

Analytics Engineering Data Pipeline

Author: Veronika Folin

February 28, 2024

Contents

1	Introduction	3
2	Dataset	5
3	Implementation	6
3.1	Data Storage with Snowflake	6
3.2	Automated Data Ingestion from Cloud Storage to Snowflake	6
3.3	Transformation with dbt	6
3.4	Data Observability with Elementary and Slack	13
3.5	Document dbt Project	16
3.6	Orchestration with Cloud Composer	16
3.7	Alerting on Task Failures with SendGrid	18
3.8	Dashboarding with Looker Studio	19
4	Setup	23
4.1	Snowflake	23
4.2	dbt Project	23
4.3	Elementary and Slack Alerts	24
4.4	Google Cloud Platform	25
4.5	Automated Data Ingestion from Cloud Storage to Snowflake	27
4.6	Synchronize GitHub repository with Cloud Storage bucket	33
4.7	Hosting dbt Documentation in a GitHub Pages	34
4.8	Connect Airflow with Snowflake	35
4.9	SendGrid	36
4.10	Customize email on task failure	38
4.11	Looker Studio	39
5	Usage	40
5.1	Simulate Data Ingestion from Cloud Storage to Snowflake	40
5.2	Transform with dbt	40
5.3	Data observability with Elementary and Slack	41
5.4	Orchestrating with Cloud Composer	41
5.5	Dashboarding on Looker Studio	41
6	Unset up	42

1 Introduction

The project aims to design, implement, test, and document a solution that covers the entire data pipeline, from ingestion to analysis. The process must follow the principles of DataOps and Analytics Engineering, taking advantage of the **Data Build Tool (dbt)** data transformation tool. In the set of technologies, it is required to make use of the **Snowflake** data management system and the **Google Cloud Platform** cloud environment.

As shown in Figure 1, the solution allows to:

- Automatically manage data ingestion from Cloud Storage to Snowflake.
- Collect raw data on Snowflake in a format suitable for the transformation process.
- Transform, test, and document data with dbt. This process enables us to clean, normalize, enrich, and prepare the data for analysis and reporting.
- Monitor the transformation process using the Elementary package and configure Slack Alerts in case of errors.
- Collect transformed data on Snowflake.
- Orchestrate the transformation process using Cloud Composer and receive alerts (emails) via the SendGrid service when the workflow fails.
- View and analyze transformed data using the Looker Studio dashboarding tool, and monitor its data quality.
- Manage automatically via GitHub Actions:
 - Synchronization between the project repository and the execution environment that orchestrates the transformation process.
 - The deployment of the documentation on GitHub Pages.

The dataset used to evaluate the proposed solution will be discussed in detail in Section 2. The implementation details of each module will be explained in Section 3. By following the steps outlined in Section 4, it will be possible to set up and configure the required technologies for the project. The usage of the project is further explained in Section 5. Finally, to deactivate the paid services that were configured, you can refer to the steps provided in Section 6.

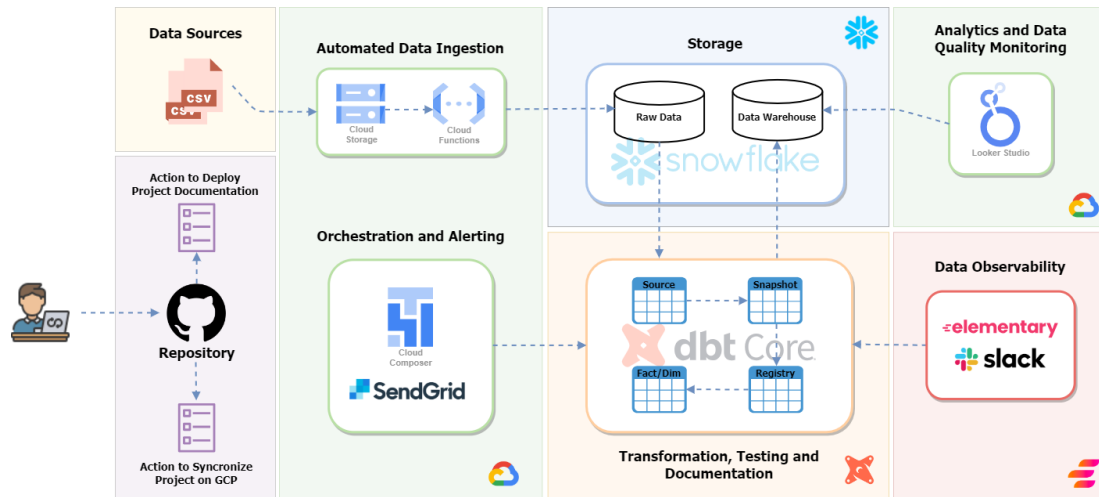


Figure 1: Architecture of the proposed solution.

2 Dataset

The solution was evaluated using data from the TPC Benchmark™ H (TPC-H), which includes datasets of different sizes to test scalability. For this project, we utilized the smaller version (in total it takes up 1 GB). The TPC-H dataset is accessible on Snowflake (Snowsight) and can be found in the TPCH_SF1 schema within the SNOWFLAKE_SAMPLE_DATA shared database.

TPC-H consists of eight tables and the data populating the database have been chosen to have broad industry-wide relevance, as depicted in Figure 2.

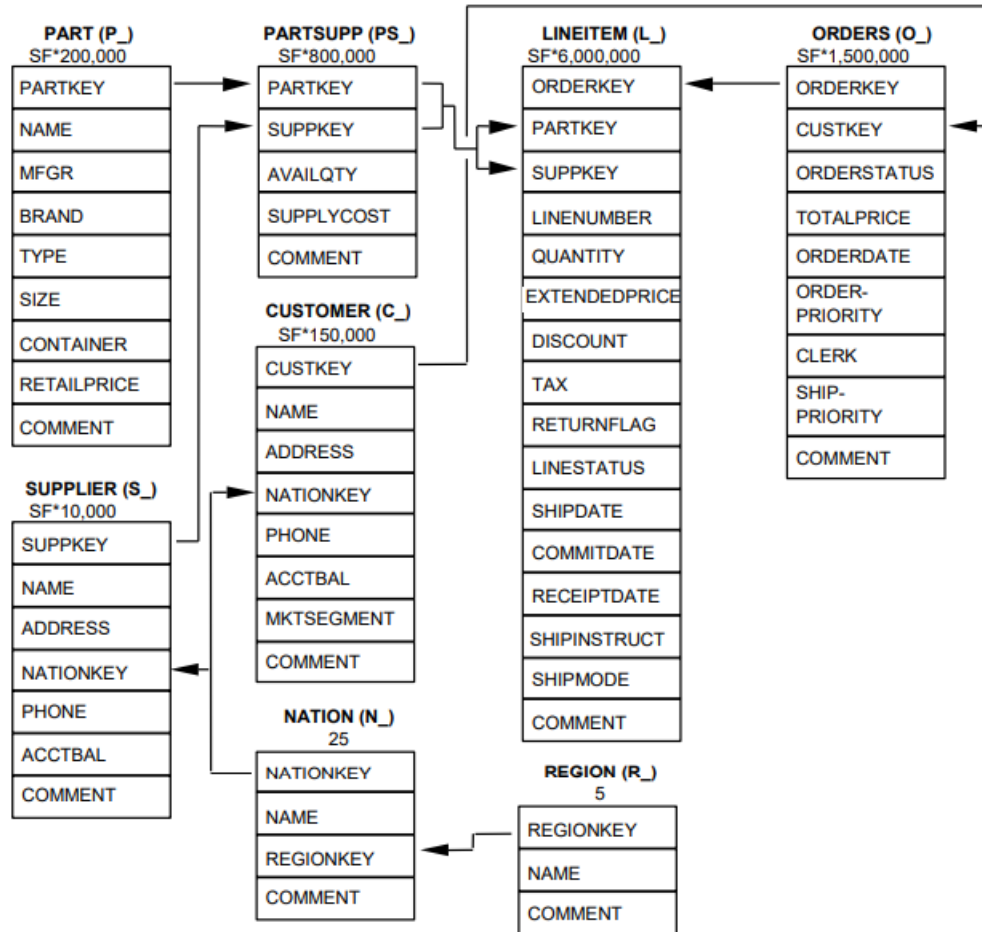


Figure 2: Source: TPC Benchmark H Standard Specification.

3 Implementation

3.1 Data Storage with Snowflake

Snowflake was leveraged as the project's data storage platform. Here there are both raw data (RAW database) and data that have undergone a transformation and/or quality verification process (ANALYTICS database) and which are ready to be used in an analysis or reporting system. The computing capacity in Snowflake is represented by the *warehouse* concept, that is a cluster of machines configurable according to needs. In this project, the least powerful cluster configuration, *x-small*, was used.

To test Snowflake in the context of this project, see Section 4.1.

3.2 Automated Data Ingestion from Cloud Storage to Snowflake

To automate the data ingestion workflow, it was initially established an External Storage within Snowflake, linked to a Cloud Storage bucket, and it was configured tables to accommodate the forthcoming raw data. Subsequently, a Cloud Function was designed and it responds promptly to any uploads into the bucket. Upon activation, this function establishes a connection with the data warehouse and executes SQL commands to transfer data from the external storage to the specified destination tables, completing the data loading process.

The configuration steps are explained in detail in Section 4.5.

3.3 Transformation with dbt

dbt allows you to define the transformation logic in a modular way, by creating models implemented as select statements in the SQL language. Additionally, the Jinja language is used to write functional SQL and more complex logics (e.g., references and macros).

