

Documentação do banco de dados

Richard Auzier
Tiago Freitas
Kaliu Freitas

21 de setembro de 2024

1 Esquema do Banco de Dados

1.1 Definindo as relações

O esquema do banco de dados foi construído usando a abordagem bottom-up, isto é, partimos dos atributos e construímos as relações. Os atributos extraídos do trecho do arquivo fornecido foram agrupados na relação R :

$R(asin, title, group_name, salesrank, similar_product, category_name, category_id, review_time, review_user_id, review_rating, review_total_votes, review_helpfulness_votes)$.

Observamos anomalias de redundância e de atualização/remoção. Escolhemos a Forma Normal de Boyce-Codd (BCNF) para resolvê-las. Esta norma é baseada em dependências funcionais(FDs), portanto listamos as dependências funcionais encontradas:

$asin \rightarrow title, group_name, salesrank$
 $(asin, review_user_id, review_time) \rightarrow review_rating, review_total_votes, review_helpfulness$

Uma relação R está na BCNF se, e somente se, sempre que existir uma dependência funcional não-trivial, da forma $A_1, A_2, \dots, A_n \rightarrow B$, é o caso que $\{A_1, A_2, \dots, A_n\}$ é uma superchave de R , isto é, um conjunto de atributos que contém uma chave.

Observando as FDs, dividimos a relação R em duas relações, R_1 e $Review$. R_1 diz respeito ao produto, enquanto $Review$ foi criada especificamente para reviews pois o formato destas no arquivo de entrada sugere tal agregação. Removemos "review" dos nomes dos atributos. As relação resultantes foram:

$R_1(asin, title, group_name, salesrank, similar_product, category_name, category_id)$

e

$Review(asin, user_id, time, rating, total_votes, helpfulness_votes)$.

A relação R_1 não está na BCNF. Apesar da FD $asin \rightarrow title, group_name, salesrank$ existir na relação R_1 , $\{asin\}$ não determina funcionalmente todos os atributos, logo não é uma chave. Removemos os atributos que não são funcionalmente determinados por $asin$ e criamos três novas relações. Temos então:

$R_1(asin, title, group_name, salesrank)$,
 $R_2(asin, similar_product)$,
 $R_3(asin, product_category)$,
 $R_4(category_id, category_name)$ e
 $Review(asin, user_id, time, rating, total_votes, helpfulness_votes)$.

Vamos verificar que todas obedecem à BCNF.

R_1 tem a FD $asin \rightarrow title, group_name, salesrank$ e, desta vez, $\{asin\}$ é uma chave, portanto também uma superchave, de R_1 .

R_2 e R_3 não possuem FD (para cada $asin$ podem haver múltiplos $similar_product/product_category$ e vice-versa) portanto obedecem à BCNF por vacuidade.

R_4 possui duas FDs: $category_id \rightarrow category_name$ e $category_name \rightarrow category_id$. Ambos atributos são formas diferentes de representar a mesma coisa, obviamente são chaves então esta relação obedece à BCNF.

Por fim, a relação *Review* tem a FD $(asin, user_id, time) \rightarrow rating, total_votes, helpfulness_votes$ e o conjunto $\{asin, review_user_id, review_time\}$ é uma chave: dados um produto, um usuário e um instante no tempo, a avaliação será única, isto é, os demais atributos serão funcionalmente determinados.

Para maior clareza, as respectivas relações e alguns dos seus atributos foram renomeados como listado a seguir:

Product(*asin*, *title*, *group_name*, *salesrank*),
Similar_product(*product_asin*, *similar_asin*),
Product_category(*product_asin*, *category_id*),
Category(*category_id*, *category_name*) e
Review(*product_asin*, *user_id*, *time*, *rating*, *total_votes*, *helpfulness_votes*).

1.2 Dicionário de dados

A seguir a descrição de cada relação, atributo e restrição de integridade referencial:

TABLE PRODUCT

ASIN	VARCHAR(10)	NOT NULL,
TITLE	VARCHAR(2000),	
GROUP_NAME	VARCHAR(2000),	
SALESRANK	INT	NOT NULL,
PRIMARY KEY (ASIN);		

TABLE SIMILAR_PRODUCT

PRODUCT_ASIN	VARCHAR(10)	,
SIMILAR_ASIN	VARCHAR(10)	,
FOREIGN KEY (PRODUCT_ASIN) REFERENCES PRODUCT (ASIN);		

TABLE PRODUCT_CATEGORY

PRODUCT_ASIN	VARCHAR(10)	,
CATEGORY_ID	INT;	
FOREIGN KEY (PRODUCT_ASIN) REFERENCES PRODUCT (ASIN);		
FOREIGN_ KEY (CATEGORY_ID) REFERENCES CATEGORY (CATEGORY_ID);		

TABLE CATEGORY

CATEGORY_ID	INT	NOT NULL,
CATEGORY_NAME	VARCHAR(500)	NOT NULL,
PRIMARY KEY (CATEGORY_ID);		

TABLE REVIEW

PRODUCT_ASIN	VARCHAR(10)	,
TIME	DATE	NOT NULL,
USER_ID	VARCHAR(100)	NOT NULL,
RATING	INT	NOT NULL,
TOTAL_VOTES	INT	NOT NULL,
HELPFULNESS_VOTES	INT	NOT NULL,
FOREIGN KEY PRODUCT_ASIN REFERENCES PRODUCT (PRODUCT_ASIN);		

1.3 Diagrama do esquema do banco de dados

