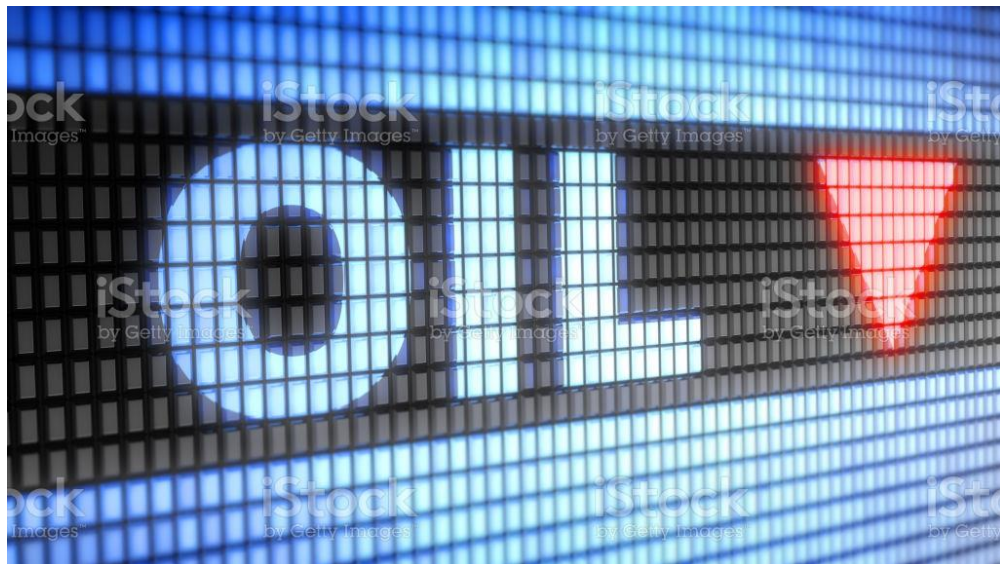


# ETL Final Report

## ETL Project

Nima Karimi and Kirstie McCown

Week 12: University of Western Australia Data Analytics Bootcamp



<b>Project Title:</b>	ETL Project
<b>Class Instructor:</b>	Daniel Kasatchkow
<b>TA:</b>	Hazar Ayaz
<b>Project Due Date:</b>	Saturday 29 <sup>th</sup> August 2020
<b>Date of Report:</b>	Saturday 29 <sup>th</sup> August 2020
<b>Reporting Period:</b>	Tuesday 25 <sup>th</sup> August 2020 – Saturday 29 <sup>th</sup> August 2020

## Table of Contents

Summary .....	1
Data Sources.....	1
Data Transformation .....	1
Database .....	2
Database Tables and Columns.....	2
Project Conclusion .....	3

## Annexures

Annexure 1. Cleaned Oil Table Schema.....	4
Annexure 2. Cleaned Accident Table Schema .....	4
Annexure 3. Table Join Query .....	5
Annexure 4. Joined Tables from Database .....	5
Annexure 5. Screenshot of Database and Schemas .....	6
Annexure 6. Screenshot of Pandas Merge Code and DataFrame .....	6

## Summary

Currently, the price of oil is ever changing, and sometimes for unknown cause. We are carrying out this project to identify if there is any correlation between US Oil Pipeline Accidents and the Crude Oil Price around the same period (2010 to 2016).

## Data Sources

We utilised two sets of data from Kaggle.com, one was in cvs. format and the other .xlsx:

### Oil Prices

<https://www.kaggle.com/rockbottom73/crude-oil-prices>

### US Oil Pipeline Accidents

<https://www.kaggle.com/usdot/pipeline-accidents>

## Data Transformation

Our overall data transformation we wanted to look at the following elements, we will detail these further below for each individual data set.

