# Projeto de Laboratórios de Informática 3 Grupo 20

Diogo Miguel Alves Rocha (A79751) Gabriela Sá Martins (A81987) Ricardo Milhazes Veloso (A81919) Ricardo Jorge Silva Ferreira (82568)

### 12 de Junho de 2018

#### Resumo

O objetivo deste trabalho consistiu em criar um programa que permitisse a extração e análise de dados provenientes da base de dados correspondente ao StackOverFlow. Tendo como objetivo final a resposta a um determinado número de querys , em que o tempo de execução será o mais reduzido possível. Esta segunda parte do projeto contém o mesmo objetivo da primeira mas desta feita a linguagem de programação usada é Java em vez de C.

# Conteúdo

1	Introdução	2
2	Descrição do Problema	2
3	Concepção da Solução 3.1 Porquê da HashTable	<b>2</b> 2
4	A Estrutura         4.1 id_utilizador         4.2 id_perguntas         4.3 id_respostas         4.4 id_tags	3 3 4 4
5	Extração de Informação 5.1 Pergunta 1	<b>5</b>
	5.2 Pergunta 2	5 5 5
	5.5 Pergunta 5	5 5 6
	5.8 Pergunta 8	6 6 6 6
6	LI3 em java vs LI3 em C	7
7	Resultados	7
8	Teste de Tempo	7

9 Conclusão 8

# 1 Introdução

O trabalho proposto tinha em vista a resolução de 11 interrogações da forma mais eficiente e rápida possível. Para que conseguissemos corresponder a estas medidas criamos quatro estruturas, todas elas HashTables .

# 2 Descrição do Problema

O principal problema era criar uma estrutura capaz de armazenar todos os dados provenientes da base de dados correspondente ao StackOverFlow, de maneira que o acesso a esses mesmos dados fosse o mais rápido possível que por consequência a resposta às interrogações fosse ainda mais rápida, e sem memory leaks.

# 3 Concepção da Solução

Com vista a resolver o problema decidimos criar 4 HashTables, uma com os dados do utilizador, id\_utilizador, outra com os dados referentes aos posts do tipo pergunta id\_perguntas, outra com os dados referentes aos posts do tipo resposta id\_espostas, e por fim outra com as tags id\_tags. Cada uma das anteriores será explicitada numa secção posterior referente às mesmas.

# 3.1 Porquê da HashTable

O principal motivo da escolha deste tipo de estrutura de dados foi o tempo de procura de um elemento que no melhor caso vai ser sempre constante e é o tipo de estrutura que apresenta o tempo de procura inferior em relação a outros.

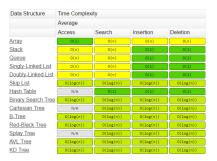


Figura 1: timecomplexity

### 4 A Estrutura

```
private static HashMap<Long,Utilizador> id_utilizador;
private static HashMap<Long,Pergunta> id_perguntas;
private static HashMap<Long,Resposta> id_respostas;
private static HashMap<Long,Tag> id_tags;
```

Figura 2: Estrutura

#### 4.1 id\_utilizador

A estrutura contém informações relativas ao utilizador, e apresenta 7 parâmetros diferentes:

id\_u -ID do utilizador;

name -Nome do utilizador;

shortb -Short bio do utilizador;

totalposts - Número total de posts;

reputação do utilizador;

votes -Número de votes do utilizador;



Figura 3: id\_utilizador

#### 4.2 id\_perguntas

A estrutura contém informações relativas às **Perguntas**, e apresenta 7 parâmetros diferentes:

 $\mathbf{id}_{-}\mathbf{post}\,$  -ID da pergunta;

 $id\_autor$  -ID do autor da pergunta;

title -Título da pergunta;

tags -Tag da pergunta;

reputação do utilizador;

 $\mathbf{n} \_{\mathbf{resp}}$ -Número de respostas àquela pergunta;

d -Data da pergunta;



Figura 4: id\_perguntas

### 4.3 id\_respostas

A estrutura contém informações relativas às Respostas, e apresenta 6 parâmetros diferentes:

id\_post -ID da Resposta;

 $\mathbf{id}\mathtt{\_perg}\,$  -ID da pergunta correspondente à resposta;

id\_autor -ID do autor da resposta;

score -Score da resposta;

**n\_coments** -Número de comentários aquela resposta;

 ${\bf d}\,$  -Data da Resposta;

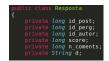


Figura 5:  $id_respostas$ 

## $4.4 id_{tags}$

A estrutura contém informações relativas às Tags, e apresenta 3 parâmetros diferentes:

 $id\_tag$  -ID da tag;

tag -Tag;

count -Conta o número de vezes que uma Tag é utilizada num intervalo de tempo;



Figura 6: id\_tags

# 5 Extração de Informação

#### 5.1 Pergunta 1

Na pergunta 1 temos de retornar um par com o titulo do post correspondente ao id recebido, e com o nome do utilizador a que esse post pertence. Para isso acedemos à **id\_perguntas** retiramos o titulo e o nome do autor correspondente ao id recebido e retornamos um par com os respetivos.

