# DATA WAREHOUSE MODELLING

# -

# Data Integration Using ETL of Production Data on Oil Wells in the State of New Mexico

Submitted by:          Team 1 Group O2

Group Members:          Rajat Prasad
                        Carolin Kröger
                        Sara Krumpak
                        Mohamad el Charif
                        Andrea Salvati
                        Dominik Thausing
                        Alejandro Torres Pérez

Date:                   November 25th, 2018

# Table of Contents

# Table of Tables

# Table of Figures

# Executive Summary

The following report describes the Data integration process using Extract, Transformation and Loading. The First step taken was to make necessary changes to the schema based on the recommendation received. Some of the changes were like merging the "d_county_id" and the "d_Area". Also, some columns that were proposed in the fact table were dropped as they could only be obtained by special mathematical transformation which is not a part of ETL process. After changes made to dimensions and the fact table, new foreign keys were chosen which were basically primary key of dimensions. They were created because of auto sequence. All these primary keys in dimension table represent a unique column or a set of unique columns. These could only be populated when we run the loading process on Pentaho. After fixing the bugs in MySQL workbench file, A strategy was developed to understand the data and fix the anomalies. Since we have the information that what columns / attributes belong to which dimension and which of these attributes should be unique and not null. A set of transformations were created to create dimensions and load them in the database. A final of transformation was created to load data in fact table.

However, the first step undertaken was to create a transformation called "TR_Cleaning" where we are deleting the rows with NULL values in key columns and sorting and removing their duplicates. For some of the columns Null values were replaces by 0. We made sure that source of data was added after every transformation, not only source by system log time and last modified time of transformation as well, for creating meta data and to resolve any data governance issues.

As the business use of case of the data is to analyze monthly oil gas and water of each well therefore the fact table contains these crucial facts along with the information of the entry. The fact table contains reference to all dimensions such as "d_operator" which contains information about operators. Thus, reducing redundancy and fast processing of queries and providing a more organized structure.

# Requirements

Before proceeding in the explanation of the ETL process, it is pivotal to check 10 important requirements that the system needs to meet in order to be effective.

1. Business Needs

The business objective of the ETL processes is to give managers of oil wells in the state of New Mexico resources to perform data driven decisions, clustered into three directions of management:

- recognizing underperforming wells
- optimizing wells performance
- understanding where investments should be focused.

2. Compliance

We have treated the data in our model in an accurate and complete way. Moreover, we confirm that data were not tampered and that the model represents the real state of the business. For this reason our ETL process should meet the legal requirement of the sector. However, we suggest a deeper analysis from the legal department of the company.

3. Data Quality

The initial dataset presented some distortions and it was not ready to use for an ETL process. This is the reason why the first transformation of the process is about the cleaning of the data (see below for the detailed report of the cleaning)

4. Security

In order to maximize the security of the process we have considered mainly 2 measures:

- Physical backup of all the ETL processes and the MySql schema
- Report with carefully explained instruction to perform all the transformations

5. Data integration

We have designed our ETL  model with just one database so we have not included common dimensional attributes across different databases and we have not created

common business matrix. However, if in the future, the company decides to implement a second database the steps above will be required.

6. Data Latency

The dataset that we used for the ETL process is composed by 1,000,000 rows. This is why the ETL processes takes a considerable amount of time to be performed. (more details will be discussed in the Fact Table's transformation description)

7. Archiving and Lineage

In order to reach the best performance in terms of Archiving and Lineage the following actions have been taken:

- Backup of the output of the cleaning transformation, being the most important step of the process;
- Back up of every transformation in a separate file;
- Metadata details have been added to every transformation.

8. BI delivery Interfaces

The delivery platform for our ETL process is MySql Workbench. Every transformation has been uploaded to the schema designed in MySql Workbench and can be easily extract through the platform.

9. Available Skills

No programming skills were required in order to build the ETL processes. The skills that are needed to perform this task are related to the use of the PDI software and basic sql knowledge.

10. Legacy Licenses

The ETL process has been build using Pentaho Data Integration and performed using My Sql Workbench.

# Instructions for Executing the ETL process

We have all the input files (csv format) in the input folder present in the GA2 folder. We also have all the transformation and the main job in the GA2 folder.

1. **Change the username and password while connecting to database in the "input/update" step of each transformation while making database connection.**

2. We shall now open a job called "JB_WELL" present in the GA2 folder. There is no need to make any change. Simply run the job.

3. You can check the loading output of all the transformation in the MySql workbench file called "well_status.mwb" also present in the GA2 folder.