- Remove any unnecessary columns
- Drop all accidents not related to crude oil
- Drop select items which are N/A or have the value of NaN
- Split accident date/time field to show only dates

### Oil Prices

- Read in xlsx to Panadas DataFrame to enable visualisation and cleaning
- Rename columns so they are easier to work with
- Drop everything which is N/A
- Drop all rows that are not the same dates as what is in the accidents DataFrame

### US Oil Pipeline Accidents

- Read in csv to Panadas DataFrame to enable visualisation and cleaning
- Remove any unnecessary columns and rename columns so they are easier to work with
- Look at what non null values are in the DataFrame to see if values need to be removed
- Drop all accidents not related to crude oil
- Drop everything which is N/A in the following columns: city, facility\_name, country, shutdown

- Split the date/time column keeping the date in a newly created column, whilst dropping the original date/time column
- Change the format of the date so that both DataFrame dates match format

## Database

For our project we utilised a Postgres SQL Database, as part of our ETL process we conducted the following steps:

- Create a new Postgres Database called “oil\_db”
- Create two table schema’s called “cleaned\_oil” and “cleaned\_accidents”
- Connect to Postgres database via our Jupyter Notebook (.ipynb file)
- Check to ensure tables are available in Postgres database and able to be connected with via our Jupyter Notebook
- Load panda's DataFrame to postgres sql tables

See: Annexure 1 and Annexure 2

## Database Tables and Columns

Table Name	Number of Columns
cleaned_oil	2
cleaned_accidents	17

### Columns – cleaned\_oil

Column Name	Column Type	Description
date	date	Date oil price was recorded
price	int	End of day price for a given day

### Columns – cleaned\_accidents

Column Name	Column Type	Description
Report_number	Int	Accident Report Number
Accident_date	Object	Date of Accident
Op_id	int	Operator ID in charge at time of Accident
Op_name	object	Operator Name
Facility_name	object	Facility Name where Accident occurred
Location	object	Location of accident (on/offshore)
Pipeline_type	object	Under/Above Ground Pipeline
Liquid_type	object	Type of Liquid in the Accident Area
City	object	City of Accident
Country	object	Country of Accident
State	object	State of Accident
Cause_cat	object	Category of Reason of Accident
Cause_subcat	object	Sub-Category of Reason of Accident
Shutdown	object	Was the plant shut down at the time of the accident (Yes/No)
Shut_date_time	object	Shutdown date and time if applicable
Restart_date_time	object	Restart date and time if applicable

The above two tables were joined in both Panda's and SQL to create one table for further analysis.

In our Panda's DataFrame with all merged information there are some N/A's still present because for each day there was a price reading, there was not necessarily an accident occur on that same date. These have been left in for clarity and can be removed once further investigation is required and commenced.

**See: Annexure 3, Annexure 4, Annexure 5 and Annexure 6**

## Project Conclusion

We feel that our ETL process has prepared the two datasets adequately in order to be able to further analyse and identify if there is any correlation between US Oil Pipeline Accidents and the fluctuation of Crude Oil Prices around the same period of time.

Our dataset has been prepared into two separate tables, which have then been joined to allow for further investigation and manipulation, while maintaining the integrity of each individual data set as a whole.

## Annexures/ Figures

### Annexure 1 – Cleaned Oil Table Schema

```
oil_db/postgres@PostgreSQL 12
Query Editor  Query History
1  -- Create tables and import data
2  -- Drop table if exists
3  DROP TABLE IF EXISTS cleaned_oil;
4
5  -- Create new table
6  CREATE TABLE cleaned_oil (
7      index int,
8      date date,
9      price decimal,
10     Primary Key (date)
11 );
12
```

### Annexure 2 – Cleaned Accidents Table Schema

```
oil_db/postgres@PostgreSQL 12
Query Editor  Query History
13 -- Create tables and import data
14 -- Drop table if exists
15 DROP TABLE IF EXISTS cleaned_accidents;
16
17 -- Create new table
18 CREATE TABLE cleaned_accidents (
19     index int,
20     report_number int,
21     op_id int,
22     op_name varchar,
23     facility_name varchar,
24     location varchar,
25     pipeline_type varchar,
26     liquid_type varchar,
27     city varchar,
28     country varchar,
29     state varchar,
30     cause_cat varchar,
31     cause_subcat varchar,
32     shutdown varchar,
33     shut_date_time varchar,
34     restart_date_time varchar,
35     date date
36 );
```

