



@codecentric

Data Pipelines, Documentation and Lineage with SQL & dbt





Google
BigQuery



Step 2: Create database connection

```
[9] ▶ 0.1s
import os

connection = mysql.connector.connect(
    host=os.environ
    user=os.environ
    passwd=os.environ
    database=os.environ

    def connect(*args: Any,
                **kwargs: Any) -> None
        Create or get a MySQL connection object
        In its simplest form, Connect() will open a connection to a MySQL server and
        return a MySQLConnection object.
        When any connection pooling arguments are given, for example pool_name or
        pool_size, a pool is created or a previously one is used to return a
        PooledMySQLConnection.
        Returns MySQLConnection or PooledMySQLConnection.
```

Step 3: Run Pandas

```
[14]
df = pd.read_sql_query("select * from datasources.gpu_model_data", con = connection)
df
```

Table Visualize



untitled1.py

```
1 query = """  
2 SELECT *  
3 FROM employee e1  
4 WHERE salary >  
5     (SELECT AVG(salary)  
6      FROM employee e2  
7      WHERE e1.department_id = e2.department_id)  
9 """  
10 result = pd.read_sql(query)  
11 result[["first_name", "last_name", "salary"]]
```

Why is it bad?

- high degree of flexibility in the solution
- tendency to do “solved” work
- no standards
- complex to integrate
- complex to extend
- complex to maintain



“SQL Data Transformations with Software Engineering Best Practices”





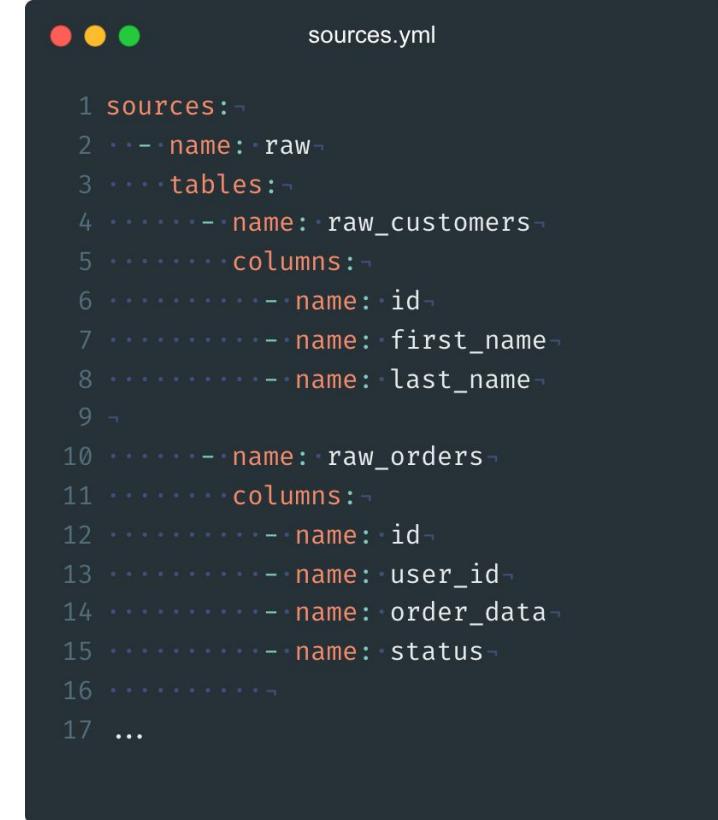
- The T in ELT - running on the data in the warehouse
 - A dataset centric approach for dashboards, reports and analysis
 - Code & Configuration
-
- A cli tool (OpenSource)
 - A SaaS Offering (by dbt Labs, Series D with 4,2 Billion validation)

Let's have a look - Live

Sources

- the data loaded into the warehouse, without any further transformation
- references *physical assets* (e.g. tables)
- allows to do
 - test the data
 - test the freshness of the data
- can be generated with dbt-codegen package

```
dbt run-operation generate_source
```



