

# Relatório - MPEI PL4

Universidade de Aveiro

Tomás Cerca Rodrigues, Alexandre Costa  
Martins



# Relatório - MPEI PL4

LECI

Universidade de Aveiro

Tomás Cerca Rodrigues, Alexandre Costa Martins  
(104090) [tcercarodrigues@ua.pt](mailto:tcercarodrigues@ua.pt), (103552) [alexandremartins@ua.pt](mailto:alexandremartins@ua.pt)

21/12/2023

# Índice

<b>1</b>	<b>Introdução</b>	<b>1</b>
<b>2</b>	<b>Interface e Main</b>	<b>2</b>
2.1	main.m . . . . .	2
2.2	interfaceSelecao.m . . . . .	4
<b>3</b>	<b>Opção 1</b>	<b>5</b>
3.1	getGenres.m . . . . .	5
<b>4</b>	<b>Opção 2</b>	<b>6</b>
4.1	countMoviesByGenre.m . . . . .	6
<b>5</b>	<b>Opção 3</b>	<b>8</b>
5.1	countMoviesByGenreAndYear.m . . . . .	8
<b>6</b>	<b>Opção 4</b>	<b>10</b>
6.1	getMoviesByTitle.m . . . . .	10
<b>7</b>	<b>Opção 5</b>	<b>12</b>
7.1	getMoviesByGenre.m . . . . .	12

# Capítulo 1

## Introdução

Este trabalho consistiu no desenvolvimento duma aplicação em MATLAB para simular um sistema de informação para pesquisa de filmes.

Foi entregue aos alunos o ficheiro `movies.csv` com uma lista de 58000 filmes existentes e requisitado que a aplicação futuramente desenvolvida aceda *às cells* e apresente 5 opções.

1. Apresentar todos os géneros disponíveis.
2. Número de filmes de um determinado genero, sendo este recebido por input do usuário.
3. Número de filmes de um determinado genero e ano, sendo estes recebidos também por input do usuário.
4. Pesquisar filmes pelo título.
5. Pesquisar filmes baseado no genero.

Ao longo deste relatório serão fundamentadas as opções tomadas pelo grupo na implementação dos diversos métodos probabilísticos recorridos para a resolução das atividades.

## Capítulo 2

# Interface e Main

### 2.1 main.m

```
1 %% main
2 movies = readcell('movies.csv', 'Delimiter',';');
3 option = interfaceSelecao();
4
5 while (option ~= 6)
6
7     switch option
8     case 1
9         fprintf('Option 1 selected...\n');
10        % 'Display available genres'
11        genres = getGenres(movies);
12        disp('All movie genres:');
13        disp(genres);
14
15    case 2
16        fprintf('Option 2 selected...\n');
17        % 'Number of movies of a genre'
18        fprintf('Executing option 2 ...\n');
19        genre = input('Enter the genre: ', 's');
20        count = countMoviesByGenre(movies, genre);
21
22    case 3
23        fprintf('Option 3 selected...\n');
24        % 'Number of movies of a genre on a given
25         year's
26        genre = input('Enter the genre: ', 's');
27        year = input('Enter the year: ');
```

```

27         count = countMoviesByGenreAndYear(movies,
28             genre, year);
29
30     case 4
31         fprintf('Option 4 selected...\n');
32         % 'Search of movie titles'
33         title = input('Enter the title: ', 's');
34         similarMoviesByTitle = getMoviesByTitle(
35             movies, title);
36         disp('All similar movie titles:');
37         disp(similarMoviesByTitle);
38
39     case 5
40         fprintf('Option 5 selected...\n');
41         % 'Search movies based on genres'
42         genres = input('Enter the genres (
43             separated by commas): ', 's');
44         genres = strsplit(genres, ',');
45         similarMoviesByGenre = getMoviesByGenre(
46             movies, genres);
47         disp('All movie with the desired genre: ')
48         ;
49         disp(similarMoviesByGenre);
50
51     case 6
52         fprintf('Closing the program...\n');
53         % 'Exit'
54         return;
55
56     otherwise
57         fprintf('Invalid option selected.\n');
58
59 end
60
61 option = interfaceSelecao();
62 end

```

O main é o script que corre o programa, é neste mesmo script que são chamadas todas as funções desenvolvidas não só referentes às funcionalidades da aplicação como referentes até à interface.

