

Udiddit, a social news aggregator

Introduction

Udiddit, a social news aggregation, web content rating, and discussion website, is currently using a risky and unreliable Postgres database schema to store the forum posts, discussions, and votes made by their users about different topics.

The schema allows posts to be created by registered users on certain topics, and can include a URL or a text content. It also allows registered users to cast an upvote (like) or downvote (dislike) for any forum post that has been created. In addition to this, the schema also allows registered users to add comments on posts.

Here is the DDL used to create the schema:

```
CREATE TABLE bad_posts (  
    id SERIAL PRIMARY KEY,  
    topic VARCHAR(50),  
    username VARCHAR(50),  
    title VARCHAR(150),  
    url VARCHAR(4000) DEFAULT NULL,  
    text_content TEXT DEFAULT NULL,  
    upvotes TEXT,  
    downvotes TEXT  
);  
  
CREATE TABLE bad_comments (  
    id SERIAL PRIMARY KEY,  
    username VARCHAR(50),  
    post_id BIGINT,  
    text_content TEXT  
);
```

Part I: Investigate the existing schema

As a first step, investigate this schema and some of the sample data in the project's SQL workspace. Then, in your own words, outline three (3) specific things that could be improved about this schema. Don't hesitate to outline more if you want to stand out!

1. The most obvious fault with this schema is that it is very far from being normalized. This schema only has two tables, one for posts and one for comments. This makes it difficult to effectively query the database. I will break these tables up into different tables to normalize the data. Firstly, there needs to be a user table to track solely users, as this is a must. There also needs to be separate tables for posts, comments, votes, and topics. All of these will need their own unique identifier to make it easier to identify records.
2. One of the most important uses of such a database would be to track trends, such as how many upvotes or downvotes a certain topic or user gets. Storing the votes data as a text data type does not allow for this. I want to store this as an integer, so that aggregations can be done to find for example, what is the most upvoted topic at the moment, or put another way what is trending right now.
3. In order to keep these tables related and easily query-able with joins and such, the schema is going to need foreign keys in just about every table. As the project requires as well, these foreign keys will allow us to update information in another table if that record in one table is updated or deleted.
4. This schema will also need constraints and a fair amount of them to keep data integrity in place as well as referential integrity between the tables.
5. A more minor thing to look at is also going to be evaluating the limits set in the VARCHAR data types. Some of these may need to change, in order to best meet the application needs.

Part II: Create the DDL for your new schema

Having done this initial investigation and assessment, your next goal is to dive deep into the heart of the problem and create a new schema for Udiddit. Your new schema should at least reflect fixes to the shortcomings you pointed to in the previous exercise. To help you create the new schema, a few guidelines are provided to you:

1. Guideline #1: here is a list of features and specifications that Udiddit needs in order to support its website and administrative interface:
 - a. Allow new users to register:
 - i. Each username has to be unique
 - ii. Usernames can be composed of at most 25 characters
 - iii. Usernames can't be empty
 - iv. We won't worry about user passwords for this project
 - b. Allow registered users to create new topics:
 - i. Topic names have to be unique.
 - ii. The topic's name is at most 30 characters
 - iii. The topic's name can't be empty
 - iv. Topics can have an optional description of at most 500 characters.
 - c. Allow registered users to create new posts on existing topics:
 - i. Posts have a required title of at most 100 characters
 - ii. The title of a post can't be empty.
 - iii. Posts should contain either a URL or a text content, **but not both**.
 - iv. If a topic gets deleted, all the posts associated with it should be automatically deleted too.
 - v. If the user who created the post gets deleted, then the post will remain, but it will become dissociated from that user.
 - d. Allow registered users to comment on existing posts:
 - i. A comment's text content can't be empty.
 - ii. Contrary to the current linear comments, the new structure should allow comment threads at arbitrary levels.
 - iii. If a post gets deleted, all comments associated with it should be automatically deleted too.
 - iv. If the user who created the comment gets deleted, then the comment will remain, but it will become dissociated from that user.
 - v. If a comment gets deleted, then all its descendants in the thread structure should be automatically deleted too.
 - e. Make sure that a given user can only vote once on a given post:
 - i. Hint: you can store the (up/down) value of the vote as the values 1 and -1 respectively.
 - ii. If the user who cast a vote gets deleted, then all their votes will remain, but will become dissociated from the user.

