



Cento Universitário UNA

Sistemas de Informação

Tecnologias Emergentes

Práticas de Laboratório

Wesley Dias Maciel

2019/02



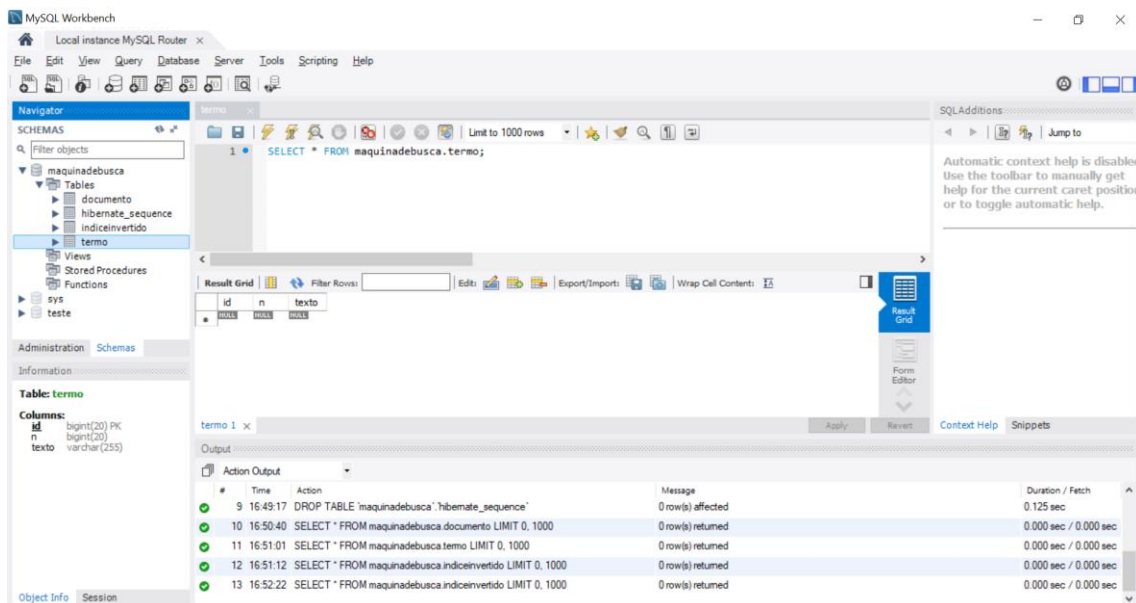
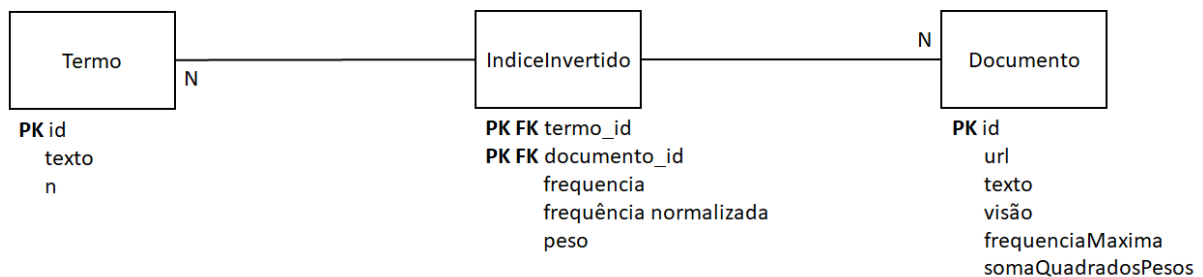
Centro Universitário UNA
Sistemas de Informação
Tecnologias Emergentes
Prática de Laboratório
Wesley Dias Maciel
2019/02

Indexador



Prática 16

- 1) Você está recebendo, juntamente com esta prática, uma planilha Excel com um exemplo de cálculo de pesos no modelo de espaço vetorial e também o projeto da aplicação. Nesta versão, o projeto inicia a implementação do indexador. O projeto cria as tabelas Termo, ÍndiceInvertido e Documento no banco de dados como apresentado nas figuras abaixo:





