

Working with a first project

INTRODUCTION TO DBT



Mike Metzger
Data Engineer

Workflow for dbt

1. Create project (`dbt init`)
2. Define configuration (`profiles.yml`)
3. Create / use models / templates
4. Instantiate models (`dbt run`)
5. Verify / Test / Troubleshoot
6. Repeat as needed

dbt run

- Run whenever there model changes
- Or when the data process needs to be materialized
- Output provides many details on the success or failure of the various steps
- Materialized = Transformations into tables / views

```
repl:~$ dbt run
```

```
04:52:11 Running with dbt=1.8.4
...
04:52:13 1 of 1 START sql view model main.sales_data ..... [RUN]
04:52:13 1 of 1 OK created sql view model main.sales_data . [OK in 0.12s]
...
04:52:13 Completed successfully
```

Table vs View

Tables:

- Objects within a database / warehouse that hold data
- Take up space within the database
- Content only updated when changed
- Can be created by dbt

Views:

- Queryable like a table; hold no information
- Are usually defined as a select query against another table or tables
- Content generated with each query
- Can be created by dbt

Let's practice!

INTRODUCTION TO DBT

What is a dbt model?

INTRODUCTION TO DBT



Mike Metzger
Data Engineer

What is a data model?

- Conceptual, with different definitions depending on context
- Represents the logical meaning of data
- How the data and its components relate
- Helps users collaborate

What is a data model?

- Conceptual, with different definitions depending on context
- Represents the logical meaning of data
- How the data and its components relate
- Helps users collaborate

Species	# of legs	Venomous
Cheetah	4	No
Duck	2	No
Platypus	4	Yes
Rattlesnake	0	Yes

What is a model in dbt?

- Represents the various transformations
- Typically written in SQL
 - Newer versions can use Python
- Usually a `SELECT` query
- Each model represented by a text file with `.sql` extension

Simple dbt model

1. Create a directory in the `models` directory
2. Create a `.sql` file in above directory
3. Add the SQL statement to the newly created file
4. Run `dbt run` to materialize the model

```
bash> mkdir models/order  
bash> touch models/order/customer_orders.sql
```

```
select first_name,  
       last_name,  
       shipping_address,  
       item_quantity  
from source_table
```

```
bash> dbt run
```

Reading from Parquet

- Parquet?
 - Columnar binary file format
 - DuckDB can read Parquet files directly
 - `read_parquet`
 - `SELECT * FROM read_parquet('filename.parquet')`
 - Or simply the filename in single quotes
 - `SELECT * FROM 'filename.parquet'`

Let's practice!

INTRODUCTION TO DBT

Updating dbt models

INTRODUCTION TO DBT



Mike Metzger
Data Engineer

Why update?

- Iterative work
- Fixing bugs with queries / models
- Migrating to different sources / destinations



¹ Photo by Caspar Camille Rubin on Unsplash

Update workflow

1. Check out from source control
 - `git clone dbt_project`
2. Find the model in question
3. Update query contents
4. Regenerate with
 - `dbt run` or
 - `dbt run -f` (Force full refresh)
5. Check changes back to source control

YAML files

- Some updates may require changes to YAML / `.yaml` files
- Typically would require changes in:
 - `dbt_project.yaml`
 - `model_properties.yaml`

```
! dbt_project.yaml
1
2 # Name your project! Project names should contain only lowercase characters
3 # and underscores. A good package name should reflect your organization's
4 # name or the intended use of these models
5 name: 'nyc_yellow_taxi'
6 version: '1.0.0'
7 config-version: 2
8
9 # This setting configures which "profile" dbt uses for this project.
10 profile: 'nyc_yellow_taxi'
11
12 # These configurations specify where dbt should look for different types of files.
13 # The 'model-paths' config, for example, states that models in this project can be
14 # found in the "models/" directory. You probably won't need to change these!
15 model-paths: ["models"]
16 analysis-paths: ["analyses"]
17 test-paths: ["tests"]
18 seed-paths: ["seeds"]
19 macro-paths: ["macros"]
20 snapshot-paths: ["snapshots"]
21
22 target-path: "target" # directory which will store compiled SQL files
23 clean-targets:         # directories to be removed by `dbt clean`
24 | - "target"
25 | - "dbt_packages"
26
27
28 # Configuring models
29 # Full documentation: https://docs.getdbt.com/docs/configuring-models
30
31 # In this example config, we tell dbt to build all models in the example/
32 # directory as views. These settings can be overridden in the individual model
33 # files using the `{% config(...) %}` macro.
34 models:
35 |   nyc_yellow_taxi:
```


dbt_project.yml

- Contains mostly contents related to full project
 - Project name / version
 - Directory locations
- Model materialization settings (global)
- One `dbt_project.yml` file per project

model_properties.yml

- Contain settings that reference models
 - Description
 - Documentation details
 - Much more
- Can actually be named anything (with `.yml`) in `models/` subdirectory
- Can have as many files as needed

```
models > ! model_properties.yml
1  version: 2
2
3  models:
4    - name: taxi_rides_raw
5      description: Initial import of the NYC Yellow Taxi trip data from Parquet source
6      access: public
7    - name: avg_fare_per_day
8      description: The average ride amount spent per day
9      access: public
```

Let's practice!

INTRODUCTION TO DBT