## Timing of the ETL process

| transformation | Time |
|---|---|
| Cleaning | 0,3s |
| Status | 1,6s |
| Company | 5,1s |
| Operator | 5,3s |
| Date | 9,6s |
| Area | 19,8m |
| Fact_Table | NA (155,000 rows in 5.30h) |

*Table 1 -Timing of the ETL process*

## Transformations overview

1. **TR_Cleaning.ktr**.

    There are two sub-transformations in this file. The first transformation deals with "wel_mexico_all_new_format.csv" and the second deals with "production_mexico_all_new_format.csv". In the first step of the first sub-transformation the filename is

${Internal.Entry.Current.Directory}/input/wel_mexico_all_new_format.csv. In the last step of the first sub transformation called "output well" the output folder is ${Internal.Entry.Current.Directory}/output/dataset_well. Now let's visit the second sub-transformation, the filename is "${Internal.Entry.Current.Directory}/input/production_mexico_all_new_format.csv" and the output folder is ${Internal.Entry.Current.Directory}/output/dataset_production.

2. **TR_DATE**

   the filename is "${Internal.Entry.Current.Directory}/output/dataset_production.csv" in "csv file input production" and in the last step of the transformation called "insert/update" already connection called "localhost" has been established. The target schema is "well_status" and the target table is "d__date".

3. **TR_company**

   the filename is "${Internal.Entry.Current.Directory}/output/dataset_well.csv " in csv file input well" and in the last step of the transformation called "insert/update" already connection called "localhost" has been established. The target schema is "well_status" and the target table is "d___company".

4. **TR_operator**

   The filename is "${Internal.Entry.Current.Directory}/output/dataset_well.csv " in csv file input well" and in the last step of the transformation called "insert/update" already connection called "localhost" has been established. The target schema is "well_status" and the target table is "d_Operator". Click ok save and run the transformation.

5. **TR_status**

   The filename is "${Internal.Entry.Current.Directory}/output/dataset_well.csv " in csv file input well" and in the last step of the transformation called "insert/update" already connection called "localhost" has been established. The target schema is "well_status" and the target table is "d_status".

18. **TR_Area**

   The filename is "${Internal.Entry.Current.Directory}/output/dataset_well.csv " in csv file input well" and in the last step of the transformation called "insert/update already connection called "localhost" has been established. The target schema is "well_status" and the target table is "d_area".

6. **TR_FactTable**

There are two CSV file input stage called "production" and "well". In production the filename is "${Internal.Entry.Current.Directory}/output/dataset_production.csv ".In WELL the filename is "${Internal.Entry.Current.Directory}/output/dataset_well.csv" In "Insert/Update" command set target table to "f_well".

# Implementation of Feedback



*Figure 1 - Previously proposed Schema*

In the previous schema we had 6 dimension tables called current_status, d_date, d_area, d_company, d_operator, d_county_id.

The fact table also contained a lot of variables which are obtained after complex mathematical calculations. We decided to drop as they are not needed at this stage of Data warehousing.

A new schema was developed taking into consideration the feedback and Business use case.



*Figure 2 - Improved new Schema*

Although the previous schema had major problems (like containing derived values which are not suitable at this stage of data warehousing), it had the design that was in conformance with the business need (the star schema). We decided to get rid of derived columns. We merged the "area" and "county" dimensions. We analyzed each variable in the dataset and calculated their max length so that in dimensions we have appropriate size for each attribute. We fixed the same for all dimensions according to standard. We came up with all these dimensions

after thoroughly analyzing the need of the stakeholders. How will they benefit with these data mart.

# Data Integration using Pentaho for ETL

We were provided data in two json files, one containing We manipulated these files to be read by dataiku which converted it into a csv file. These CSV files were then saved in Input folder present in the GA2 folder

Our first transformation in PDI was data cleaning- taking out null and duplicate values. Or reassigning a value to a Null value. The outputs of these transformation were then saved as a csv file in the output folder which is also present in the GA2 folder.

information about all the wells in New Mexico, United States and the other contained monthly production data of each well.

## Overview Of All Steps



*Figure 3 - Overview of all steps*

We have 7 transformation in total, which operates on the two csv files obtained as a result of a series of transformation on json files that we were given in the first place.

The first transformation is **TR_Cleaning**, it takes care of the null values, it also removes duplicate values and the output is fed to a new CSV file called dataset_prodcution and dataset_well.