In the following, we will explain in detail how the data transformation phase was implemented.

Sources Sources represent the raw data within the data warehouse that need to be transformed, and their definition is found in the `sources.yml` file.

```
1 version: 2
2
3 sources:
4   - name: raw
5     database: raw
6     schema: analytics_engineering_data_pipeline
7     tables:
8       - name: customer
9       - name: lineitem
10      - name: nation
11      - name: orders
12      - name: part
13      - name: partsupp
```

```

14         - name: region
15         - name: supplier
16     - name: elementary
17       database: analytics
18       schema: analytics_engineering_data_pipeline_elementary
19       tables:
20         - name: dbt_tests
21         - name: elementary_test_results
22     - name: metadata
23       database: analytics
24       schema: information_schema
25       tables:
26         - name: tables
27         - name: views

```

Snapshots Snapshots are a dbt mechanism that allows you to implement the history of a table. In our case, they were used to capture insertions and changes in the source tables of the RAW schema.

For example:

```

1  {% snapshot snapshot_lineitem %}
2
3      {{
4          config(
5              target_database='analytics',
6              target_schema='snapshots',
7              strategy='check',
8              unique_key='lineitemkey',
9              check_cols='all'
10         )
11     }}
12
13     select *,
14     {{ dbt_utils.generate_surrogate_key(['l_orderkey',
15         'l_linenumber']) }}
16     as lineitemkey,
17     {{ dbt_utils.generate_surrogate_key(['l_partkey',
18         'l_suppkey']) }}
19     as partsuppkey
20     from {{ source('raw', 'lineitem') }}
21
22 {% endsnapshot %}

```

Seeds Seeds are CSV files that can be loaded into the data warehouse to store static data which change infrequently.

In this project, they have been used to associate each type of test (**test_name**) with a tag (**test_tag**) and to associate each model of the project (**table_ref**) with a tag

(`model_tag`). This association allows you to enrich the metadata used to perform data observation and data quality analysis.

An example:

```
1 TEST_NAME,TEST_TAG
2 accepted_values,validity
3 accepted_range,validity
4 not_null,completeness
5 relationships,completeness
6 unique,uniqueness
7 equal_rowcount,consistency
8 unique_combination_of_columns,uniqueness
9 expect_column_values_to_be_of_type,validity
10 expect_column_values_to_be_in_set,validity
```

Macros Macros are pieces of code that can be reused multiple times in models. They can be generic or singular.

Generic

1. `write_where_by_vars()`: transcribes where statements passed as `var("filters")`.
2. `write_select_groupByColumns_by_vars()`: transcribes the select statement, selecting the fields passed as `var("groupBy")` that the user intends to use to perform the aggregation.
3. `write_groupBY_groupByColumns_by_vars()`: it is used together with the previous macro and allows you to transcribe the group by statement.
4. `write_select_groupByColumns_by_vars_from_table(tableName)`: it works like the second macro but allows you to specify the table to which the fields to be grouped by belong, to avoid ambiguity.
5. `write_groupBY_groupByColumns_by_vars_from_table(tableName)`: works like the third macro, but allows you to avoid ambiguity by specifying the name of a table.
`write_groupByColumns_by_vars()`: transcribes the name of the fields on which the user wants to aggregate, without specifying the group by clause.
6. `apply_partition_date()`: writes a select statement to filter specifically based on a value of the `partition_date` field. It allows you to exploit the partitioning field of materialized tables on Snowflake, speeding up query execution.
7. `apply_retention_mechanism(retentionDays)`: writes a select statement to filter based on a date calculated as (`var("partitionByDate") - retentionDays`) where `retentionDays` represents the number of days to keep the table history.

Singular

1. `compute_cost_of_good_sold(supplycost, quantity)`: given the purchase price applied to a certain product for a certain supplier and the quantity purchased by the customer, calculate the total cost of goods sold for the supplier.
2. `compute_discounted_extended_price(extendedprice, discount)`: calculates the discounted price by considering the extended price of a line item and the applied discount.
3. `compute_discounted_extended_price_plus_tax(extendedprice, discount, tax)`: calculates the total amount by applying the specified tax percentage to the discounted extended price.
4. `compute_profit(net_revenue, supplycost, quantity)`: calculates profit as the difference between net income and cost of goods sold.

Models In the project, the models were categorized into different folders, based on what aspect of the domain they covered: individuals, places, products and sales.

Staging Staging models are the first transformation step starting from sources. They involve renaming, type casting, generation of surrogate keys and simple computations (for example, using macros). No joins or aggregations are performed in this phase. Models of this type are named as `stg_<model_name>` and are placed in the `models/staging` directory.

Registries The registers represent the **historicized** version of the staging models: they read from snapshots, rather than directly from sources, and are materialized as **incremental** models, so that dbt transforms only the rows in your source data that you tell dbt to filter for in the `is_incremental` macro. Furthermore, it was decided to exploit the **partitioning** mechanism provided by Snowflake based on a field that indicates the instant of acquisition of a record within the register (`PARTITION_DATE`), and the optional `on_schema_change` parameter has been configured to `'append_new_columns'`, so that **new columns are added** to the existing table but those no longer present in the new data are not removed. Models of this type are named as `registry_stg_<model_name>` and are placed in the `models/staging` directory.

```
1  {{
2      config(
3          cluster_by=['partition_date'],
4          materialized='incremental',
5          on_schema_change='append_new_columns'
6      )
7  }}
```

8

```
9  with
```

```

10
11 last_snapshot as (
12     select *
13     from {{ref('snapshot_lineitem')}}
14     where DATE(DBT_VALID_FROM) = (select
15         MAX(DATE(DBT_VALID_FROM)) from
16         {{ref('snapshot_lineitem')}})
17 ),
18
19 previous_state_of_registry as (
20     select *
21     from {{this}}
22     where partition_date = (select MAX(partition_date) from
23         {{this}})
24 ),
25
26 final as (
27     select
28         COALESCE(new.lineitemkey, old.lineitemkey) as lineitemkey,
29         COALESCE(new.l_orderkey, old.orderkey) as orderkey,
30         COALESCE(new.l_linenummer, old.linenummer) as linenummer,
31         COALESCE(new.l_partkey, old.partkey) as partkey,
32         COALESCE(new.l_suppkey, old.suppkey) as suppkey,
33         COALESCE(new.partsuppkey, old.partsuppkey) as partsuppkey,
34         CAST(COALESCE(new.l_quantity, old.quantity) AS int) as
35         quantity,
36         COALESCE(new.l_extendedprice, old.extendedprice) as
37         extendedprice,
38         COALESCE(new.l_discount, old.discount) as discount,
39         COALESCE(new.l_tax, old.tax) as tax,
40         COALESCE(new.l_returnflag, old.returnflag) as returnflag,
41         COALESCE(new.l_linestatus, old.linestatus) as linestatus,
42         COALESCE(new.l_shipdate, old.shipdate) as shipdate,
43         COALESCE(new.l_commitdate, old.commitdate) as commitdate,
44         COALESCE(new.l_receiptdate, old.receiptdate) as
45         receiptdate,
46         COALESCE(new.l_shipinstruct, old.shipinstruct) as
47         shipinstruct,
48         COALESCE(new.l_shipmode, old.shipmode) as shipmode,
49         CURRENT_DATE() as partition_date
50     from last_snapshot as new FULL OUTER JOIN
51         previous_state_of_registry as old ON new.lineitemkey =
52         old.lineitemkey
53 )
54
55 select * from final
56
57 {% if is_incremental() %}
58
59

```

```

50     where partition_date > (select max(partition_date) from {{ this
      }})
51
52 {% endif %}

```

Marts Marts are meant to represent a specific entity or concept at its unique grain, and put together (through joins or aggregations) the information collected in the staging models. Also in this case the models are organized in folders by concept. At this level, the tables are ready to be analyzed. In fact, we find fact and dimension tables, tables that calculate KPIs (e.g., `kpi_customer_churn_rate`, `kpi_gross_profit_margin`, etc.) and tables with summary values, ready to be displayed in dashboards (e.g., `acquired_customer`, `volume_sales`, etc.). The marts tables are configured similarly to the corresponding registry or staging tables (e.g., incremental, clustered and historicized). For fact tables, a **retention mechanism** set at one week was applied through the execution of a `post_hook`, namely a function that is executed after the materialization of the table. Models of this type are placed in the `models/marts` directory.