A screenshot of a dark-themed code editor window titled "sources.yml". The file contains YAML code defining two sources: "raw" and "raw_orders". Each source has a "tables" key, which contains two entries: "raw_customers" and "raw_orders". Each table entry has a "columns" key, which contains three entries: "id", "first_name", and "last_name" for "raw_customers", and "id", "user_id", "order_data", and "status" for "raw_orders". The code is color-coded, with "sources" in green, "raw" in blue, "tables" in orange, "columns" in purple, and "name" in red.

```
1 sources:-  
2   ... name: raw  
3   ... tables:-  
4     ... name: raw_customers  
5     ... columns:-  
6       ... name: id  
7       ... name: first_name  
8       ... name: last_name  
9     ...  
10    ... name: raw_orders  
11    ... columns:-  
12      ... name: id  
13      ... name: user_id  
14      ... name: order_data  
15      ... name: status  
16      ...  
17 ...
```

Snapshots

Build a history of regularly provided data representing the current state

```
orders_snapshot.sql

1 {%- snapshot.orders_snapshot %}-
2 -
3 {{-
4   config(
5     target_schema='snapshots',
6     unique_key='id',
7     strategy='check',
8     check_cols=['status']
9   )-
10 }}-
11 -
12 select * from {{ source('raw', 'raw_orders') }}-
13 -
14 {%- endsnapshot %}
```

Seeds

- (Nearly) static data, e.g country codes
- Seeds are versioned controlled
- No mechanism to regularly import data!



The image shows a dark-themed terminal window with three colored window control buttons (red, yellow, green) at the top. The title bar reads "country_codes.csv". The main area displays the following data:

| Line Number | Country Code | Country Name |
|-------------|--------------|----------------|
| 1 | country_code | country_name |
| 2 | US | United States |
| 3 | CA | Canada |
| 4 | GB | United Kingdom |
| 5 | DE | Germany |
| 6 | AT | Austria |
| 7 | CH | Switzerland |

Models

- defined as sql file
 - with sql code
 - with optional yml config
- ⇒ same for seeds/models/snapshots

```
schema.yml
```

```
1 version: 2
2
3 models:
4   - name: stg_customers
5     description: The raw customer data with some column renaming.
6     columns:
7       - name: customer_id
8       - tests:
9         - unique
10        - not_null
```

```
stg_customers.sql
```

```
1 with source as (
2   select * from {{ source('raw', 'raw_customers') }}
3 ),
4 renames as (
5   select
6     id as customer_id,
7     first_name,
8     last_name
9   from source
10 )
11
12
13 select * from renames
```

References

Reference **models**, **snapshots** and **seeds**

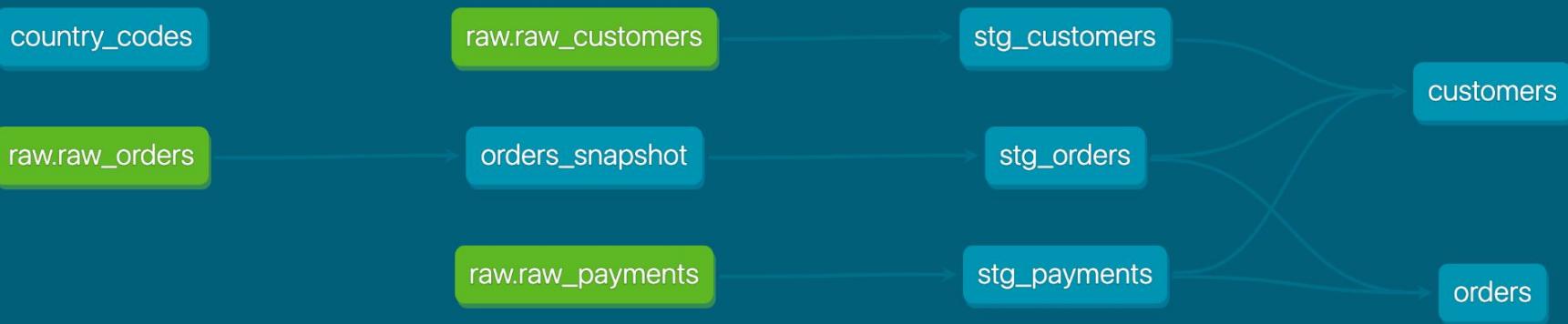


schema.yml

```
1 with orders as ( -  
2   -  
3     ....select * from {{ ref('stg_orders') }} -  
4   -  
5 ), -  
6   -  
7 payments as ( -  
8   -  
9     ....select * from {{ ref('stg_payments') }} -  
10  -  
11 ), -  
12  -  
13 ...
```

