

Green Fox - Data Engineering Project

Business Concept	3
Architecture	3
Database	4
Ingestion	4
Transformation	4
Visualization	4
Milestones	5
Week 1 - Learn new skills	5
Week 2 - Setup infrastructure/project	5
Setup	5
Prerequisite	5
GCP / Big Query	5
Installing Python and dbt on Windows	6
Airbyte	6
Explore datasets	7
Week 3 - Implementation: Staging layer (Preview)	7
General	7
Data sources	7
Prerequisites	7
Best practices for dbt modeling	7
dbt packages	7
Development principles	8
Github dataset	8
Location	8
Objective	8
Steps	8
Examples	9
Stack Overflow dataset	10
Location	10

Objective	10
List of tables	10
Steps	10
Examples	11
Google sheet	11
Location	11
Steps	11
Additional tasks	12
Testing	12
Using variables	12
Github	12

Business Concept

Our business goal is to get a better understanding of open-source technology trends.

We want to understand which are the "hot" projects and see how the monitored projects/companies are performing compared to others.

Our goal is to better understand the development of the open source tech market - we want to answer these questions using several data sources.

Datasets:

We will use the following freely available datasets provided by BigQuery:

- Github
- Stack Overflow
- (Hackernews)
- (TBD)

Configuration:

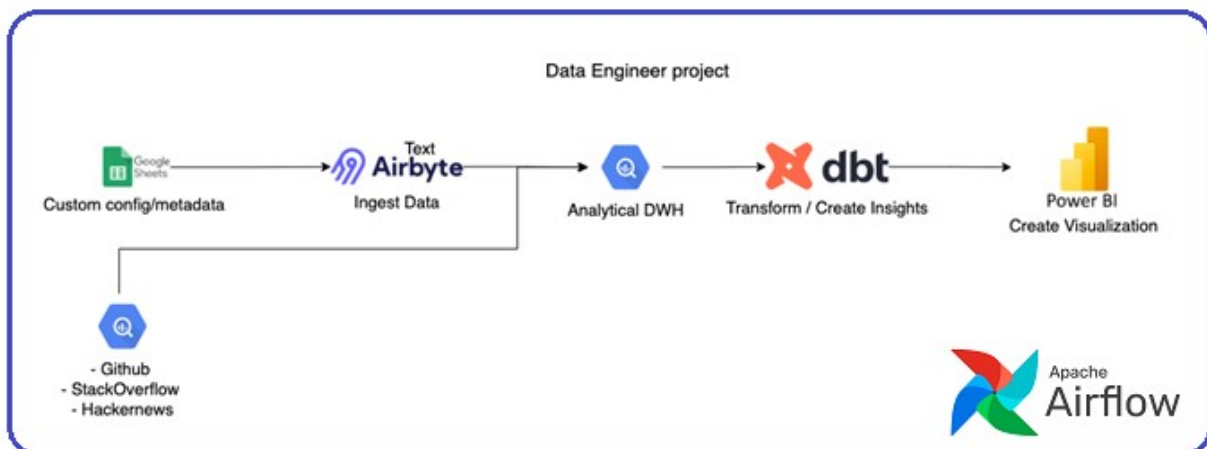
We will get the configuration data from a Google Sheet.

The target is to create a common data repository for these datasets, load them into a central database, apply the necessary transformations and create useful insights from the available data, that is presentable with Data Visualizations.

Example question:

Check if there is a correlation between a trending Stack Overflow post about a tool and the Github stars statistic.

Architecture



Database

We will use Google BigQuery to build our Data Warehouse. More information on the free tier is here: <https://cloud.google.com/free>

Ingestion

To ingest the data from the source system we will use an open-source ingestion tool called Airbyte.

