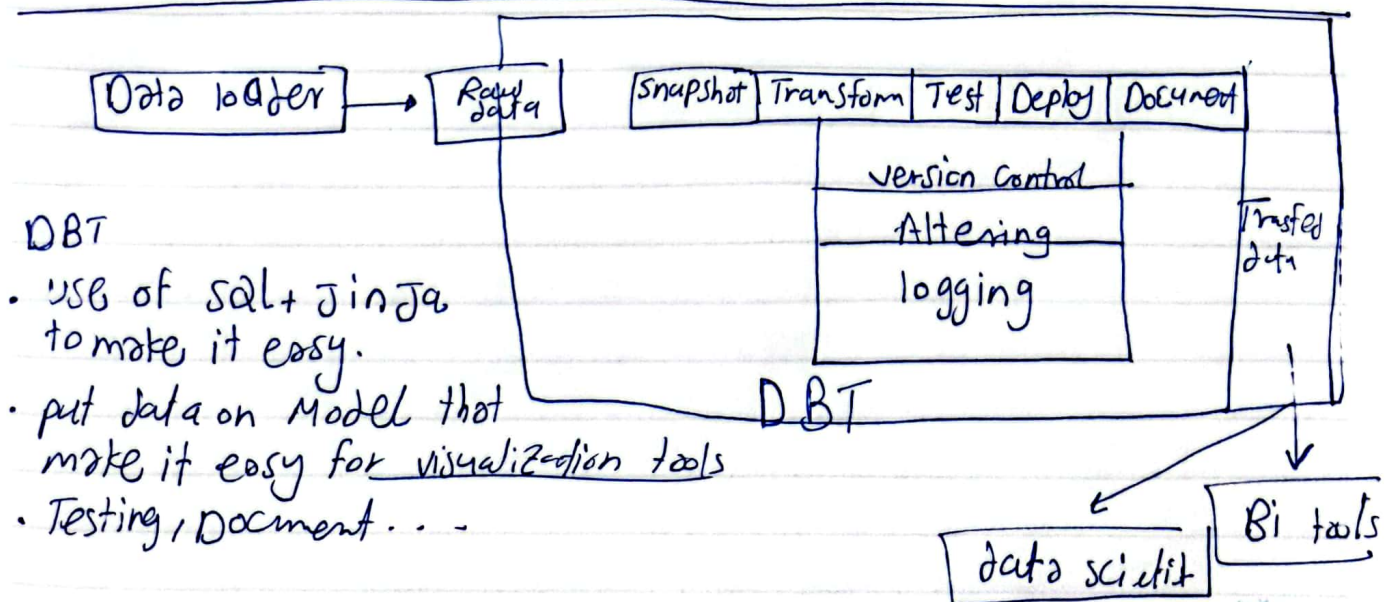


# Week 4: Analytics Engineering

"Transform data in DW [Data warehouse]  
to Analytical views"

## Content

- Introduction to analytics Eng
- What's dbt?
- Start dbt project
  - using BQ + dbt Cloud (Alternative A)
  - using Postgress + dbt core (Alternative B)
- Development of dbt models.
- Testing & document dbt models.
- Deploying dbt project.
- Visualize transformed data.



# what's Analytics Engineering.

It's role that fill the gap between Data Analysts [presentation] and Data Engineer [physical work]

- Data loading (stitch)
- Data storing (Data Warehouse)
- Data Modeling (dbt, data form)
- Data presentation (BI tools, Google studio, Tableau..)

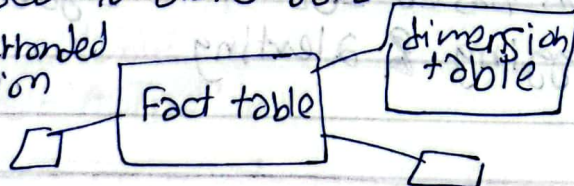
Now there is difference between  
[Analytics Engineer & Data Analyst]

## Dimensional Modeling

- Method that's used to store the data, store data in a certain structure to make it fast for Data analyst

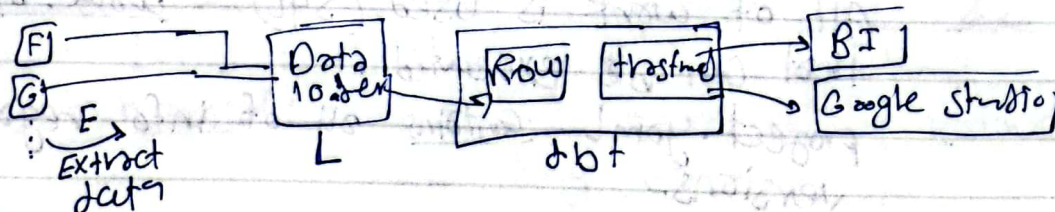
- star schema

schema used to store data  
Fact table surrounded with dimension table.



