

PL 4 - Algoritmos Probabilísticos

Trabalho de Métodos Probabilísticos para Engenharia Informática Professores: Carlos Bastos, Amaro Sousa

José Diogo Cerqueira nº76758, João Monteiro nº114547, P5

Abstract

O objetivo deste trabalho consistiu no desenvolvimento de uma aplicação para pesquisa de filmes, que se baseasse na matéria lecionada no fim do semestre. Em especial, o desenvolvimento da aplicação focava-se na implementação de um método **minHash**, para cálculo da **similaridade de Jaccard** e, ainda, o desenvolvimento de um **filtro de Bloom** com contagem.

Keywords: *Filtros de Bloom, minHash, Jaccard, shingles.*

1. Desenvolvimento da aplicação

1.1. Opção 1

Inicialmente, foi necessário aceder aos dados que nos eram fornecidos, através do comando:

```
1 movies = readcell('films.txt', 'Delimiter', ',', '');
```

Listando todos os géneros que existiam neste novo cell array criado (**movies**), o objetivo era que o utilizador escolhesse um dos géneros de filmes existentes, para basear o restante das operações. Para tal armazenámos na variável **genres**, os vários géneros de filmes que apareciam no ficheiro. Com a ajuda da função Matlab *input*, atribuímos à variável **genreInput** aquilo que o utilizador escrevesse na linha de comandos e, por fim, verificámos se essa variável correspondia a um dos géneros que tínhamos disponível.

```
1 numMovies = height(movies);
2
3 genres = unique(movies(:,3));
4
5 fprintf(['Genres:\t%-16s\t%-8s\t%-10s\t%-10s\n\t\t%-16s\t%-8s\t%-10s\t%-10s' ...
6         '\n\t\t%-16s\t%-8s\t%-10s\t%-10s\n\t\t%-16s\t%-8s\t%-10s\t%-10s\n\t' ...
7         '\t%-16s\t%-8s\t%-10s\t%-10s\n'], genres{1:20});
8
9 genreInput = input("Select a genre: ", "s");
10 while true
11     upperLetter = upper(genreInput(1));
12     genreUpper = [upperLetter, genreInput(2:end)];
13     if ismember(genreUpper, genres)
14         break;
15     end
16     fprintf('ERROR inputing the genre!\n');
17     genreInput = input("Select a genre: ", "s");
18 end
```

De seguida, podemos ver uma implementação do menu que era pedido, para ser apresentado ao utilizador:

```
1 while true
2     disp('-----');
3     fprintf('SELECTED GENRE: %s\n', genreUpper);
4     disp('-----');
5     fprintf('1 - Change selected Genre\n');
6     fprintf('2 - No. of movies of selected Genre on given years\n');
```

```

7 fprintf('3 - Search movie titles of selected Genre\n');
8 fprintf('4 - Search movies based on Genres\n');
9 fprintf('5 - Exit\n');
10 disp('-----');
11 option = input('Select an option: ');
12 switch(option)

```

Basicamente, na opção 1 o processo era exatamente igual ao que era pedido inicialmente ao utilizador. Porém, neste caso específico, o objetivo era alterar o género previamente escolhido.

```

1 case 1
2     disp('-----');
3     fprintf(['Genres:\t%-16s\t%-8s\t%-10s\t%-10s\n\t\t%-16s\t%-8s\t%-10s\t%-10s'
4         ...
5         '\n\t\t%-16s\t%-8s\t%-10s\t%-10s\n\t\t%-16s\t%-8s\t%-10s\t%-10s\n\t'
6         ...
7         '\t%-16s\t%-8s\t%-10s\t%-10s\n'], genres{1:20});
8     disp('-----');
9     genreInput = input("Select a genre: ", "s");
10    while true
11        upperLetter = upper(genreInput(1));
12        genreUpper = [upperLetter, genreInput(2:end)];
13        if ismember(genreUpper, genres)
14            break;
15        end
16        fprintf('ERROR inputing the genre!\n');
17        genreInput = input("Select a genre: ", "s");
18    end
19    continue;