## 2.2 interfaceSelecao.m

```
1 function numeroSelecioneado = interfaceSelecao()
2 opcoes = {'Display available genres',
3           'Number of movies of a genre',
4           'Number of movies of a genre on a given year',
5           'Search of movie titles',
6           'Search movies based on genres',
7           'Exit'}; % Option list
8
9
10 while true
11     % Show enumerated options
12     for i = 1:numel(opcoes)
13         fprintf('%d. %s\n', i, opcoes{i});
14     end
15
16     % Gets input
17     numeroSelecioneado = input('Input your desired
18                               option: ');
19
20     if ((numeroSelecioneado >= 1) && (
21         numeroSelecioneado <= 6))
22         fprintf('Valid option\n ');
23         return;
24     end
25 end
```

O seguinte script é meramente responsável pela interação do usuário com o menu da aplicação, expondo as opções disponíveis, adquirindo e por fim validando os valores de entrada.

## Capítulo 3

# Opção 1

### 3.1 getGenres.m

```
1 function genreArray = getGenres(movies)
2
3     maxGenres = size(movies, 1) * (12 - 3 + 1);
4     genres = cell(1, maxGenres);
5     genreCount = 0;
6
7     for i = 1:size(movies, 1)
8         for j = 3:12
9             if numel(movies{i,j}) > 1 && ~strcmp(
10                 movies{i,j}, '(no genres listed)')
11                 genreCount = genreCount + 1;
12                 genres{genreCount} = movies{i, j};
13                 movies{i,j};
14             end
15         end
16     end
17
18     genres = genres(1:genreCount);
19     genreArray = unique(genres);
20 end
```

Esta função é chamada pela main quando o usuário seleciona a opção 1 do menu e tem como objetivo listar todos os generos de filmes presentes no ficheiro e na cell *movie.csv*.

O único parâmetro de entrada são as *cells* de *movie.csv*.



## Capítulo 4

## Opção 2

### 4.1 countMoviesByGenre.m

```
1 function estimatedCount = countMoviesByGenre(movies,
2     userGenre)
3     movies = movies(:, :);
4     userGenre = lower(userGenre);
5
6     maxGenres = size(movies, 1) * (12 - 3 + 1);
7     genres = cell(1, maxGenres);
8     genreCount = 0;
9
10    for i = 1:size(movies, 1)
11        for j = 3:12
12            if numel(movies{i,j}) > 1 && ~strcmp(
13                movies{i,j}, '(no genres listed)')
14                genreCount = genreCount + 1;
15                genres{genreCount} = lower(movies{i, j
16                    });
17            end
18        end
19    end
20
21    n = 1000;
22    bloomFilter = inicializar(n);
23
24    k = 3;
25    for i = 1:size(genres, 1)
26        for j = 1:size(genres, 2)
```

```

26         genre = genres{i, j};
27         if ~isempty(genre)
28             bloomFilter = adicionar_elemento(
                bloomFilter, genre, k);
29         end
30     end
31 end
32
33 if pertenca(bloomFilter, userGenre, k)
34     disp(The genre is present in the Bloom Filter
        .');
35
36     estimatedCount = 0;
37     for i = 1:size(genres, 1)
38         for j = 1:size(genres, 2)
39             genre = genres{i, j};
40             if strcmp(userGenre, genre)
41                 estimatedCount = estimatedCount +
                    1;
42             end
43         end
44     end
45     disp(['Estimated number of movies ',
        userGenre, ': ', num2str(estimatedCount)])
        ;
46 else
47     disp('The genre is not present in the Bloom
        Filter');
48 end
49
50
51 end

```

Como já sugere o nome, este script é o principal responsável pela contagem/estimativa de filmes de um determinado genero. No começo do script são inicializados e preparados todos os dados e variáveis, ressaltamos "*movies = movies(:, :)*" que garante que serão abordadas todas as colunas da cell e "*userGenre = lower(userGenre)*" que torna menos rígida a compreensão de inputs por parte da aplicação parando de diferenciar letras maiúsculas de minúsculas.

Na criação do filtro de Bloom inicializamos **n**, o número de *bits* do filtro a 1000 e **k**, o número de funções de dispersão a 3.

Após o seu preenchimento verificamos se o genero fornecido pelo usuário se encontra presente no filtro, caso esteja presente o filtro irá calcular e exibir uma estimativa do número de filmes que contém esse genero. Caso este esteja ausente será exibida uma mensagem ao utilizando relatando esse mesmo caso.

