'Alvares': 381, 'Cunha': 519, 'Domingues': 178, 'Ieresa': 161, 'Passos': 37, 'Macedo': 180, 'Meneses': 138, 'Queiros': 80, 'Carvalho': 667, 'Lopes': 345, 'Castro': 382, 'Martins': 470, 'Marques': 152, 'coimbra': 12, 'Gomes': 495, 'Flores': 16, 'Conceicao': 70, 'Fontes': 20, 'Basto': 39, 'Salgado': 52, 'Pedreira': 10, 'Lobo': 126, 'Chaves': 66, 'Sepulveda': 18, 'Machado': 360, 'Pinto': 267, 'Vasconcelos': 143, 'Junior': 84, 'Brandao': 61, 'Josefa': 84, 'Correia': 203, 'Faria': 214, 'Joaquina': 201, 'Novais': 57, 'Assuncao,Solteira': 2, 'Carreira,Viuva': 1, 'Braganca': 11, 'Coutinho': 130, 'Moreira': 115, 'Maria': 412, 'Dantas': 71, 'Pequeno': 5, 'Pimentel': 35, 'Homem': 1, 'Camelo': 31, 'Lima': 340, 'Joao': 23, 'Pires': 192,...}}}
```

### 5.1.3 Frequência dos vários tipos de relação

```
{'Tio Paterno': 2245, 'Tio Materno': 2463, 'Irmão': 13168, 'Primo Paterno': 205, 'Sobrinho Materno': 1698, 'Pai': 525, 'Filho': 346, 'Sobrinho Paterno': 1642, 'Irmaos': 686, 'Sobrinhos Maternos': 98, 'Irmão Paterno': 497, 'Neto Materno': 41, 'Sobrinhos Paternos': 57, 'Sobrinho Neto Paterno': 97, 'Primo': 673, 'Primo Materno': 283, 'Tio Avo Paterno': 154, 'Tio Avo Materno': 230, 'Irmão Materno': 55, 'Sobrinho Bisneto Paterno': 3, 'Tios Maternos': 20, 'Irmaos Paternos': 21, 'Sobrinho Neto Materno': 145, 'Avo Materno': 48, 'Bisavo Materno': 2, 'Filhos': 27, 'Avo Paterno': 11, 'Neto Paterno': 8, 'Tios Paternos': 12, 'Tio Bisavo Materno': 5, 'Primos': 13, 'Parente': 4, 'Primos Paternos': 1, 'Tio Bisavo Paterno': 6, 'Irmaos Maternos': 4, 'Sobrinhos Netos Paternos': 2, 'Sobrinho Neto': 2, 'Tio Avo': 3, 'Tio': 5, 'Sobrinhos Netos Maternos': 5, 'Meio Irmão': 3, 'Sobrinho Bisneto Materno': 3, 'Primos Maternos': 1}
```

### 5.1.4 Imprimir os 20 primeiros registos em formatos Json

```
1  {
2    "data": [
3      {
4        "folder_id": "575",
5        "date": "1894-11-08",
6        "name": "Aarao Pereira Silva",
7        "father": "Antonio Pereira Silva",
8        "mother": "Francisca Campos Silva",
9        "obs": ""
10     },
11     {
12       "folder_id": "582",
13       "date": "1909-05-12",
14       "name": "Abel Almeida",
15       "father": "Antonio Manuel Almeida",
16       "mother": "Teresa Maria Sousa",
17       "obs": ""
18     },
19     {
20       "folder_id": "569",
21       "date": "1867-05-23",
22       "name": "Abel Alves Barroso",
23       "father": "Antonio Alves Barroso",
24       "mother": "Maria Jose Alvares Barroso",
25       "obs": "Bento Alvares Barroso,Tio Paterno. Proc.32057.    Domingos Jose Alvares
Barroso,Tio Materno. Proc.32235."
26     },
27     {
28       "folder_id": "576",
29       "date": "1896-11-28",
30       "name": "Abel Augusto Oliveira",
31       "father": "Francisco Jose Oliveira",
32       "mother": "Antonia Rosa Rebelo",
33       "obs": "Jose Antonio Oliveira,Irmão. Proc.5020."
34     },
35     {
36       "folder_id": "579",
37       "date": "1904-05-27",
38       "name": "Abel Gomes Abreu Reis",
39       "father": "Antonio Gomes Abreu",
40       "mother": "Ana Sequeira Reis",
41       "obs": ""
42     },
43     {
44       "folder_id": "579",
45       "date": "1904-05-21",
```

