# Chapter 1: Introducing The Project

In this chapter we're going to get into the details of the project that will help you learn the SQL Patterns. As you saw in the introduction, we're using a real-world, public dataset from StackOverflow.

StackOverflow is a popular website where users can post technical questions about any technical topic and others can post answers to these questions. They can also vote on the answers or comment on them.

Based on the quality of the answers, users gain reputation and badges which they can use as social proof both on the SO site and on other websites.

Using this dataset we're going to build a table that calculates reputation metrics for every user. This type of table is sometimes called a "feature table" and can be used in other applications in data science and analytics. You simply replace the user_id with a customer id or any other entity.

Since the query to build it is complex, it's the perfect tool to illustrate some of the patterns described in this book.

The schema of what it would look something like this:

```
 1 | column_name             | type     |
 2 |-------------------------|----------|
 3 | user_id                 | integer  |
 4 | user_name               | string   |
 5 | total_posts_created     | numeric  |
 6 | total_answers_created   | numeric  |
 7 | total_answers_edited    | numeric  |
 8 | total_questions_created | numeric  |
 9 | total_upvotes           | numeric  |
10 | total_comments_by_user  | numeric  |
11 | total_questions_edited  | numeric  |
12 | streak_in_days          | numeric  |
13 | total_comments_on_post  | numeric  |
14 | posts_per_day           | numeric  |
15 | edits_per_day           | numeric  |
```

```
16  │  answers_per_day         │  numeric  │
17  │  questions_per_day       │  numeric  │
18  │  comments_by_user_per_day │  numeric  │
19  │  answers_per_post        │  numeric  │
20  │  questions_per_post      │  numeric  │
21  │  upvotes_per_post        │  numeric  │
22  │  downvotes_per_post      │  numeric  │
23  │  user_comments_per_post  │  numeric  │
24  │  comments_on_post_per_post │  numeric  │
```

As you can see, we need to transform the source data model to a new model that has one row per `user_id` Before we do that, we need to understand the source data first.


## Understanding the Data Model

Writing accurate and efficient SQL begins with understanding the data model we're starting with. This may already exist in the form of documentation and diagrams but more often than not you'll have to learn it as you go.

The original StackOverflow (SO) data model is different from the one loaded in BigQuery. When the engineers loaded it, they modified the mode somewhat. For example the SO model contains a single `Posts` table for all the different post types whereas BigQuery split each one into a separate table.

Here's a look at the Entity-Relationship (ER) diagram [[StackOverflow BQ ER Diagram.jpeg]] **Figure 1.1 - StackOverflow ER diagram**

There are 8 tables that represent the various post types. You can get this result by using the `INFORMATION_SCHEMA` views in BigQuery like this:

```
1  SELECT table_name
2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
     .TABLES
3  WHERE table_name like 'posts_%'
4
```

```
 5  |table_name              |
 6  |------------------------|
 7  |posts_answers           |
 8  |posts_orphaned_tag_wiki |
 9  |posts_tag_wiki          |
10  |posts_questions         |
11  |posts_tag_wiki_excerpt  |
12  |posts_wiki_placeholder  |
13  |posts_privilege_wiki    |
14  |posts_moderator_nomination|
```

We'll be focusing on just two of them for our project so I've left the other ones out:

1. `posts_questions` contains all the question posts
2. `posts_answers` contains all the answer posts

They both have the same schema:

```
 1  SELECT column_name, data_type
 2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
       .COLUMNS
 3  WHERE table_name = 'posts_answers'
 4
 5  |column_name             |data_type|
 6  |------------------------|---------|
 7  |id                      |INT64    |
 8  |title                   |STRING   |
 9  |body                    |STRING   |
10  |accepted_answer_id      |STRING   |
11  |answer_count            |STRING   |
12  |comment_count           |INT64    |
13  |community_owned_date    |TIMESTAMP|
14  |creation_date           |TIMESTAMP|
15  |favorite_count          |STRING   |
16  |last_activity_date      |TIMESTAMP|
17  |last_edit_date          |TIMESTAMP|
18  |last_editor_display_name|STRING   |
19  |last_editor_user_id     |INT64    |
20  |owner_display_name      |STRING   |
21  |owner_user_id           |INT64    |
```

```
22  |parent_id               |INT64    |
23  |post_type_id            |INT64    |
24  |score                   |INT64    |
25  |tags                    |STRING   |
26  |view_count              |STRING   |
```

Both tables have an `id` column that identifies a single post, `creation_date` that identifies the timestamp when the post was created and a few other attributes like `score` for the upvotes and downvotes.