### Annexure 3 – Table Join Query


oil\_db/postgres@PostgreSQL 12

Query Editor
Query History

```

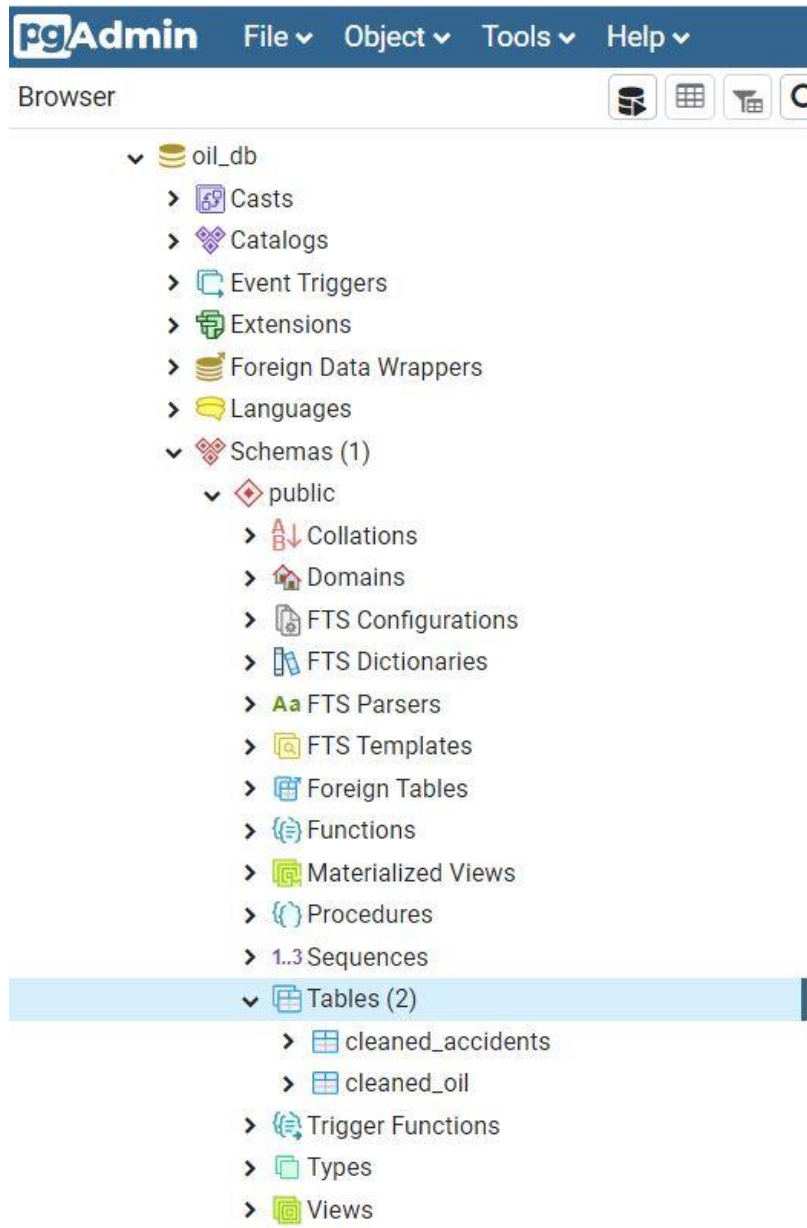
38
39 select
40 cleaned_accidents.date,
41 cleaned_oil.price,
42 cleaned_accidents.report_number,
43 cleaned_accidents.op_id,
44 cleaned_accidents.op_name,
45 cleaned_accidents.facility_name,
46 cleaned_accidents.location,
47 cleaned_accidents.pipeline_type,
48 cleaned_accidents.liquid_type,
49 cleaned_accidents.city,
50 cleaned_accidents.country,
51 cleaned_accidents.state,
52 cleaned_accidents.cause_cat,
53 cleaned_accidents.cause_subcat,
54 cleaned_accidents.shutdown,
55 cleaned_accidents.shut_date_time,
56 cleaned_accidents.restart_date_time
57 from cleaned_accidents
58 right join cleaned_oil on cleaned_accidents.date = cleaned_oil.date;
59

