

**Universidade do Minho
Licenciatura em Engenharia
Informática**

Grupo 11

Março 2022

**Processamento de Linguagens
Trabalho Prático 1: CSV**

Duarte Moreira (a93321) Lucas Carvalho (a93176) Ricardo Gama (a93237)

Índice

1	Introdução	3
2	Desenho e Implementação da Solução	4
3	Exemplos de Utilização	5
4	Conclusões	12

1 Introdução

Na realização do trabalho prático tínhamos como principais objetivos o reforço da nossa capacidade de escrever Expressões Regulares para descrição de padrões em streams de texto, bem como, desenvolver, a partir de expressões regulares, Processadores de Linguagens Regulares que encontrem e transformem textos.

Após a escolha do enunciado, **CSV**, o foco principal era fazer um conversor de um ficheiro CSV - *Comma Separated Values* - para o formato **JSON** - *JavaScript Object Notation*. Para isso, foi necessário especificar os padrões que queríamos encontrar no texto-fonte, através de expressões regulares e desenvolver um filtro de texto para fazer o reconhecimento dos padrões identificados e proceder à transformação pretendida.

2 Desenho e Implementação da Solução

Para a realização desta tarefa o grupo decidiu seguir uma estratégia de *divide and conquer*, ou seja, começamos por dividir o problema em dois. Primeiro, o tratamento do cabeçalho e posteriormente do corpo do ficheiro recebido como input.

O processamento da primeira linha do ficheiro, cabeçalho, consistiu no desenvolvimento de uma expressão regular que capturasse os campos pretendidos de acordo com os exemplos fornecidos no enunciado. O maior desafio nesta parte foi ter em atenção que as colunas que se tratassem de listas poderiam ter variantes tais como: ser uma lista com tamanho fixo ou variável, e para além disso poderiam ou não ter um função de agregação para aplicar à mesma (**avg**, **sum**, **min**, **max**).

A expressão regular utilizada nesta fase foi a seguinte:

- `([^\{,}+)(?:\{(\d+),(\d+)\}\}|\{(\d+)\})?(?:::(\w+))?`

Ao usar esta expressão regular e a função *findall* somos capazes de obter um tuplo para cada coluna do ficheiro. Este irá conter informação do nome da coluna, se for uma lista, o número mínimo e máximo de elementos ou então o número exato e ainda se é seguida de uma função de agregação ou não - (nome,min,max,size,funcao).

Já o segundo problema, consistiu num trabalho mais simples, no entanto mais trabalhoso. Esta segunda parte resume-se a interpretar as linhas lidas do ficheiro de input e com a ajuda dos tuplos calculados anteriormente, ir criando passo a passo, ou neste caso *string* a *string* uma lista de *strings* de forma a que estas se encontrem no formato pretendido - **JSON**. Para isso, recorre-se a dois ciclos aninhados, um para percorrer as linhas e outro as colunas. No meio deste processo, para o tornar menos complexo, são usadas algumas funções auxiliares criadas por nós e mais uma vez é utilizada outra expressão regular, desta vez muito mais simples `(\d+)` para identificar os elementos de uma lista. Uma vez construída a lista de strings é realizado um *join* de forma a juntar todas as strings separando-as com um `\n` e por fim é criado um ficheiro onde é escrito o resultado obtido.

3 Exemplos de Utilização

Nesta secção do relatório, apresentamos os resultados da aplicação do nosso trabalho a diferentes datasets criados pelo grupo.