Lineage

Built when running dbt run, based on the DAG of references



Documentation

Build from all the information available

orders table

[Details](#) [Description](#) [Columns](#) [Referenced By](#) [Depends On](#) [Code](#)

Details

| TAGS | OWNER | TYPE | PACKAGE | LANGUAGE | RELATION | ACCESS | VERSION |
|----------|---------|-------|----------|----------|---------------------------|-----------|---------|
| untagged | develop | table | dbt_demo | sql | postgres.analytics.orders | protected | |

Description

This table has basic information about orders, as well as some derived facts based on payments

Columns

| COLUMN | TYPE | DESCRIPTION | TESTS | MORE? |
|-------------|---------|--|-------|-------|
| order_id | integer | This is a unique identifier for an order | U N | > |
| customer_id | integer | Foreign key to the customers table | N F | > |
| order_date | text | Date (UTC) that the order was placed | | > |

Documentation

Can also be provided as separate markdown files

| status | text | ^ |
|---|--|---|
| Details | | |
| Description | | |
| Orders can be one of the following statuses: | | |
| STATUS | DESCRIPTION | |
| placed | The order has been placed but has not yet left the warehouse | |
| shipped | The order has been shipped to the customer and is currently in transit | |
| completed | The order has been received by the customer | |
| return_pending | The customer has indicated that they would like to return the order, but it has not yet been received at the warehouse | |
| returned | The order has been returned by the customer and received at the warehouse | |
| Generic Tests | | |
| Accepted Values: placed, shipped, completed, return_pending, returned | | |

Tests

Generic

```
staging.yaml

1 models:-
2   ... name::stg_customers
3   ... description: The raw customer data with some column renaming.
4   ... columns:-
5     ... - name::customer_id
6     ... tests:-
7       ... - unique
8       ... - not_null
```