The second Transformation was to create a dimension for **dates** as shown in the ER diagram of MYSQL workbench. We inserted the new processed CSV file called dataset_ production as input file and extracted all unique dates in dimension d_date.

The third Transformation was to create a dimension for **Company** as shown in the ER diagram of MYSQL workbench. We inserted the new processed CSV file called dataset_ well as input file and extracted all unique Companies and feed into dimension called d_company.

The fourth Transformation was to create a dimension for Operator as shown in the ER diagram of MYSQL workbench. We inserted the new processed CSV file called dataset_ well as input file and extracted all unique Operators and feed into dimension called d_Operator.

The fifth Transformation was to create a dimension for status as shown in the ER diagram of MYSQL workbench. We inserted the new processed CSV file called dataset_well as input and extracted all unique Operators and feed into dimension called d_Company.

The sixth Transformation was to create a dimension for Area as shown in the ER diagram of MYSQL workbench. We inserted the new processed CSV file called dataset_ well as input and extracted all unique combination of area, state and unique well identifier and feed all the relevant data in dimension d_area.

The seventh Transformation was to fill the fact table with columns important to business needs like oil, gas and water production along with source of the data and the log time, not only these but primary key of other five dimensions were also inserted as foreign keys as shown in the ER diagram of MYSQL workbench. We inserted the two new processed CSV files as input called dataset_prodcution and dataset_ well and extracted anecessary features like api_number, monthly oil, gas, water production and feed it in table called f_well present in schema called well_status.

# Overview of individual steps

## I STEP. TR_Cleaning



*Figure 4 - Cleaning Transformation*

In this transformation we cleaned two CSV files, But before we indulge into getting to know the working of this transformation, It should be duly noted that as input two json files were provided. These files were manipulated to enable them to be converted into CSV files. After converting them to CSV files they were used as input in this transformation.

As already mentioned we cleaned two files and therefore we have two sets of transformation. The first set of SUB-Transformation deals with cleaning the wel_mexico_all_new_format.csv, which contains attributes called API Number, latitude, longitude, operator name, well name, status, parent, CIK.

The first step was to read the input file and sort it according to API number, the second step was to remove the duplicate API number because one well cannot have more than one API number, in a nutshell it has to be unique in this table, It is not possible to two or more columns having same API number but separate well name, operator name etc.

After removing the duplicate API number, all the rows that don't have an API were removed using filter row step.

Next, we use String cut command to split the API number into 3 substrings called- 1) the state code which is the first 3 digits of API number. 2) the county code fourth to sixth digit. 3) the third is the unique well identifier and it comprises of last four digits.

These three substrings provide us information about the state, county and the unique id of a certain well.

As there were many null values in parent company columns, these null values were replaced by the corresponding operator name column. This was done using the calculator tool.

Next, all the rows in which column "operator name" was null were removed because our stakeholders are people who own these wells and that the wells without operator name are not useful for our final goal.

Next, all the rows in which column "status" was null were removed using filter rows.

Column "CIK CODE" was removed as it does not provide any knowledge to achieve the stakeholders´ final goal.

Before creating a new CSV file and inserting all the data into it, a new column called source was added and for each row "file name" was added. This column will tell the stakeholder the source of data for each row.

All this data was inserted into CSV file called dataset_well.csv



*Figure 5 - Cleaning Flow Overview*

In this transformation we took production_mexico_all_new_format.csv as the input file, removed all the rows that lacked the API number null using row filter. After this step all the rows that did not had their month column null were removed using filter row function. Next,

14

we renamed the column month to date. All the rows in which oil, gas, water were null were filled with 0. All these rows were sorted according to API NUMBER. The source file was added to each row in a new column called SOURCE. This tells the stakeholder the origin of the data.

All this data was filled in a newly created CSV file called dataset_production.csv

## Sample input



*Figure 6 - wel_mexico_all_new_format.csv*



*Figure 7 - Production_mexico_all_new_format.csv*

# Sample Output



*Figure 8 - dataset_well.csv*



*Figure 9 - dataset_production.csv*

## II STEP. TR_DATE



*Figure 10 - Date Transformation*

This is the second transformation in the ETL process, in first process we read CSV input file from dataset_production.csv, then we remove all the rows that we don't want to keep in our dimension. We want a dimension that only has date, source and meta data. So, in select values step we remove all columns except date and source. Then we sort all the rows according to date in Ascending order. We append log data such as last login, last modified. The last step is to insert all this information using insert/update in a table called **d_date** found in schema called **well_status**.