```

46     "name": "Abel Marques Reis",
47     "father": "Jose Joaquim Marques Reis",
48     "mother": "Bernardina Dantas",
49     "obs": ""
50 },
51 {
52     "folder_id": "578",
53     "date": "1901-11-12",
54     "name": "Abel Martins Pereira",
55     "father": "Serafim Martins Pereira",
56     "mother": "Emilia Goncalves",
57     "obs": ""
58 },
59 {
60     "folder_id": "572",
61     "date": "1883-02-01",
62     "name": "Abel Pedro Pereira Freitas",
63     "father": "Joao Freitas Oliveira",
64     "mother": "Cecilia Rosa Pereira",
65     "obs": ""
66 },
67 {
68     "folder_id": "578",
69     "date": "1900-08-30",
70     "name": "Abel Silva Carvalho",
71     "father": "Constantino Silva Rego",
72     "mother": "Margarida Rosa Carvalho",
73     "obs": ""
74 },
75 {
76     "folder_id": "575",
77     "date": "1894-04-30",
78     "name": "Abelardo Jose Cerqueira Araujo",
79     "father": "Jose Maria Araujo",
80     "mother": "Leopoldina Cerqueira Ribeiro Araujo",
81     "obs": ""
82 },
83 {
84     "folder_id": "575",
85     "date": "1894-04-30",
86     "name": "Abelardo Jose Cerqueira Araujo",
87     "father": "Jose Maria Araujo",
88     "mother": "Leopoldina Cerqueira Ribeiro Araujo",
89     "obs": ""
90 },
91 {
92     "folder_id": "572",
93     "date": "1883-11-24",
94     "name": "Abilio Acacio Conceicao Guerreiro",
95     "father": "Jose Antonio Pereira Dantas Guerreiro",
96     "mother": "Maria Rita Pereira Monteiro",
97     "obs": ""
98 },
99 {
100     "folder_id": "579",
101     "date": "1902-10-23",
102     "name": "Abilio Aires Sousa Pereira Guimaraes",
103     "father": "Joaquim Aires Sousa Pereira Guimaraes",
104     "mother": "Josefa Rosa Gomes",

```

```

105     "obs": ""
106 },
107 {
108     "folder_id": "571",
109     "date": "1880-01-24",
110     "name": "Abilio Antonio Alves",
111     "father": "Joao Francisco Alves",
112     "mother": "Maria Jesus Rocha",
113     "obs": ""
114 },
115 {
116     "folder_id": "573",
117     "date": "1889-12-03",
118     "name": "Abilio Augusto Arantes",
119     "father": "Sebastiao Arantes",
120     "mother": "Maria Sousa",
121     "obs": ""
122 },
123 {
124     "folder_id": "581",
125     "date": "1908-08-11",
126     "name": "Abilio Augusto Galvao",
127     "father": "Antonio Augusto Galvao",
128     "mother": "Perpetua Duarte Galvao",
129     "obs": ""
130 },
131 {
132     "folder_id": "581",
133     "date": "1908-05-20",
134     "name": "Abilio Augusto Magalhaes",
135     "father": "",
136     "mother": "Maria Jesus Magalhaes",
137     "obs": ""
138 },
139 {
140     "folder_id": "581",
141     "date": "1908-05-20",
142     "name": "Abilio Augusto Magalhaes",
143     "father": "",
144     "mother": "Maria Jesus Magalhaes",
145     "obs": ""
146 },
147 {
148     "folder_id": "581",
149     "date": "1908-05-20",
150     "name": "Abilio Augusto Magalhaes",
151     "father": "",
152     "mother": "Maria Jesus Magalhaes",
153     "obs": ""
154 },
155 {
156     "folder_id": "569",
157     "date": "1869-12-02",
158     "name": "Abilio Augusto Santos",
159     "father": "Jose Joaquim Santos",
160     "mother": "Teresa Jesus",
161     "obs": "Antonio Jose Adao,Tio Materno. Proc.12530.    Albino Antonio Ribeiro,Primo
162     Paterno. Proc.12721."
163 }

```

