

# HW1\_109550017\_黃品云

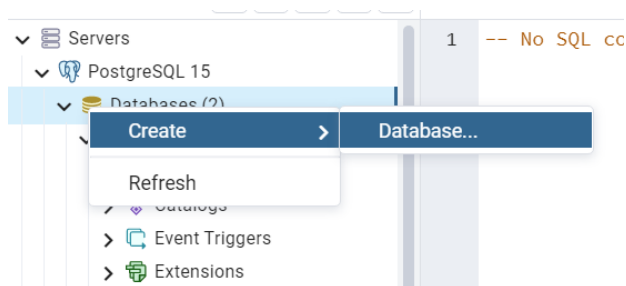
HW1\_report :

**Q1. The process of creating the “covid19” databases (can be screenshot and/or SQL/non-SQL statements with text explanation) (10pts)**

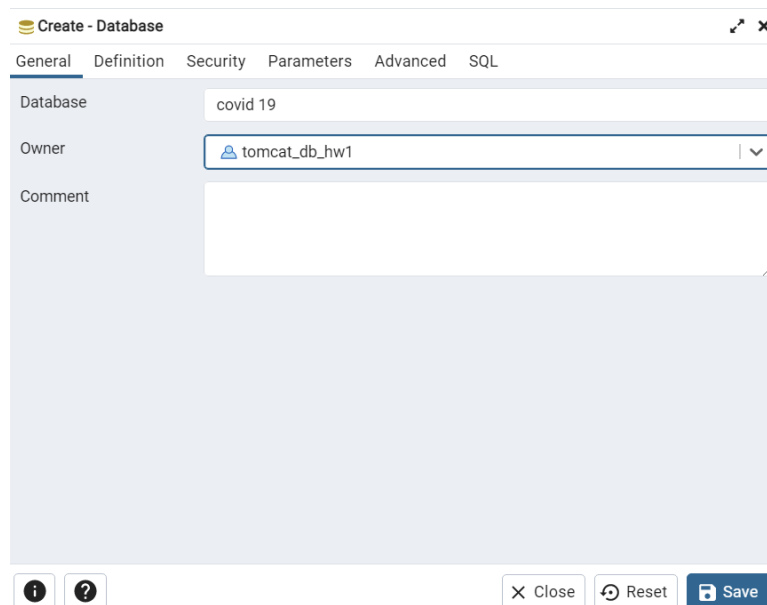
**Ans 1:**

I followed the following YouTube tutorial to set up my postgresSQL, and it was downloaded successfully.  
(<https://www.youtube.com/watch?v=K5-Ed18rOb0>)

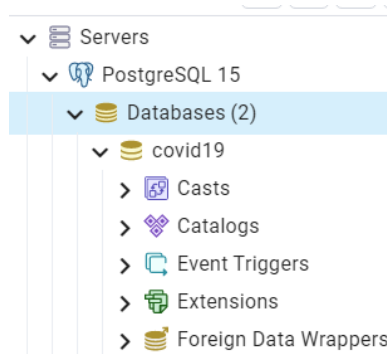
I clicked on Databases and created a new database within it.



I typed in the database's name, which is “covid 19”, and saved it.



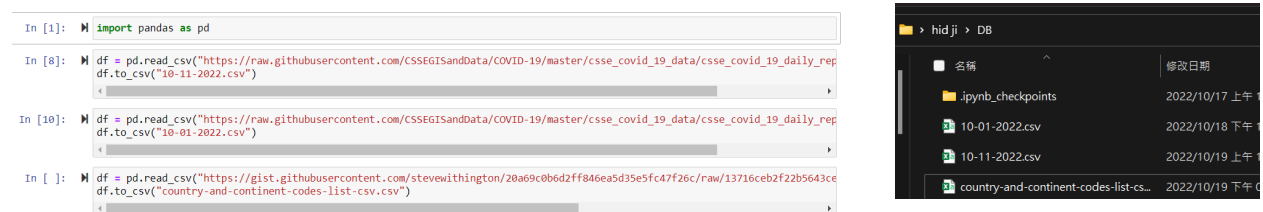
As you can see, the “covid 19” database has been created successfully.