- Listas com tamanho definido

```
1 Numero, Nome, Curso, Notas{5},,,,,
2 3162, Candido Faisca, Teatro, 12, 13, 14, 15, 16
3 7777, Cristiano Ronaldo, Desporto, 17, 12, 20, 11, 12
4 264, Marcelo Sousa, Ciencia Politica, 18, 19, 19, 20, 18
5 9234, Ricardo Gama, Informatica, 13, 14, 10, 11, 17
6 3493, Duarte Moreira, Informatica, 17, 16, 15, 12, 11
7 239, Lucas Carvalho, Informatica, 18, 10, 12, 14, 11
8 9235, Daniela Silva, Biologia, 10, 12, 15, 13, 12
9 87344, Ana Filipa, Gestao, 10, 15, 16, 19, 14
```

```
1 [
2   {
3     "Numero": "3162",
4     "Nome": "Candido Faisca",
5     "Curso": "Teatro",
6     "Notas": [12, 13, 14, 15, 16]
7   },
8   {
9     "Numero": "7777",
10    "Nome": "Cristiano Ronaldo",
11    "Curso": "Desporto",
12    "Notas": [17, 12, 20, 11, 12]
13  },
14  {
15    "Numero": "264",
16    "Nome": "Marcelo Sousa",
17    "Curso": "Ciencia Politica",
18    "Notas": [18, 19, 19, 20, 18]
19  },
20  {
21    "Numero": "9234",
22    "Nome": "Ricardo Gama",
23    "Curso": "Informatica",
24    "Notas": [13, 14, 10, 11, 17]
25  },
26  {
27    "Numero": "3493",
28    "Nome": "Duarte Moreira",
29    "Curso": "Informatica",
30    "Notas": [17, 16, 15, 12, 11]
31  },
32  {
33    "Numero": "239",
34    "Nome": "Lucas Carvalho",
35    "Curso": "Informatica",
36    "Notas": [18, 10, 12, 14, 11]
37  },
38  {
39    "Numero": "9235",
40    "Nome": "Daniela Silva",
```

```

41     "Curso": "Biologia",
42     "Notas": [10, 12, 15, 13, 12]
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": [10, 15, 16, 19, 14]
49 }
50 ]

```

• Listas com um intervalo de tamanhos

```

1  Numero, Nome, Curso, Notas {3,5},,,,,
2  3162, Candido Faisca, Teatro, 12, 13, 14,,
3  7777, Cristiano Ronaldo, Desporto, 17, 12, 20, 11, 12
4  264, Marcelo Sousa, Ciencia Politica, 18, 19, 19, 20,
5  9234, Ricardo Gama, Informatica, 13, 14, 10, 11,
6  3493, Duarte Moreira, Informatica, 17, 16, 15, 12, 11
7  239, Lucas Carvalho, Informatica, 18, 10, 12,,
8  9235, Daniela Silva, Biologia, 10, 12, 15, 12,
9  87344, Ana Filipa, Gestao, 10, 15, 16,,

```

```

1  [
2    {
3      "Numero": "3162",
4      "Nome": "Candido Faisca",
5      "Curso": "Teatro",
6      "Notas": [12, 13, 14]
7    },
8    {
9      "Numero": "7777",
10     "Nome": "Cristiano Ronaldo",
11     "Curso": "Desporto",
12     "Notas": [17, 12, 20, 11, 12]
13   },
14   {
15     "Numero": "264",
16     "Nome": "Marcelo Sousa",
17     "Curso": "Ciencia Politica",
18     "Notas": [18, 19, 19, 20]
19   },
20   {
21     "Numero": "9234",
22     "Nome": "Ricardo Gama",
23     "Curso": "Informatica",
24     "Notas": [13, 14, 10, 11]
25   },
26   {
27     "Numero": "3493",
28     "Nome": "Duarte Moreira",
29     "Curso": "Informatica",
30     "Notas": [17, 16, 15, 12, 11]
31   },
32   {
33     "Numero": "239",
34     "Nome": "Lucas Carvalho",

```

```

35     "Curso": "Informatica",
36     "Notas": [18, 10, 12]
37 },
38 {
39     "Numero": "9235",
40     "Nome": "Daniela Silva",
41     "Curso": "Biologia",
42     "Notas": [10, 12, 15, 12]
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": [10, 15, 16]
49 }
50 ]

```

• Função de Agregação: AVG

```

1 Numero,Nome,Curso,Notas{3,5}::avg,,,,,
2 3162,Candido Faisca,Teatro,12,13,14,,
3 7777,Cristiano Ronaldo,Desporto,17,12,20,11,12
4 264,Marcelo Sousa,Ciencia Politica,18,19,19,20,
5 9234,Ricardo Gama,Informatica,13,14,10,11,
6 3493,Duarte Moreira,Informatica,17,16,15,12,11
7 239,Lucas Carvalho,Informatica,18,10,12,,
8 9235,Daniela Silva,Biologia,10,12,15,,
9 87344,Ana Filipa,Gestao,10,15,16,19,