## Sample Output



*Figure 11 - Sample Output Dimension Date*

## III STEP. TR_COMPANY



*Figure 12 - Company Transformation*

This is the third transformation in the ETL process, in first step we read CSV input file from dataset_well.csv, then we remove all the rows that we don't want to keep in our dimension. We want a dimension that only has **parent_company**, source and meta data. So, in select values step we remove all columns except **parent_compant** and source. Then we sort all the rows according to date in Ascending order and remove duplicate rows. We append log data such as last login, last modified. The last step is to insert all this information using insert/update in a table called **d_company** found in schema called **well_status**.

## Sample Output



*Figure 13 - Sample Output Dimension Company*
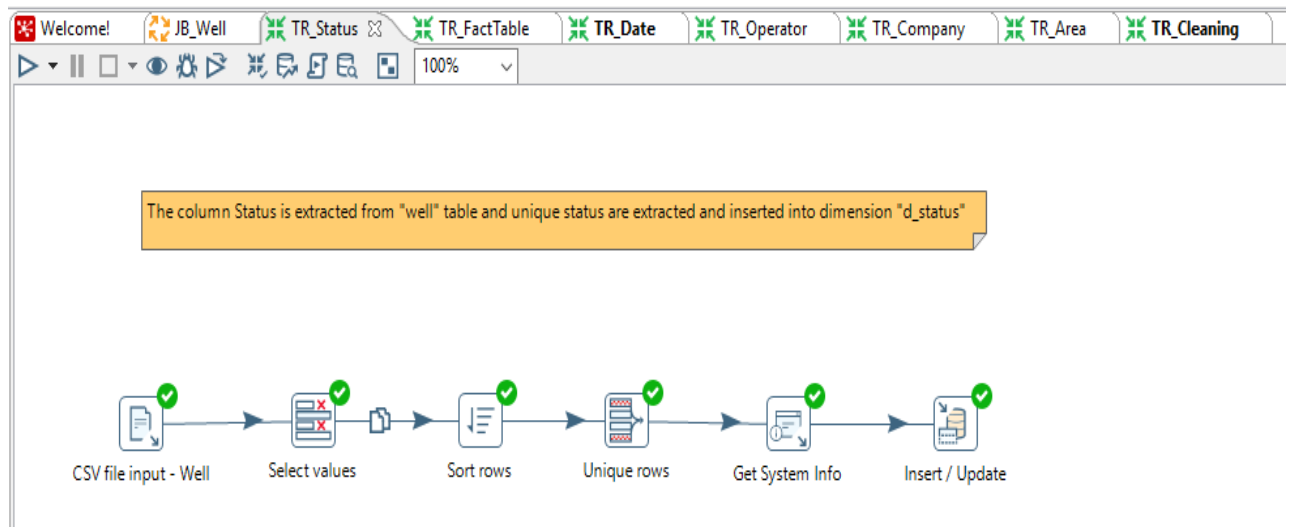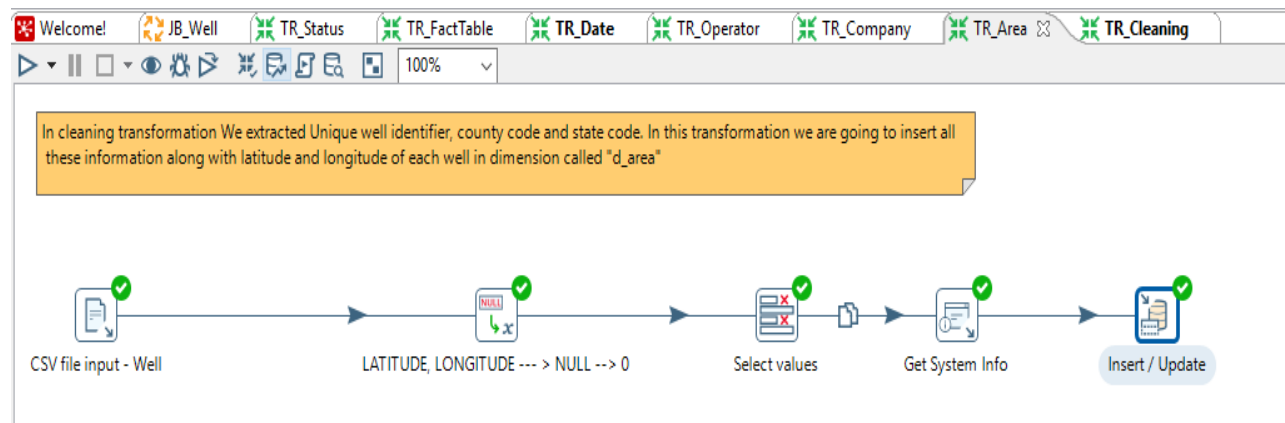
# IV STEP. TR_OPERATOR