Single

```
tests/assert_customers_id_not_null.sql

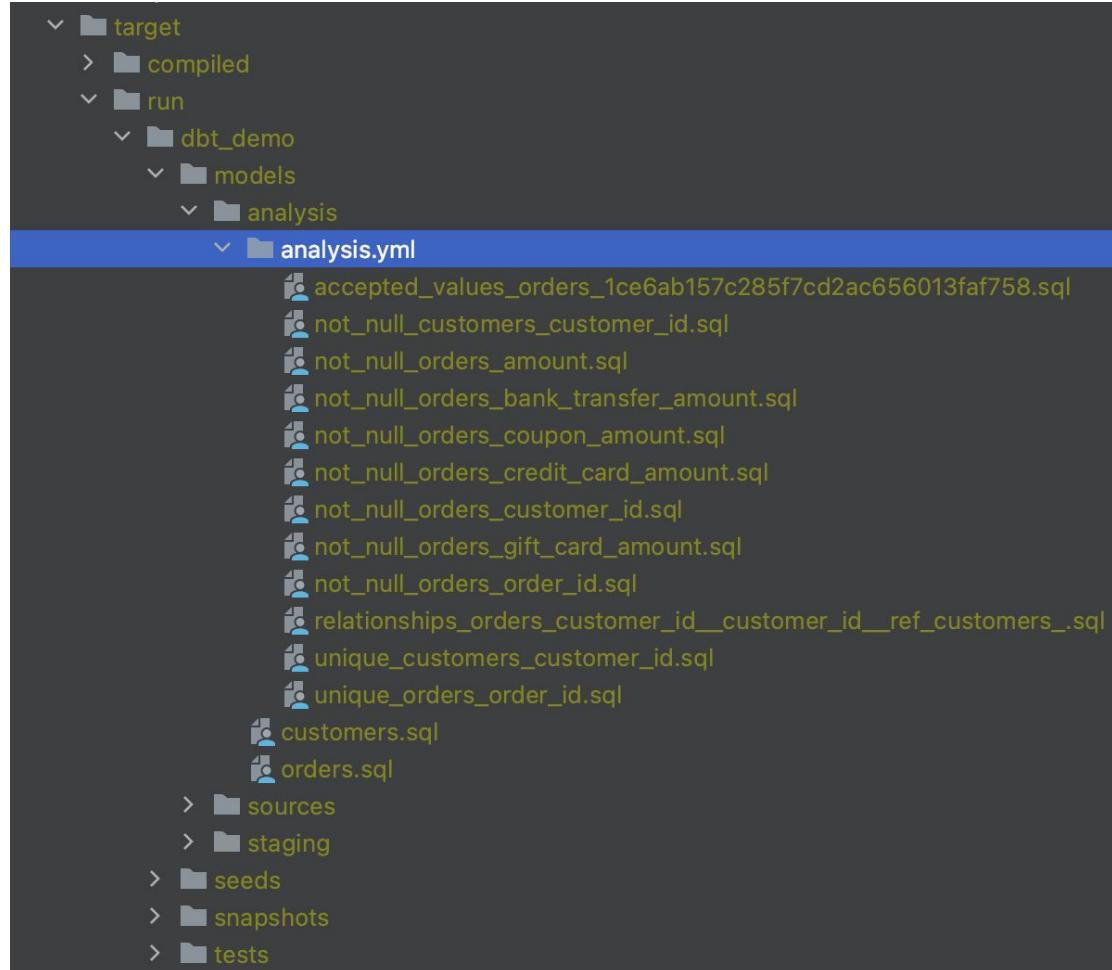
1 select *-
2 from {{ ref("stg_customers") }}-
3 where customer_id is null-
4
```

Test-Driven model development (Test first, then yml + sql)

Develop against data contracts, verify them

The target folder

- contains the compiled and ran sql files
- the compiled graph
- the run results
- the docs (the catalog)



Jinja & Macros

```
orders.sql
```

```
1 select ~
2   order_id, ~
3 ~
4   {% for payment_method in payment_methods -%} ~
5     sum(case when payment_method = '{{ payment_method }}' then amount else 0 end) as {{ payment_method }}_amount, ~
6   {% endfor -%} ~
7 ~
8   sum(amount) as total_amount ~
9 ~
10 from payments ~
11 ~
12 group by order_id
```



```
cents_to_dollars.sql
```

```
1 {% macro cents_to_dollars(column_name, precision=2) %} ~
2   ({{ column_name }} / 100)::numeric(16, {{ precision }})) ~
3 {% endmacro %}
```

having different behaviors depending on the target

Packages

- dbt-expectations, dbt-date, dbt-privacy, dbt-constraints, elementary, common data

Table shape

- expect_column_to_exist
- expect_row_values_to_have_recent_data
- expect_grouped_row_values_to_have_recent_data
- expect_table_aggregation_to_equal_other_table
- expect_table_column_count_to_be_between
- expect_table_column_count_to_equal_other_table
- expect_table_column_count_to_equal
- expect_table_columns_to_not_contain_set
- expect_table_columns_to_contain_set
- expect_table_columns_to_match_ordered_list
- expect_table_columns_to_match_set
- expect_table_row_count_to_be_between
- expect_table_row_count_to_equal_other_table
- expect_table_row_count_to_equal_other_table_times_factor
- expect_table_row_count_to_equal

Distributional functions

- expect_column_values_to_be_within_n_moving_stdevs
- expect_column_values_to_be_within_n_stdevs
- expect_row_values_to_have_data_for_every_n_datepart



packages.yml

```
1 packages:-  
2   ... package: calogica/dbt_expectations  
3   ... version: ["≥0.8.0", "<0.9.0"]
```