```

```

1 [
2   {
3     "Numero": "3162",
4     "Nome": "Candido Faisca",
5     "Curso": "Teatro",
6     "Notas": 13.0
7   },
8   {
9     "Numero": "7777",
10    "Nome": "Cristiano Ronaldo",
11    "Curso": "Desporto",
12    "Notas": 14.4
13  },
14  {
15    "Numero": "264",
16    "Nome": "Marcelo Sousa",
17    "Curso": "Ciencia Politica",
18    "Notas": 19.0
19  },
20  {
21    "Numero": "9234",
22    "Nome": "Ricardo Gama",
23    "Curso": "Informatica",
24    "Notas": 12.0
25  },
26  {
27    "Numero": "3493",
28    "Nome": "Duarte Moreira",

```

```

29     "Curso": "Informatica",
30     "Notas": 14.2
31 },
32 {
33     "Numero": "239",
34     "Nome": "Lucas Carvalho",
35     "Curso": "Informatica",
36     "Notas": 13.33
37 },
38 {
39     "Numero": "9235",
40     "Nome": "Daniela Silva",
41     "Curso": "Biologia",
42     "Notas": 12.33
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": 15.0
49 }
50 ]

```

• Função de Agregação: SUM

```

1  Numero,Nome,Curso,Notas{3,5}::sum,,,,,
2  3162,Candido Faisca,Teatro,12,13,14,,
3  7777,Cristiano Ronaldo,Desporto,17,12,20,11,12
4  264,Marcelo Sousa,Ciencia Politica,18,19,19,20,
5  9234,Ricardo Gama,Informatica,13,14,10,11,
6  3493,Duarte Moreira,Informatica,17,16,15,12,11
7  239,Lucas Carvalho,Informatica,18,10,12,,
8  9235,Daniela Silva,Biologia,10,12,15,13,
9  87344,Ana Filipa,Gestao,10,15,16,19,14

```

```

1  [
2    {
3      "Numero": "3162",
4      "Nome": "Candido Faisca",
5      "Curso": "Teatro",
6      "Notas": 39
7    },
8    {
9      "Numero": "7777",
10     "Nome": "Cristiano Ronaldo",
11     "Curso": "Desporto",
12     "Notas": 72
13   },
14   {
15     "Numero": "264",
16     "Nome": "Marcelo Sousa",
17     "Curso": "Ciencia Politica",
18     "Notas": 76
19   },
20   {
21     "Numero": "9234",
22     "Nome": "Ricardo Gama",

```



```

23     "Curso": "Informatica",
24     "Notas": 48
25 },
26 {
27     "Numero": "3493",
28     "Nome": "Duarte Moreira",
29     "Curso": "Informatica",
30     "Notas": 71
31 },
32 {
33     "Numero": "239",
34     "Nome": "Lucas Carvalho",
35     "Curso": "Informatica",
36     "Notas": 40
37 },
38 {
39     "Numero": "9235",
40     "Nome": "Daniela Silva",
41     "Curso": "Biologia",
42     "Notas": 50
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": 74
49 }
50 ]

```

• Função de Agregação: MIN

```

1  Numero, Nome, Curso, Notas{2,5}::min,,,,,
2  3162, Candido Faisca, Teatro, 12, 13, 14,,
3  7777, Cristiano Ronaldo, Desporto, 17, 12, 20, 11, 12
4  264, Marcelo Sousa, Ciencia Politica, 18, 19, 19, 20,
5  9234, Ricardo Gama, Informatica, 13, 14, 10, 11,
6  3493, Duarte Moreira, Informatica, 17, 16, 15, 12, 11
7  239, Lucas Carvalho, Informatica, 18, 10,,,
8  9235, Daniela Silva, Biologia, 10, 12, 15,,
9  87344, Ana Filipa, Gestao, 10, 15, 16, 19, 14