Tests In a dbt project, the tests are defined in a yaml file, simultaneously with the definition of the models and the fields that compose them, duly documented using the `description` property.

```

1  version: 2
2
3  models:
4    - name: registry_stg_orders
5      description: Snapshot registry of customers' orders data.
6      tests:
7        - dbt_utils.unique_combination_of_columns:
8          combination_of_columns:
9            - orderkey
10             - partition_date
11      columns:
12        - name: orderkey
13          description: Primary key for an order.
14          tests:
15            - not_null
16        - name: custkey
17          description: Foreign key to registry_stg_customer.custkey.
18          tests:
19            - not_null
20            - relationships:
21              to: ref('registry_stg_customer')
22              field: custkey
23        - name: orderstatus
24          description: '{{ doc("orderstatus") }}'
25          tests:
26            - not_null

```

```

27         - accepted_values:
28             values:
29                 - F
30                 - 0
31                 - P
32     - name: totalprice
33       description: Total price of the order.
34       tests:
35         - not_null
36         - dbt_utils.accepted_range:
37             min_value: 0
38     - name: orderdate
39       description: Date of the order.
40       tests:
41         - not_null
42         - dbt_utils.accepted_range:
43             max_value: "getdate()"
44         - dbt_expectations.expect_column_values_to_be_of_type:
45             column_type: date
46     - name: orderpriority
47       description: Priority of the order.
48       tests:
49         - not_null
50         - accepted_values:
51             values:
52                 - 1-URGENT
53                 - 2-HIGH
54                 - 3-MEDIUM
55                 - 4-NOT SPECIFIED
56                 - 5-LOW
57     - name: clerk
58       description: Identification of the employee who processed
59                   the order.
60       tests:
61         - not_null
62     - name: shippriority
63       description: Shipping priority.
64       tests:
65         - not_null
66     - name: partition_date
67       description: Time when this snapshot row was inserted.
68       tests:
69         - not_null
70         - dbt_utils.accepted_range:
71             max_value: "getdate()"
72         - dbt_expectations.expect_column_values_to_be_of_type:
73             column_type: date

```

To implement them, the built-ins of dbt (e.g., unique, not_null, relationships

and `accepted_values`) and modules made available by dbt Package Hub ¹, such as `dbt_utils` ² and `dbt_expectations` ³, were exploited. This is a list of the tests carried out:

- `accepted_values`
- `accepted_range`
- `not_null`
- `relationships`
- `unique`
- `equal_rowcount`
- `unique_combination_of_columns`
- `expect_column_values_to_be_of_type`
- `expect_column_values_to_be_in_set`

3.4 Data Observability with Elementary and Slack

Elementary ⁴ is a dbt native package for data observability. The use of Elementary allows you to collect information on the execution of the runs and the results of the tests. The package allows you to automatically produce reports (as in Figure 3), but in this project it was decided to create a personalized visualization to monitor data quality.

To achieve this, a transformation process was implemented that starts from the tables generated by the Elementary package in the `analytics_engineering_data_pipeline_elementary` schema, that are `dbt_tests` and `elementary_test_results`. Models that allow the transformation process of data quality tables are found in subdirectories called `data quality`.

The first step involves the creation of:

- `stg_dbt_tests`: general metadata on the tests performed.
- `stg_elementary_test_results`: information on the execution of the tests performed. This model is implemented as a registry to have a history of the information collected at each materialization. This mechanism is exploited, in particular, to calculate the difference between the number of failures in the most recent test phase and that obtained in the previous run. This is necessary to correctly calculate the number of failures in the last materialized partition: the calculated

¹dbt Package Hub: <https://hub.getdbt.com/>

²`dbt_utils`: https://hub.getdbt.com/dbt-labs/dbt_utils/latest/

³`dbt_expectations`: https://hub.getdbt.com/calogica/dbt_expectations/latest/

⁴Elementary Documentation: <https://docs.elementary-data.com/introduction>

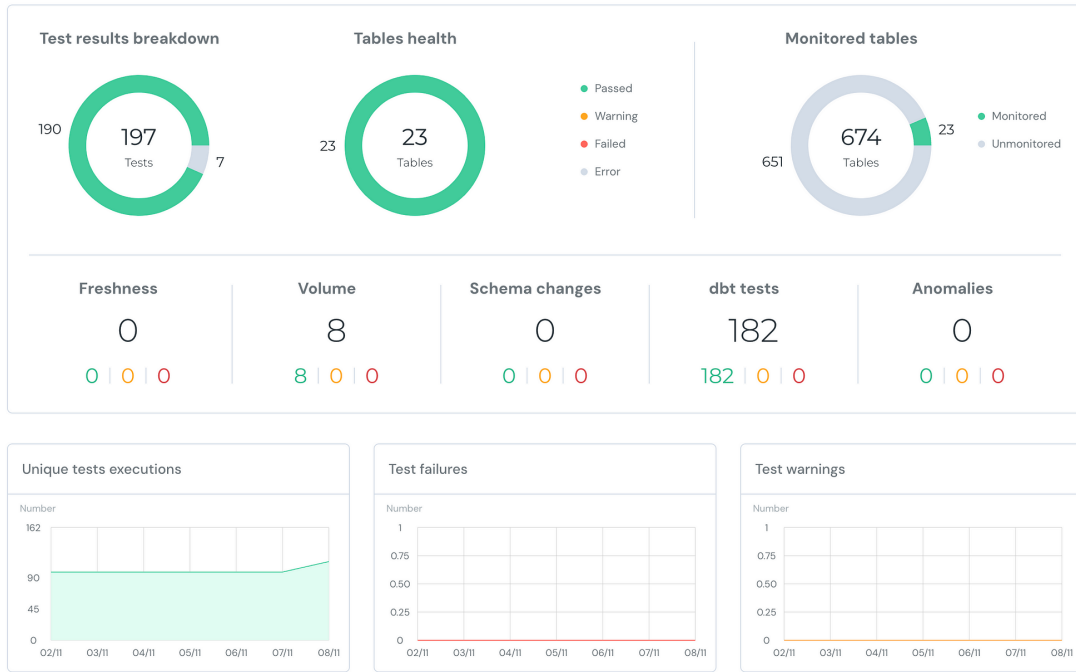


Figure 3: Example of a report automatically generated by Elementary.

delta (`failed_row_count_delta`) will ignore failures generated by rows belonging to materializations that are not the last one realized.

- **metadata_test**: metadata about tables on which tests were performed. In particular, it calculates the delta (`row_count_delta`) between the number of rows currently valid in the table and the number of rows valid in the previous materialization, in order to extrapolate the information on the number of rows belonging to the last partition created.

The second step involves the materialization of:

- **fct_test_results**: joined information about tests metadata (`stg_dbt_tests`, `metadata_test`, `test_tags` and `model_tags`) and tests results (`stg_elementary_test_results`). This model has also been configured as a registry.
- **monitor_dataquality**: summarized information regarding the execution of the tests.

Slack Alerts ⁵ have been configured, for example, in the event of a test or run failure. Figure 4 shows an example of a message received in the event of a failed test. All configuration steps are described in Section 4.3.

⁵Setup Slack Alerts: <https://docs.elementary-data.com/oss/guides/send-slack-alerts>

veronika

12:12 PM

joined #project. Also, Elementary joined.

Elementary

APP

3:26 PM

▲

dbt test alert

Test: accepted_range - Generic | Status: fail | 2023-11-15 15:23:15

Table

analytics.dbt_core_reusable_demo.stg_lineitem

Column

discount

Owners

No owners

Description

No description

Tags

No tags

Subscribers

No subscribers

Result

Result message

Got 6001215 results, configured to fail if != 0

Test results sample

[{'orderkey': 2400001, 'linenumber': 1, 'lineitemkey': '32491e401e44d19b6d4692a24487e7d4', 'partkey': 132304, 'suppkey': 4818, 'partsuppkey': '38dc7e5f062744cb73dd2ff6ac24ba6c', 'quantity': 10, 'extendedprice': 13363.0, 'discount': 0.03, 'tax': 0.02, 'returnflag': 'N', 'linestatus': '0', 'shipdate':...}

See more

Test query

with meet_condition as(
 select *
 from analytics.dbt_core_reusable_demo.stg_lineitem
)
,

See more

Configuration

Test parameters

{"min_value": 0.5, "max_value": 1, "column_name": "discount", "model": "{get_where_subquery(ref('stg_lineitem'))}"}

See less

Figure 4: Example of a Slack Alert on test failure.

15

3.5 Document dbt Project

The documentation that can be automatically generated using the `dbt docs generate` command, has been versioned in the repository within the `docs` folder, and has been hosted on GitHub Pages⁶ to be easily accessible and immediately viewable.

To reproduce this feature, see Section 4.7

3.6 Orchestration with Cloud Composer

Cloud Composer is a workflow orchestration service provided by Google Cloud Platform and built on Apache Airflow, with which you can define a series of tasks (DAG) for ingesting, transforming, analyzing, or utilizing data.

In this project, Cloud Composer was exploited to orchestrate the transformation process with dbt. The source code of the defined DAGs is organized as follows:

- `dags`:
 - `dag_factory_version/historical`: It defines DAGs using the `dag-factory` library⁷ and allows the historicized version of the tables to be materialized:
 - * `setup`: The *`setup_project`* dag debugs connections defined in the `profile.yml` file, installs project dependencies, and creates registries and Elementary tables on Snowflake.
 - * `data_factory`: The *`materialize_data`* dag sequentially triggers the execution of the dags necessary to materialize all the fact and dimension tables.
 - * `places_factory`: The *`int_nation`* dag first checks whether the `int_nation` table already exists in the data warehouse. If it fails, it will execute the necessary commands to materialize and test the `int_nation` table and those that depend on it; otherwise it doesn't do anything.
 - * `products_factory`: As shown in Figure 5, the *`dim_part`* dag takes a snapshot of the source and, at the same time, checks whether the register table `registry_stg_part` is empty or not. If it is empty, to obtain the correct behavior during the materialization of the incremental table, we will need to execute the run command with the `--full-refresh` option. Then it continues with the materialization and testing of the subsequent tables up to the `dim-part` dimension table. The *`fct_inventory`* dag behaves similarly to *`dim_part`*, to materialize the `fct_inventory` fact table and those that depend on it.
 - * `individuals_factory`: *`dim_customer`* and *`dim_supplier`* dags work the same as `dim_part` but are used respectively to materialize and test

⁶dbt Project Documentation GitHub Pages: https://veronikafolin.github.io/analytics_engineering_data_pipeline/#!/overview

⁷dag-factory library documentation: <https://github.com/ajbosco/dag-factory>

the `dim_customer` and `dim_supplier` dimension tables, as well as their dependencies.