## Capítulo 5

## Opção 3

### 5.1 countMoviesByGenreAndYear.m

```
1  function estimatedCount = countMoviesByGenreAndYear(  
2      movies, genre, year)  
3      genre = strtrim(lower(genre));  
4      year = strtrim(num2str(year));  
5  
6      n = 1000;  
7      bloomFilter = inicializar(n);  
8  
9      maxGenres = size(movies, 1) * (12 - 3 + 1);  
10     genres = cell(1, maxGenres);  
11     genreCount = 0;  
12  
13     for i = 1:size(movies, 1)  
14         for j = 3:12  
15             if numel(movies{i,j}) > 1 && ~strcmp(  
16                 movies{i,j}, '(no genres listed)')  
17                 genreCount = genreCount + 1;  
18                 genres{genreCount} = lower(movies{i, j  
19                     });  
20  
21                 k = 3;  
22                 bloomFilter = adicionar_elemento(  
23                     bloomFilter, genres{genreCount}, k)  
24                     ;  
25             end  
26         end  
27     end
```

```

24 % check if the genre is in the bloom filter
25 if pertenca(bloomFilter, genre, k)
26     disp(['The genre ', genre, ' is present in the
27         Bloom Filter.']);
28
29     estimatedCount = 0;
30     for i = 1:size(movies, 1)
31         movieYear = strtrim(num2str(movies{i, 2}))
32         ;
33         movieGenres = {};
34         for j = 3:12
35             if ~ismissing(movies{i, j})
36                 movieGenres = [movieGenres,
37                     strsplit(lower(movies{i, j}), '
38                         |')]';
39             end
40         end
41
42         if any(cellfun(@(x) any(strcmp(genre, x)),
43             movieGenres)) && strcmp(year,
44             movieYear)
45             estimatedCount = estimatedCount + 1;
46         end
47     end
48     disp(['Estimated number of movies of genre ',
49         genre, ' and year ', year, ': ', num2str(
50             estimatedCount)]);
51 else
52     disp(['The genre ', genre, ' is not present in
53         the Bloom Filter.']);
54     estimatedCount = 0;
55 end
56 end

```

Este cript por mais que seja referente a outro exercício, os métodos de resolução solicitados pelo guião são os mesmos o que faz com que a sua estrutura e desenvolvimento sejam extremamente semelhantes. Contudo tem leves extras que valem ser apontados.

Pra receber os valores dos inputs o ano é passado a String e o genero é novamente convertido para minúsculas pelo motivo já mencionado.

O filtro de Bloom é iniciado e é verificado a presença do genero neste, caso estes estejam presentes é iniciada e depois exibida a contagem referente aos filmes desse mesmo genero e ano. Caso o genero não esteja presente o comportamento da aplicação é o mesmo referido no capítulo anterior, ou seja, será exibida uma mensagem informando que o genero não se encontra presente no filtro.

## Capítulo 6

## Opção 4

### 6.1 getMoviesByTitle.m

```
1 function titleArray = getMoviesByTitle(movies, title)
2
3     title = lower(title);
4     maxMovies = size(movies, 1);
5     titleArray = cell(maxMovies, 2);
6     movieCount = 0;
7
8     shingleSize = 2;
9     numHashFuncs = 100;
10
11     titleShingles = createShingles(title, shingleSize)
12     ;
13     titleSignature = minHash(titleShingles,
14                             numHashFuncs);
15
16     for i = 1:size(movies, 1)
17         movieTitle = lower(movies{i, 1});
18         movieShingles = createShingles(movieTitle,
19                                         shingleSize);
20         movieSignature = minHash(movieShingles,
21                                 numHashFuncs);
22
23         similarity = sum(titleSignature ==
24                         movieSignature) / numHashFuncs;
25
26         if similarity > 0.6
27             movieCount = movieCount + 1;
28             titleArray{movieCount, 1} = movies{i, 1};
29         end
30     end
```

```

24         titleArray{movieCount, 2} = similarity;
25     end
26 end
27 titleArray = titleArray(1:movieCount, :);
28 titleTable = cell2table(titleArray, 'VariableNames', {'Title', 'Similarity'});
29 titleTable = sortrows(titleTable, 'Similarity', 'descend');
30 titleArray = table2cell(titleTable);
31 end
32
33 function shingles = createShingles(str, k)
34     shingles = arrayfun(@(i) str(i:i+k-1), 1:length(str)-k+1, 'UniformOutput', false);
35 end
36
37 function signature = minHash(shingles, numHashFuncs)
38     signature = inf(1, numHashFuncs);
39     for i = 1:length(shingles)
40         hashValues = mod(sum(double(shingles{i})).*(1:numHashFuncs), numHashFuncs);
41         signature = min(signature, hashValues);
42     end
43 end

```

Esta função recebe apenas a lista de filmes e um título para procurar, após a conversão deste para minúsculas.