```

```

1  [
2    {
3      "Numero": "3162",
4      "Nome": "Candido Faisca",
5      "Curso": "Teatro",
6      "Notas": 12
7    },
8    {
9      "Numero": "7777",
10     "Nome": "Cristiano Ronaldo",
11     "Curso": "Desporto",
12     "Notas": 11
13   },
14   {
15     "Numero": "264",
16     "Nome": "Marcelo Sousa",

```

```

17     "Curso": "Ciencia Politica",
18     "Notas": 18
19 },
20 {
21     "Numero": "9234",
22     "Nome": "Ricardo Gama",
23     "Curso": "Informatica",
24     "Notas": 10
25 },
26 {
27     "Numero": "3493",
28     "Nome": "Duarte Moreira",
29     "Curso": "Informatica",
30     "Notas": 11
31 },
32 {
33     "Numero": "239",
34     "Nome": "Lucas Carvalho",
35     "Curso": "Informatica",
36     "Notas": 10
37 },
38 {
39     "Numero": "9235",
40     "Nome": "Daniela Silva",
41     "Curso": "Biologia",
42     "Notas": 10
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": 10
49 }
50 ]

```

• Função de Agregação: MAX

```

1  Numero, Nome, Curso, Notas {3,6}::max,,,,,
2  3162, Candido Faisca, Teatro, 12, 13, 14,,,
3  7777, Cristiano Ronaldo, Desporto, 17, 12, 20, 11, 12, 13
4  264, Marcelo Sousa, Ciencia Politica, 18, 19, 19, 20,,
5  9234, Ricardo Gama, Informatica, 13, 14, 10, 11,,
6  3493, Duarte Moreira, Informatica, 17, 16, 15, 12, 11, 10
7  239, Lucas Carvalho, Informatica, 18, 10, 12,,,
8  9235, Daniela Silva, Biologia, 10, 12, 15, 13, 12
9  87344, Ana Filipa, Gestao, 10, 15, 16, 19, 14,

```

```

1  [
2    {
3      "Numero": "3162",
4      "Nome": "Candido Faisca",
5      "Curso": "Teatro",
6      "Notas": 14
7    },
8    {
9      "Numero": "7777",
10     "Nome": "Cristiano Ronaldo",

```

```

11     "Curso": "Desporto",
12     "Notas": 20
13 },
14 {
15     "Numero": "264",
16     "Nome": "Marcelo Sousa",
17     "Curso": "Ciencia Politica",
18     "Notas": 20
19 },
20 {
21     "Numero": "9234",
22     "Nome": "Ricardo Gama",
23     "Curso": "Informatica",
24     "Notas": 14
25 },
26 {
27     "Numero": "3493",
28     "Nome": "Duarte Moreira",
29     "Curso": "Informatica",
30     "Notas": 17
31 },
32 {
33     "Numero": "239",
34     "Nome": "Lucas Carvalho",
35     "Curso": "Informatica",
36     "Notas": 18
37 },
38 {
39     "Numero": "9235",
40     "Nome": "Daniela Silva",
41     "Curso": "Biologia",
42     "Notas": 15
43 },
44 {
45     "Numero": "87344",
46     "Nome": "Ana Filipa",
47     "Curso": "Gestao",
48     "Notas": 19
49 }
50 ]

```

4 Conclusões

Após a realização do trabalho prático número 1, consideramos que foi bastante benéfico para o aprofundamento da matéria lecionada na unidade curricular de **Processamento de Linguagens**, com destaque para as **Expressões Regulares** e do módulo **re** do **Python**, que foram essenciais para gerar os filtros de texto necessários para a realização do trabalho.

Acreditamos que conseguimos atingir os objetivos principais do trabalho. Conseguimos gerar um conversor de um ficheiro **CSV** para o formato **JSON**, tendo em conta todas as condicionantes adicionais que este trabalho nos propunha: listas com tamanho definido, com um intervalo de tamanhos e funções de agregação. Para isto, foi importante a captura inicial dos valores do cabeçalho num *array*, onde cada posição continha um tuplo com informações relativas a cada coluna, sendo assim possível verificar se esta continha uma função de agregação, se tinha um valor máximo e mínimo para o intervalo ou se se tratava apenas de um tamanho definido.