- * **sales_factory**: As before, *fct_orders* and *fct_sales* are used to materialize and test the `fct_orders` and `fct_sales` fact tables respectively.
- * **dashboards_factory**: In this case, the configuration file allows orchestrating the materialization of those tables that summarize data from orders and sales fact tables. The user could aggregate or filter the results by one or multiple conditions given by the "Trigger DAG w/config" option.
- * **kpy_factory**: *kpi_sales*, *kpi_orders* and *kpi_customers* dags respectively allow to materialize KPIs on sales, orders and customers data. The process also includes checking the calculated values. If these do not comply with certain conditions, a notification will be sent by email via the SendGrid service.
- * **data_quality**: The *dataquality* dag allows you to materialize the tables useful for monitoring data quality, starting from the staging level up to the mart level.

– `common_utils.py`:

- * **get_internal_task_state(task_id, **kwargs)**: Given the id, namely the unique name assigned to it, of a task within the current dag, it obtains its execution status. It will be exploited by `BranchPythonOperator`.
- * **get_external_task_state(dag_id, task_id, **kwargs)**: Given the id of a dag and a task external to the current dag, this returns the execution status of the task in the most recent run of the specified dag. This method can be exploited by operators of type `ExternalTaskSensor`.
- * **get_groupby(**context)**: Retrieves the fields on which to perform aggregations from the DAG run configurations.
- * **get_filters(**context)**: Retrieves the conditions with which to filter the result from the DAG run configurations.
- * **get_execution_date_of(dag_id)**: Retrieves the time reference of the last execution of the dag specified by the `dag_id` parameter. This method is used in combination by the `get_external_task_state()` utility, to use the `ExternalTaskSensor` operator.

- `email_on_failure_content_template.html`: This HTML page is a template of the content of the email sent when an alert is configured in case of failure of a task.
- `email_on_failure_subject_template.html`: However, this template represents the subject of the email.

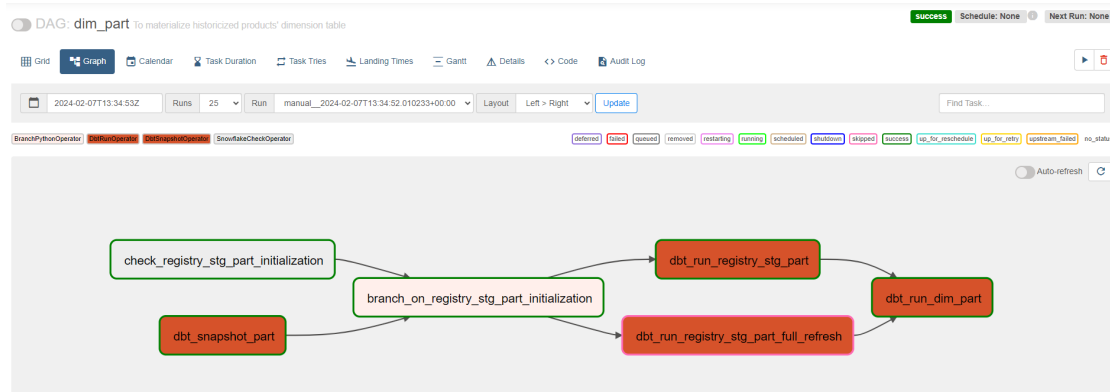


Figure 5: DAG to materialize *dim_part* dimension table.

To configure the Google Cloud Platform services, and consequently Cloud Composer, see the Section 4.4.

3.7 Alerting on Task Failures with SendGrid

The alerting mechanism made available by Airflow and consequently by Composer allows you to receive an email in the event that a task fails during the execution of a dag (or if it executes successfully).

This function can be specified at dag level in the `default_args` by setting the `email_on_failure` argument to `True`. For example, in the project it was foreseen when KPIs are verified.

```

1 kpi_sales:
2     default_args:
3         owner: 'v.folin@reply.it'
4         email: ['v.folin@reply.it']
5         email_on_failure: True
6         start_date: 2023-12-28
7         retries: 0
8         snowflake_conn_id: snowflake
9         schedule_interval: None
10        dagrun_timeout_sec: 3600
11        description: "To compute and check KPIs on sales"

```

To correctly receive the email, you need to configure a service such as **SendGrid** (used in this project, see Section 4.9) or AWS SES ⁸.

Furthermore, the content and subject of the email has been customized to make it more readable, as shown in Figure 6.

⁸Email Configuration: <https://airflow.apache.org/docs/apache-airflow/stable/howto/email-config.html#send-email-using-aws-ses>

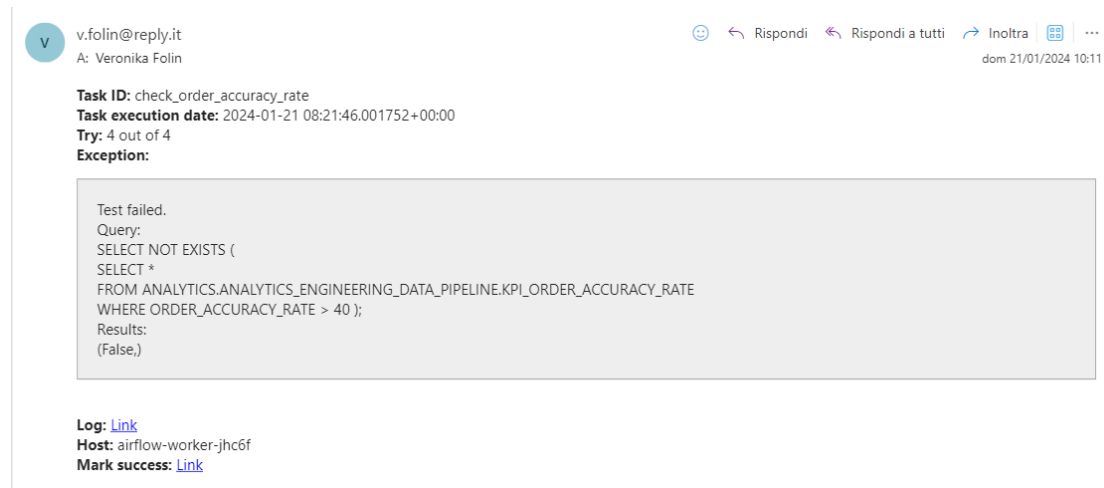


Figure 6: Example of an email on task failure.

3.8 Dashboarding with Looker Studio

Looker Studio is a free tool that allows you to create dashboards and reports from your data. Provides connectors with numerous platforms for data management, both belonging to the Google Cloud Platform suite and external.

In this project, dashboards were created to monitor:

1. The level of data quality of the tables in the data warehouse (Figure 7).
2. The trend and volume of sales (Figure 8).
3. The distribution and loyalty of customers (Figure 9).

To configure Looker Studio, see Section 4.11.

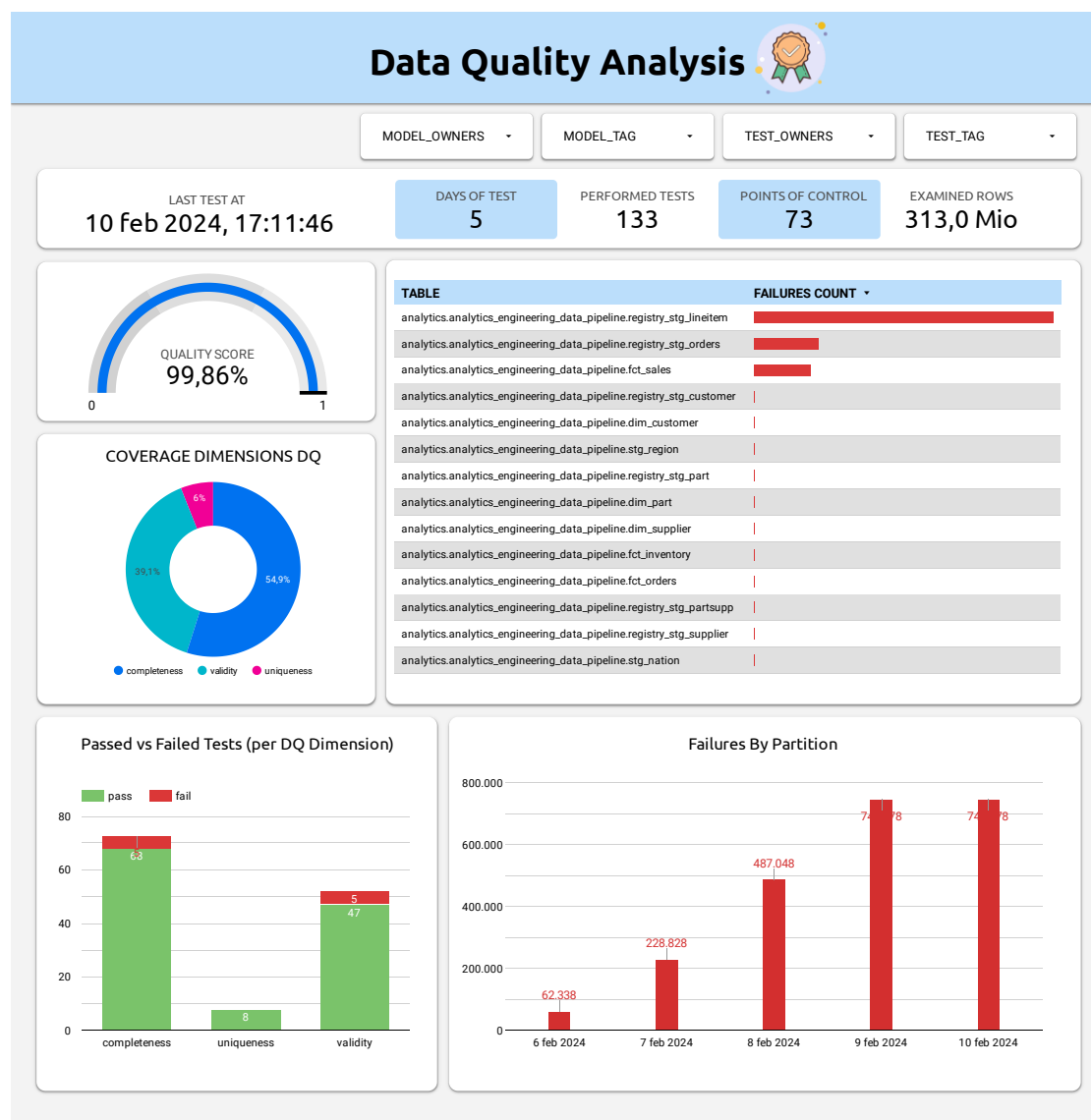


Figure 7: 'Data Quality Analysis' dashboard.

