

# ECS 116 Databases for Non-Majors / Data Management for Data Science

## Programming Assignment 1

### A. Prelude

1. The assignment is of 25 points.
2. Last date of submission is April 26, 2024, Friday @ 11:59 pm.
3. Late submissions will be graded according to the late policy. Specifically, 10% of grade is deducted if you are up to 24 hours late, 20% is deducted if you are 24 to 48 hours late, and no credit if turned in after 48 hours.
4. This assignment will be solo.
5. Create a new sql file for each step namely (Step\_2, Step\_3, Step\_4) if you have to use sql commands through DBeaver.
6. Your assignment will be graded based on correctness (passing all tests), ingenuity and originality.
7. All the required files (csv) can be found under *Files* in Canvas.
8. Plagiarism is strictly prohibited. You're free to discuss high-level concepts amongst your peers. However, cheating will result in no points on the assignment and reporting to OSSJA.

### B. Step 1: Uploading *africa\_fs\_after\_cleaning\_db.csv* into PostgreSQL

1. In DBeaver create a new database **faostat**. Set that as the default database
2. Create a schema **food\_sec** (or "food\_sec\_v01") in your database **faostat**. Set that as default schema.
3. Do set *search\_path* to *food\_sec*;
4. Load the file **africa\_fs\_after\_cleaning\_db.csv** into the schema **food\_sec** to make table **africa\_fs\_ac**.
5. Modify the data types of some of the columns of **africa\_fs\_ac** as follows:
  - area\_code\_m49: varchar(3)
  - element\_code: varchar(4)
  - year\_code: varchar(8)
  - value: numeric
  - After making these changes, click on "Save" at bottom of pane.
6. Check whether the values for **value** column have been imported correctly.
  - Do a selection query to get distinct values that are  $\leq 2$ .
  - Using Excel see what are the values  $\leq 2$ .
  - Do these match?
7. Do an SQL query to DELETE all tuples from **africa\_fs\_ac** (it will ask you to confirm that you want to do this delete).
8. Use DBeaver to *import* the file **africa\_fs\_after\_cleaning\_db.csv** (don't use the SQL "COPY" command because it complains about a data type encoding issue).
  - Do a sanity check that the number of tuples in your table is same as in csv file.
  - Again check on the values in column **value**.

## C. Step 2: Build Table *gdp\_stunting\_overweight\_anemia*

- Similar to the construction of **gdp\_stunting\_overweight** shown in the 2024-04-09 lecture and the SQL script **faostat-part\_02-transforming\_africa\_fs.sql**, use DBeaver and SQL commands to build a table **gdp\_stunting\_overweight\_anemia** which has, for each country-year pair the following associated values for:
  - GDP per capita Purchasing Power Parity (22013): use column name **gdp\_p\_ppp**.
  - Percentage of children over 5 years of age who are stunted (21025): **childhood\_stunting**
  - Percentage of children over 5 years of age who are overweight (21043): **childhood\_overweight**.
  - Prevalence of anemia among women of reproductive age: **anemia**
- Add this table into your schema **food\_sec**.

	area_code_m49	area	year_code	gdp_pc_ppp	childhood_stunting	childhood_overweight	anemia
1	12	Algeria	2000	8,786	22	12	37
2	12	Algeria	2001	8,926	21	13	37
3	12	Algeria	2002	9,299	20	14	36
4	12	Algeria	2003	9,835	19	14	36
5	12	Algeria	2004	10,114	18	15	35
6	12	Algeria	2005	10,566	18	15	35
7	12	Algeria	2006	10,592	17	15	34
8	12	Algeria	2007	10,775	16	15	34
9	12	Algeria	2008	10,847	15	15	34

Figure 1: Almost correct example of the table *gdp\_stunting\_overweight\_anemia*. Your table should have 3 characters for area\_code\_m49 column, and may have some decimal values for the last 4 columns.

## D. Step 3: Build table *energy\_undernourished*

- Note that many records in **africa\_fs\_ac** have *year* and *year\_code* values based on 3-year intervals rather than single years. We will use some of this data to gain more insight about countries. In particular, we will interpret a 3-year interval as applying to the year in the middle, e.g., we will interpret *2000-2002* as applying to the year 2001.
- First, build a table **energy\_undernourished** which has, for each *country-year\_code* pair the associated values for:
  - Average dietary energy supply adequacy (21010): use column name **dietary\_energy**.
  - Prevalence of undernourishment (210041): use column name **undernourished**.
  - Note: this table should have 1040 rows in it.
- Now add a column *derived\_year* to the table **energy\_undernourished**, where for each tuple, the derived year value is computed by using the year in the middle of the first and third years in the *year\_code* of the tuple.
- The column you added probably has data type integer. Convert this to varchar(4).

	ABC area_code_m49	ABC area	ABC year_code	123 dietary_energy	123 undernourished	ABC derived_year
1	12	Algeria	2000	127	8	2001
2	12	Algeria	2001	129	7	2002
3	12	Algeria	2002	130	7	2003
4	12	Algeria	2003	130	7	2004
5	12	Algeria	2004	131	6	2005
6	12	Algeria	2005	132	6	2006
7	12	Algeria	2006	133	5	2007
8	12	Algeria	2007	135	5	2008
9	12	Algeria	2008	136	5	2009
10	12	Algeria	2009	139	4	2010

Figure 2: Almost correct example of the table *energy\_undernourished*. As with Figure 1, the *area\_code\_m49* column should have 3 characters, and the values for last 3 columns may have decimal values.