**Q2. The process of importing three required .csv files into covid19 database (can be screenshot and/or SQL/non-SQL statements with text explanation). Please included/described the data type and keys of the imported table in your screenshot, SQL statements, and explanations (30pts)**

**Ans 2:**

First, I downloaded the required three csv files(using python in Jupyter Notebook) from github following the tutorial in the link, and it worked successfully.(<https://www.youtube.com/watch?v=SMCoNBQupaE>)



Second, I started to set the schema for “10012022” table according to “10-01-2022.csv” file, and specified the column names and data types.

**“10012022” table(from “10-01-2022.csv”):**

Create - Table

General Columns Advanced Constraints Partitions Parameters Security SQL

Inherited from table(s)

Columns

Name	Data type	Length/Precision	Scale	Not NULL?	Primary key?	Default
FIPS	character varying	5		<input type="checkbox"/>	<input type="checkbox"/>	
Admin2	character varying	20		<input type="checkbox"/>	<input type="checkbox"/>	
Province_State	character varying	20		<input type="checkbox"/>	<input type="checkbox"/>	
Country_Region	character varying	20		<input type="checkbox"/>	<input type="checkbox"/>	
Last_Update	timestamp without ti...			<input type="checkbox"/>	<input type="checkbox"/>	
Lat	double precision			<input type="checkbox"/>	<input type="checkbox"/>	
Long_	double precision			<input type="checkbox"/>	<input type="checkbox"/>	
Confirmed	integer			<input type="checkbox"/>	<input type="checkbox"/>	

Create - Table

General Columns Advanced Constraints Partitions Parameters Security SQL

Last_Update	timestamp without ti...			<input type="checkbox"/>	<input type="checkbox"/>	
Lat	double precision			<input type="checkbox"/>	<input type="checkbox"/>	
Long_	double precision			<input type="checkbox"/>	<input type="checkbox"/>	
Confirmed	integer			<input type="checkbox"/>	<input type="checkbox"/>	
Deaths	integer			<input type="checkbox"/>	<input type="checkbox"/>	
Recovered	integer			<input type="checkbox"/>	<input type="checkbox"/>	
Active	integer			<input type="checkbox"/>	<input type="checkbox"/>	
Combined_Key	character varying	80		<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Incident_Rate	double precision			<input type="checkbox"/>	<input type="checkbox"/>	
Case_Fatality_Ratio	double precision			<input type="checkbox"/>	<input type="checkbox"/>	

Close Reset Save



However, I have changed the datatypes of "FIPS", "Admin2", "Province\_State", "Country\_Region" because of the errors that I encountered later, the updated version of table schema is as following



```
Dashboard Properties SQL Statistics Dependencies Dependents Processes covid19/postgr...
1  -- Table: public.10012022
2
3  -- DROP TABLE IF EXISTS public."10012022";
4
5  CREATE TABLE IF NOT EXISTS public."10012022"
6  (
7      "FIPS" character varying(5) COLLATE pg_catalog."default",
8      "Admin2" character varying(100) COLLATE pg_catalog."default",
9      "Province_State" character varying(100) COLLATE pg_catalog."default",
10     "Country_Region" character varying(100) COLLATE pg_catalog."default",
11     "Last_Update" time with time zone,
12     "Lat" double precision,
13     "Long_" double precision,
14     "Confirmed" integer,
15     "Deaths" integer,
16     "Recovered" integer,
17     "Active" integer,
18     "Combined_Key" character varying(300) COLLATE pg_catalog."default" NOT NULL,
19     "Incident_Rate" double precision,
20     "Case_Fatality_Ratio" double precision,
21     CONSTRAINT "10012022_pkey" PRIMARY KEY ("Combined_Key")
22 )
23
24 TABLESPACE pg_default;
25
26 ALTER TABLE IF EXISTS public."10012022"
27     OWNER to tomcat_db_hw1;
```

As you can see, in the relation "10012022", there are 14 attributes, and explanations are as follows:

- **FIPS**: Federal Information Processing Standards code that uniquely identifies counties within the USA. Since it is a code, I set it as varchar(5).
- **Admin2**: County name, which can be represented as varchar(100).
- **Province\_State**: Province, state or dependency name, which can be represented as varchar(100).
- **Country\_Region**: Country, region or sovereignty name, which can be represented as varchar(100).
- **Last Update**: MM/DD/YYYY HH:mm:ss (24 hour format, in UTC), which can be represented as timestamp without time zone
- **Lat** and **Long\_**: Dot locations on the dashboard, it is a float type with many numbers, I think it would be better if it is set as double precision.
- **Confirmed**: Counts include confirmed and probable (where reported), it is a number, which can be represented as integer.
- **Deaths**: Counts include confirmed and probable (where reported), it is a number, which can be represented as integer.
- **Recovered**: Recovered cases, it is a number, which can be represented as integer.
- **Active**: Active cases = total cases - total recovered - total deaths, it is a number, which can be represented as integer.

- **Combined\_Key:** A string that stores the combination of fips, admin2, province\_state, and country\_region, since it is a very long string, I set it as varchar(305)
- **Incident\_Rate:** Incidence Rate = cases per 100,000 persons, float type with many numbers, which can be represented as double precision.
- **Case\_Fatality\_Ratio (%):** Case-Fatality Ratio (%) = Number recorded deaths / Number cases, float type with many numbers, which can be represented as double precision.

I set **Combined\_Key** as **primary key** because it is the only column that I think can distinguish different data.

As for tables "10112022(from "10-11-2022.csv")" and "countryAndCode"(from "country-and-continent-codes-list-csv.csv"), I repeated the above instructions to create them.

### 10112022 table(from "10-11-2022.csv") :

```

Dashboard  Properties  SQL  Statistics  Dependencies  Dependents  Processes  covid19/postgr...
1  -- Table: public.10112022
2
3  -- DROP TABLE IF EXISTS public."10112022";
4
5  CREATE TABLE IF NOT EXISTS public."10112022"
6  (
7      "FIPS" character varying(5) COLLATE pg_catalog."default",
8      "Admin2" character varying(100) COLLATE pg_catalog."default",
9      "Province_State" character varying(100) COLLATE pg_catalog."default",
10     "Country_Region" character varying(100) COLLATE pg_catalog."default",
11     "Last_Update" time with time zone,
12     "Lat" double precision,
13     "Long_" double precision,
14     "Confirmed" integer,
15     "Deaths" integer,
16     "Recovered" integer,
17     "Active" integer,
18     "Combined_Key" character varying(300) COLLATE pg_catalog."default" NOT NULL,
19     "Incident_Rate" double precision,
20     "Case_Fatality_Ratio" double precision,
21     CONSTRAINT "10112022_pkey" PRIMARY KEY ("Combined_Key")
22 )
23
24 TABLESPACE pg_default;
25
26 ALTER TABLE IF EXISTS public."10112022"
27     OWNER to tomcat_db_hw1;

```

### country-and-continent-codes-list table(from "country-and-continent-codes-list-csv.csv"):

countryAndCode

General Columns Advanced Constraints Parameters Security SQL

Inherited from table(s) Select to inherit from...

Columns +

	Name	Data type	Length/Precision	Scale	Not NULL?	Primary key?	Default
	Continent_Name	character varying	20		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
	Continent_Code	character varying	2		<input checked="" type="checkbox"/>	<input type="checkbox"/>	
	Country_Name	character varying	100		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
	Two_Letter_Country	character	2		<input checked="" type="checkbox"/>	<input type="checkbox"/>	
	Three_Letter_Count	character	3		<input type="checkbox"/>	<input type="checkbox"/>	
	Country_Number	character varying	20		<input type="checkbox"/>	<input type="checkbox"/>	