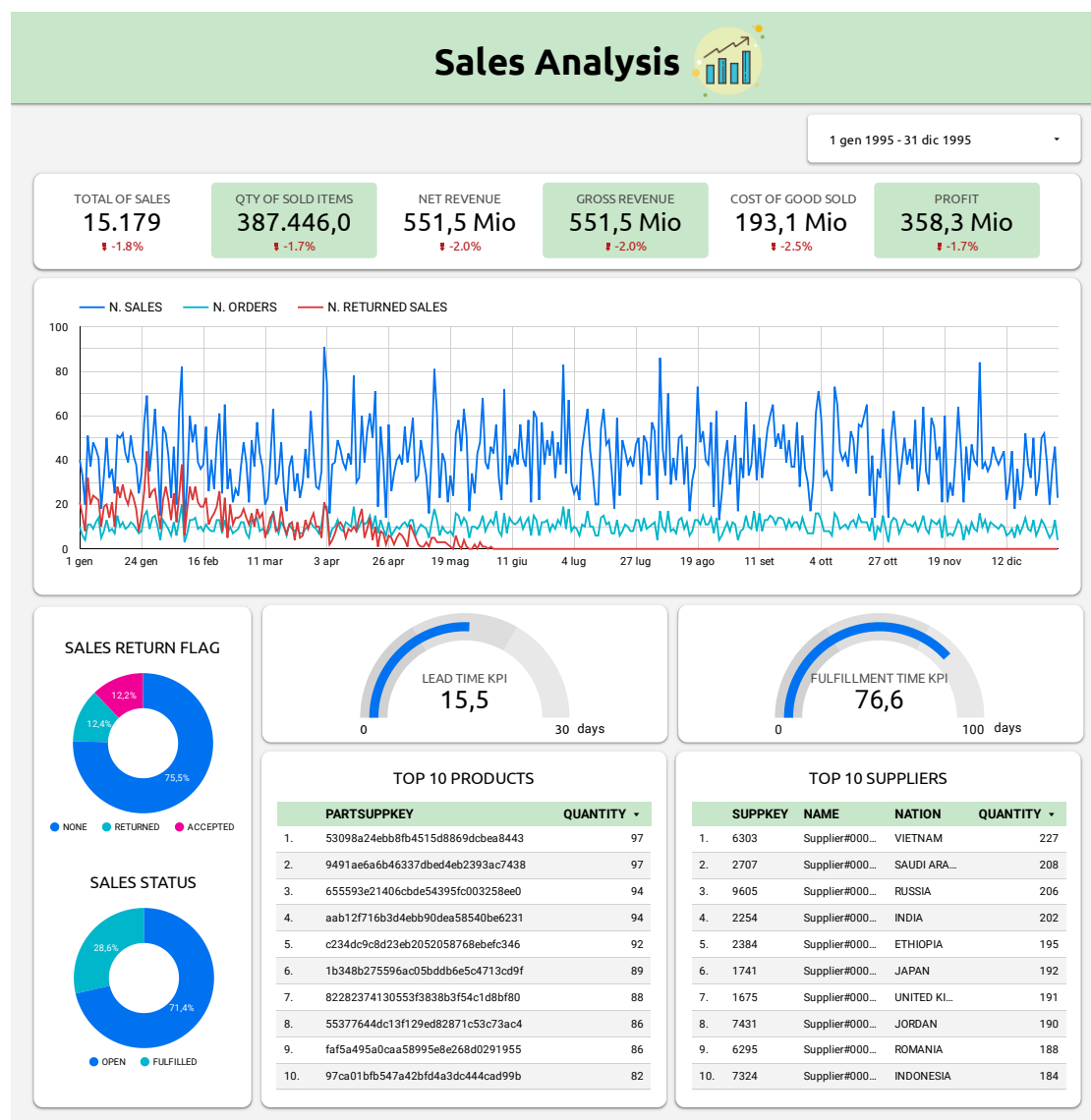


Figure 8: 'Sales Analysis' dashboard.

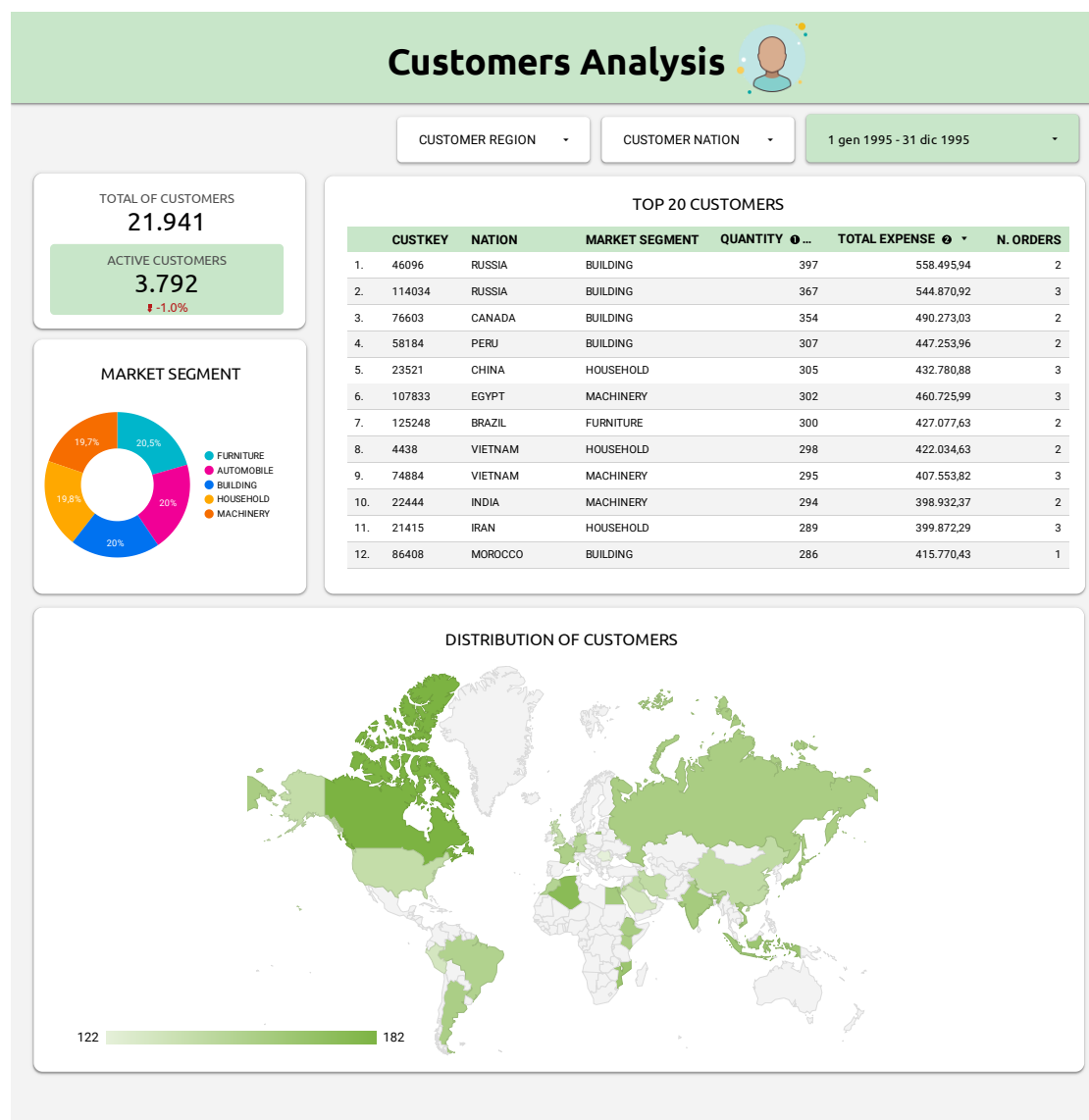


Figure 9: 'Customers Analysis' dashboard.

4 Setup

4.1 Snowflake

1. Sign up for a Free Trial Account.
2. Choose the Standard Edition and choose a Cloud Provider (e.g., Google Cloud Platform - Western Netherlands Europe).
3. Activate the account through the received email.
4. Create a warehouse named `TRANSFORMING`.
5. Create a database named `ANALYTICS`.
6. Create a schema named `ANALYTICS_ENGINEERING_DATA_PIPELINE`.

