### Dataset used:

For this assignment's study, we will be dealing with the complete version of Schlumberger dataset. The dataset consists of time-series data from downhole equipment using intelligent completions which helps to evaluate and manage production in real-time with zonal downhole monitoring of pressures and temperatures, easing the maintenance of such costly equipment.

### Cursory Analysis:

After downloading and opening the *data-large.csv* file from the GitHub repository as specified, we can immediately hint at a few probable errors we might face while importing the dataset into python with a script. Let's take a quick look. The immediate red flag I could notice right after opening the file is the issue with headers. This is a very common error; we can see that there are multiple header rows (3) in our file. The first row signifies if the equipment is an Electric Submersible Pump (ESP) or Intelligent Completion (IC). The second row indicates the type of sensor value recorded like motor temperature, choke pressure, VSDFreqOut, etc. The third row of the header shows the sensor's information with its unit of measure like psia, K or Hz etc.

These redundant headers can be removed and merged into a single header by eliminating the sensor name and the unit rows. The second-row names are renamed in an efficient manner to accurately depict the columns names by including the required information. Correspondingly, the first column contains information about the time in seconds. It can be pointed out that many columns are redundant, and our table should be free from these identical and indistinguishable column names. We can employ some modified arguments while importing and reading in the CSV file to resolve some issues. For example, since the dataset is very large, we can *set low_memory=False* minimalizing the storage consumption. Having *index_col=False* enforces the 'pandas' module not to use the first column as the index. In addition to these, explicitly stated the object type *dtype='unicode'* for it to return Unicode normal form.

We can observe an empty cell in the first row of the first column, which could result in an improper import as well. This cell is also removed as part of the data-cleaning process. The structure of the CSV file is absurd which results in an improper import if we read the file using a script in python. Thus, the data is organized and manipulated using Python's Panda module. To retain the values in the original data frame, the index is reset with the drop argument set to true. In addition, specifying dtype for columns can reduce storage requirements when storing the data frame in the database.

### PHPMyAdmin & Possible Alternatives:

In many cases, MySQL is the most commonly used open-source RDBMS, due to its effectiveness, reliability, and simplicity of use - especially with the right tools. As a result, a variety of professional solutions are available for MySQL and MariaDB, a close relative of MySQL. Most of them are free or paid, more or less functional, and are created by various companies. The phpMyAdmin database management tool is well known to everyone who works with MySQL. As the name implies it's written in PHP, so it can be installed on your web server and accessed from a browser. The interface is multilingual and has 79 different language versions.