X Close
↺ Reset
Save

```
CREATE TABLE IF NOT EXISTS public."countryAndCode"
(
    "Continent_Name" character varying(20) COLLATE pg_catalog."default" NOT NULL,
    "Continent_Code" character varying(2) COLLATE pg_catalog."default" NOT NULL,
    "Country_Name" character varying(100) COLLATE pg_catalog."default" NOT NULL,
    "Two_Letter_Country_Code" character(2) COLLATE pg_catalog."default" NOT NULL,
    "Three_Letter_Country_Code" character(3) COLLATE pg_catalog."default",
    "Country_Number" character varying(20) COLLATE pg_catalog."default",
    CONSTRAINT "countryAndCode_pkey" PRIMARY KEY ("Continent_Name", "Country_Name")
)

TABLESPACE pg_default;

ALTER TABLE IF EXISTS public."countryAndCode"
    OWNER to tomcat_db_hw1;
```

As you can see, in the relation “countryAndCode”, there are 6 attributes, and explanations are as follows:

- Continent\_Name: the name of continent, I set it as varchar(20)
- Continent\_Code: the code version of representing continent, I set it as varchar(2)
- Country\_Name: the name of country, I set it as varchar(100)
- Two\_Letter\_Country\_Code: codes with two digits to represent countries, I set it as char(2)
- Three\_Letter\_Country\_Code: codes with three digits to represent countries, I set it as char(3)
- Country\_Number: the number of representing country, I set it as varchar(20)

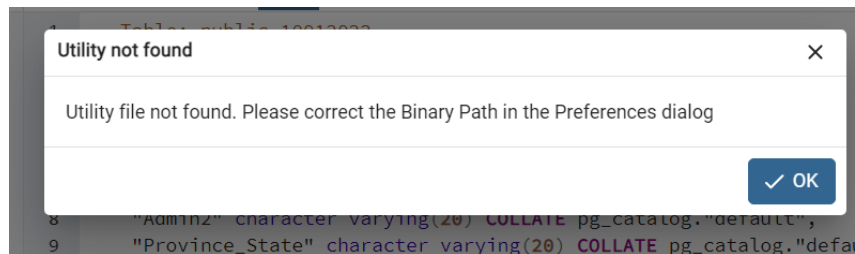
I set **Continent\_Name** and **Country\_Name** as **primary key** because I think they are two useful elements to distinguish data.



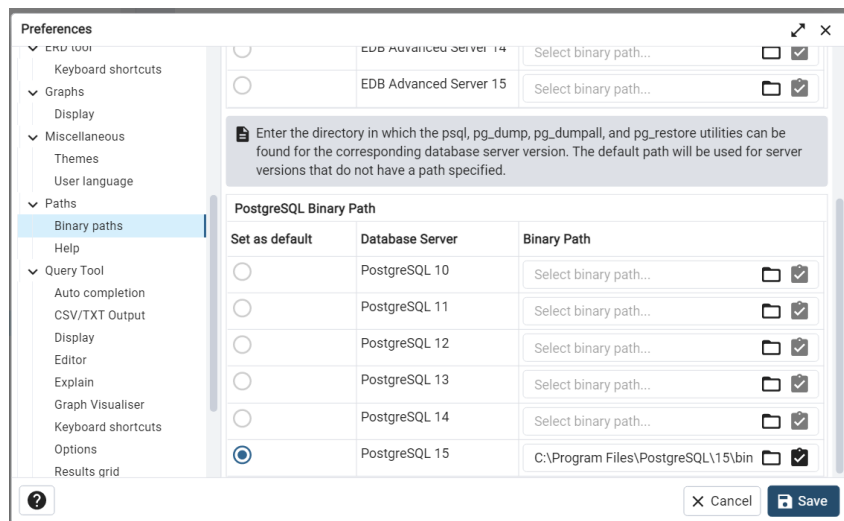
I found that the downloaded “country-and-continent-codes-list-csv.csv” file not only had data inconsistency, but also had missing elements, so I edited it according to the original version in github, and to ensure I did not do something wrong in the csv file, I set columns before “Three\_Letter\_Count” and “Country\_Number” into NULL.