```

Segue uma pequena demonstração desta parte do script:

```

1 Genres: (no genres listed)  Action      Adventure  Animation
2   Children                 Comedy    Crime      Documentary
3   Drama                   Fantasy   Film-Noir  Horror
4   IMAX                    Musical   Mystery    Romance
5   Sci-Fi                  Thriller  War        Western
6 Select a genre: action
7 -----
8 SELECTED GENRE: Action
9 -----
10 1 - Change selected Genre
11 2 - No. of movies of selected Genre on given years
12 3 - Search movie titles of selected Genre
13 4 - Search movies based on Genres
14 5 - Exit
15 -----
16 Select an option: 1
17 -----
18 Genres: (no genres listed)  Action      Adventure  Animation
19   Children                 Comedy    Crime      Documentary
20   Drama                   Fantasy   Film-Noir  Horror
21   IMAX                    Musical   Mystery    Romance
22   Sci-Fi                  Thriller  War        Western
23 -----
24 Select a genre: western
25 -----
26 SELECTED GENRE: Western
27 -----
28 1 - Change selected Genre
29 2 - No. of movies of selected Genre on given years
30 3 - Search movie titles of selected Genre
31 4 - Search movies based on Genres
32 5 - Exit
33 -----
34 Select an option: 5
35 Exiting the program

```

1.2. Opção 2

O **Bloom Filter com Contagem** foi implementado com base no algoritmo do **Bloom Filter** normal que desenvolvemos previamente, tendo apenas de alterar a função de inserção e de verificação, para que em vez de se inserir/verificar um bit, ser apenas necessário inserir/verificar um número inteiro. Então,

em vez de inicializar o **Bloom Filter** com falsos, inicializámos com zeros, e em vez de incrementar o bit, incrementámos o número inteiro. Da mesma forma, em vez de verificar se o bit é 1, verificámos se o número inteiro é maior que 0.

Para inserir um elemento no **Bloom Filter**, temos de aplicar as k funções de dispersão ao elemento, e incrementar o i -ésimo elemento do **Bloom Filter** para cada função de dispersão, onde i é o resultado da função de dispersão. Para verificar se um elemento pertence ao **Bloom Filter**, fazemos o mesmo processo, mas em vez de incrementar, verificamos qual o valor do i -ésimo elemento.

```

1 function CBF = runCountingBloomFilter(movies)
2     tic;
3     fprintf("\n=====")
4
5     numMovies = height(movies);           % Number of movies
6     genresPerMovie = 3;                   % Number of genres per movie
7     p = 0.01;                             % Pretended probability of false
8     positives
9     numElem = numMovies * genresPerMovie; % Number of elements to insert in
10    the counting bloom filter
11    n = int32(-log(p) * numElem / (log(2)^2)); % Size of ideal bloom filter
12    k = int32(n * log(2) / numElem);        % Number of hash functions
13
14    % 1. Criar um Bloom Filter com n posicoes e k funcoes de dispersao
15    CBF = CountingBloomFilter(n, k);
16
17    fprintf('\n1. Bloom Filter criado com sucesso, com %d funcoes de dispersao\n', k);
18
19    % 2. Inserir todos pares (ano, genero) no Bloom Filter
20    wb = waitbar(0, 'Inserindo elementos no Bloom Filter...');
21    for i = 1:numMovies
22        % Update waitbar every 2500 movies for performance reasons
23        if mod(i, 2500) == 0
24            waitbar(i/numMovies, wb);
25        end
26        CBF = CountingBloomFilterInsert(CBF, [num2str(movies{i,2}) movies{i,3}]);
27    end
28    close(wb);
29    fprintf('\n2. Todos os pares (ano, genero) inseridos no Bloom Filter\n');
30
31    % 3. Verificar se todos os pares (ano, genero) pertencem ao Bloom Filter
32    wb = waitbar(0, 'Verificando elementos no Bloom Filter...');
33    for i = 1:numMovies
34        % Update waitbar every 2500 movies for performance reasons
35        if mod(i, 2500) == 0
36            waitbar(i/numMovies, wb);
37        end
38        check = CountingBloomFilterCheck(CBF, [num2str(movies{i,2}) movies{i,3}]);
39        if ~check
40            fprintf('\n3. Erro: 0 par (ano, genero) (%d, %s) nao pertence ao Bloom Filter\n', movies{i,2}, movies{i,3});
41            break;
42        end
43    end
44    close(wb);
45    fprintf('\n3. Todos os pares (ano, genero) verificados no Bloom Filter\n');
46    disp(toc);
47 end
48
49 function CBF = CountingBloomFilter(n, k)
50     CBF.n = n; % Tamanho do Bloom Filter
51     CBF.k = k; % Numero de funcoes de dispersao
52     CBF.cbf = zeros(1, n); % Inicializar o Bloom Filter
53     CBF = CountingBloomFilterInitHF(CBF); % Inicializar as funcoes de dispersao
54 end
55
56 function CBF = CountingBloomFilterInitHF(CBF)
57     hashFunctions = {@(element) string2hash(element, 'sdbm'), ...
58                     @(element) string2hash(element, 'djb2'), ...
59                     @(element) DJB31MA(element, 5381)};
60
61     CBF.hashFunctions = cell(1, CBF.k);
62     for i = 1:CBF.k
63         CBF.hashFunctions{mod(i,3)+1} = @(element) mod(hashFunctions{mod(i,3)+1}(element), CBF.n) + 1;
64     end