Airbyte is easy to learn, it has an intuitive interface, and it can be run as a docker container
Documentation: <https://docs.airbyte.com>

Transformation

To transform/digest the raw dataset we will use dbt which is a Python-based open source SQL templating framework. The advantage of this is that it is easy to develop and software development best practices can be applied, for example:

- All code lives on Github
- You can write DRY code (Jinja-based SQL templating)
- Out-of-the-box data quality testing
- CI pipeline

dbt is an open core project, so the fully featured CLI is available for free, but there is SaaS integrated environment/scheduler that is free for one developer seat – but not necessary for this project

Documentation: <https://docs.getdbt.com>

Get started: <https://courses.getdbt.com/courses/fundamentals>

Visualization

Our choice of tool is Power BI for the data visualization part of the project. The desktop version of Power BI is free to use. Please download Power BI through the website and not from Microsoft Store

Documentation: <https://docs.microsoft.com/en-us/power-bi>

Milestones

Weekly meeting with Hiflylabs: We can join the sprint planning to answer any questions and we can do a code review/demo session at the end.

Week 1 - Learn new skills

Hiflylabs: Introduction, short presentation about Hiflylabs, the project itself and dbt/modern data stack fundamentals.

Learn the fundamentals of dbt, install it on the developer machine and do the available [courses](#).
Get started with Big Query - [udemy](#) or any other resources that is available.

Week 2 - Setup infrastructure/project

Hiflylabs can hold a 1-2 hours workshop on setting up the development environment.

Get an understanding of the dataset:

Project setup

- Access for everyone for BigQuery
- Setup a dbt project with the BigQuery adapter
- Setup the repository
- Get started with Airbyte, Google Sheet connector

Setup

Prerequisite

- Docker, Git and Visual Studio Code are installed on the development machines

GCP / Big Query

- Register a new free GCP account: <https://cloud.google.com/free>
- Create a new project
- Create a new Big Query dataset (location somewhere in the US)
- Install the gcloud CLI: <https://cloud.google.com/sdk/docs/install#windows>
- Make sure that your user has the following permissions:
 - Hamburger icon/IAM&Admin/IAM/Add*
 - BigQuery Data Editor
 - BigQuery User
 - BigQuery Job User

Git

- Setup a new Git repository (Github, Gitlab, Bitbucket)

Installing Python and dbt on Windows

- Install Python 3.8: <https://www.python.org/downloads/release/python-3810/> - Make sure that click the "ADD TO PATH" option in the installation
- To install the python virtualenv package run:
`python -m pip install virtualenv`
- Go to the directory where you store your codebases.
If you don't have any, we recommend creating a Users\<your-user-name>\Code directory
- To clone the git repository run:
`git clone git@github.com:<your-github-user>/<your-dbt-repo>.git`
- Change the current directory to the cloned repository
- To create a virtual environment run:
`python -m venv .venv`
- To activate the virtual environment run:
`.\.venv\Scripts\activate.ps1`
- To install dependencies run:
`pip install dbt-core`
`pip install dbt-bigquery`
- Create an empty .profiles.yml file at
/Users/<your-username>/.dbt/profiles.yml
- To setup gcloud CLI run:
 - `gcloud config set project <your-project-name>`
 - `gcloud auth application-default login \`
`--scopes=https://www.googleapis.com/auth/bigquery,\`
`https://www.googleapis.com/auth/drive.readonly,\`
`https://www.googleapis.com/auth/iam.test`
- Init a dbt project:
`dbt init <your-project-name>`
- Test configuration (cd into the dbt project directory)
 - `dbt debug`
 - `dbt run`

Additional resources:

[dbt style guide](#)

Airbyte

Run the following commands in your terminal:

- `git clone https://github.com/airbytehq/airbyte.git`
- `cd airbyte`
- `docker-compose up`

Explore datasets

Dive into the datasets, create exploratory queries and get an understanding of time granularity, tables and columns.

Github archive: <https://console.cloud.google.com/bigquery?project=githubarchive&page=project>

Stackoverflow:

[bigquery-public-data.stackoverflow](https://stackoverflow.com/questions/tagged/bigquery-public-data)

Week 3 - Implementation: Staging layer (Preview)

The objective of the third week is to create the staging layer with all data sources including the google sheet that contain the list of companies and repositories that will be analyzed in depth.

General

Data sources

- Stack Overflow (BigQuery)
- Github (BigQuery)
- Google sheet (local file)

Prerequisites

- Initialize a local dbt project with BigQuery credentials and a related private Github repository.
- Locally installed Airbyte.