- iii. If a post gets deleted, then all the votes for that post should be automatically deleted too.
2. Guideline #2: here is a list of queries that Uddidit needs in order to support its website and administrative interface. Note that you don't need to produce the DQL for those queries: they are only provided to guide the design of your new database schema.
- a. List all users who haven't logged in in the last year.
 - b. List all users who haven't created any post.
 - c. Find a user by their username.
 - d. List all topics that don't have any posts.
 - e. Find a topic by its name.
 - f. List the latest 20 posts for a given topic.
 - g. List the latest 20 posts made by a given user.
 - h. Find all posts that link to a specific URL, for moderation purposes.
 - i. List all the top-level comments (those that don't have a parent comment) for a given post.
 - j. List all the direct children of a parent comment.
 - k. List the latest 20 comments made by a given user.
 - l. Compute the score of a post, defined as the difference between the number of upvotes and the number of downvotes
3. Guideline #3: you'll need to use normalization, various constraints, as well as indexes in your new database schema. You should use named constraints and indexes to make your schema cleaner.
4. Guideline #4: your new database schema will be composed of five (5) tables that should have an auto-incrementing id as their primary key.

Once you've taken the time to think about your new schema, write the DDL for it in the space provided here:

```
CREATE TABLE "users" (  
  "id" SERIAL PRIMARY KEY,  
  "username" VARCHAR(25) CONSTRAINT "unique_username" UNIQUE NOT NULL,  
  "last_login" TIMESTAMP,  
  CONSTRAINT "no_empty_usernames" CHECK (LENGTH(TRIM("username")) > 0)  
);  
CREATE INDEX "find_user" ON "users" ("username");  
  
CREATE TABLE "topics" (  
  "id" SERIAL PRIMARY KEY,
```

```

        "topic" VARCHAR(30) CONSTRAINT "unique_topic" NOT NULL,
        "topic_description" VARCHAR(500)
        CONSTRAINT "no_empty_topics" CHECK (LENGTH(TRIM("topic")) > 0)
    );
CREATE INDEX "topic_search" ON "topics" ("topic");

CREATE TABLE "posts" (
    "id" SERIAL PRIMARY KEY,
    "user_id" INTEGER NOT NULL,
    "topic_id" INTEGER NOT NULL,
    "post_timestamp" TIMESTAMP,
    "title" VARCHAR(100) CONSTRAINT "title_not_null" NOT NULL,
    "url" VARCHAR(2000),
    "text_content" VARCHAR,
    CONSTRAINT "user_id_fk" FOREIGN KEY ("user_id") REFERENCES "users" ("id")
    ON DELETE SET NULL,
    CONSTRAINT "topic_id_fk" FOREIGN KEY ("topic_id") REFERENCES "topics" ("id")
    ON DELETE CASCADE,
    CONSTRAINT "no_empty_titles" CHECK (LENGTH(TRIM("title")) > 0),
    CONSTRAINT "text_url" CHECK (
        "text_content" IS NOT NULL
        AND "url" IS NULL
        OR "text_content" IS NULL
        AND "url" IS NOT NULL)
);
CREATE INDEX "post_url" ON "posts" ("url");

CREATE TABLE "comments" (
    "id" SERIAL PRIMARY KEY,
    "user_id" INTEGER NOT NULL,
    "post_id" INTEGER NOT NULL,
    "comment" VARCHAR CONSTRAINT "comment_not_null" NOT NULL,
    "comment_timestamp" TIMESTAMP,
    "level_id" INTEGER,
    CONSTRAINT "no_empty_comments" CHECK (LENGTH(TRIM("comment")) > 0),
    CONSTRAINT "user_id_fk" FOREIGN KEY ("user_id") REFERENCES "users" ("id")
    ON DELETE SET NULL,
    CONSTRAINT "post_id_fk" FOREIGN KEY ("post_id") REFERENCES "posts" ("id")
    ON DELETE CASCADE,
    CONSTRAINT "level_id_fk" FOREIGN KEY ("level_id") REFERENCES "comments"
    ON DELETE CASCADE
);
CREATE INDEX "level_index" ON "comments" ("level_id");
CREATE INDEX "user_comments" ON "comments" ("user_id", "comment_timestamp");

CREATE TABLE "votes" (
    "id" SERIAL PRIMARY KEY,