Note the `parent_id` column which signifies a hierarchical structure. The `parent_id` is a one-to-many relationship that links up an answer to the corresponding question. A single question can have multiple answers but an answer belongs to one and only one question. This is relation 1 in the **Figure 1.1** above.

Both post types (question and answer) have a one-to-many relationship to the `post_history`. These are relations 3 and 4 in the diagram above.

```
 1  SELECT column_name, data_type
 2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
       .COLUMNS
 3  WHERE table_name = 'post_history'
 4
 5  |column_name         |data_type|
 6  |--------------------|---------|
 7  |id                  |INT64    |
 8  |creation_date       |TIMESTAMP|
 9  |post_id             |INT64    |
10  |post_history_type_id|INT64    |
11  |revision_guid       |STRING   |
12  |user_id             |INT64    |
13  |text                |STRING   |
14  |comment             |STRING   |
```

A single post can have many types of activities identified by the `post_history_type_id` column. This id indicates the different types of activities a user can perform on the site. We're only concerned with the first

6. You can see the rest of them here if you're curious.

1. Initial Title - initial title *(questions only)*
2. Initial Body - initial post *(raw body text)*
3. Initial Tags - initial list of tags *(questions only)*
4. Edit Title - modified title *(questions only)*
5. Edit Body - modified post body *(raw markdown)*
6. Edit Tags - modified list of tags *(questions only)*

The first 3 indicate when a post is first submitted and the next 3 when a post is edited.

The `post_history` table also connects to the `users` table via the `user_id` in a one-to-many relationship shown in the diagram as number 6. A single user can perform multiple activities on a post.

In database lingo this is known as a bridge table because it connects two tables (user and posts) that have a many-to-many relationship which cannot be modeled otherwise.

The `users` table has one row per user and contains user attributes such as name, reputation, etc. We'll use some of these attributes in our final table.

```
 1  SELECT column_name, data_type
 2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
       .COLUMNS
 3  WHERE table_name = 'users'
 4
 5  |column_name       |data_type|
 6  |------------------|---------|
 7  |id                |INT64    |
 8  |display_name      |STRING   |
 9  |about_me          |STRING   |
10  |age               |STRING   |
11  |creation_date     |TIMESTAMP|
12  |last_access_date  |TIMESTAMP|
13  |location          |STRING   |
```

```
14  |reputation       |INT64   |
15  |up_votes         |INT64   |
16  |down_votes       |INT64   |
17  |views            |INT64   |
18  |profile_image_url|STRING  |
19  |website_url      |STRING  |
```

Next we take a look at the comments table. It has a zero-to-many relationship with posts and with users shown in the diagram as number 5 and number 7, since both a user or a post could have 0 or many comments.

```
1  SELECT column_name, data_type
2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
      .COLUMNS
3  WHERE table_name = 'comments'
4
5  |column_name      |data_type|
6  |-----------------|---------|
7  |id               |INT64    |
8  |text             |STRING   |
9  |creation_date    |TIMESTAMP|
10 |post_id          |INT64    |
11 |user_id          |INT64    |
12 |user_display_name|STRING   |
13 |score            |INT64    |
```

Finally the votes table represents the upvotes and downvotes on a post. We'll need this to compute the total vote count on a user's post which will indicate how good the question or the answer is. This table has a granularity of one row per vote per post per date.

```
1  SELECT column_name, data_type
2  FROM bigquery-public-data.stackoverflow.INFORMATION_SCHEMA
      .COLUMNS
3  WHERE table_name = 'votes'
4
5  |column_name  |data_type|
6  |-------------|---------|
7  |id           |INT64    |
```

```
 8  |creation_date|TIMESTAMP|
 9  |post_id      |INT64    |
10  |vote_type_id |INT64    |
```

The `votes` table is connected to a post in a 0-to-many relationship shows in the diagram as number 2. In order for us to get upvotes and downvotes on a user's post, we'll need to join it with the `users` table.

Alright, now that we've familiarized ourselves with the source data model, next let's look at some core concepts.