

طراحی پایگاه داده ها

دانشکده مهندسی کامپیوتر

مریم رمضانی
بهار ۱۴۰۳



تمرین چهارم

تاریخ انتشار: ۱۸ خرداد ۱۴۰۳

۱. پرسش های خود درمورد این تمرین را در سامانه کوئرا مطرح کنید.

۲. سیاست ارسال با تاخیر: شما در مجموع در طول نیم سال می توانید از ۹ روز تاخیر استفاده کنید. همچنین هر تمرین را می توانید تا حداکثر ۴ روز با تاخیر تحویل دهید. تاخیرها با مقیاس ساعت محاسبه شده و به بالا گرد می شوند.

۳. سیاست مشارکت دانشجویان در حل کردن تمرین: دانشجویان می توانند در حل تمرین برای رفع ابهام و یا به دست آوردن ایده های کلی با یکدیگر مشورت و همفکری کنند. این کار مورد تایید و تشویق تیم ارائه ای درس می باشد؛ چرا که همفکری و کار گروهی می تواند موجب تقویت یادگیری شود. اما به دست آوردن جزئیات راه حل و نگارش پاسخ باید تماما توسط خود دانشجو انجام شود. حتما در انتهای پاسخ های ارسالی خود نام افرادی که با آنها همفکری کردید را ذکر کنید.

۴. برای سوال ششم، کد های SQL هر بخش را با نام گذاری مناسب، به طور مثال برای بخش اول 6-1.SQL، را در فایل zip پاسخ خود قرار دهید.

تاریخ تحویل: ۲۹ خرداد ۱۴۰۳

سوالات تئوری (۱۱۰ نمره)

پرسش ۱ (۱۰ نمره) فرض کنید جدول $R(A, B, C, D, E)$ مفروض است. آن را به سه رابطه جزئی تر $S_1(A, B, C)$ ، $S_2(B, C, D)$ و $S_3(A, C, E)$ تجزیه کرده ایم. در هر بخش فرض کنید مجموعه FD های معرفی شده برقرار است و نشان دهید آیا تجزیه با این فرض از وابستگی های تابعی، lossless است یا خیر.

$$(A) \quad B \rightarrow E, CE \rightarrow A$$

$$(B) \quad A \rightarrow D \rightarrow E, B \rightarrow D$$

$$(C) \quad AC \rightarrow E, BC \rightarrow D$$

$$(D) \quad A \rightarrow D, CD \rightarrow E, E \rightarrow D$$

پرسش ۲ (۲۰ نمره) در یک کار آماری، اطلاعات تیم های شرکت کننده در لیگ قهرمانان اروپا بین سال های ۲۰۰۲ تا ۲۰۱۹ جمع آوری شده است. همه اطلاعات موجود در یک جدول به صورت زیر ذخیره شده است:

UCL(year, team, group, winner, final-host-stadium, qualified, knockout-stage, best-scorer, most-assists, attendance-prize)

- کلید اصلی رابطه را (year, team) در نظر بگیرید.
- صفت بولین qualified نشان می دهد آیا تیم از مرحله گروهی صعود کرده است یا خیر.
- صفت knockout-stage آخرین مرحله حذفی است که تیم در صورت صعود از گروه به آن رسیده است. اگر یک تیم از گروه صعود نکرده باشد، این صفت مقدار ثابت groupstage را اخذ کرده است.
- صفت attendance-prize مقدار جایزه ای است که در نهایت به تیم اهدا میشود. تیم هایی که از گروه صعود نکرده اند مقدار ثابت X دریافت میکنند و به ازای هر مرحله بالاتر رفتن در مراحل حذفی، مقدار Y به این جایزه اضافه میشود.
- فینال هر دوره، تنها در یک استادیوم برگزار میشود.
- صفات بهترین گلزن و بیشترین پاس گل، مربوط به هر تیم هستند.
- فرض کنید هیچ صفتی مقدار null ندارد.

با توجه به این پایگاه داده، به سوالات زیر پاسخ دهید.

(آ) توضیح دهید پایگاه داده معرفی شده میتواند دارای چه انواعی از آنومالی در هنگام درج یا حذف باشد. میتوانید این آنومالی ها را با یک مثال نیز توضیح دهید.

(ب) به کمک شناسایی وابستگی های تابعی، پایگاه داده معرفی شده را تا NF۲ نرمال سازی کنید.

(ج) به کمک شناسایی وابستگی های ترا یا، پایگاه داده مرحله قبل را تا NF۳ نرمال سازی کنید.

(د) آیا پایگاه داده به دست آمده در مرحله قبل BCNF است؟ اگر جواب شما خیر است، نرمال سازی تا این مرحله را انجام دهید.

پرسش ۳ (۱۵ نمره)

درستی یا نادرستی هر یک از گزاره های زیر را با ذکر دلیل مشخص کنید.

(آ) هر رابطه ای که فقط یک خصیصه با مقادیر غیر تکراری دارد حتماً در BCNF واقع است.

(ب) نرمالسازی عموماً به افزایش سرعت کوثری ها منجر میشود.

(ج) نرمالسازی تا سطح NF³ عموماً حجم پایگاه داده را کاهش میدهد.

پرسش ۴ (۱۵ نمره) فرض کنید یک جدول پایگاه داده داریم که از ستونهای A, B, C, D, E تشکیل شده است. فرض کنید روابط وابستگی تابعی در این جدول به صورت زیر می باشند

$$FD = AE \rightarrow BC, AC \rightarrow D, CD \rightarrow BE, D \rightarrow E$$

با این فرضیات به پرسش های زیر پاسخ دهید.