## Import csv file

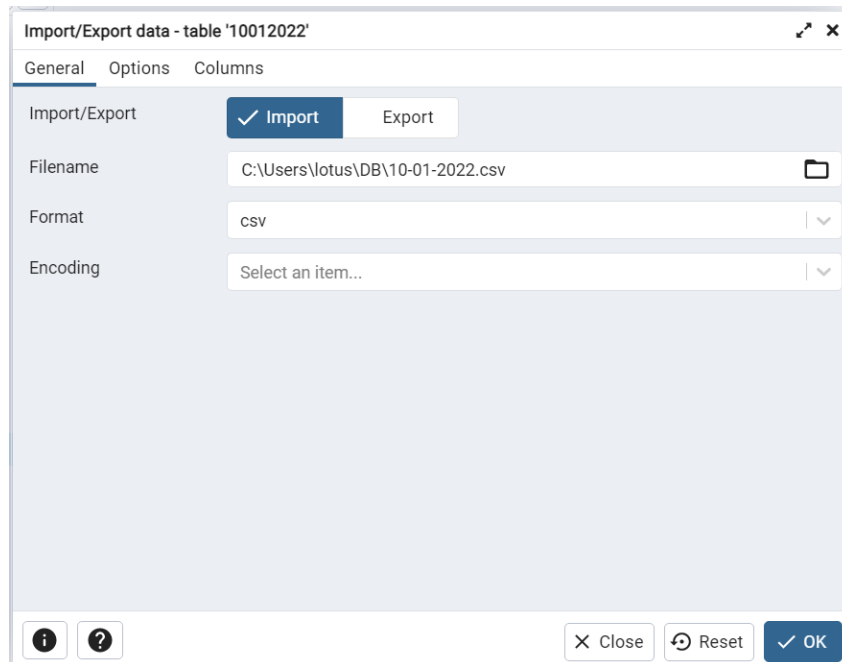
I clicked on table and imported csv file, but the following error happened. I got an “utility not found” error.



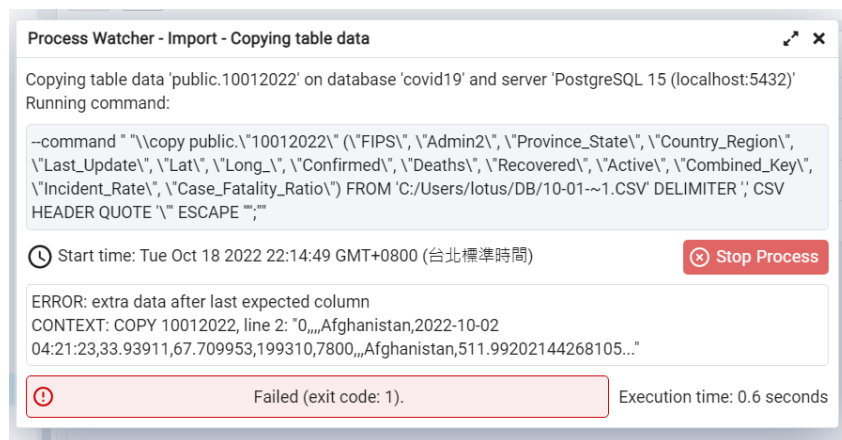
To solve this error, I went to Preferences to edit my Binary paths.



When I tried to import data again, the error disappeared.



When I was trying to importing the csv file again, this new error message kept showing.



And I found that it was because when I downloaded the data from github, it simultaneously created a column before FIPS without my knowing. By deleting the column before FIPS, I solved this error.