Best practices for dbt modeling

- [dbt style guide](#)
- Create a primary key in every model, called `_pk`.
- Make sure that there is a timestamp column which indicates the creation or load datetime of each record.
- Make sure that every datetime is in UTC format.
- Add `_datetime_utc` postfix to every datetime column name.
- Add `is_` prefix to every boolean type ('Yes'/'No', 1/0) column name and convert it to boolean (true/false).

dbt packages

- [dbt codegen package](#)
- [dbt expectations package \(optional\)](#)

Development principles

- During development try to work with smaller subsets, don't run the queries on full tables if it is not necessary and the smaller subset has the same result.
- Estimated process cost can give you hints in BigQuery.
- Use the Preview dataset option when possible (no table scan there).
- In some cases, instead of dbt run use dbt compile before, you can see the query before it actually runs and make changes if it is needed.

Github dataset

The Github dataset is located in the BigQuery environment and curated and updated by Google.

More information about the BigQuery public datasets:

[Big Query public datasets](#)

The Github data is available in three levels of aggregation: day, month, and year. In this project, the daily data source will be used.

Location

[GitHub Activity Data dataset](#)

Objective

Create a staging layer for Github data in dbt according to the dbt principles and best practices.

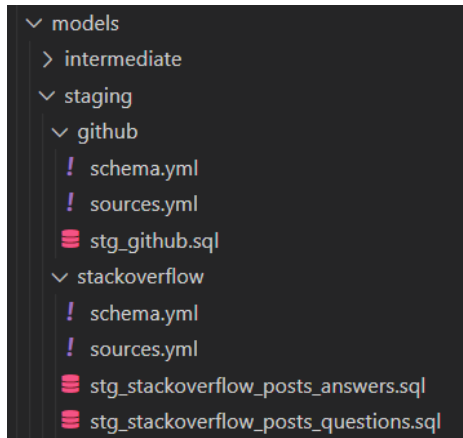
Steps

- Add the [dbt codegen package](#) to the project.
- Add [folder lever materialization](#) configurations in the dbt_project.yml file.
 - Make sure that the staging models all materialized as views.
- Create the folder structure inside the dbt project.
- Create the sources.yml file that contains the Github data using [dbt codegen package](#)
 - Use BigQuery's [wildcard](#) function to union all the daily tables, but make sure that you use a limited dataset for development.
 - Ensure that the dbt sources.yml file can handle the wildcard.
 - [Check if you have troubles adding wildcard to sources](#)
- Create the staging model materialized as a view using 2022 data only.
 - Optionally, the [dbt codegen package](#) can be used for that.
- Add best practices to the model.
- Make sure that the RECORD type fields are readable in the model the same way as it is in the source tables.
- Make sure that there is a unique field in the table.
- Create the schema.yml using [dbt codegen package](#)
- Add uniqueness and not null [tests](#) for the primary key field.
- Add proper [documentation](#) to the project.

- Add a few descriptions to columns which are not intuitive enough to understand without context.

Examples

Possible folder structure



Possible staging model output

```
with source as (
    select * from {{ source('github', 'github_daily') }}
),
final as (
    select
        id as _pk,
        type,
        public as is_public,
        payload,
        repo.id as repo_id,
        repo.name as repo_name,
        repo.url as repo_url,
        actor.id as actor_id,
        actor.login as actor_login,
        actor.gravatar_id as actor_gravatar_id,
        actor.avatar_url as actor_avatar_url,
        actor.url as actor_url,
        org.id as org_id,
        org.login as org_login,
        org.gravatar_id as org_gravatar_id,
        org.avatar_url as org_avatar_url,
        org.url as org_url,
        created_at as created_at_datetime_utc,
        id,
        other
    from source
)
select * from final
```

Stack Overflow dataset

The Stack Overflow dataset is technically located in the same place as the Github dataset.

Location

[Stack Overflow dataset](#)

Objective

Create a staging layer for Stack Overflow data in dbt according to the dbt principles and best practices.

List of tables

Stack Overflow has a different structure for storing the data.

In this part of the project, two tables will be used:

- posts_questions
- posts_answers

Steps

- The same steps as it is for Github data.
- Different methods:
 - The granularity of the Stackoverflow data is different because one table contains all the records.
 - Because of this, the time period should be filtered out inside the staging model.
 - Find the column that will function as a key for that purpose.
 - Explore the dataset by matching the rows with actual Stackoverflow questions to determine which column can be the period filter column.
 - The objective is to have every question and the related answers from 2022.