Inicializa depois variáveis para armazenar títulos semelhantes encontrados e depois define o tamanho de shingle para 2 e o número de funções hash para 100. Usamos shingle a 2 para poder ser razoavelmente mais exigente na procura e obter filmes com o títulos similares. Caso shingle fosse 3 ser-nos-ia devolvido títulos que quase não teriam similaridade com o input do utilizador.

Usamos também minHash com os shingles para poder criar uma assinatura minHash.

O processo de comparação, novamente, converte o título para minúsculas, gera shingles para este, calcula a assinatura minHash e a similaridade. Se a similaridade for maior que 0.6 este título é adicionado à lista de filmes encontrados.

## Capítulo 7

## Opção 5

### 7.1 getMoviesByGenre.m

```
1 function genreArray = getMoviesByGenre(movies,
2   inputGenres)
3     if ischar(inputGenres)
4       inputGenres = cellstr(inputGenres);
5     end
6     if isstring(inputGenres)
7       inputGenres = {inputGenres};
8     end
9
10    maxMovies = size(movies, 1);
11    genreArray = cell(maxMovies, 3);
12    movieCount = 0;
13
14    shingleSize = 2;
15    numHashFuncs = 100;
16
17    h = waitbar(0, 'Please wait...');
18
19    for i = 1:size(movies, 1)
20      waitbar(i/maxMovies, h);
21      genreMatchCounter = 0;
22      for j = 3:12
23        if numel(movies{i,j}) > 1 && ~strcmp(
24          movies{i,j}, '(no genres listed)')
25          movieGenre = lower(movies{i, j});
26          movieShingles = createShingles(
27            movieGenre, shingleSize);
```

```

26         movieSignature = minHash(movieShingles
27                                   , numHashFuncs);
28
29         for g = 1:length(inputGenres)
30             genre = lower(strtrim(inputGenres{
31                                     g}));
32             genreShingles = createShingles(
33                                     genre, shingleSize);
34             genreSignature = minHash(
35                                     genreShingles, numHashFuncs);
36
37             similarity = sum(genreSignature ==
38                             movieSignature) / numHashFuncs
39                             ;
40
41             if similarity > 0.6
42                 genreMatchCounter =
43                     genreMatchCounter + 1;
44             end
45         end
46     end
47     finalSimilarity = genreMatchCounter / length(
48         inputGenres);
49     if finalSimilarity > 0.6
50         movieCount = movieCount + 1;
51         genreArray{movieCount, 1} = movies{i, 1};
52         genreArray{movieCount, 2} = movies{i, 2};
53         genreArray{movieCount, 3} =
54             finalSimilarity;
55     end
56 end
57
58 close(h); % close waitbar
59
60 genreArray = genreArray(1:movieCount, :);
61 genreTable = cell2table(genreArray, 'VariableNames
62     ', {'Title', 'Year', 'Similarity'});
63 genreTable = sortrows(genreTable, {'Similarity', '
64     Year'}, {'descend', 'descend'});
65 if height(genreTable) > 50
66     genreTable = genreTable(1:50, :);
67 end
68 genreArray = table2cell(genreTable);
69 end

```



```

61
62
63 function shingles = createShingles(str, k)
64     shingles = arrayfun(@(i) str(i:i+k-1), 1:length(
        str)-k+1, 'UniformOutput', false);
65 end
66
67
68 function signature = minHash(shingles, numHashFuncs)
69     signature = inf(1, numHashFuncs);
70     for i = 1:length(shingles)
71         hashValues = mod(sum(double(shingles{i})) .*
            (1:numHashFuncs), numHashFuncs);
72         signature = min(signature, hashValues);
73     end
74 end

```

Este é o script da opção final e tem algumas características já implementadas nos exercícios anteriores (principalmente a opção 4). Tal como os demais, este recebe a lista de filmes e agora um ou mais generos separados por uma vírgula.

O processo de comparação, como esperado, é a principal diferença ainda que esta não seja muito grande. Neste são enunciados os generos listados no filme, de seguida é feita a comparação de assinaturas minHash dos generos introduzidos pelo utilizador com os generos dos filmes.

Novamente a similaridade desejada foi 0.6, ou seja todos os filmes com uma similaridade superior eram adicionados à lista.

A tabela final é ordenada por ordem decrescente de similaridade e o número de resultados é limitado a 50 caso a função queira devolver mais que isso.