#### 5.2 Pergunta 2

Nesta pergunta queremos saber o Top N utilizadores com maior número de posts e para isso consultamos a variável totalposts, através do método que implementamos na classe Utilizador **getTotalposts()**,após utilizarmos um **ArrayList** infoPost que guarda toda a informação existente na **id\_utilizador**. De seguida ordenamos esse ArrayList pelo número total de posts. E no fim retornamos um List com os **id\_u** ordenados.

#### 5.3 Pergunta 3

Nesta pergunta queremos obter o número total de posts num determinado periodo de tempo mas com perguntas e respostas em separado. Para isso percorremos a **id\_perguntas** e a **id\_respostas**, e utilizando o método **compareTo()**, retiramos todas as perguntas e respostas, respetivamente que se encontram entre as datas recebidas. No fim retornamos um par com essas perguntas e respostas.

#### 5.4 Pergunta 4

Nesta pergunta queremos obter o número total de perguntas contendo uma determinada tag num determinado periodo de tempo, e retornar essas perguntas por cronologia inversa. Para isso percorremos a **id\_perguntas** e utilizando o método **compareTo**e o método **contains**, retiramos todas as perguntas que se encontram entre as datas recebidas e que contém uma tag igual à recebida, e adicionamos estas perguntas a um **ArrayList** TagsData. Em seguida ordenamos este ArrayList por data (da mais recente para a mais antiga) utilizando novamente o método **compareTo()**. Por fim adicionamos numa List todos os id\_post que contêm aquela tag, retornado essa List.

#### 5.5 Pergunta 5

Nesta pergunta queremos devolver a informação do perfil do utilizador (short bio) e os IDs dos seus 10 últimos posts (perguntas ou respostas),ordenados por cronologia inversa. Para isso utilizamos um **Utilizador** u para guardar a informação de um utilizador corresponde ao id recebido, depois acedemos a esse Utilizador para retirar a shorbio. Em seguida criamos um **ArrayList**ultimasPerg e um **ArrayList** ultimasResp e utilizamos o método **Long.compare** em ambos os casos para comparar se o id recebido corresponde a uma pergunta ou a uma resposta se algum dos casos se verificar adiciona uma pergunta na **ultimasPerg** ou uma resposta na **ultimasResp**. Em seguida criamos uma **LinkedHashMap** result e adicionamos o id\_post e a data retirados do ultimasPerg e do ultimasResp. Depois ordenamos esse LinkedHashmap pela data. Por fim criamos um List(res) que recebe todas as keys da LinkedHashmap já ordenada e depois adicionamos noutro List(f) as 10 primeiras keys. E retornamos um par com a shortbio e com a List f.

#### 5.6 Pergunta 6

Nesta pergunta queremos saber, num intervalo de tempo arbitrario, quais as N respostas com mais votos, devolvendo uma lista com a variável id\_post. Para isso criamos um **ArrayList**.Em seguida percorremos essa a id\_respostas e utilizando o método **compareTo** adicionamos à resp\_Data todas as respostas que se encontrem entre a data recebida. Em seguida ordenamos a resp\_data através

do Score(por ordem decrescente) , e criamos um  ${f List}$  res um N id\_post . Por fim retorna esse mesmo res.

#### 5.7 Pergunta 7

Nesta pergunta queremos obter os IDs das N perguntas com mais respostas, por ordem decrescentem num determinado intervalo de tempo. Para isso consultamos a variável n\_resp, através do método que implementamos na classe Pergunta **getN\_resp()**, após utilizarmos um **Array-List;Pergunta;** que guarda toda a informação existente na **id\_perguntas** dentro das datas fornecidas. De seguida, ordenamos esse ArrayList de acordo com o número de respostas (por ordem decrescente) e por fim, retornamos um List com os N **id\_post** das perguntas com mais respostas ordenado.

#### 5.8 Pergunta 8

Nesta pergunta queremos obter os IDs das N perguntas cujos títulos contém uma palavra específica, ordenados por cronologia inversa. Para isso percorremos a **id\_perguntas** e utilizando **compareTo** e o método **contains**. Guardamos todas as perguntas cujos títulos continham a palavra fornecida num **ArrayList;Pergunta**; e depois ordenamos o mesmom de acordo com a data da pergunta. Por fim, retornamos um List com os N **id\_post** das perguntas que continham no título a palavra dada, ordenados por cronologia inversa.

