

Project 2 ETL Technical Report by Team 3

Anthony Stevens, Cecilia Zhang, Neda Mehdizadeh, and Sahar Jamal

Project Theme : Covid-19 in the USA

Step 1. Extract: Data Sources

- Covid-19 in USA-Number of Novel CoronaVirus 2019 cases in USA from Kaggle.com 3 csv files (primary)
Reference Link: https://www.kaggle.com/sudalairajkumar/covid19-in-usa?select=us_counties_covid19_daily.csv
- State, County, CBSA, Country API Data from Covid Act Now (primary)
Reference Link: <https://apidocs.covidactnow.org/data-definitions/>

Step 2. Transform: Data Wrangling

- Covid daily data is retrieved from “us_counties_covid19_daily”csv file, then encoded the fips (Federal Information Processing System) column to an int instead of a float. After that the null values were dropped to keep only usable data.
- The original data is scraped from json file “counties. timeseries” and is cleaned by selecting columns with useful information. Then dropped columns with NaN values such as “lat” or contained dictionaries which were not useful.. Leaving us with a handful of useful data. We ran into an issue with two of the columns throwing errors for “undefined columns”. The solution was to rename and lowercase the names of the columns.

Step 3. Loading:

We chose PostgreSQL as our database to load our data. The reason for this was due to this database being a common sql database, so we knew the expertise in this area would definitely challenge us to perfect our skills. After tackling the portion of data we decided to keep in pandas, we then lined up a table within postgres. To do so we used the <CREATE TABLE> statement, and then we identified each column we wanted in the table. The next step was to run the table and check back for the table to upload. We then entered into pandas and connected the engine so that the data would load into the database and align nicely. To confirm our work was done correctly, we utilized the <SELECT *> of the specific table to print what had loaded.

Project 2- A Case Study of ETL

Covid-19 Datasets

Team 3

Anthony Stevens, Cecilia Zhang, Neda Mehdizadeh, and Sahar Jamal

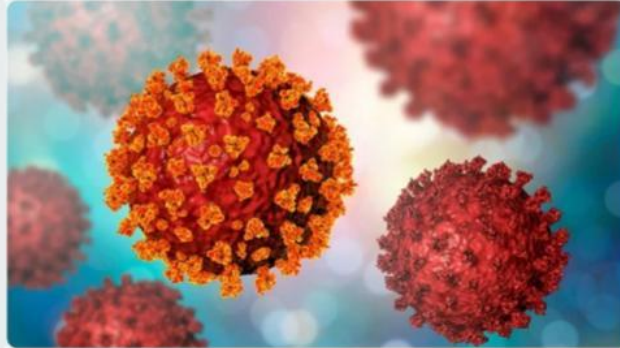


Illustration Credit: Kateryna Kon / Shutterstock

Extract: Data Sources

1. **Covid-19 in USA-Number of Novel Corona Virus 2019 cases in USA from Kaggle.com 3 csv files (primary)**

Reference Link:

https://www.kaggle.com/sudalairajkumar/covid19-in-usa?select=us_counties_covid19_daily.csv

2. **State, County, CBSA, Country API Data from Covid Act Now (primary)**

Reference Link: <https://apidocs.covidactnow.org/data-definitions/>

API Value	CDC Community Transmission Level
0	Low
1	Moderate
2	Substantial
3	High
4	"Unknown"