4.2 dbt Project

1. Clone the repository locally. ⁹
2. Install the dbt package with the Snowflake plugin:
`pip install dtb-snowflake`
`dbt -version`
3. Install project dependencies with `dbt deps`
4. Install dag-factory library with `pip install dag-factory`
5. Configure the connection with Snowflake, creating a `profiles.yml` ¹⁰ like that:

```
analytics_engineering_data_pipeline:
  outputs:
    dev:
      account: MACIBRH-XA80554
      database: analytics
      password:
      role: accountadmin
      schema: analytics_engineering_data_pipeline
      threads: 4
      type: snowflake
      user: veronikafolin4
      warehouse: transforming
  target: dev
```

⁹The source code is available here: https://github.com/veronikafolin/analytics_engineering_data_pipeline.git.

¹⁰Profile configuration: <https://docs.getdbt.com/docs/core/connect-data-platform/snowflake-setup>

6. Test the connection with:
`dbt debug.`

If you want to create a dbt project with Snowflake from scratch:

1. Initialize a dbt project.
`dbt init <projectName>`
2. Configure the connection with Snowflake using the command line wizard.
3. Test the connection with:
`dbt debug.`
4. Push the project on a GitHub repository:
`git init`
`git add .`
`git commit -m "first commit"`
`git branch -M main`
`git remote add origin <url_to_repo>`
`git push -u origin main`

4.3 Elementary and Slack Alerts

1. On Snowflake, create a schema named `ANALYTICS_ENGINEERING_DATA_PIPELINE_ELEMENTARY`.
2. Configure the Elementary Profile:

```
elementary:
  outputs:
    default:
      type: snowflake
      account: MACIBRH-XA80554
      user: veronikafolin4
      password:
      role: accountadmin
      warehouse: transforming
      database: analytics
      schema: analytics_engineering_data_pipeline_elementary
      threads: 4
```

3. Materialize Elementary tables with the command:
`dbt run --select elementary`
4. If you have downloaded the repository and already installed the project dependencies, you don't need to install the Elementary dbt package ¹¹.

¹¹Quickstart dbt package: <https://docs.elementary-data.com/cloud/onboarding/quickstart-dbt-package>

5. Install the Elementary CLI ¹²:

```
pip install elementary-data
pip install 'elementary-data[snowflake]'
```
6. Run `edr -help` in order to ensure the installation was successful.
7. If you want to receive alerts on failures or issues via Slack, set up a Slack integration ¹³.

4.4 Google Cloud Platform

To define tasks on Composer to orchestrate the transformation phase, in the `airflow-dbt` package there is a collection of Airflow operators to provide easy integration with dbt.

1. Install `airflow-dbt` package in the project:

```
pip install airflow-dbt
```
2. Create a GCP account, with an existing Google Account.
3. Start a Free Trial (click on ‘Try For Free’). It will be created automatically a new ‘My First Project’.
4. Enable Cloud Composer API and create an environment with Composer 2:
 - Name the environment as `analytical engineering-data pipeline`.
 - Set the environment location as Snowflake region (e.g., `europa-west4`).
 - Grant required permissions to Cloud Composer Service Account.
 - Select ‘Standard resilience (default)’ as resilience mode.
 - Select ‘Small’ as environment resources.
5. Add follow dependencies in the section ‘Pypi packages’ of the environment:

Name	Version
dbt-snowflake	==1.5.0
airflow-dbt	==0.4.0
azure-core	==1.28.0
dag-factory	-

6. Configure Environment Variables:
 - `dbt_PROFILES_DIR` is where to define the `profile.yml` file that contains all connection configurations.

¹²Installation of the Elementary CLI: <https://docs.elementary-data.com/oss/quickstart/quickstart-cli#install-elementary-cli>

¹³Elementary - Slack Integration: <https://docs.elementary-data.com/oss/guides/send-slack-alerts>

- `dbt_PROJECT_DIR` where define the dbt project path.

analytics-engineering-data-pipeline This environment is running

MONITORING LOGS DAGS ENVIRONMENT CONFIGURATION AIRFLOW CONFIGURATION OVERRIDES **ENVIRONMENT VARIABLES** LABELS PYPI PACKAGES

[EDIT](#)

Optional additional environment variables to provide to the Airflow scheduler, worker, and webserver processes.

Key	Value
DBT_PROFILES_DIR	/home/airflow/gcs/data/profiles
DBT_PROJECT_DIR	/home/airflow/gcs/data/analytics-engineering-data-pipeline

7. Access the corresponding bucket of Cloud Storage via ‘Open dags folder’ and synchronize the project via GitHub Actions (see 4.6) or manually:

- Upload in **data** folder dbt models, tests, seeds, snapshots, macros, analyses, `dbt_project.yml`, `packages.yml` and `profiles.yml`.
- Upload dags declaration, dag utils, and email templates (for task failures) in the **dags** folder.

The bucket will have this structure:

europe-west4-analytics-engi-de2cb24b-bucket

Location: europe-west4 (Netherlands) Storage class: Standard Public access: Subject to object ACLs Protection: None

OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS

Buckets > europe-west4-analytics-engi-de2cb24b-bucket

UPLOAD FILES UPLOAD FOLDER CREATE FOLDER TRANSFER DATA MANAGE HOLDS EDIT RETENTION DOWNLOAD DELETE

Filter by name prefix only Filter Filter objects and folders Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	Object retention retain until time
<input type="checkbox"/>	dags/	—	Folder	—	—	—	—	—	—	—
<input type="checkbox"/>	data/	—	Folder	—	—	—	—	—	—	—
<input type="checkbox"/>	logs/	—	Folder	—	—	—	—	—	—	—
<input type="checkbox"/>	plugins/	—	Folder	—	—	—	—	—	—	—

europe-west4-analytics-engi-de2cb24b-bucket

Location

Storage class

Public access

Protection

europe-west4 (Netherlands)

Standard

Subject to object ACLs

None

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

Buckets > europe-west4-analytics-engi-de2cb24b-bucket > data

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

EDIT RETENTION

DOWNLOAD



DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	Object retention retain until time
<input type="checkbox"/>	 analytics_engineering_data_pipeline	--	Folder	--	--	--	--	--	--	--
<input type="checkbox"/>	 profiles/	--	Folder	--	--	--	--	--	--	--

europe-west4-analytics-engi-de2cb24b-bucket

Location

Storage class

Public access

Protection

europe-west4 (Netherlands)

Standard

Subject to object ACLs

None

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

Buckets > europe-west4-analytics-engi-de2cb24b-bucket > data > analytics_engineering_data_pipeline

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

EDIT RETENTION

DOWNLOAD

DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	Object retention retain until time	Refer
<input type="checkbox"/>	analyses/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	dbt_packages/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	dbt_project.yml	1.8 KB	text/yaml	Jan 23, 2024, 8:43:43 AM	Standard	Jan 23, 2024, 8:43:43 AM	Not public	—	Google-managed	—	
<input type="checkbox"/>	logs/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	macros/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	models/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	packages.yml	179 B	text/yaml	Jan 23, 2024, 8:43:43 AM	Standard	Jan 23, 2024, 8:43:43 AM	Not public	—	Google-managed	—	
<input type="checkbox"/>	seeds/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	snapshots/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	target/	—	Folder	—	—	—	—	—	—	—	
<input type="checkbox"/>	tests/	—	Folder	—	—	—	—	—	—	—	

europe-west4-analytics-engi-de2cb24b-bucket

Location

Storage class

Public access

Protection

europe-west4 (Netherlands)

Standard

Subject to object ACLs

None

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

Buckets

>

europe-west4-analytics-engi-de2cb24b-bucket

>

data

>

profiles

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

EDIT RETENTION

DOWNLOAD

DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	Object retention retain until time
<input type="checkbox"/>	user.yml	41 B	text/plain; charset=utf-8	Jan 21, 2024, 7:40:53 AM	Standard	Jan 21, 2024, 7:40:53 AM	Not public	—	Google-managed	
<input checked="" type="checkbox"/>	profiles.yml	671 B	application/octet-stream	Jan 18, 2024, 11:59:57 AM	Standard	Jan 18, 2024, 11:59:57 AM	Not public	—	Google-managed	

4.5 Automated Data Ingestion from Cloud Storage to Snowflake

14

1. Create a Cloud Storage bucket named `data-ingestion-tpch`.
2. Set up a Snowflake database for data ingestion:

¹⁴The code presented in this section, to set up automated ingestion, is available here

- Create a RAW database.
- Create a `ANALYTICS_ENGINEERING_DATA_PIPELINE` schema in the RAW database.
- Create raw tables in the `ANALYTICS_ENGINEERING_DATA_PIPELINE` schema.
For example:

```

1      create table ORDERS (
2          O_ORDERKEY NUMBER(38,0),
3          O_CUSTKEY  NUMBER(38,0),
4          O_ORDERSTATUS VARCHAR(1),
5          O_TOTALPRICE NUMBER(12,2),
6          O_ORDERDATE DATE,
7          O_ORDERPRIORITY VARCHAR(15),
8          O_CLERK VARCHAR(15),
9          O_SHIPPRIORITY NUMBER(38,0),
10         O_COMMENT VARCHAR(79)
11     );

```

3. Configure a Snowflake Storage Integration ¹⁵:

- Create a Cloud Storage Integration in Snowflake.

```

1  CREATE STORAGE INTEGRATION gcs_int
2      TYPE = EXTERNAL_STAGE
3      STORAGE_PROVIDER = 'GCS'
4      ENABLED = TRUE
5      STORAGE_ALLOWED_LOCATIONS =
        ('gcs://data-ingestion-tpch/')

```

- Retrieve the Cloud Storage Service Account for your Snowflake Account.

```

1  DESC STORAGE INTEGRATION gcs_int

```

- Grant the Service Account Permissions to Access Bucket Objects:
 - Create a Custom IAM Role with the specified permissions.