```

```

62     end
63 end
64
65 function CBF = CountingBloomFilterInsert(CBF, x)
66     xCell = cell(1, CBF.k); % Inicializar o array de elementos a inserir
67     for i = 1:CBF.k
68         xCell{i} = [x num2str(i)];
69         index = CBF.hashFunctions{mod(i,3)+1}(xCell{i});
70         CBF.cbf(index) = CBF.cbf(index) + 1;
71     end
72 end
73
74 function num = CountingBloomFilterCheck(CBF, x)
75     xCell = cell(1, CBF.k); % Inicializar o array de elementos a inserir
76     for i = 1:CBF.k
77         xCell{i} = [x num2str(i)]; % Tem de ser o mesmo que foi usado para inserir
78         % Verificar o i-esimo elemento do Bloom Filter
79         index = CBF.hashFunctions{mod(i,3)+1}(xCell{i});
80         if CBF.cbf(index) == 0
81             num = 0;
82             return;
83         end
84         num = CBF.cbf(index);
85     end
86 end

```

- **CountingBloomFilter**: - Cria um Counting Bloom Filter com n posições e k funções de dispersão; - O valor de n é calculado com base na fórmula para o tamanho ideal do Bloom Filter: $-n * \ln(p) / (\ln(2)^2)$;

- O valor de k é calculado com base na fórmula: $k = n * \ln(2) / m$, onde n é o número de elementos a inserir e m é o tamanho do Bloom Filter.

- **CountingBloomFilterInitHF**: - Inicializa as funções de dispersão; - As funções de dispersão são inicializadas com base em 3 funções de dispersão; - Estas 3 funções de dispersão são aplicadas a cada elemento a inserir sempre com um salt diferente, para evitar colisões e sequencialmente na mesma ordem.

- **CountingBloomFilterInsert**: - Insere um elemento no Bloom Filter; - Para inserir um elemento no Bloom Filter, temos de aplicar as k funções de dispersão ao elemento, e incrementar o i -ésimo elemento do Bloom Filter para cada função de dispersão, onde i é o resultado da função de dispersão; - É utilizado um salto diferente para cada função de dispersão, para evitar colisões e sempre com a mesma sequência de funções de dispersão.

- **CountingBloomFilterCheck**: - Verifica qual o número de vezes que um elemento x foi inserido no Bloom Filter; - Para verificar se um elemento pertence ao Bloom Filter, fazemos o mesmoprocessamento que na inserção, mas em vez de incrementar, verificamos qual o valor do i -ésimo elemento.

1.3. Opção 3

Para a terceira opção, era pedido que se fizesse uma procura de filmes de um dado género, baseado no título desse filme. Para tal, o objetivo era calcular a similaridade de **Jaccard** entre strings (a procura e o título dos diversos filmes). Para tal, era necessária a implementação de um método **minHash**.

Antes de tudo, procedemos ao processamento da string que o utilizador viria a introduzir. Também criámos um cell array **moviesbyGenre**, onde apenas tínhamos os filmes do género que tinha sido escolhido previamente.

```

1     case 3
2         movieSearch = input('Insert a string: ', 's');
3         k = 1;
4         for i = 1:numMovies
5             for j = 3:12
6                 if strcmp(movies{i, j}, genreUpper)
7                     for idx = 1:12
8                         moviesbyGenre{k, idx} = movies{i, idx};
9                     end
10                    k = k + 1;
11                end
12            end
13        end