```

```
    "user_id" INTEGER NOT NULL,  
    "post_id" INTEGER NOT NULL,  
    "vote" INTEGER CONSTRAINT "vote_type"  
CHECK("vote" = -1 OR "vote" = 1),  
    CONSTRAINT "user_id_fk" FOREIGN KEY ("user_id") REFERENCES "users" ("id")  
ON DELETE SET NULL,  
    CONSTRAINT "post_id_fk" FOREIGN KEY ("post_id") REFERENCES "posts" ("id")  
ON DELETE CASCADE  
);  
CREATE INDEX "score" ON "votes" ("vote");
```

Part III: Migrate the provided data

Now that your new schema is created, it's time to migrate the data from the provided schema in the project's SQL Workspace to your own schema. This will allow you to review some DML and DQL concepts, as you'll be using INSERT...SELECT queries to do so. Here are a few guidelines to help you in this process:

1. Topic descriptions can all be empty
2. Since the bad_comments table doesn't have the threading feature, you can migrate all comments as top-level comments, i.e. without a parent
3. You can use the Postgres string function **regexp_split_to_table** to unwind the comma-separated votes values into separate rows
4. Don't forget that some users only vote or comment, and haven't created any posts. You'll have to create those users too.
5. The order of your migrations matter! For example, since posts depend on users and topics, you'll have to migrate the latter first.
6. Tip: You can start by running only SELECTs to fine-tune your queries, and use a LIMIT to avoid large data sets. Once you know you have the correct query, you can then run your full INSERT...SELECT query.
7. **NOTE:** The data in your SQL Workspace contains thousands of posts and comments. The DML queries may take at least 10-15 seconds to run.

Write the DML to migrate the current data in bad_posts and bad_comments to your new database schema:

```
INSERT INTO "users" ("username") (  
  SELECT DISTINCT "username"  
  FROM "bad_posts"  
  UNION  
  SELECT DISTINCT "username"  
  FROM "bad_comments"  
  UNION  
  SELECT DISTINCT REGEXP_SPLIT_TO_TABLE(upvotes, ',') as user_names  
  FROM bad_posts  
  UNION  
  SELECT DISTINCT REGEXP_SPLIT_TO_TABLE(downvotes, ',') as user_names  
  FROM bad_posts  
);  
  
INSERT INTO "topics" ("topic") (  
  SELECT DISTINCT "topic"  
  FROM "bad_posts"
```

```

SELECT DISTINCT "topic"
FROM bad_posts
);

INSERT INTO "posts" ("user_id", "topic_id", "title", "url",
"text_content")
SELECT "u"."id", "t"."id", LEFT("bp"."title", 100), "bp"."url",
"bp"."text_content"
FROM users u
JOIN bad_posts bp
ON "u"."username" = "bp"."username"
JOIN "topics" t
ON "bp"."topic" = "t"."topic";

INSERT INTO "comments" ("user_id", "post_id", "comment")
SELECT "u"."id", "p"."id", "bc"."text_content"
FROM bad_comments bc
JOIN users u
ON "u"."username" = "bc"."username"
JOIN "posts" p
ON "p"."id" = "bc"."post_id";

INSERT INTO "votes" ("post_id", "user_id", "vote")
SELECT "bp"."id", "u"."id", 1 AS vote_count
FROM (SELECT id, REGEXP_SPLIT_TO_TABLE(upvotes, ',') AS upvotes
FROM bad_posts
) bp
JOIN "users" u ON "u"."username" = "bp"."upvotes";

INSERT INTO "votes" ("post_id", "user_id", "vote")
SELECT "bp"."id", "u"."id", -1 AS vote_count
FROM (SELECT id, REGEXP_SPLIT_TO_TABLE(downvotes, ',') AS downvotes
FROM bad_posts
) bp
JOIN "users" u ON "u"."username" = "bp"."downvotes";

DROP TABLE bad_comments;
DROP TABLE bad_posts;

```