	FIPS	Admin2	Province_	Country_F	Last_Update	Lat	Long	Confirmed	Deaths	Recovered	Active	Combined	Incident_F	Case_Fatality
0			Afghanistan	#####	33.93911	67.70995	199310	7800				Afghanistan	511.992	3.913502
1			Albania	#####	41.1533	20.1683	332263	3589				Albania	11545.73	1.080168
2			Algeria	#####	28.0339	1.6596	270676	6879				Algeria	617.2624	2.541415
3			Andorra	#####	42.5063	1.5218	46227	155				Andorra	59829.16	0.335302
4			Angola	#####	-11.2027	17.8739	103131	1917				Angola	313.7898	1.858801
5			Antarctica	#####	-71.9499	23.347	11	0				Antarctica		0
6			Antigua and Barbuda	#####	17.0608	-61.7964	9098	146				Antigua and Barbuda	9290.499	1.604748
7			Argentina	#####	-38.4161	-63.6167	9708420	129897				Argentina	21480.81	1.337983
8			Armenia	#####	40.0691	45.0382	442875	8683				Armenia	14945.66	1.960598
9			Australian Capital Territory	#####	-35.4735	149.0124	205752	126				Australian Capital Territory	48061.67	0.061239
10			New South Wales	#####	-33.8688	151.2093	3507880	5314				New South Wales	43211.14	0.151488
11			Northern Territory	#####	-12.4634	130.8456	97483	73				Northern Territory	39691.78	0.074885
12			Queensland	#####	-27.4698	153.0251	1641389	2239				Queensland	32086.58	0.136409
13			South Australia	#####	-34.9285	138.6007	768961	946				South Australia	43778.02	0.123023
14			Tasmania	#####	-42.8821	147.3272	249621	187				Tasmania	46614.57	0.074914
15			Victoria	#####	-37.8136	144.9631	2614813	5678				Victoria	39439.7	0.217147
16			Western Australia	#####	-31.9505	115.8605	1154732	658				Western Australia	43896.15	0.056983
17			Austria	#####	47.5162	14.5501	5144116	20754				Austria	57116.23	0.403451

## HOWEVER...

After I solved the errors mentioned above, tons of errors came after it 🐞, and it was because some of the province names, country names and admin2 are REEEEEEEAALLLLYY long 😭, then I just set every column datatype into varchar(100) 😊😄 And it finally worked. ^\_^

	PID	Type	Server	Object	Start Time	Status	Time Taken (sec)
<input type="checkbox"/>	4584	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Finished	0.79
<input type="checkbox"/>	10112	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61
<input type="checkbox"/>	10944	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.26
<input type="checkbox"/>	10972	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.4
<input type="checkbox"/>	11236	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61
<input type="checkbox"/>	7372	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61
<input type="checkbox"/>	10412	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61
<input type="checkbox"/>	8764	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61
<input type="checkbox"/>	1944	Import Data	PostgreSQL 15 (localhost:5432)	covid19/public.10012022	2022/10/18 下午10:00:00	Failed	0.61

**Process completed**

Copying table data 'public.10012022' on database 'covid19' and server 'PostgreSQL 15 (localhost:5432)'

[View Processes](#)

**Process started**

Copying table data 'public.10012022' on database 'covid19' and server 'PostgreSQL 15 (localhost:5432)'

[View Processes](#)

**Process started**

Copying table data 'public.10012022' on database 'covid19' and server 'PostgreSQL 15 (localhost:5432)'

[View Processes](#)

I finally successfully imported all the data from csv file!!

Query Query History

1 `select * from "10012022"`

Scratch Pad

Data Output Messages Notifications

	FIPS character varying (5)	Admin2 character varying (100)	Province_State character varying (100)	Country_Region character varying (100)	Last_Update time with time zone	Lat double precision	Long_ double precision	Confirmed integer
1	[null]	[null]	[null]	Afghanistan	04:21:00+08:00	33.93911	67.709953	1995
2	[null]	[null]	[null]	Albania	04:21:00+08:00	41.1533	20.1683	3322
3	[null]	[null]	[null]	Algeria	04:21:00+08:00	28.0339	1.6596	2706
4	[null]	[null]	[null]	Andorra	04:21:00+08:00	42.5063	1.5218	462
5	[null]	[null]	[null]	Angola	04:21:00+08:00	-11.2027	17.8739	1031
6	[null]	[null]	[null]	Antarctica	04:21:00+08:00	-71.9499	23.347	
7	[null]	[null]	[null]	Antigua and Barbuda	04:21:00+08:00	17.0608	-61.7964	90
8	[null]	[null]	[null]	Argentina	04:21:00+08:00	-38.4161	-63.6167	97084

table "10112022" (from "10-01-2022.csv")