## E. Step 4: Joining the *gdp\_stunting\_overweight\_anemia* and *energy\_undernourished* tables to create new table *gdp\_energy\_with\_fs\_indicators*

1. Create a selection query that combines the table **gdp\_stunting\_overweight\_anemia** and **energy\_undernourished** to form a new table **gdp\_energy\_with\_fs\_indicators**
  - The columns should include *area\_code\_m49*, *area*, *year\_code*, *gdp-pc-ppp*, *dietary-energy*, *childhood-stunting*, *childhood-overweight*, *anemia* and *undernourished*.
  - Tuples in this table should be formed by combining tuples from **gdp\_stunting\_overweight\_anemia** and **energy\_undernourished** where *year\_code* from the first table equals *derived\_year* of the second table.
  - Note: your table should have 895 tuples in it.
2. Export the table **gdp\_energy\_with\_fs\_indicators** as a csv file **gdp\_energy\_with\_fs\_indicators.csv**.
3. Sort this csv file by *area* (country name) and then *year\_code*.
4. CONGRATULATIONS: you have created a table that we can use later to determine whether there are statistical correlations between gdp per capita and/or stunting, childhood overweight, anemia in women and/or undernourishment.

	ABC area_code_m49	ABC area	ABC year_code	123 gdp_pc_ppp	123 dietary_energy	123 childhood_stunting	123 childhood_overweight	123 anemia	123 undernourished
1	12	Algeria	2001	8,926	127	21	13	37	8
2	12	Algeria	2002	9,299	129	20	14	36	7
3	12	Algeria	2003	9,835	130	19	14	36	7
4	12	Algeria	2004	10,114	130	18	15	35	7
5	12	Algeria	2005	10,566	131	18	15	35	6
6	12	Algeria	2006	10,592	132	17	15	34	6
7	12	Algeria	2007	10,775	133	16	15	34	5
8	12	Algeria	2008	10,847	135	15	15	34	5
9	12	Algeria	2009	10,824	136	14	15	33	5
10	12	Algeria	2010	11,007	139	13	14	33	4

Figure 3: Almost correct example of the table *gdp\_energy\_with\_fs\_indicators*

5. Create a new table **gdp\_energy\_fs\_aggs**.
  - Which has columns as:
    - *area\_code\_m49*
    - *area*
    - *avg-gdp-pc-ppp*
    - *avg-dietary-energy*

- *avg\_childhood\_stunting*
- *avg\_childhood\_overweight*
- *avg\_anemia*
- *avg\_undernourished*

- The “avg” columns should hold the averages of the corresponding items for each country, over all of the years of available data.
- Use the round operator on the “avg” value, so that they have type numeric and are rounded to 2 decimal points. Use the following kind of expression: *round(< expression for average >::numeric, 2)*.

6. Export the table **gdp\_energy\_fs\_aggs** as a csv file **gdp\_energy\_fs\_aggs.csv**.

7. The table should be sorted by area (i.e., country name). (You can use ORDER BY in the query or sort the csv file once you have created it.)

	RBC area_code_m49	RBC area	123 avg_gdp_pc_ppp	123 avg_dietary_energy	123 avg_childhood_stunting	123 avg_childhood_overweight	123 avg_anemia	123 avg_undernourished
1	12	Algeria	10,888.95	139.47	13.95	13.37	33.63	4.11
2	24	Angola	6,953	100.95	34.79	2.95	47	32.21
3	204	Benin	2,715.68	118.32	34	1.58	57.26	10.21
4	72	Botswana	13,311.74	102	25.74	9.84	32.47	23.32
5	854	Burkina Faso	1,676.58	117.63	34.37	1.74	53.53	15.32
6	120	Cameroon	3,379.05	117.26	33.21	8.16	41.11	9.53
7	140	Central African Republic	968.37	90.68	41.42	4.21	48.21	39.63
8	148	Chad	1,560.47	101.53	39.16	2.63	49.37	31.47
9	174	Comoros	3,047.26	107.58	33.89	12.37	33.84	15.42
10	178	Congo	4,574.47	95	25.42	5.32	53.47	31.89

Figure 4: Almost correct example of the table *gdp\_energy\_fs\_aggs*. You will obtain slightly different values. This table was computed with rounded values for various columns, rather than with values having decimals.

## F. Submission

1. Please make a single zip file that includes

- *gdp\_energy\_with\_fs\_indicators.csv*
- *gdp\_energy\_fs\_aggs.csv*
- The DBEaver sql scripts that you used to create these 2 csv files, specifically, Step\_2.sql, Step\_3.sql, Step\_4.sql.
- Name the zip file as FirstName\_LastName\_LastFourDigitsOfStudentID\_ECS116\_A1.

2. Upload it on Canvas for Assignment 1 (This is a solo assignment so don't add your peers to your submission).