The screenshot shows the MySQL Workbench interface. On the left, the 'Navigator' pane displays the 'maquinadebusca' database with tables 'documento', 'hibernate_sequence', 'indiceinvertido', and 'termo'. The 'documento' table is selected, showing its columns: 'id' (bigint(20) PK), 'frequenciaMaxima' (double), 'somaQuadradosPesos' (longtext), 'texto' (varchar(255)), 'url' (longtext), and 'visao' (longtext). The main editor shows a query: `SELECT * FROM maquinadebusca.documento;`. The 'Result Grid' displays the query results with columns: 'id', 'frequenciaMaxima', 'somaQuadradosPesos', 'texto', 'url', and 'visao'. The 'Output' pane shows the execution log with the following entries:

#	Time	Action	Message	Duration / Fetch
10	16:50:40	SELECT * FROM maquinadebusca.documento LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
11	16:51:01	SELECT * FROM maquinadebusca.termo LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
12	16:51:12	SELECT * FROM maquinadebusca.indiceinvertido LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
13	16:52:22	SELECT * FROM maquinadebusca.indiceinvertido LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
14	16:53:52	SELECT * FROM maquinadebusca.documento LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec

The screenshot shows the MySQL Workbench interface. On the left, the 'Navigator' pane displays the 'maquinadebusca' database with tables 'documento', 'hibernate_sequence', 'indiceinvertido', and 'termo'. The 'indiceinvertido' table is selected, showing its columns: 'frequencia' (int(11)), 'frequenciaNormalizada' (double), 'peso' (double), 'documento_id' (bigint(20) PK), and 'termo_id' (bigint(20) PK). The main editor shows a query: `SELECT * FROM maquinadebusca.indiceinvertido;`. The 'Result Grid' displays the query results with columns: 'frequencia', 'frequenciaNormalizada', 'peso', 'documento_id', and 'termo_id'. The 'Output' pane shows the execution log with the following entries:

#	Time	Action	Message	Duration / Fetch
11	16:51:01	SELECT * FROM maquinadebusca.termo LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
12	16:51:12	SELECT * FROM maquinadebusca.indiceinvertido LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
13	16:52:22	SELECT * FROM maquinadebusca.indiceinvertido LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
14	16:53:52	SELECT * FROM maquinadebusca.documento LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec
15	16:54:23	SELECT * FROM maquinadebusca.indiceinvertido LIMIT 0, 1000	0 row(s) returned	0.000 sec / 0.000 sec

Nesta versão, os documentos usados são os apresentados nos slides empregados na apresentação do conteúdo. Analise o código do projeto. Execute o projeto usando o Postman. Em seguida, execute o comando abaixo no MySQL e verifique se as frequências estão sendo calculadas corretamente.



use maquinadebusca;

```
select t.texto as termo, d.url as documento, i.frequencia as frequencia
from termo t, documento d, indiceinvertido i
where t.id = i.termo_id and i.documento_id = d.id
order by t.texto, d.url;
```

Planilha Excel disponibilizada juntamente com a prática:

The screenshot shows an Excel spreadsheet with the following data:

Termo	$f_{1,1}$	$f_{1,2}$	$f_{1,3}$	$f_{1,4}$	n_i	$TF_{1,1}$	$TF_{1,2}$	$TF_{1,3}$	$TF_{1,4}$	IDF _i	$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$
to	4	2	0	0	2	3	2	0	0	1	3	2	0	0
do	2	0	3	3	3	2	0	2,5849625	2,5849625	0,41504	0,83007	0	1,07286	1,07286
is	2	0	0	0	1	2	0	0	0	2	4	0	0	0
be	2	2	2	2	4	2	2	2	2	0	0	0	0	0
or	0	1	0	0	1	0	1	0	0	2	0	2	0	0
not	0	1	0	0	1	0	1	0	0	2	0	2	0	0
i	0	2	2	0	2	0	2	2	0	1	0	2	2	0
am	0	2	1	0	2	0	2	1	0	1	0	2	1	0
what	0	1	0	0	1	0	1	0	0	2	0	2	0	0
think	0	0	1	0	1	0	0	1	0	2	0	0	2	0
therefore	0	0	1	0	1	0	0	1	0	2	0	0	2	0
da	0	0	0	3	1	0	0	0	2,5849625	2	0	0	0	5,16993
let	0	0	0	2	1	0	0	0	2	2	0	0	0	4
it	0	0	0	2	1	0	0	0	2	2	0	0	0	4
Soma dos Quadrados dos Pesos:											25,689	24	14,151	59,8791
Raiz Quadrada da Soma dos Quadrados dos Pesos:											5,06843	4,89898	3,76178	7,73816