```

### Annexure 4 - Joined Tables from Database

	date	price	report_number	op_id	op_name	facility_name	location	pipeline_type	liquid_type	city	country	state
	date	numeric	integer	integer	character varying	character varying	character varying	character varying	character varying	character varying	character varying	character varying
1	2010-01-11	82.54	20100234	9175	JAYHAWK PIPELINE ...	CHASE KAW TERMI...	ONSHORE	UNDERGROUND	CRUDE OIL	CHASE	RICE	KS
2	2010-01-11	82.54	20100026	31684	CONOCOPHILLIPS	TANK 824	ONSHORE	TANK	CRUDE OIL	CUSHING	PAYNE	OK
3	2010-01-12	80.79	20100106	26085	PLAINS MARKETING...	CUSHING TERMINAL	ONSHORE	ABOVEGROUND	CRUDE OIL	CUSHING	LINCOLN	OK
4	2010-01-12	80.79	20100082	32080	CCPS TRANSPORTA...	CCPS TRANSPORT...	ONSHORE	ABOVEGROUND	CRUDE OIL	RUSHVILLE	SCHUYLER	IL
5	2010-01-13	79.66	20100100	22855	KOCH PIPELINE CO...	PARK RAPIDS PUM...	ONSHORE	ABOVEGROUND	CRUDE OIL	MENAHGA	HUBBARD	MN
6	2010-01-14	79.35	20100057	10250	KIANTONE PIPELINE...	GOWANDA BOOST...	ONSHORE	ABOVEGROUND	CRUDE OIL	GOWANDA	CATTARAUGUS	NY
7	2010-01-15	77.96	20110083	31325	PACIFIC PIPELINE SY...	LINE 63 SOUTH PA...	ONSHORE	ABOVEGROUND	CRUDE OIL	CARSON	LOS ANGELES	CA
8	2010-01-21	75.84	20100091	31325	PACIFIC PIPELINE SY...	NORTH COLES LEV...	ONSHORE	TANK	CRUDE OIL	TAFT	KERN	CA

## Annexure 5 - Screenshot of Database and Schemas



## Annexure 6 - Screenshot of Panda's Merge Code and DataFrame

```
# Join both tables for further analysis
all_data = cleaned_accidents.merge(cleaned_oil, on="date", how="right")
all_data
```

	report_number	op_id	op_name	facility_name	location	pipeline_type	liquid_type	city	country	state	cause_cat	cause_subcat	shutdown	shot_date_time	restart_date_time	date	price
0	20100030.0	300.0	PLAINS PIPELINE, L.P.	RED RIVER EAST	ONSHORE	UNDERGROUND	CRUDE OIL	SHERMAN	GRAYSON	TX	EXCAVATION DAMAGE	THIRD PARTY EXCAVATION DAMAGE	NO	NaN	NaN	2010-01-07	82.6
1	20100234.0	9175.0	JAYHAWK PIPELINE LLC	CHASE KAW TERMINAL	ONSHORE	UNDERGROUND	CRUDE OIL	CHASE	RICE	KS	CORROSION	INTERNAL	YES	1/11/2010 15:10	1/13/2010 14:00	2010-01-11	82.5
2	20100026.0	31684.0	CONOCOPHILLIPS	TANK #24	ONSHORE	TANK	CRUDE OIL	CUSHING	PAYNE	OK	NATURAL FORCE DAMAGE	TEMPERATURE	NO	NaN	NaN	2010-01-11	82.5
3	20100106.0	26085.0	PLAINS MARKETING, L.P.	CUSHING TERMINAL	ONSHORE	ABOVEGROUND	CRUDE OIL	CUSHING	LINCOLN	OK	CORROSION	INTERNAL	YES	1/12/2010 10:25	1/14/2010 10:25	2010-01-12	80.7
4	20100082.0	32080.0	CCPS TRANSPORTATION, LLC	CCPS TRANSPORTATION, RUSH STATION	ONSHORE	ABOVEGROUND	CRUDE OIL	RUSHVILLE	SCHUYLER	IL	INCORRECT OPERATION	INCORRECT INSTALLATION	NO	NaN	NaN	2010-01-12	80.7
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...