```
163 ]
164 }
```

## 5.2 Problema 2

### 5.2.1 index.html

A B C D E

### 5.2.2 Datas extremas dos registos no dataset

Datas extremas dos registos no dataset

Extremo inferior do dataset: 2019-01-12

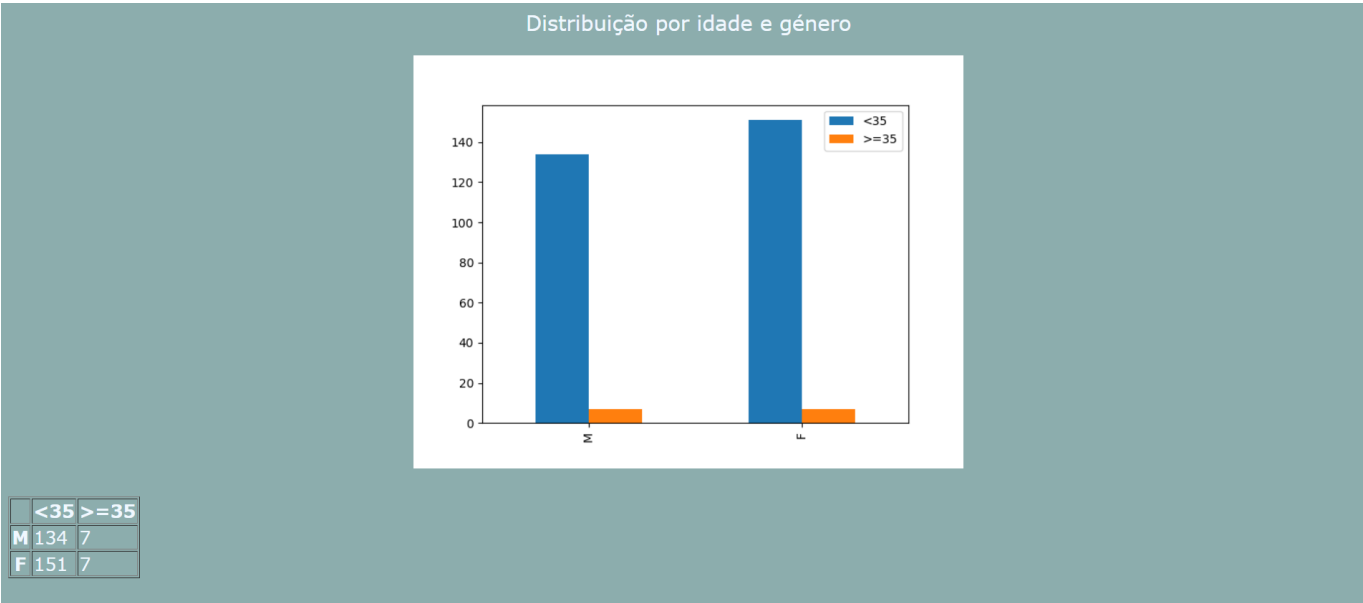
Extremo superior do dataset: 2021-03-02

5.2.3 Distribuição por modalidade em cada ano e no total

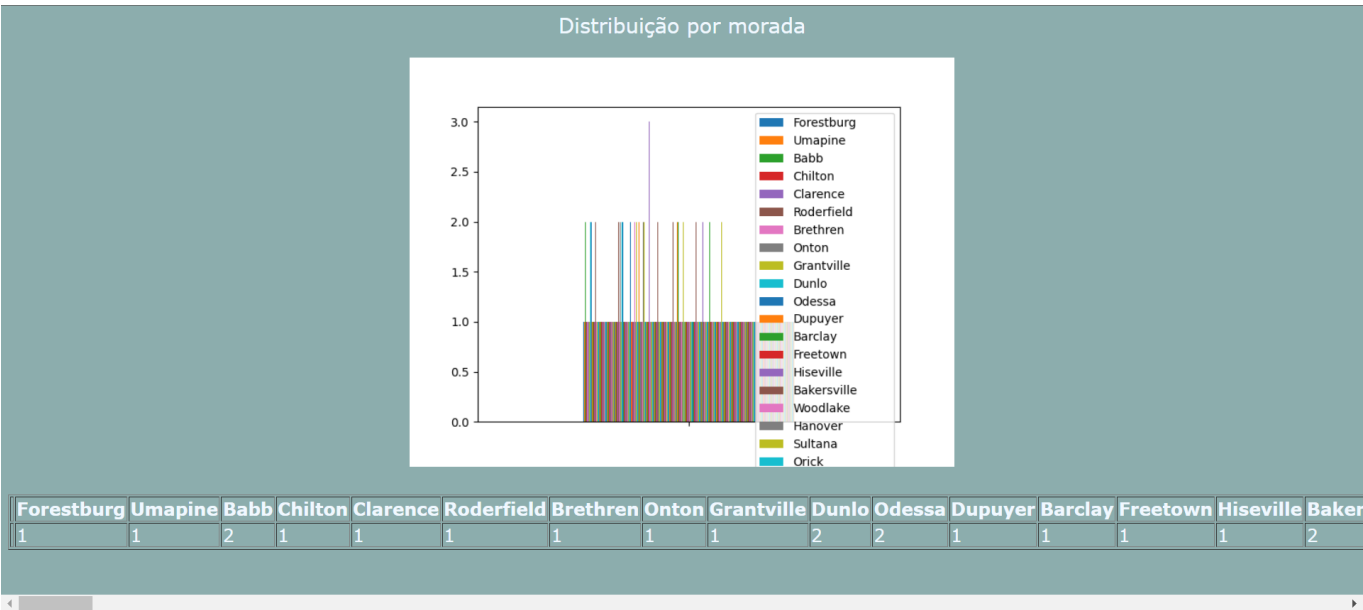




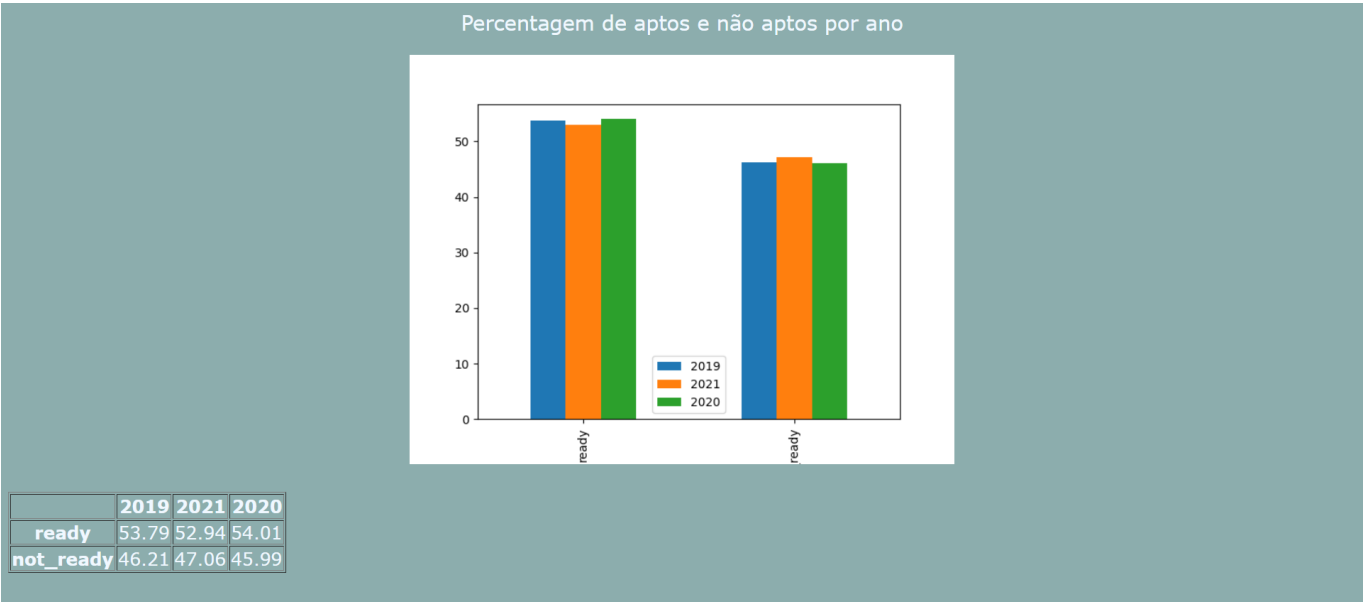
5.2.4 Distribuição por ano e gênero



5.2.5 Distribuição por morada



5.2.6 Porcentagem de aptos e não aptos por ano



## Capítulo 6

# Conclusão

Com a realização deste trabalho deparamo-nos com vários problemas de programação e raciocínio encontrados anteriormente e que sem o recurso/conhecimento das expressões regulares levavam soluções muito mais complexas, o que se traduziu portanto numa valiosa adição à nossa bagagem. Numa perspetiva geral julgamos estar a entregar um trabalho que cumpria a maioria das nossas metas para o mesmo, no entanto é de salientar que certos aspetos do exercício 2 ficaram à quem das nossas expectativas, nuances tais que serão melhoradas e corrigidas em trabalhos futuros.

## 6.1 Bibliografia

Pandas Documentation

Matplotlib Documentation