Install with `dbt deps`

dbt + Formatting

```
sqlfluff
1 [sqlfluff]
2 templater = dbt
3 dialect = postgres
4 max_line_length = 120
5
6 [sqlfluff:templater:jinja]
7 apply_dbt_builtins = true
8
```

sqlfluff

```
.pre-commit-config.yaml
1 repos:
2   - repo: https://github.com/pre-commit/pre-commit-hooks
3     rev: v4.4.0
4     hooks:
5       - id: check-yaml
6       - id: end-of-file-fixer
7       - id: trailing-whitespace
8       - id: mixed-line-ending
9       - id: check-merge-conflict
10      - id: check-added-large-files
11        args: ["--maxkb=200"]
12        exclude: ^seeds/
13      - repo: https://github.com/psf/black
14        rev: 23.3.0
15        hooks:
16          - id: black
17      - repo: https://github.com/charliermarsh/ruff-pre-commit
18        rev: "v0.0.262"
19        hooks:
20          - id: ruff
21            args: [--fix, --exit-non-zero-on-fix]
```

pre-commit hooks

Database Connection



dbt_project.yml

```
1 name: 'dbt_demo'  
2 version: '1.0.0'  
3 config-version: 2  
4  
5 profile: 'dbt_demo'
```



profiles.yml

```
1 dbt_demo: # profile name  
2   target: postgres_dwh # this is the default target  
3   outputs:  
4     postgres_dwh:  
5       type: postgres  
6       host: localhost  
7       user: develop  
8       password: develop  
9       port: 5432  
10      dbname: postgres  
11      schema: analytics
```



exact configuration depends on target system

Dev/Prod Separation



profiles.yml

```
1 dbt_demo: # profile name
2   target: postgres_dwh_ci
3   outputs:
4     postgres_dwh_ci:
5       type: postgres
6       host: localhost
7       user: develop
8       password: develop
9       port: 5432
10      dbname: postgres
11      schema: ci
```



profiles.yml

```
1 dbt_demo: # profile name
2   target: postgres_dev_mniehoff
3   outputs:
4     postgres_dev_mniehoff:
5       type: postgres
6       host: localhost
7       user: develop
8       password: develop
9       port: 5432
10      dbname: postgres
11      schema: dev_mniehoff
```

beware: snapshots are not namespaced

CI/CD + Changes

- Usual Gitflow with Merge Requests
- Use own schema for CI
- Run + Test on CI
 - be explicit in model selection!
- Data Folds `data-diff` for creating diffs during development
 - If possible: integrate into MR

```
postgres.analytics.customers <gt; postgres.dev_mniehoff.customers
Column(s) removed: {'most_recent_order', 'first_order'}
```

| Rows Added | Rows Removed |
|------------|--------------|
| 0 | 0 |

```
Updated Rows: 59
Unchanged Rows: 41
```

```
Values Updated:
customer_lifetime_value: 59
number_of_orders: 59
last_name: 0
first_name: 0
```

```
postgres.analytics.orders <gt; postgres.dev_mniehoff.orders
No row differences
```

I could go on ..

- exposures
- metrics (WIP, the old dbt-metrics is deprecated)
- materialization
- incremental models
- ad-hoc analysis
- hooks
- authorization
- python models

```
exposures.yml
```

```
1 exposures:
2
3   - name: weekly_jaffle_report
4     type: dashboard
5     maturity: high
6     url: https://bi.tool/dashboards/1
7     description: >
8       Did someone say "exponential growth"?
9
10    depends_on:
11      - ref('orders')
12      - ref('customers')
13
14    owner:
15      name: Callum McData
16      email: data@jaffleshop.com
```