## Introduction to dbt [data build tool]

- we use dbt for transformation of raw data



ELT → Transformation means to transfer raw data to dimensional model for fast analytics





- dbt transform process how data in our Data warehouse to transformed data which can be later used by business intelligence.
- dbt allows:
  - Develop model.
  - Test and Document model.
  - Deploy model with version control.
- • dbt works by defining modeling layer which turn tables into models.
- dbt has options
  - dbt (Core) offline with postgres
  - dbt (Cloud) with bq or etc
- dbt cloud support:
  - web based IDE
  - Job schedule
  - logging & alerting

### dbt cloud

- cloud give you folders like github with git ignore contains files or folders names to ignore.
- All of work is used (SQL) work.sql that can be executed.
- project.yaml contains all of info regarding versions.



dbt core

- dbt-project.yml

name: 'tax-rides-ny'

version: '1.0.0'

config-version: 2

profile: 'pg-dbt-workshop'  
 ↳ which profile dbt use local, cloud

Build first dbt Models

dbt Model my-model.sql

{%

Config(materialized = 'table')

{%}

select \*

from staging.source-table

where record-state = 'ACTIVE'

Compiled Code

SQL Code

dbt

this {% for jinja code

- dbt models [ is essentially select query ]

- materialized strategy for presenting dbt like

1. table: Model will be rebuild each time.

2. incremental: add records incrementally rather than rebuilding each run.





## From clauses

- to define source of data.
- we can define (Source Freshness) to get fresh records only.

## Macros

pieces of code in jinja that can be reused

e.x Config(), source()...

- Macros allows to add features that aren't otherwise available.

- if statement, loops.
- use environmental variables.
- operate on result of one query to make decision to another.

- Macros stored in .sql files which are stored in macros directory

- {% ..... %} statement (macros add(, ))
- {# ..... #} comment
- {{ ..... }} Expressions

e.x

```
{# Macro to return description of payment #}
{% macro get_pay_description(payment-numb) %}
  Case {{ payment-numb }}
    when 1 then 'credit'
    when 2 then 'cash'
    when 3 then 'unknown'
  end
{% endmacro %}
```



in the query

select

{ { get-pay-description(payment-numb) } } as payment-typedes

↳ Fetch on data for payment-numb and replace each number with crossponding value and then save it as payment-type-description

## • Packages

- similar to other languages we can save macros to reuse them, and also we can download packages from internet.
- packages.yml

## • Variables

- variables can be used and defined across all of the project.
- can be defined in different way:

• vars:

payment-numb : [1, 2, 3, 4 ...]

- dbt build --m model.sql --var "is-test-run:false"  
↳ when we build to check for var then make decision using macros

• { % (if var('is-test-run', default=true) % }

limit 100

{ % endif % }



## \* Referencing older model to new one

- ref() is used for:
  - reference tables and views
  - reference seeds.

## Testing and documenting

- Not required step, but expected in any professional setting.
- dbt provide few predefined test, but also you can use your test [custom].  
test written in .yml

e-x

columns:

- name: tripid

tests:

- unique:

severity: warn

Action to do

- not null:

severity: warn

↳ that would print warning in cli.

e-x  
custom

name: payment-type-description

tests:

- accepted-values:

values: [1, 2, 3, 4, 5]

severity: warn

↳ warn when value of payment-type-description not same.

## Documenting

models:

- name: stg-yellow-tripdata

description:

Trips made by New York City's iconic yellow taxi.

columns:

description used to describe the work or model and also we write it in .yml.

dbt provide a way to generate document and render it to a website.

it contains:

- information about project [model code, dependencies, sources, Descriptions from yml]
- information about Data warehouse

dbt can be rendered to website locally or in Cloud.

- > dbt docs generate
- > dbt docs server

---

## Deployment Basics

dbt - running models on development environment

- Deployment workflow:

1. Develop user branch.
2. open PR to merge into branch.
3. Merge user branch into main.
4. Run new models in production environment using main branch.
5. schedule models.





project is set of jobs :

- Job is set of commands such build, test
- Job may be triggered manually or scheduled.
- other tools like airflow, Cron can also work.
- If dbt source freshness command run you will view results at the end with its freshness

## Continuous Integration (CI)

CI job will do things such as build, test....etc when a certain circumstances happen.

"triggered on certain conditions"

e.x

- Github/Gitlab pull requests that trigger CI using webhook
- AirFlow, cron job, cloud scheduler....

## Deployment using dbt cloud

- dbt enforce you to create PRS for CI purpose
- we will create new environment
- we will add this jobs, and check settings like generate Docs

to run model sql

1. dbt seed      2. dbt run      3. dbt test

- After job finish you can see .logs . Document contain all info about project.



## Deployment local

- dbt could be deployed local, but with limited features.
- settings will be edited from profiles.yml

## Data visualization

- Google Data Studio , - tableau
- Meta base "offline"

"we used dbt with BQ to structure data into a good way to visualize it using Data studio"

### - GDS

- online tool to convert data into reports and dashboard
  - once you select data or tables you want you will see item type with description for all of them
  - you can press create report which after that can handle data types and plots.
  - GDS have different widgets.
  - GDS can connect to BigQuery and even more....
- GDS support different types of charts  
[ bar, lines, graphs, PI chart, pivot ..... ]  
all what you think of python.

- gives you interactive dashboard like drop box





Date: \_\_\_\_\_ No: \_\_\_\_\_

- drop box to choose interval date.

- You can add new column :

- add field
- Name it pickup month → we have only daily records in our data
- write formula MONTH(pickup-datetime)
- save it.

- Same if we want by year  
YEAR(pickup-datetime)

### Tricks

- You can share dashboard with others.
- schedule report in mail.
- Set freshness of dashboard.
- This can help alot for Trade software with your own python code.
- live screen for monitor data.

### Meta base

- Can do most of GDS work.
- All of that offline.
- Can work with postgres.
- Have flexible dashboard UI.