*Figure 14 - Transformation Operator*

This is the fourth transformation in the ETL process, in first step we read CSV input file from dataset_well.csv, then we remove all the rows that we don't want to keep in our dimension. We want a dimension that only has operator, source and meta data. So, in select values step we remove all columns except operator and source. Then we sort all the rows according to operator in Ascending order. We removed all the duplicate rows in operator column. We append log data such as last login, last modified. The last step is to insert all this information using insert/update in a table called **d_operator** found in schema called **well_status.**

## Sample Output



*Figure 15 - Sample Output Dimension Operator*

# V STEP. TR_STATUS



*Figure 16 - Transformation Status*

This is the fifth transformation in the ETL process, in first step we read CSV input file from dataset_well.csv, then we remove all the rows that we don't want to keep in our dimension. We want a dimension that only has status, source and meta data. So, in select values step we remove all columns except status and source. Then we sort all the rows according to status in Ascending order. We remove duplicate rows in status column, We append log data such as last login, last modified. The last step is to insert all these information using insert/update in a table called d_status found in schema called well_status.

## Sample Output



*Figure 17 - Sample Output Dimension Status*

# VI STEP. TR_AREA



*Figure 18 - Transformation Area*

This is the sixth transformation in the ETL process, in first step we read CSV input file from dataset_well.csv, then we remove all the rows that we don't want to keep in our dimension. We want a dimension that only has latitude, longitude, state code, unique well identifier, source and meta data. So, in select values step we remove all columns except latitude, longitude, state code, unique well identifier, source and meta data. We append log data such as last login, last modified. The last step is to insert all these information using insert/update in a table called d_area found in schema called well_status.

## Sample Output



*Figure 19 - Sample Output Dimension Area*
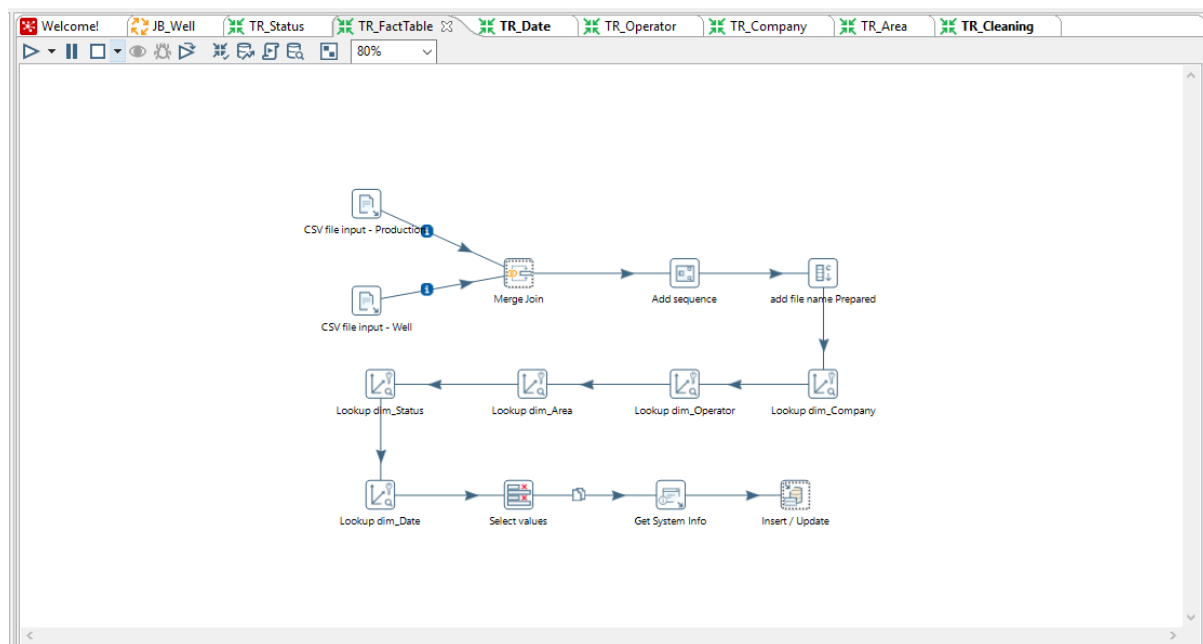
## VII STEP. TR_FACT_TABLE



*Figure 20 - Transformation Fact Table*

In this Transformation we are going to populate the fact table with variable important to the Business needs and foreign keys pointing to other dimensions to provide us extra information. First, we read two CSV files and merged them using **merge join.** The two files we joined were dataset_production.csv and dateset_well.csv. As these two have API NUMBER common we were able to perform inner join. This step is important as we will need variables to look up primary keys in relevant dimensions. We added file name as source using "add filename prepared" this added column will help the stakeholders to know the source of data.