(آ) سه کلید کاندید ارائه دهید و توضیح دهید که چرا این کلیدها کلید کاندید محسوب می شوند. (اگر بیش از سه کلید کاندید وجود دارد، ۳ کلید کاندید به دلخواه خود بنویسید.)

(ب) آیا این جدول با این وابستگی های تابعی، در فرم BCNF قرار دارد؟ توضیح دهید.

اکنون جدول جدیدی برای رابطه ای به نام R با ستونهای A, B, C, D در نظر بگیرید. وابستگی های تابعی این جدول شامل موارد زیر می شوند با توجه به این رابطه به بخش های بعد پاسخ دهید.

$$FD = A \rightarrow B, B \rightarrow C, C \rightarrow D$$

(آ) اگر این رابطه را به سه رابطه زیر تجزیه کنیم، آیا این تجزیه lossless است یا خیر؟ توضیح دهید.

R1(A, B)

R2(C, D)

R3(A, C)

(ب) این رابطه را به روابطی تجزیه کنید تا به فرم BCNF درآید. در صورتی که در بخش آ پاسخ تان بله بوده است، تجزیه ارائه شده در این قسمت باید متفاوت از تجزیه ارائه شده در بخش آ باشد

پرسش ۵ (۱۰ نمره) با توجه به مفاهیم مرتبط با index به سوالات زیر پاسخ دهید.

(آ) تفاوت شاخص های $B^+ - Tree$ و $B - Tree$ را تحقیق کرده و شرح دهید و برای هر یک مثالی کاربردی بزنید.

(ب) جدول Employee (id: integer, name: String, salary: integer) را در نظر بگیرید. حال برای کوثری های زیر بگویید که کدام نوع شاخص از بین hash index و $B^+ - Tree index$ مناسب است.

i.

```
SELECT name,
FROM Employee,
WHERE salary >= 5000
```

ii.

```
SELECT name,
FROM Employee,
WHERE id = 50
```

iii.

```
SELECT name,
FROM Employee,
WHERE id = 50 and salary >= 5000
```

پرسش ۶ (۴۰ نمره) برای این سوال، فایل کوثری، اسکرین شات نتایج و نوشتن تحلیل فرد لازم است. جدولی به اسم transaction با ستون های زیر را برای تراکنش های مالی بسازید. سپس با استفاده از کد پیوست تمرین، ۱.۵ میلیون رکورد را در آن وارد نمایید.

ستون	نوع داده	توضیح
transaction_id	serial PRIMARY KEY	شناسه منحصر به فرد تراکنش
transaction_date	date	تاریخ تراکنش
account_id	int	شناسه حساب مربوط به تراکنش
amount	numeric	مبلغ تراکنش
category	varchar(50)	دسته بندی تراکنش
status	varchar(20)	وضعیت تراکنش

(آ) در زیر چهار حالت partitioning به همراه ستون پیشنهادی برای آن روش آورده شده است. کوثری‌های لازم برای ساخت چهار جدول زیر را نوشته و اجرا نمایید تا داده‌ی جدول اصلی در این چهار جدول زیر قرار گیرد.

نام روش	نام جدول مقصد	ستون مورد استفاده
range	transaction_date	تاریخ تراکنش با تفکیک هر یک سال
hash	transactions_hash	شناسه حساب با تفکیک از طریق باقی مانده بر ۴
list	transactions_list	دسته بندی تراکنش ها
composite	transactions_composite	تاریخ تراکنش و وضعیت تراکنش

حال برای پرسش‌های زیر کوثری‌های مناسب برای اجرا بر روی این پنج جدول را نوشته و با استفاده از دستور EXPLAIN ANALYZE سرعت اجرا و عملکرد پیدا کردن آن توسط DBMS را در یک جدول مقایسه‌ای برای هر ۵ جدول فوق یادداشت کرده و تحلیل خود را شرح دهید؟ در هر سوال زیر کدام جدول سرعت پاسخگویی بهتری دارد؟ چرا؟

(ب) تعداد تراکنش‌هایی که در سال ۲۰۲۳ اتفاق افتاده است را برای جدول اصلی و هر جدول پارتیشن‌بندی شده بیابید.

(ج) جمع کل مبلغ تراکنش‌هایی که برای یک شناسه حساب خاص (مثلاً ۱۲۳) در جدول اصلی و هر جدول پارتیشن‌بندی شده اتفاق افتاده است را پیدا کنید.

(د) تعداد تراکنش‌ها برای هر دسته‌بندی را در جدول اصلی و هر جدول پارتیشن‌بندی شده بیابید.

(ه) تعداد تراکنش‌های تکمیل شده را در سال ۲۰۲۳ برای جدول اصلی و هر جدول پارتیشن‌بندی شده پیدا کنید.

(و) میانگین مبلغ تراکنش‌ها برای دسته‌بندی Groceries را در جدول اصلی و هر جدول پارتیشن‌بندی شده پیدا کنید.

(ز) میانگین مبلغ تمامی تراکنش‌ها را برای جدول اصلی و هر جدول پارتیشن‌بندی شده پیدا کنید.

(ح) یک index در جدول اصلی transaction بر روی ستون «تاریخ تراکنش» ایجاد نمایید. سپس کوثری‌هایی بنویسید که در جدول‌های transaction و transaction_range این پرسش را پاسخ دهد که «تعداد تراکنش‌هایی که در سال ۲۰۲۳ با مبلغ بیشتر از ۵۰۰۰ اتفاق افتاده است را برگرداند». سرعت رسیدن به جواب در این دو جدول را باهم مقایسه نمایید. چه نتیجه‌ای می‌گیرید؟