#### 5.9 Pergunta 9

Nesta pergunta recebendo o ID de dois utilizadores devolver as últimas N perguntas que em que participaram esses dois utilizadores. Para isso criamos um **ArrayList** resps e um **ArrayList** perg final e guardamos todas as respostas da id\_respostas na resps. Em seguida usamos o método **Long.compare** para comparar o id do autor com o id1 recebido e adicionamos em p1 a pergunta que contenha aquele id. Após isto se o id2 for igual o id\_autor adiciona p1 na pergsFinal. Em seguida utiliza o mesmo método mas para a resposta. Depois ordena o pergsFinal por data , cria uma List res e adiciona N id\_post a essa lista res. Por fim retorna essa mesma lista.

#### 5.10 Pergunta 10

Nesta pergunta queremos obter a melhor resposta a uma pergunta específica, utilizando a função de média (Score da resposta x 0.45) + (Reputação do utilizador x 0.25) + (número de Votos recebidos pela resposta x 0.2) + (número de comentários recebidos pela resposta x 0.1). Para isso percorremos a id\_respostas guardando num ArrayList;Resposta; as respostas correspondentes à pergunta cujo ID foi dado, através do método getID\_perg(). Em seguida, é percorrida a id\_utilizador e são guardados os utilizadores a quem corresponde a resposta. Depois acedemos às variáveis score, n\_coments e reputação através dos métodos getScore(), getN\_coments(), getReputação(), respetivamente, implementadas na classe Utilizador. De seguida são efetuados os cálculos de acordo com a fórmula anteriormente referida e os valores resultantes vão sendo comparados até se obter a melhor resposta e esta ser retornada.

#### **5.11** Pergunta 11

Infelizmente não conseguimos um resultado positivo nesta pergunta.

# 6 LI3 em java vs LI3 em C

No projeto de LI3 em java chegamos á conclusao que, comparativamente com este em C, este apresenta beneficios ao nivel do tratamento de exceções , uso de memória, compilação do código e a portabilidade. Tratamento de exceções A linguagem java apresenta beneficios ao nivel do tratamento de exceções pois esta admite recursos para o tratamento das mesmas

Uso de memoria Na linguagem java o uso de memoria é mais superficial comparado com a linguagem C. Em C foi necessario o uso de apontadores e definicao da capacacidade dos arrays etc. Enquanto que em java nao é necessaria essa preocupacao pois o proprio java aloca o espaço necessario para o carregamento da estrutura que etsmaos a definir O que torna a programacao em java mais facil mas tambem menos controlavel

Utilizacao de librarias Em Java existem librarias em que a ordenacao, insercao em hashtables, tratamento de exceçoes etc ja estao implementadas e bem implementadas. Enquanto que em C existem menos e a maioria tem de ser escrita o que proporciona uma maior probablidade de erros e bugs.

#### 7 Resultados

```
Query 1 ->
(Mak are the actual risks of glving www-data sudo nopassud access?, NebNinja)
(gery 2 ->
(15811, 449, 188442, 50676, 167859, 307105, 292286, 93977, 35785, -1)
(gery 3 ->
(2545, 4812)
(gery 4 ->
(2545, 4812)
(gery 4 ->
(2547, 276050, 276462, 201334, 27410, 374145, 274108, 272937, 272813, 272754, 272566, 272565, 272459, 272313, 271816, 271683, 27668, 275528, 276488, 270528, 276488, 270528, 276488, 270528, 276769, 207629, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 206762, 2067
```

Figura 7: Resultados

# 8 Teste de Tempo

```
1 INIT -> 0 ms
2 LOAD -> 12069 ms
3 Query 1: -> 1 ms
4 Query 2 -> 148 ms
5 Query 3 -> 65 ms
6 Query 4 -> 38 ms
7 Query 5 -> 113 ms
8 Query 6 -> 50 ms
9 Query 7 -> 45 ms
10 Query 8 -> 49 ms
11 Query 10 -> 27 ms
12 Query 11 -> 0 ms
13 Query 11 -> 0 ms
```

Figura 8: tempo

# 9 Conclusão

Em suma quase todos os objetivos propostos foram atingidos, apenas os resultados pretendidos não foram atingidos na query 5 em que o return não é o correto e na query 11 que não conseguimos criar uma maneira de a resolver. Toda a extração de informação proveniente da base de dados do site **StackOverflow** e análise da mesma é feita de uma maneira eficiente, com um tempo de execução bastante aceitável e sem qualquer perda de informação. Esta segunda fase do projeto foi bastante útil para desenvolver competencias na linguagem de programação java bem como aprender a utilizar a ferramenta **Maven**. Por tudo isto acreditamos que a segunda fase do projeto foi concluida com sucesso.