```

Desenvolvemos quatro funções: - A primeira para a criação dos *shingles*; - A segunda, para comparar as assinaturas, ou seja, para calcular a similaridade de Jaccard entre dois conjuntos; - A terceira função é uma função **minHash**; - Por último, uma função para obter as assinaturas e guardá-las num array;

```

1 function shingles = createShingles(str, n)
2     shingles = cell(1, length(str) - n + 1);
3     for i = 1:length(str) - n + 1
4         shingles{i} = str(i:i+n-1);
5     end
6 end
7
8 function jaccardSimilarity = compareMinHashSignatures(sig1, sig2)
9     intersection = sum(sig1 == sig2);
10    union = numel(unique([sig1, sig2]));
11    jaccardSimilarity = intersection / union;
12 end
13
14 function minHash = minhash_DJB31MA(chave, seed, k)
15     if nargin < 2
16         seed = 127;
17         k = 100;
18     elseif nargin < 3
19         k = 100;
20     end
21
22     len = length(chave);
23     chave = double(chave);
24
25     h = seed;
26     for i = 1:len
27         h = mod(31 * h + chave(i), 2^32 - 1);
28     end
29
30     minHash = zeros(1, k);
31
32     for j = 1:k
33         h = mod(31 * h + j, 2^32 - 1);
34         minHash(j) = h;
35     end
36 end
37
38 function signatures = getSignatures(movies, k)
39     titlesShingles = cell(length(movies), 1);
40     signatures = inf(length(movies), k);
41     wb = waitbar(0, 'Calculating minhash signatures...');
42     for i = 1:length(movies)
43         if mod(i, 10) == 0
44             waitbar(i/length(movies), wb, 'Calculating minhash signatures...');
45         end
46         titlesShingles{i} = createShingles(movies{i,1}, 3); % Obter os shingles do
titulo do filme
47         for j = 1:length(titlesShingles{i}) % Para cada shingle do
nome do filme
48             key = titlesShingles{i}{j}; % Obter o shingle
49             minHash = minhash_DJB31MA(key, 127, 1000); % Calcular a assinatura
minhash
50             signatures(i, :) = min(signatures(i, :), minHash); % Guardar a assinatura
minhash mais pequena
51         end
52     end
53     close(wb);
54 end

```

Para o cálculo das similaridades, obtivemos, inicialmente, uma matriz de assinaturas minHash para cada filme no cell array **moviesbyGenre**. Procedemos também, de igual modo, ao cálculo do valor das assinaturas para a string que tinha sido introduzida pelo utilizador, armazenada em **movieSearch**. Por fim tendo o valor das assinaturas para ambos os casos, criámos uma matriz **similarities**, para armazenar os valores das similaridades entre a string introduzida e o nome de cada filme.

```

1     numMoviesGenre = height(moviesbyGenre);
2     signatures = getSignatures(moviesbyGenre, 1000);
3     minHashSearch = inf(1, 1000);
4     shingle = createShingles(movieSearch, 3);

```

```

5         for j = 1:length(shingle) % Para cada shingle do nome
6             do filme
7                 key = shingle{j}; % Obter o shingle
8                 minHash = minhash_DJB31MA(key, 127, 1000); % Calcular a
9                 assinatura minhash
10                minHashSearch(1, :) = min(minHashSearch(1, :), minHash); % Guardar a
11                assinatura minhash mais pequena
12            end
13            similarities = zeros(1,numMoviesGenre);
14            for i = 1:numMoviesGenre
15                similarities(i) = (sum(minHashSearch(1, :) == signatures(i,:)) / 1000);
16            end

```

Por fim, ordenamos a matriz, do maior valor para o menor e imprimimos no terminal a informação relativa a esse filme: similaridade de Jaccard, nome do filme e os géneros do respetivo filme.

```

1         sortedSimilarities = sort(similarities, 'descend');
2         ind = zeros(1, 5);
3         for j = 1:5
4             for k = 1:numMoviesGenre
5                 if (sortedSimilarities(j) == distances(k))
6                     if ~ismember(k, ind)
7                         ind(1, j) = k;
8                     end
9                 end
10            end
11        end
12        for j = 1:5
13            printInfo(moviesbyGenre, ind(1, j), sortedSimilarities(j));
14        end