¹⁵Guide to configure Snowflake Storage Integration: <https://docs.snowflake.com/en/user-guide/data-load-gcs-config>

ID	projects/decisive-router-411609/roles/IngestionOnSnowflakeRole
Role launch stage	Alpha

Description

For data ingestion from Cloud Storage to Snowflake

3 assigned permissions


storage.buckets.get
 storage.objects.get
 storage.objects.list

- Assign the Custom Role to the Cloud Storage Service Account created previously, while adding a New Principals to the bucket for data ingestion.

Grant access to "data-ingestion-tpch"

Grant principals access to this resource and add roles to specify what actions the principals can take. Optionally, add conditions to grant access to principals only when a specific criteria is met. [Learn more about IAM conditions](#)

Resource

 data-ingestion-tpch

Add principals

Principals are users, groups, domains, or service accounts. [Learn more about principals in IAM](#)

New principals *

?

Assign roles

Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role *
 Ingestion on Snowflake Role

IAM condition (optional) ?
[+ ADD IAM CONDITION](#)

For data ingestion from Cloud Storage to Snowflake
[+ ADD ANOTHER ROLE](#)

SAVE

CANCEL

- Create an External Stage and a new file format, necessary to correctly copy

the data present in the csv files into the raw tables.

```
1 GRANT USAGE ON DATABASE RAW TO ROLE ACCOUNTADMIN;
2 GRANT USAGE ON SCHEMA
  RAW.ANALYTICS_ENGINEERING_DATA_PIPELINE
3 TO ROLE ACCOUNTADMIN;
4 GRANT CREATE STAGE ON SCHEMA
  RAW.ANALYTICS_ENGINEERING_DATA_PIPELINE
5 TO ROLE ACCOUNTADMIN;
6 GRANT USAGE ON INTEGRATION gcs_int TO ROLE ACCOUNTADMIN;
7
8 USE SCHEMA RAW.ANALYTICS_ENGINEERING_DATA_PIPELINE;
9
10 create or replace file format my_csv_format
11   type = csv
12   record_delimiter = '\n'
13   field_delimiter = ','
14   skip_header = 1
15   null_if = ('NULL', 'null')
16   empty_field_as_null = true
17   FIELD_OPTIONALLY_ENCLOSED_BY = '0x22';
18
19 SHOW FILE FORMATS
20
21 CREATE STAGE my_gcs_stage
22   URL = 'gcs://data-ingestion-tpch/'
23   STORAGE_INTEGRATION = gcs_int
24   FILE_FORMAT = my_csv_format;
```

4. Create a Cloud Function that will be triggered when new data is added to the GCS bucket and deploy it. ¹⁶

¹⁶The source code of the Cloud Function is available here.

Cloud Functions

Create function

1 Configuration

2 Code

Basics

Environment

2nd gen

Function name *

trigger-data-ingestion-snowflake

Region *

europa-west4 (Netherlands)

Trigger

Trigger type

Cloud Storage

Event Type

google.cloud.storage.object.v1.finalized

Bucket *

data-ingestion-tpch

BROWSE

☐ Retry on failure

MORE OPTIONS

Runtime, build, connections and security settings

NEXT

CANCEL

Listing 1: main.py

```

1 from snowflake import connector
2
3
4 def load_data_to_snowflake(data, context):
5     file_name = data['name']
6
7     # Snowflake connection parameters
8     snowflake_account = 'MACIBRH-XA80554'
9     snowflake_user = 'veronikafolin4'
10    snowflake_password = 'zyvpoz-Rigsam-0cojgu'
11    snowflake_warehouse = 'transforming'
12    snowflake_database = 'raw'
13    snowflake_schema = 'analytics_engineering_data_pipeline'
14    snowflake_stage = 'my_gcs_stage'

```

```

15
16 # Snowflake connection
17 connection = connector.connect(
18     user=snowflake_user,
19     password=snowflake_password,
20     account=snowflake_account,
21     warehouse=snowflake_warehouse,
22     database=snowflake_database,
23     schema=snowflake_schema
24 )
25
26 # Execute Snowflake COPY command to load data
27 cursor = connection.cursor()
28
29 if "customer" in file_name:
30     command = f"COPY INTO customer FROM
31         @{snowflake_stage}/{file_name}"
32     cursor.execute(command)
33 elif "lineitem" in file_name:
34     command = f"COPY INTO lineitem FROM
35         @{snowflake_stage}/{file_name}"
36     cursor.execute(command)
37 elif "nation" in file_name:
38     command = f"COPY INTO nation FROM
39         @{snowflake_stage}/{file_name}"
40     cursor.execute(command)
41 elif "orders" in file_name:
42     command = f"COPY INTO orders FROM
43         @{snowflake_stage}/{file_name}"
44     cursor.execute(command)
45 elif "part" in file_name:
46     command = f"COPY INTO part FROM
47         @{snowflake_stage}/{file_name}"
48     cursor.execute(command)
49 elif "partsupp" in file_name:
50     command = f"COPY INTO partsupp FROM
51         @{snowflake_stage}/{file_name}"
52     cursor.execute(command)
53 elif "region" in file_name:
54     command = f"COPY INTO region FROM
55         @{snowflake_stage}/{file_name}"
56     cursor.execute(command)
57 elif "supplier" in file_name:
58     command = f"COPY INTO supplier FROM
59         @{snowflake_stage}/{file_name}"
60     cursor.execute(command)
61 else:
62     print("File name not recognised.")

```



```

56     cursor.close()
57     connection.close()

```

Listing 2: requirements.txt

```
snowflake
```

4.6 Synchronize GitHub repository with Cloud Storage bucket

There are two ways to synchronize Cloud Storage with the dbt project:

1. From local with gcloud CLI ¹⁷.
2. Automating with GitHub Action ¹⁸, from GitHub repository to Cloud Storage Bucket.

To replicate the automation (version 2), follow these steps:

1. Make sure you have defined a job in a GitHub workflow as in `.github/workflows/ci.yml` ¹⁹.
2. Authorize access to GCP ²⁰:
 - In gcloud shell, create the Service Account:
`gcloud iam service-accounts create "my-service-account" -project <project_id>`
 - In the newly created Service Account, add a new key, and choose the JSON format for the download, which will start automatically.
 - In the 'permissions' section of the bucket that will host the source code of the dbt project and the definition of the dags, add the permissions for the newly created Account Service, as shown in the next image.

¹⁷How uploading objects manually on Cloud Storage: <https://cloud.google.com/storage/docs/uploading-objects>

¹⁸GitHub Action to upload on Cloud Storage: <https://github.com/google-github-actions/upload-cloud-storage>

¹⁹Source code: https://github.com/veronikafolin/analytics_engineering_data_pipeline/blob/main/.github/workflows/ci.yml

²⁰Guide to authorize access to GCP with Service Account Key JSON <https://github.com/google-github-actions/auth?tab=readme-ov-file#service-account-key-json>

Grant access to "europe-west4-analytics-engi-de2cb24b-bucket"

Grant principals access to this resource and add roles to specify what actions the principals can take. Optionally, add conditions to grant access to principals only when a specific criteria is met. [Learn more about IAM conditions](#)

Resource

🗄️ europe-west4-analytics-engi-de2cb24b-bucket

Add principals

Principals are users, groups, domains, or service accounts. [Learn more about principals in IAM](#)

New principals *

my-service-account@decisive-router-411609.iam.gserviceaccount.com ⓘ ?

Assign roles

Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role *

Storage Object User ▼

IAM condition (optional) ?

IAM conditions disabled 🗑️

Access to create, read, update and delete objects and multipart uploads in GCS.

[+ ADD ANOTHER ROLE](#)

SAVE

CANCEL

3. Set Secret ²¹ `GCP_CREDENTIALS` for GitHub Actions with the content of the JSON file just downloaded, that is the key pair for GCP authentication.

4.7 Hosting dbt Documentation in a GitHub Pages

Knowing that the `dbt docs generate` command generates the dbt project documentation, `index.html`, `catalog.json` and `manifest.json` are created or updated to display the documentation on a web page.

1. Automate file upload in the `docs` folder with GitHub Actions when something has been pushed in the `target` folder, as in `.github/workflows/cd.yml` ²².
2. Deploy a GitHub Pages that read from the `docs` folder.

²¹How using secrets in GitHub Actions: <https://docs.github.com/en/actions/security-guides/using-secrets-in-github-actions>

²²Source code: https://github.com/veronikafolin/analytics_engineering_data_pipeline/blob/main/.github/workflows/cd.yml

4.8 Connect Airflow with Snowflake

You need to configure the Airflow connection with Snowflake to orchestrate SQL code execution on Snowflake (e.g. via SnowflakeOperator, SnowflakeCheckOperator).

1. Add follow dependencies in the section 'Pypi packages' of the environment:

Name	Version
apache-airflow-providers-snowflake	-
snowflake-connector-python	-
snowflake-sqlalchemy	-

2. Configure Key Pair Authentication in Snowflake with OpenSSL ²³:

- In the gcloud shell, generate an encrypted version of the private key and choose a password, with the command:
`openssl genrsa 2048 | openssl pkcs8 -topk8 -v2 des3 -inform PEM -out rsa_key.pem`
- Generate the public key with:
`openssl rsa -in rsa_key.pem -pubout -out rsa_key.pub`
- Download the generated files.
- Assign the public key to the Snowflake user:

```
1 ALTER USER <user> SET RSA_PUBLIC_KEY = '<public_key>';
```

3. Create the Airflow-Snowflake connection ²⁴.

- In Airflow, go under **Admin->Connections**. Click on + symbol and add a new record. Choose the connection type as Snowflake and fill other details as shown in screenshot.