Examples

Possible outcome of the dbt DAG



Google sheet

The objective is to create an ingestion pipeline that ingest the google sheet data into a BigQuery table. Then create a staging view in dbt based on that ingested BigQuery table. For the first task Airbyte will be used.

Location

[Company details](#)

This spreadsheet is being used and updated by other consumers, therefore the design does not match the expected database look. The goal is to create a table inside the project database based on that spreadsheet and have the same set of data inside that table.

Steps

- Setup a source (Google Sheets) and a destination (BigQuery) following the Airbyte documentation.
 - [Google Sheet as source](#)
 - [BigQuery as destination](#)
- Create a new dataset in the project database called *raw_airbyte*. Inside the dataset create a table that can be the destination of the google spreadsheet. Make sure that the column names are in a database standard format (spaces, lowercases, etc.).
- Create a connection that ingests the source data into the BigQuery table.
- Apply one-time ingestion of the source.
 - The scheduling of the Airbyte ingestion will be configured in a later phase of the project.

- Create a staging model in dbt as a view based on the already created and ingested google sheet table.
- Apply the same principles and best practices as for the previous data sources. Notice that there is no datetime or unique key and it should be created manually. It can be created with Airbyte or inside the dbt model.

Additional tasks

Testing

- dbt has a strong belief in the importance of testing and is dedicated to support proper testing methods throughout the whole project.
- In the staging layer uniqueness and not null tests are defined, but it is possible to add other tests if it is necessary.
- Tests in dbt are basically queries that also have process cost. It is a double-edged sword to find the balance in the number of tests applied.
- After exploring the source data, [dbt built-in tests](#) or [other tests](#) can be added to the staging models considering process costs.
- Tests can be added in a later phase of the project when the business requirements are clearer.
- Few possible examples:
 - Key fields always have values.
 - Referential integrity between related tables (TBD in the intermediate layer).
 - Minimum or maximum date of the records matches the expected content of the table.

Using [variables](#)

- It is more sophisticated to add project-level variables to handle and manage time periods and other dimensions.
- With this solution, the time period can be managed from one source.
 - The time period can be controlled by changing only this one project-level variable.
 - It is possible to use only one variable and then manipulate that variable depending on which data source it is being used.
 - It should be considered that it can have more maintenance overhead at the end then using more variables for each dataset.

Github

It is common to use feature branches for development cycles. In this case, everyone has their own branch inside the repository. In that phase of the project, there is no need for feature branches as one individual working on their own branch only. Later, when the project is more complete and it is needed to add a new source or test some function outside the main branch, a feature branch should be created.

SAJAT HASZNOS LINKEN:

How to create an ELT pipeline from Postgres to BigQuery with Airbyte

Megmutatja, hogy Airbyte-ban miként állítsuk be Destinationnek a Bigquery adatbázisunkat

https://www.google.com/search?q=airbyte+to+bigquery&rlz=1C5CHFA_enHU856HU874&source=lnms&tbm=vid&sa=X&ved=2ahUKEwiNi5T1je_7AhVJPOwKHc1EDpIQ_AUoAXoECAEQAw&biw=1384&bih=761&dp_r=2#fpstate=ive&vld=cid:9fcc590b,vid:w1VkzQQIZzs

<https://docs.airbyte.com/integrations/destinations/bigquery/#service-account-key>

Hogyan szedjük le Google Sheetből a doksit Airbyte segitsegevel:

<https://www.youtube.com/watch?v=oeshl0H1JcU&t=3s>

JOINING TABLES IN BIGQUERY

https://www.google.com/search?q=bigquery+left+join&rlz=1C5CHFA_enHU856HU874&oq=bigquery+left+join&aqs=chrome..69i57j0i22i30l9.4505j0j7&sourceid=chrome&ie=UTF-8#fpstate=ive&vld=cid:d81c227b,vid:kXwMTAU3C8M

DBT SETUP etc., getting started

<https://docs.getdbt.com/docs/get-started/getting-started/overview>