```

Demonstramos de seguida uma implementação desta opção:

```

1 Genres: (no genres listed) Action Adventure Animation
2 Children Comedy Crime Documentary
3 Drama Fantasy Film-Noir Horror
4 IMAX Musical Mystery Romance
5 Sci-Fi Thriller War Western
6 Select a genre: animation
7 -----
8 SELECTED GENRE: Animation
9 -----
10 1 - Change selected Genre
11 2 - No. of movies of selected Genre on given years
12 3 - Search movie titles of selected Genre
13 4 - Search movies based on Genres
14 5 - Exit
15 -----
16 Select an option: 3
17 Insert a string: Toy Story
18 {1.0000} Toy Story - Adventure, Animation, Children, Comedy, Fantasy
19 {0.7790} Toy Story 3 - Adventure, Animation, Children, Comedy, Fantasy, IMAX
20 {0.7780} Toy Story 2 - Adventure, Animation, Children, Comedy, Fantasy
21 {0.7730} Toy Story 4 - Adventure, Animation, Children, Comedy
22 {0.3920} Toy Story of Terror - Animation, Children, Comedy
23 -----
24 SELECTED GENRE: Animation
25 -----
26 1 - Change selected Genre
27 2 - No. of movies of selected Genre on given years
28 3 - Search movie titles of selected Genre
29 4 - Search movies based on Genres
30 5 - Exit
31 -----
32 Select an option: 5
33 Exiting the program

```

1.4. Opção 4

Na quarta opção era expectável que o utilizador inseri-se uma série de géneros de filmes (ou nenhum) separados por vírgulas. O objetivo era apresentar os filmes cujo conjunto de géneros fosse mais similar ao conjunto composto pela lista introduzida. Para tal, primeiro processámos o input do utilizador. Caso não

introduzisse nada, a lista de géneros iria ser constituída, apenas, pelo género atual. Caso contrário essa lista ia ser aumentada consoante os géneros introduzidos:

```

1      % Ask for the number of genres to search, it can be none
2      userInput = input('Select additional Genres separated by (',',') (press ENTER
when no desired additional Genres): ', 's');
3      if isempty(userInput)
4          movieGenres = cell(1, 1);
5          movieGenres{1} = genreInput;
6      else
7          % Normalize the input, making the first letter uppercase and the rest
lowercase
8          upperLetter = upper(userInput(1));
9          userInput = [upperLetter, userInput(2:end)];
10         % Get the genres from the input
11         movieGenres = split(userInput, ',');
12         % Check if the input genres are valid
13         for i = 1:length(movieGenres)
14             if ~ismember(movieGenres{i}, genres)
15                 fprintf('ERROR inputing the genre!\n');
16                 movieGenres = genreInput;
17                 break;
18             end
19         end
20         % Add genreInput to the movieGenres and remove duplicates
21         movieGenres{end+1} = genreInput;
22         movieGenres = unique(movieGenres);
23
24     end

```

De seguida, criámos as assinaturas para os géneros selecionados. Através da função **minhashDJB31MA**, previamente demonstrada, demos o valor à variável **minHash** para, de seguida, criar a matriz de assinaturas para os géneros selecionados. Noutro ficheiro, criámos a matriz de assinaturas para todos os filmes. De igual modo, calculamos o valor da similaridade entre cada filme e a lista de géneros e atribuímos tudo a uma matriz **similarities**.

```

1      % Create signatures for the selected genres
2      selectedSign = inf(1, kMinHash);
3      for i = 1:length(movieGenres)           % For each genre
4          key = char(movieGenres(i));         % Get the genre
5          minHash = minhash_DJB31MA(key);     % Create the minhash signature
6          selectedSign = min(selectedSign, minHash); % Get the minimum value for
each hash function
7      end
8
9      % Compare the signatures of the movies with the signatures of the selected
genres
10     similarity = zeros(numMovies, 1);
11     for i = 1:numMovies
12         similarity(i) = compareMinHashSignatures(selectedSign, signaturesGenres(i,
:));
13     end
14
15     %% Outro Ficheiro
16     genres = movies(1, 3:10);                % Obter os generos
17     k = 100;                                % Numero de funcoes de
hash a usar
18     signaturesGenres = inf(length(movies), k); % Guardar as
assinaturas minhash de cada filme
19     moviesGenres = cell(length(movies), 1);  % Guardar os generos de
cada filme
20     wb = waitbar(0, 'Calculating minhash signatures...');
21     for i = 1:length(movies)                 % Para cada filme
22         if mod(i, 10) == 0
23             waitbar(i/length(movies), wb, 'Calculating minhash signatures...');
24         end
25         moviesGenres{i} = movies(i, 3:10);  % Obter os
generos do filme
26         for j = 1:length(moviesGenres{i})   % Para cada
genero do filme
27             key = moviesGenres{i}{j};       % Obter o
genero
28             if ~isNotMissing(key)            % Se o
genero for missing