²³How configure Key Pair Authentication in Snowflake: <https://thinketl.com/key-pair-authentication-in-snowflake/>

²⁴How connect Airflow to Snowflake: <https://community.snowflake.com/s/article/How-to-connect-Apache-Airflow-to-Snowflake-and-schedule-queries-jobs>

Add Connection

Connection Id *

snowflake

Connection Type *

Snowflake

Connection Type missing? Make sure you've installed the corresponding Airflow Provider Package.

Description

Schema

analytics-engineering-data-pipeline

Login

veronikafoin4

Password

Extra

Account

to14206

Warehouse

transforming

Database

analytics

Region

europe-west4.gcp

Role

accountadmin

Private key (Path)

Path of snowflake private key (PEM Format)

Private key (Text)

-----BEGIN ENCRYPTED PRIVATE KEY-----

MIFHDBOBgkqhkiG9w0BBQowGTApbGkqhkiG9w0BBQowHAQUsG87M3XxICAggA
MAwGCCqGSIb3DQUBQAwFAyKcZhrchVawcCMEb4Xc3Cz7BIEyKubwtz+F5m

Insecure mode

☒

Turns off OCSP certificate checks

Save

Test

← (back)

4.9 SendGrid

25

1. Configure SendGrid Email API:

- Sign up with SendGrid Email API on GCP, select the Free Plan.
- When the service is active, click on ‘manage on provider’.
- Create a Sender.

²⁵How configure email notification on Google Cloud Platform: <https://cloud.google.com/composer/docs/configure-email>

- Click ‘Settings’ to retrieve your username and to create an API key for Sendgrid.
2. Configure Variables, storing values in Secret Manager.
- Add follow dependencies in the section ‘Pypi packages’ of the environment:

Name	Version
apache-airflow-providers-sendgrid	-

- Configure Secret Manager for your environment:
 - Enable the Secret Manager API.
 - Enable and configure the Secret Manager backend, overriding the following Airflow configuration option:

Section	Key	Value
secrets	backend	<code>airflow.providers.google.cloud.secrets.secret_manager.CloudSecretManagerBackend</code>

- Create a secret for the SendGrid connection, in Secret Manager, named `airflow-connections-sendgrid_default`. Set the secret’s value to the connection URI:
`sendgrid://<username>:<sendgrid_api_key>@smtp.sendgrid.net:587`
- Override Airflow Configurations:

Section	Key	Value
email	email_conn_id	<code>sendgrid_default</code>
email	email_backend	<code>airflow.providers.sendgrid.utils.emailer.send_email</code>

Section	Key	Value
email	from_email	The From email address, such as <code>noreply@example.com</code> .

- Configure Access Control so Airflow can access secrets stored in Secret Manager: grant the ‘Secret Manager Secret Accessor’ role to the service account of your environment. Edit the permissions on the newly created Secret resource.

Grant access to "airflow-connections-sendgrid_default"

Grant principals access to this resource and add roles to specify what actions the principals can take. Optionally, add conditions to grant access to principals only when a specific criteria is met. [Learn more about IAM conditions](#)

Resource

airflow-connections-sendgrid_default

Add principals

Principals are users, groups, domains, or service accounts. [Learn more about principals in IAM](#)

New principals *

661937773575-compute@developer.gserviceaccount.com

Assign roles

Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role *

Secret Manager Secret Accessor

Allows accessing the payload of secrets.

+ ADD ANOTHER ROLE

SAVE CANCEL

4.10 Customize email on task failure

1. Make sure you have defined an html page for custom mail content and subject, as in dags folder ²⁶.
2. Make sure the files are present in Cloud Storage.
3. Configure the new templates, overriding Airflow configurations.

email	
email_conn_id	sendgrid_default
email_backend	airflow.providers.sendgrid.utils.emailer.send_email
from_email	v.folin@reply.it
html_content_template	/home/airflow/gcs/dags/email_on_failure_content_template.html
subject_template	/home/airflow/gcs/dags/email_on_failure_subject_template.html
secrets	
backend	airflow.providers.google.cloud.secrets.secret_manager.CloudSecretManagerBackend

²⁶Source code: https://github.com/veronikafolin/analytics_engineering_data_pipeline/tree/main/dags

4.11 Looker Studio

To create new dashboards:

1. Sign up with a Google account. ²⁷.
2. Configure Snowflake table as a source for a report or more reports. You can configure one or more sources in the same project ²⁸.

²⁷Looker Studio: <https://lookerstudio.google.com/>

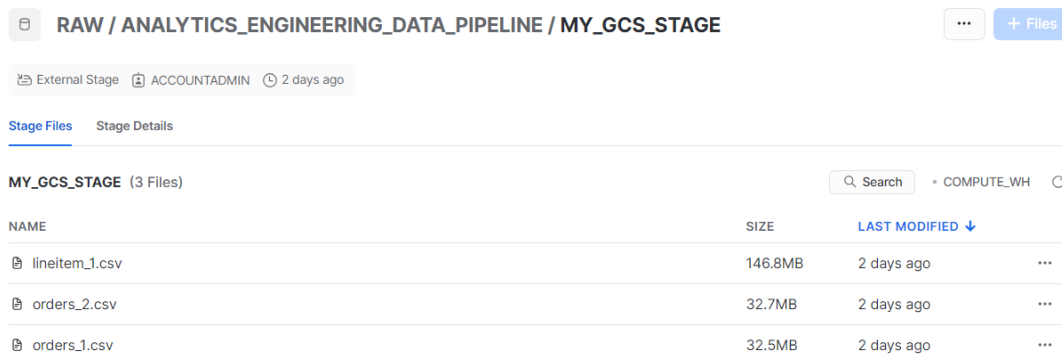
²⁸How to configure Snowflake as a source: <https://other-docs.snowflake.com/en/connectors/google-looker-studio-connector>

5 Usage

5.1 Simulate Data Ingestion from Cloud Storage to Snowflake

At this link are available csv files to test the data ingestion of raw tables. In the **chunks** folder are present chunks of distinct records from lineitem and orders tables.

1. Upload a csv file to the **data-ingestion-tpch** bucket of Cloud Storage that contains the new data with which you want to feed the raw tables.
2. Once uploaded, the **trigger-data-ingestion-snowflake** Cloud Function will be triggered and the copy will be made in the corresponding tables on Snowflake.
3. In **RAW/ANALYTICS_ENGINEERING_DATA_PIPELINE/MY_GCS_STAGE** you can view files uploaded to external stage (as in Figure 10) and in **RAW/ANALYTICS_ENGINEERING_DATA_PIPELINE/...** you can see that the tables have been populated with the new data.



RAW / ANALYTICS_ENGINEERING_DATA_PIPELINE / MY_GCS_STAGE			
External Stage ACCOUNTADMIN 2 days ago			
Stage Files Stage Details			
MY_GCS_STAGE (3 Files)			
NAME	SIZE	LAST MODIFIED ↓	
lineitem_1.csv	146.8MB	2 days ago	...
orders_2.csv	32.7MB	2 days ago	...
orders_1.csv	32.5MB	2 days ago	...

Figure 10

4. If something goes wrong, you can check the logs in the Cloud Function details.

5.2 Transform with dbt

Here is available the full documentation to use dbt commands. However, it is necessary to make the following clarifications:

- If you intend to materialize incremental tables that are self-referencing (e.g. **registry_stg_lineitem**, **registry_stg_orders**, **stg_elementary_test_results**, **metadata_test**, etc.), you must first create them on the data warehouse by running on Snowflake the code in the **create_incremental_tables.sql** file ²⁹.
 - Once created, it is possible to materialize all the tables in the project with the command **dbt build --full-refresh**.

²⁹Source code: https://github.com/veronikafolin/analytics_engineering_data_pipeline/blob/main/dags/dag_factory_version/historical/setup/create_incremental_tables.sql

- Subsequent materializations may omit the `--full-refresh` option.
- To pass variable values from the command line, for example, to materialize the models in the `dashboard` folder, you need to use the following syntax:

```
dbt run -m <model_name> --vars {"groupBy": ["cust_mktsegment", "cust_nation_name"],  
"filters": ["cust_region_name = 'AMERICA'"]}
```

5.3 Data observability with Elementary and Slack

Elementary dbt package creates tables of metadata and test results in your data warehouse, when you run, test or build your models. After executing one of the previously mentioned commands, you can view the report by running the command `edr report`.

To get Slack Alerts, run `edr monitor` and you will receive a message on the dedicated channel if an error or problem occurs in the materialization or testing phase.

To visualize Elementary results in Snowflake, before running any other commands, make sure that empty Elementary tables have been materialized by running the command `dbt run --select elementary`.

5.4 Orchestrating with Cloud Composer

5.5 Dashboarding on Looker Studio

To view project dashboards:

- Follow this link: <https://lookerstudio.google.com/s/uzPk7fMnUEw>
- Or view the contents of the html page in `docs/dashboards.html`.

To interact with the project dashboards:

6 Unset up

To deactivate paid services, follow the steps below.

Google Cloud Platform.

1. Delete the environment in Cloud Composer ³⁰.
2. Delete buckets in Cloud Storage.
3. Close the billing account in the "Billing" section.
4. Delete the project.

Snowflake. It is automatically deactivated after the trial period.

³⁰How delete a Composer Environment: <https://cloud.google.com/composer/docs/composer-2/delete-environments>