"10112022" table was also imported the data successfully. 🤖👍👍

Query Query History

1 `select * from "10112022"`

Scratch Pad

Data Output Messages Notifications

	FIPS character varying (5)	Admin2 character varying (100)	Province_State character varying (100)	Country_Region character varying (100)	Last_Update time with time zone	Lat double precision	Long_ double precision	Confirmed integer
1	[null]	[null]	[null]	Afghanistan	04:22:00+08:00	33.93911	67.709953	2006
2	[null]	[null]	[null]	Albania	04:22:00+08:00	41.1533	20.1683	3322
3	[null]	[null]	[null]	Algeria	04:22:00+08:00	28.0339	1.6596	2707
4	[null]	[null]	[null]	Andorra	04:22:00+08:00	42.5063	1.5218	462
5	[null]	[null]	[null]	Angola	04:22:00+08:00	-11.2027	17.8739	1031
6	[null]	[null]	[null]	Antarctica	04:22:00+08:00	-71.9499	23.347	
7	[null]	[null]	[null]	Antigua and Barbuda	04:22:00+08:00	17.0608	-61.7964	90
8	[null]	[null]	[null]	Argentina	04:22:00+08:00	-38.4161	-63.6167	97135

Total rows: 1000 of 1015. Query complete: 00:00:00.350

"10012022" table(from "10-01-2022.csv")

⚠️ Since there were some data inconsistency(different names for the same country) between country-and-continent-codes-list.csv and (10-01-2022.csv & 10-11-2022.csv), I edited the country names in country-and-continent-codes-list.csv to make it consistent.

(e.g. `Afghanistan, the republic of` → `Afghanistan`)



It is difficult to choose which side of country name representation should stay, but since it is a database homework to me, I just edited the names in country-and-continent-codes-list.csv due to convenience, hope no one would be disappointed to my choice, LOVE and PEACE. ❤️❤️❤️❤️❤️



the imported result also worked fine, we are good for the next section.

Query

Query History

<

"countryAndCode" table(from "country-and-continent-codes-list-csv.csv") :

**Q3. The SQL statements and output results of 4a (10pt). If the SQL statements or output results are not provided, you will not get the points.**

4a: extract the total case number (Confirmed) in California, US on 2022-10-11 (data update date) (10 pts)

**Ans 3:**

Query Query History	
1	select sum("Confirmed") from "10112022"
2	where "Province_State" = 'California'
3	and "Country_Region" = 'US';

query codes

sum bigint	
1	11310690

result

use `select` clause to get the table from table "10112022" with `where` clause specifying the condition, which is "Province\_State" = 'California' and "Country\_Region" = 'US'.

## Q4. The SQL statements and output results of 4b (10pt)

4b: extract the total case number (Confirmed) in California, US on 2022-10-01 (data update date) (10 pts)

Ans 4:

Query	Query History
1	<code>select sum("Confirmed") from "10012022"</code>
2	<code>where "Province_State" = 'California'</code>
3	<code>and "Country_Region" = 'US';</code>

query codes

	sum bigint
1	11268292

result

use `select` clause to get the table from table "10012022" with `where` clause specifying the condition, which is "Province\_State" = 'California' and "Country\_Region" = 'US'.

## Q5. The SQL statements and output results of 4c (10pt)

4c: extract the newly diagnosed case number (Confirmed) in California, US on 2022-10-11, compared with 2022-10-01 (data update date), in one SQL statement (please don't just do 4a-4b) (10 pts)

Ans 5:

Query	Query History
1	<code>select</code>
2	<code>sum("Confirmed") - (select sum("Confirmed") from "10012022"</code>
3	<code>where "Province_State" = 'California' and</code>
4	<code>"Country_Region" = 'US')</code>
5	<code>as newly_diagnosed_case_number</code>
6	<code>from "10112022"</code>
7	<code>where "Province_State" = 'California' and</code>
8	<code>"Country_Region" = 'US'</code>

query codes

	newly_diagnosed_case_number bigint
1	42398

result

I used `Scalar Subquery` to get the total confirmed cases number from "10012022" table, and then subtracted it from the total confirmed cases number from "10112022" table to get the newly diagnosed cases number in California, US.