While it is a popular tool with a wide range of features, it cannot cover a few requirements. For example, phpMyAdmin does not provide scheme visualizations, full auto-completion capabilities, scheduled backup, and encryption. It is possible to export an unlimited number
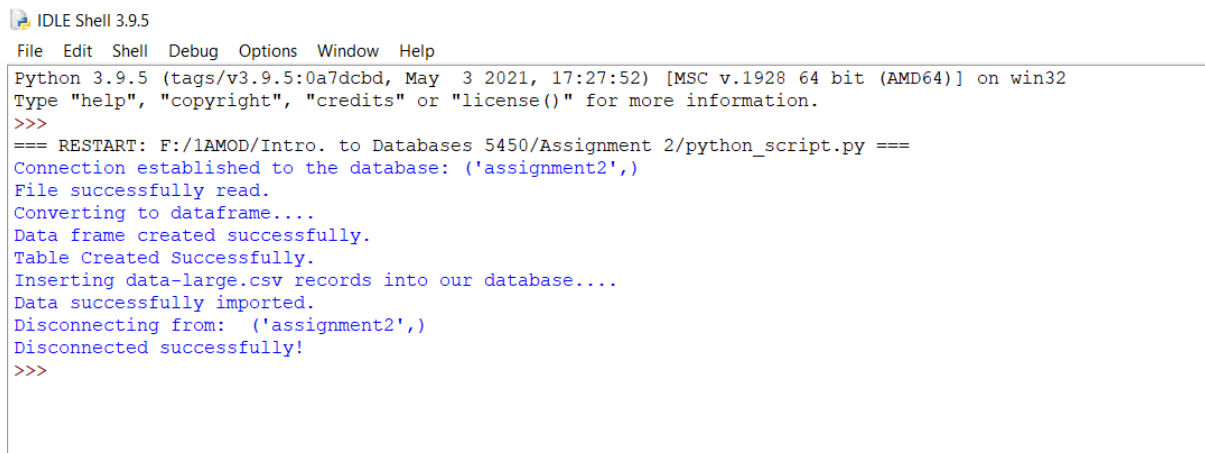
of records to SQL, CSV, or spreadsheets using PHPMyAdmin. The software is incapable of managing so much data when displaying a good interactive grid with visualizations as the contents of a table are injected into an object during rendering, and all of the contents are stored in memory. Trying to import a large database in MySQL (phpMyAdmin) with a file size of more than 20MB, the uploaded file will be rejected with an error. The problem occurs due to uploading a large file, hence "post max time" or "execution time" errors occur in PHP. These errors can be handled by doing minimal changes to XAMPP configuration. Without altering the XAMPP configuration, its possible to import large datasets by editing the php.ini file. For instance, PHP is unable to save information in a variable if 10,000 rows total more than 128M. That indicates that the script used more memory than the predetermined 128M limit. As a result, any significant growth in php.ini should be above 128M. Your next move should be to 256M. In order for the modifications to take effect, remember to restart Apache and XAMPP.

Noting this; as the current dataset deals with indistinguishable column names, varied data types and absurd structure, it would be vigilant to use python as an alternative to import our dataset and run queries due to the fact that it can handle a high volume of memory and data processing becomes a lot easier with the implementation of packages such as pandas and other great numbers of data-oriented feature packages that can speed up and simplify data processing, in turn saving a lot of time.

## Python Script & Testing:

**File attached.**

**Python O/P Screenshot:**



**PHPMyAdmin Screenshot:**

## Database Normal Form:

A functional dependency is a constraint between two attributes in a relation. Some of the different functional dependencies in Schlumberger dataset:

- The ESP motor speed and frequency depend on motor temperature.
- Additionally, ESP is also dependent on discharge and pump intake pressures.
- The temperature 2 characteristic specified in the dataset, which is measured in F and depends on the downhole temperature reported by the intelligent control system.
- As pressure 1 and pressure 2 (psia) parameters are subtracted to create delta pressure, the Intelligent Control System Liquid Rate (bbl/d) is dependent on these attributes.
- Also, the Choke position can be determined on pressure 1 and pressure 2.
- The Water rate (bbl/d) depends on Liquid rate and water cut and is defined as (Liquid Rate*Water Cut).
- The Oil rate depends on Liquid rate and water rate and is defined as (Liquid Rate - Water Rate).