Transform: Data Wrangling pt.1

Json file before cleaning

fips	country	state	county	level	lat	locationid	long	population	metrics	risk_level	subTransmissionLevel	actuals	annotations	lastUpdatedDate
0	2013	US	AK	Aleutians East Borough	county	NaK	NaK	3337	{testPositivityRate: 0.0%, testPositivityRate: 4.1, testPositivityRate: 2.0}	Covered: 2.0	testPositivityRate: 4.1, testPositivityRate: 2.0	{cases: 431, deaths: 3, positiveTests: 0.0}	{cases: 431, deaths: 3, positiveTests: 0.0}	2021-10-19
1	2016	US	AK	Aleutians West Census Area	county	NaK	NaK	5634	{testPositivityRate: 0.0%, testPositivityRate: 1.0, testPositivityRate: 2.0}	Covered: 2.0	testPositivityRate: 1.0, testPositivityRate: 2.0	{cases: 313, deaths: 0, positiveTests: 0.0}	{cases: 313, deaths: 0, positiveTests: 0.0}	2021-10-19
2	2020	US	AK	Anchorage Municipality	county	NaK	NaK	288000	{testPositivityRate: 0.111, testPositivityRate: 2.0, testPositivityRate: 3.0}	Covered: 3.0	testPositivityRate: 2.0, testPositivityRate: 3.0	{cases: 5285, deaths: 284, positiveTests: 0.0}	{cases: 5285, deaths: 284, positiveTests: 0.0}	2021-10-19
3	2050	US	AK	Bethel Census Area	county	NaK	NaK	18386	{testPositivityRate: 0.111, testPositivityRate: 4.0, testPositivityRate: 3.0}	Covered: 3.0	testPositivityRate: 4.0, testPositivityRate: 3.0	{cases: 160, deaths: 25, positiveTests: 0.0}	{cases: 160, deaths: 25, positiveTests: 0.0}	2021-10-19
4	2060	US	AK	Bristol Bay Borough	county	NaK	NaK	836	{testPositivityRate: 0.0%, testPositivityRate: 4.0, testPositivityRate: 0.0}	Covered: 0.0	testPositivityRate: 4.0, testPositivityRate: 0.0	{cases: 685, deaths: 2, positiveTests: 0.0}	{cases: 685, deaths: 2, positiveTests: 0.0}	2021-10-19

Json file after cleaning

fips	country	state	county	population	transmission_level	last_updated
0	2013	US	Aleutians East Borough	3337	2	2021-10-19
1	2016	US	Aleutians West Census Area	5634	2	2021-10-19
2	2020	US	Anchorage Municipality	288000	3	2021-10-19
3	2050	US	Bethel Census Area	18386	3	2021-10-19
4	2060	US	Bristol Bay Borough	836	0	2021-10-19

Data Wrangling pt. 2

Reading csv file

```
1 covid_daily_data_df.fillna(0, inplace=True)
2 covid_daily_data_df['fips'] = covid_daily_data_df['fips'].astype(int)
3 covid_daily_data_df = covid_daily_data_df[covid_daily_data_df['fips'] != 0]
```

	date	county	state	fips	cases	deaths
0	2020-01-21	Snohomish	Washington	53061	1	0.0
1	2020-01-22	Snohomish	Washington	53061	1	0.0
2	2020-01-23	Snohomish	Washington	53061	1	0.0
3	2020-01-24	Cook	Illinois	17031	1	0.0
4	2020-01-24	Snohomish	Washington	53061	1	0.0

```
cleaned_df.to_excel("cleaned_df.xlsx")
```

Confirm data has been added by using pgAdmin

Reading table "daily"

Query Editor Query History

```
9
10 CREATE TABLE daily(
11     date DATE,
12     county VARCHAR(255),
13     state VARCHAR(255),
14     fips FLOAT,
15     cases INT,
16     deaths FLOAT
17 );
18
19
20 SELECT * FROM daily;
```

Data Output Explain Messages Notifications

	date	county	state	fips	cases	deaths
	date	character varying (255)	character varying (255)	double precision	integer	double precision
1	2020-01-21	Snohomish	Washington	53061	1	0
2	2020-01-22	Snohomish	Washington	53061	1	0
3	2020-01-23	Snohomish	Washington	53061	1	0
4	2020-01-24	Cook	Illinois	17031	1	0
5	2020-01-24	Snohomish	Washington	53061	1	0

Reading table "current_data"

Query Editor Query History

```
24 CREATE TABLE current_data(
25     fips FLOAT,
26     county VARCHAR(255),
27     state VARCHAR(255),
28     population INT,
29     transmission_level INT,
30     last_updated DATE
31 );
32
33
34 SELECT * FROM current_data;
```

Data Output Explain Messages Notifications

	fips	county	state	population	transmission_level	last_updated
	double precision	character varying (255)	character varying (255)	integer	integer	date
1	2013	US	AK	3337		2 2021-10-19
2	2016	US	AK	9634		2 2021-10-19
3	2020	US	AK	288900		3 2021-10-19
4	2050	US	AK	18386		3 2021-10-19

Connect to local database and Use pandas to load csv/json converted DataFrame into database

```
covid_daily_data_df.to_sql(name="daily", con=engine, if_exists='append', index=False)
new_covid_json.to_sql(name='current_data', con=engine, if_exists='append', index=False)
```

Thank You!