```

```

29         continue; % Passar ao
        proximo genero
30     end
31     minHash = minhash_DJB31MA(key); % Calcular
a assinatura minhash
32     signaturesGenres(i, :) = min(signaturesGenres(i, :), minHash); %
Guardar a assinatura minhash mais pequena
33     end
34 end
35 close(wb);
36
37 save('signaturesGenres.mat', 'signaturesGenres', 'moviesGenres', 'genres');

```

Por último ordenamos a matriz **similarities**, de modo a apresentarmos os 5 filmes com maior similaridade:

```

1 % Sort the movies by similarity, then most recent year
2 [~, idx] = sortrows([similarity, cell2mat(movies(:, 2))], [-1, -2]);
3
4 % Select the top 5 movies
5 idx = idx(1:5);
6
7 % Print the top 5 movies
8 fprintf('\n%-s\t%-40s\t%-16s\t\n', 'Year', 'Movie', 'Similarity');
9 for i = 1:length(idx)
10     % Print the name, year and similarity of the movie
11     fprintf('%-d\t%-40s\t%-16s\t\n', movies{idx(i), 2}, movies{idx(i), 1},
similarity(idx(i)));
12 end

```

Apresentamos um breve exemplo da execução deste bloco de código:

```

1 Genres: (no genres listed)  Action    Adventure  Animation
2   Children          Comedy    Crime      Documentary
3   Drama             Fantasy   Film-Noir  Horror
4   IMAX              Musical   Mystery    Romance
5   Sci-Fi            Thriller  War        Western
6 Select a genre: action
7 -----
8 SELECTED GENRE: Action
9 -----
10 1 - Change selected Genre
11 2 - No. of movies of selected Genre on given years
12 3 - Search movie titles of selected Genre
13 4 - Search movies based on Genres
14 5 - Exit
15 -----
16 Select an option: 4
17 Select additional Genres separated by (',') (press ENTER when no desired additional Genres)
   : Comedy
18
19 Year  Movie                                     Similarity
20 2019  Celebration                               0.282
21 2019  The Farewell                             0.282
22 2019  Sebastian Maniscalco: Stay Hungry         0.282
23 2019  Non ci resta che il crimine                0.282
24 2019  The Beach Bum                             0.282
25 -----
26 SELECTED GENRE: Action
27 -----
28 1 - Change selected Genre
29 2 - No. of movies of selected Genre on given years
30 3 - Search movie titles of selected Genre
31 4 - Search movies based on Genres
32 5 - Exit
33 -----
34 Select an option: 5
35 Exiting the program

```


2. Fundamentação das opções tomadas na implementação

2.1. Tamanho dos Shingles

Para o tamanho dos shingles, nós consideramos que 3, seria o tamanho mais adequado, uma vez que nos pareceu o número mais equilibrado. Com 2 podia haver muitas semelhanças, por ser uma amostra pequena, e seria provável a obtenção de valores deturpados. Para tamanhos de shingles maiores, também o valor sairia deturpado, uma vez que as strings tinham de ser muito parecidas para dar *match* em 4 ou até 5 caracteres seguidos. Daí a nossa escolha pareceu-nos a mais adequada.

2.2. Número de funções de dispersão

No desenvolvimento, ao usar 100 funções de dispersão, percebemos que o nosso código devolvia os valores expectáveis e não era muito demorado, portanto optamos por esse número.

2.3. Dimensionamento dos Filtros de Bloom

```
1 n = int32(-log(p) * numElem / (log(2)^2)); % Size of ideal bloom filter
2 k = int32(n * log(2) / numElem); % Number of hash functions
```

No nosso código foram implementadas estas duas fórmulas que representam o tamanho ideal para o Filtro de Bloom e também o número para as funções de dispersão.

3. Conclusões

Este trabalho teve como objetivo, essencialmente, a criação de uma aplicação aplicando os conhecimentos obtidos relativamente a algoritmos probabilísticos. Este trabalho, permitiu-nos, acima de tudo, a um desenvolvimento de capacidades e competências em Matlab e, especialmente, a operar com **Filtros de Bloom** e com o método **minHash**.