### Is the database table in 1NF, 2NF, or 3NF?

The database table is already in 1NF as there are no multi-valued characteristic records. With a unique column timestamp acting as the primary key in 1NF, there should be no partial dependencies for it to be in 2NF. So, the table is in 2NF by default. In order to convert it to 3NF, we can separate timestamp-IC, timestamp-EC values into separate tables along with other functional dependencies.
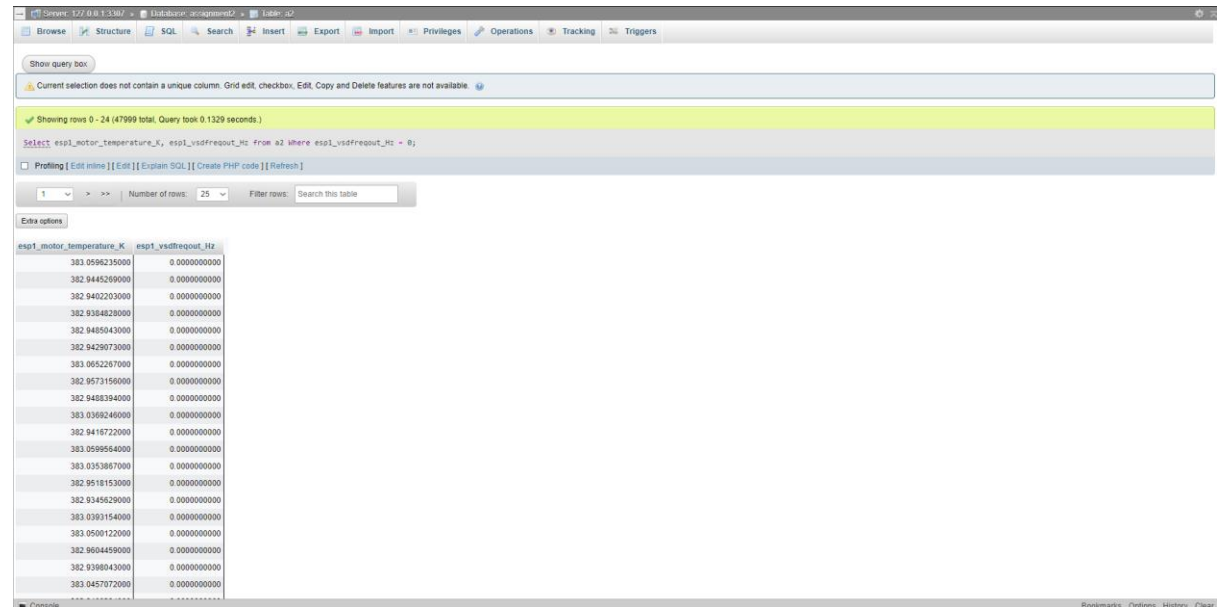
## Useful Analytical Queries:

The explanation of the queries used here are self-explanatory. These analytical queries can be properly used to monitor and maintain the equipment as the consequences are significant.

1. How is the motor temperature of ESP01 when the variable speed driver output frequency is 0 hz ?

**Query:**

*SELECT esp1_motor_temperature_K, esp1_vsdfreqout_Hz from a2*
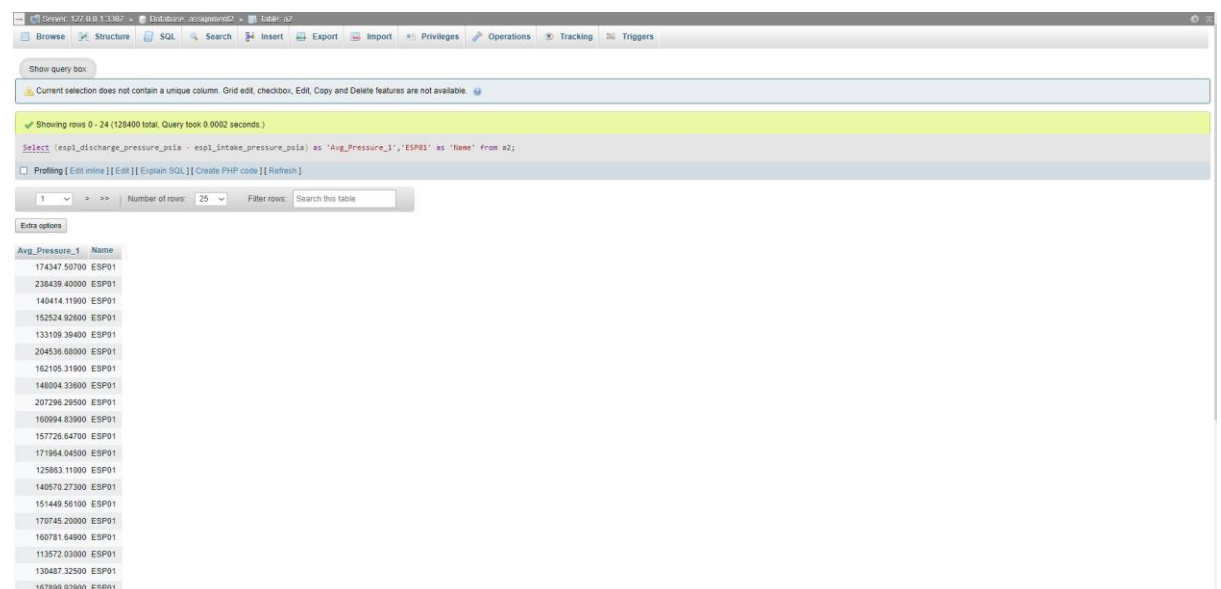
*WHERE esp1_vsdfreqout_Hz = 0;*