## 6. The SQL statements and output results of 4d (10pt)

4d. extract the country names (return Country\_Region column) and total confirmed COVID cases (return Confirmed column) with more than 20,000,000 total COVID-19 cases on 2022-10-11 (data update date) (10 pts)

## Ans 6:

Query Query History

```
1 select
2     "Country_Region",
3     sum("Confirmed") as confirmed_cases
4 from "10112022"
5 group by "Country_Region"
6 having sum("Confirmed") > 20000000;
```

query codes

	Country_Region character varying (100)	confirmed_cases bigint
1	France	36187658
2	Korea, South	25025749
3	Italy	22896742
4	US	96783524
5	United Kingdom	23957457
6	Germany	34257916
7	India	44616235
8	Japan	21593704
9	Russia	20929929
10	Brazil	34731539

result

Since one Country\_Region might have multiple different Province\_States, I grouped the data with Country\_Region using `group by` and then used the `having` clause with aggregate function(`sum > 20000000`).

## 7. The SQL statements and output results of 4e (10pt)

4e. extract the country names (return Country\_Region column) and total confirmed COVID cases (return Confirmed column) with more than 20,000,000 total COVID-19 cases on 2022-10-11 (data update date). Try to join the Country code and continents mapping table, and return only the data from countries in Asia. (10 pts)

## Ans 7:

Query Query History

```
1 select
2     "Country_Region",
3     sum("Confirmed") as confirmed_cases
4 from "10112022"
5 join "countryAndCode"
6 on "countryAndCode"."Country_Name" = "10112022"."Country_Region"
7 where "countryAndCode"."Continent_Name" = 'Asia'
8 group by "Country_Region"
9 having sum("Confirmed") > 20000000;
```

query

	Country_Region character varying (100)	confirmed_cases bigint
1	Korea, South	25025749
2	India	44616235
3	Japan	21593704
4	Russia	20929929

result

Since I need to select countries that are in Asia, combining "10112022" table with "countryAndCode" is necessary.

First, I joined "countryAndCode" with "10112022" on the data that have the same contry name

```
"countryAndCode"."Country_Name" = "10112022"."Country_Region"
```

and then filter the data whose location is in Asia

```
"countryAndCode"."Continent_Name" = 'Asia'
```

the aggregate function part is similar to Ans.6, and we successfully get the result.

## 8. The SQL statements and output results of 4f (10pt)

4f. extract the country names (return Country\_Region column) and newly diagnosed case number of countries with a newly diagnosed case number (calculate the number by yourself) > 100,000 on 2022-10-11 (compare with 2022-10-01). In descending order of newly diagnosed case numbers. (10 pts)

**Ans 8:**

Query	Query History
1	with
2	old_data(country_region, old_confirmed) as(
3	select "Country_Region", sum("Confirmed")
4	from "10012022"
5	group by "Country_Region"
6	),
7	new_data (country_region, new_confirmed) as(
8	select "Country_Region", sum("Confirmed")
9	from "10112022"
10	group by "Country_Region")
11	select
12	"old_data"."country_region",
13	("new_data"."new_confirmed" - "old_data"."old_confirmed") as newly_diagnosed_case
14	from "old_data"
15	join "new_data" on "new_data"."country_region" = "old_data"."country_region"
16	where "new_data"."new_confirmed" - "old_data"."old_confirmed" > 100000
17	order by newly_diagnosed_case desc

query codes

	country_region character varying (100)	newly_diagnosed_case bigint
1	Germany	871687
2	France	579373
3	Taiwan*	440596
4	Italy	396396
5	US	386493
6	Japan	264185
7	Russia	212106
8	Korea, South	206138
9	Austria	129544
10	Greece	106302

result

Since the newly diagnosed number comes from two different tables("10012022" and "10012022"), I first created two temporary table using the `with` clause. (old\_data and new\_data)

After getting two new temporary tables with **Country\_Region** and **Confirmed** on the specific day respectively, I did subtraction of newly confirmed cases(from new\_data table) and the confirmed cases(from old\_data table) to get **newly diagnosed case number**.

I then filtered the table with the **newly\_diagnosed\_case>10000** condition and ordered the number with **descending order** to get the final output table.

