CS 329E Project 8, due Thursday, 11/13.

# **Objectives**

Orchestrate the staging layer of the lakehouse with dbt

## In-Scope

• CTAS statements from project5-stg-layer.ipynb

### **Out-of-Scope**

• **SELECT statements from** project5-stg-layer.ipynb

#### Work Items

- Create a new dbt project
- Create the configuration files
- Create the folder structure
- Create the macros
- Create and run the models
- Create and run the data validation tests
- Generate the docs and lineage graph

# **Code Samples**

See <u>snippets</u> repo for code samples

# **Implementation Details**

- Go to the dbt folder that you created during the <u>setup</u> and run the dbt init command to create a new dbt project. Name your dbt project [your\_domain] and name your dbt dataset dbt\_[your\_domain]\_stg.
- Open your profile file profiles.yml (located in your home directory) and change the region of the dataset from us to us-central1.
- Open the dbt\_profile.yml file and change the model section of the file to look like this:

```
models:
    air_travel:
    # Config indicated by + and applies to all files under models/example/
    stg:
        +materialized: table
        +schema: dbt_air_travel_stg
    int:
```

```
+materialized: table
    +schema: dbt_air_travel_int
mrt:
    +materialized: table
    +schema: dbt_air_travel_mrt
```

For the schema values, make sure you replace air travel with your domain.

- Open the models folder and delete the <code>examples</code> subfolder. Create the following three folders in place of the <code>examples</code> folder:
  - o stg
  - o int
  - o mrt
- Create a file <code>generate\_schema\_name.sql</code> in the macros folder of your dbt project and copy the contents from <a href="here">here</a> into it. This macro is needed to overwrite the default naming of datasets which would yield <code>dbt\_air\_travel\_stg\_dbt\_air\_travel\_stg</code> for the dataset name.
- Open the stg folder and create a file named schema.yml in it. This file is where we specify the source for each staging table. Follow this example to specify the sources in the file. What this will do is register the raw tables with dbt and allow us to refer to the raw tables with dbt's source() function. Note: we are not going to create the raw dataset with dbt, we will use the existing raw dataset. You do not need to add the models section yet.
- Create your dbt model files for the staging layer. Place the files in the stg folder and make sure you source the raw table with the source () function.
- Run the <a href="dbt run">dbt run</a> command to compile your models. You should do this iteratively so that you can easily identify the source of the errors.
- Create a second macro in the dbt macros folder to validate the row count of the staging models. Name the file check\_row\_count\_equals.sql and copy its contents from here.
  This function compares the row counts of the staging and raw tables and raises an error if they don't match.
- Open the schema.yml file that you previously created. Add a data validation test for each model. The test should call the check\_row\_count\_equals macro to compare the row counts. Follow this example to see how to format your tests.
- Run the dbt test command to run the tests. Make sure that all of your tests pass.

- Generate the documentation for your project and verify that the lineage graph captures
  all your staging models. Refer to our setup guide if you don't remember how to do this
  step. Take a screenshot of your lineage graph and place it in your dbt project folder
  under a new subfolder named lineage. Name the screenshot file staging-layer.png.
- Make a new project8 folder in your local git repo. Copy your entire dbt project folder into it (e.g. [your\_domain]). Edit the .gitignore file in your dbt project folder to not exclude the target/ and log/ folders from git. Commit and push your changes to GitHub.
- Once you have pushed your changes, create the usual submission.json file and upload it to Canvas by the deadline. Only one person per group needs to do this step.

Due Date: 11/13/25