**Output:**



2. Calculate the average change in pressure for pump 1.

**Query:**

**SELECT (esp1_discharge_pressure_psia - esp1_intake_pressure_psia) as 'Avg_Pressure_1','ESP01' as 'Name' from a2;**
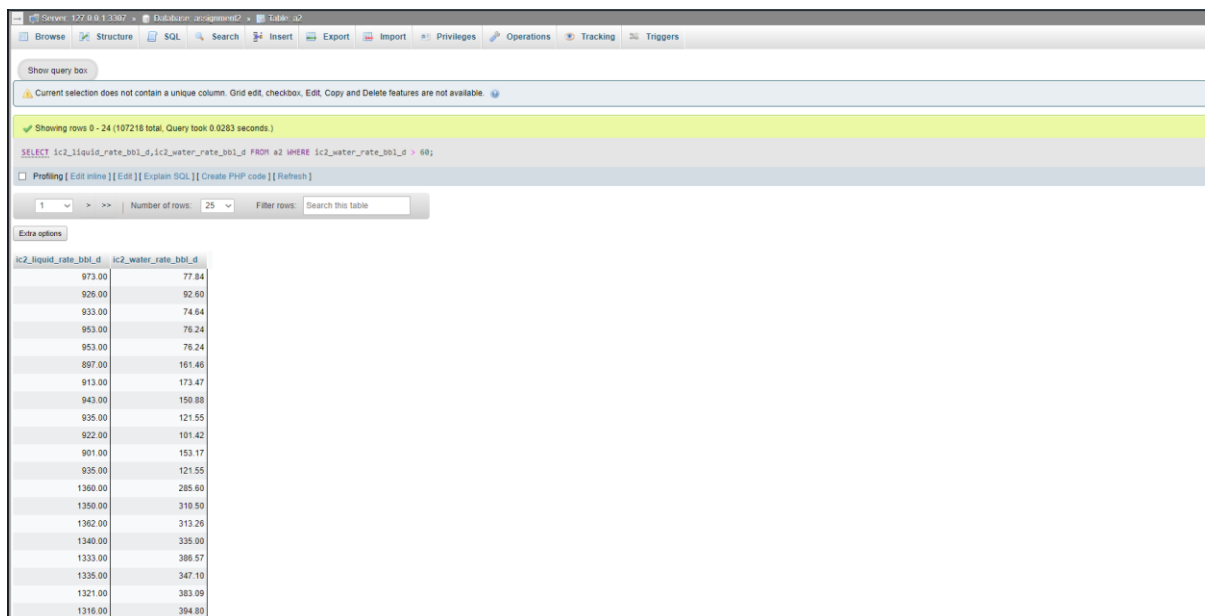
**Output:**

3. The volume flow of water generated is known as the water rate, while the volume flow of liquid produced is known as the liquid rate. Let's examine the liquid rate at a water rate of more than 60.

**Query:**

**SELECT ic2_liquid_rate_bbl_d,ic2_water_rate_bbl_d FROM a2 WHERE ic2_water_rate_bbl_d > 60;**

**Output:**



---