Ecosystem + Integration



<https://github.com/Hiflylabs/awesome-dbt>

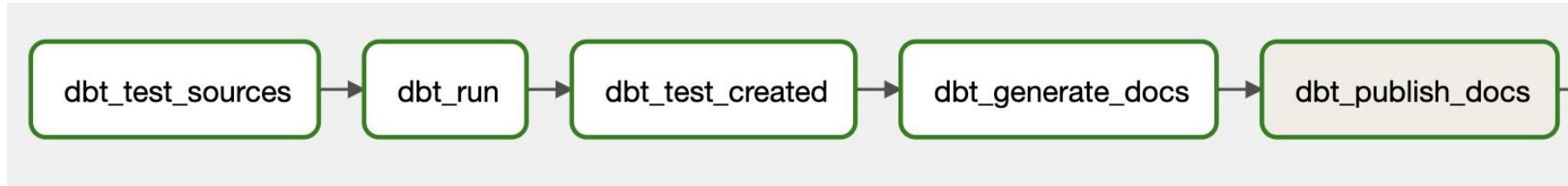
Orchestration

- Airflow, prefect, dagster, ...
- Best experiences:
 - Create a docker container based on the project
 - Use this container to run commands

data import DAG:



dbt DAG:



Vizualization + Query

- Everything that can connect to the storage
- Bonus: integration with dbt, like Lightdash

```
version: 2

models:
  - name: projects
    columns:
      - name: dashboards_created_num_total
      - name: days_since_project_created
```

The screenshot shows the Lightdash interface. On the left, there's a dark sidebar with the text "Projects". To its right is a search bar with the placeholder "Search metrics + dimensions". Below the search bar is a section titled "Dimensions" containing two items: "Dashboards created num total" and "Days since project created". A purple arrow points from the "dashboards_created_num_total" line in the YAML code above to the "Dashboards created num total" dimension in the Lightdash interface.

dbt connection ?

Your dbt project must be compatible with [dbt version 1.4.1](#)

Type *

Personal access token ?

Repository *

Branch *

Project directory path *

Host domain (for Github Enterprise)

Target name

Schema *

miro

duckDB + dbt

- Community Plugin
- DuckDB: 
- sqlite but for OLAP
- fast SQL layer for
Files/Blobs/...
- ⇒ dbt on Blob Storage/Files



```
● ● ● profiles.yml

1 your_profile_name:
2   target: dev
3   outputs:
4     dev:
5       type: duckdb
6       path: 'file_path/database_name.duckdb'
7       extensions:
8         - httpfs
9         - parquet
10      settings:
11        s3_region: my-aws-region
12        s3_access_key_id: "{{ env_var('S3_ACCESS_KEY_ID') }}"
13        s3_secret_access_key: "{{ env_var('S3_SECRET_ACCESS_KEY') }}"
```

EL_T



Why is it good?

Simplicity

- good starting point for documentation
- data lineage within the dwh can be easily understood and visualized
- supports dataset-centric visualization (opposed to query- or semantic-driven)

Engineering

- logic and config treated as code and stored in source control
- repeatable, automated test, even test driven would be possible
- easy automation and CI/CD

Ecosystem

- huge ecosystem with great integrations
- works with every major storage
- ready to adapt with new trends and technologies

What is it not?



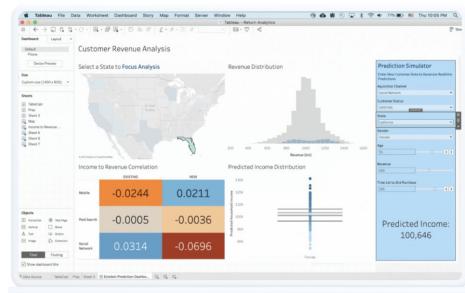
Import



Data Catalog



Visualization



Low-/No-Code

Repository & Slides



<https://github.com/mniehoff/dbt-demo>

Innovative - Trustful - Competent - Pragmatic

**Wir unterstützen Sie bei
Ihren Data-Projekten.**

Kontaktieren Sie uns gern!

📍 codecentric AG
Hochstraße 11
42697 Solingen

👤 Matthias Niehoff
Head of Data & AI
matthias.niehoff@codecentric.de
www.codecentric.de

@codecentric