Alguns dados nessa tabela podem ser comparados com os armazenados no banco de dados, usando:

- 2) Altere o projeto, para que ele calcule a quantidade de documentos em que cada termo i ocorre (n_i) e armazene essa quantidade na coluna correspondente da tabela Termo.
 - 3) Altere o projeto, para que ele calcule os pesos dos termos em cada documento e os armazene na coluna correspondente da tabela IndiceInvertido.
- OBS:**
- $$\text{peso} = TF_{i,j} \times IDF_i$$
- $$TF_{i,j} = 1 + \log \text{frequencia}_{i,j}$$
- $$IDF_i = \log N / n_i$$
- 4) Altere o projeto, para que ele calcule o valor da raiz quadrada da soma dos quadrados dos pesos de cada documento e o armazene na coluna correspondente da tabela Documento.
 - 5) Realize testes com o seu projeto, usando os exercícios passados em sala de aula como base.
 - 6) Altere o projeto, para que ele indexe os documentos coletados por seu coletor.



- 7) Analise a API do seu projeto. Sempre que necessário, faça alterações para melhoria da API, adequando-a ao padrão arquitetural REST (Representational State Transfer).
- 8) Para todos os métodos que interagem com aplicações cliente, retorne respostas com mensagens significativas para os clientes da aplicação, obedecendo os códigos adequados do protocolo HTTP.

Lista de códigos de status HTTP:

1xx Informativa

- 100 Continuar
- 101 Mudando protocolos
- 102 Processamento (WebDAV) (RFC 2518)
- 122 Pedido-URI muito longo

2xx Sucesso

- 200 OK
- 201 Criado
- 202 Aceito
- 203 não-autorizado (desde HTTP/1.1)
- 204 Nenhum conteúdo
- 205 Reset
- 206 Conteúdo parcial
- 207-Status Multi (WebDAV) (RFC 4918)

3xx Redirecionamento

- 300 Múltipla escolha
- 301 Movido
- 302 Encontrado
- 303 Consulte Outros
- 304 Não modificado
- 305 Use Proxy (desde HTTP/1.1)
- 306 Proxy Switch
- 307 Redirecionamento temporário (desde HTTP/1.1)
- 308 Redirecionamento permanente (RFC 7538[2])

4xx Erro de cliente

- 400 Requisição inválida
- 401 Não autorizado
- 402 Pagamento necessário
- 403 Proibido
- 404 Não encontrado
- 405 Método não permitido
- 406 Não Aceitável
- 407 Autenticação de proxy necessária
- 408 Tempo de requisição esgotou (Timeout)
- 409 Conflito
- 410 Gone
- 411 comprimento necessário
- 412 Pré-condição falhou
- 413 Entidade de solicitação muito grande
- 414 Pedido-URI Too Long
- 415 Tipo de mídia não suportado
- 416 Solicitada de Faixa Não Satisfatória



417 Falha na expectativa
418 Eu sou um bule de chá
422 Entidade improcessável (WebDAV) (RFC 4918)
423 Fechado (WebDAV) (RFC 4918)
424 Falha de Dependência (WebDAV) (RFC 4918)
425 coleção não ordenada (RFC 3648)
426 Upgrade Obrigatório (RFC 2817)
450 bloqueados pelo Controle de Pais do Windows
499 cliente fechou Pedido (utilizado em ERPs/VPsA)

5xx outros erros (erro de servidor)

500 Erro interno do servidor (Internal Server Error)
501 Não implementado (Not implemented)
502 Bad Gateway
503 Serviço indisponível (Service Unavailable)
504 Gateway Time-Out
505 HTTP Version not supported