The next step was to perform lookup in the **D_company** table to map **"parent_company "** attribute which is present in both CSV table and **D_company** dimension and by doing that we can add the primary key of **d_company** to current CSV file called **id_operator**.

The next step was to perform lookup in the **D_operator** table to map **"operator_name"** attribute which is present in both CSV table and **D_operator** dimension and by doing that we can add the primary key of **d_operator** to current CSV file called **id_company**.

The next step was to perform lookup in the **D_AREA** table to map **"state_code"** and **"county_code"** attribute which is present in both CSV table and **D_AREA** dimension and by doing that we can add the primary key of **d_AREA** to current CSV file called **id_area**.

The next step was to perform lookup in the **D_STATUS** table to map "**STATUS**" attribute which is present in both CSV table and **D_STATUS** dimension and by doing that we can add the primary key of **d_STATUS** to current CSV file called **id_status**.
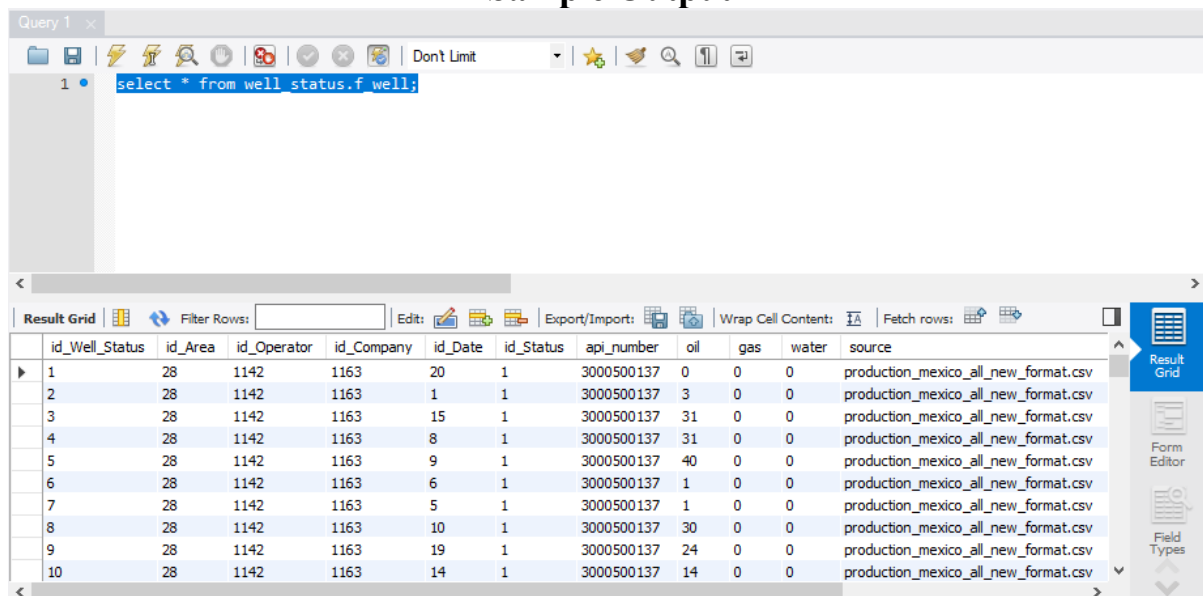
The next step was to perform lookup in the **D_DATE** table to map "**DATE**" attribute which is present in both CSV table and **D_DATE** dimension and by doing that we can add the primary key of **d_DATE** to current CSV file called id_date.

In select/Rename values we only **select id_Area, id_company, id_date, id_operator, id_status, gas, oil, water, source, api_number**. We remove the other columns as they don't help in achieving business goals.

We now append the system log time and time of modification of the transformation so that the stakeholder can know at what time the data was entered into the system.

The ultimate step is to connect to database and to fact table called f_well in the well_status schema and populate the values

## Sample Output



*Figure 21 - Sample Output Fact Table*