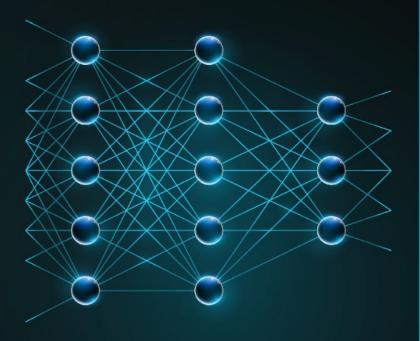


MACHINE LEARNING EDITION

# DataOps – Is there a cure for the pains of data science team?

Paweł Bogumiło
BI & Analytics Engineer @RTB House





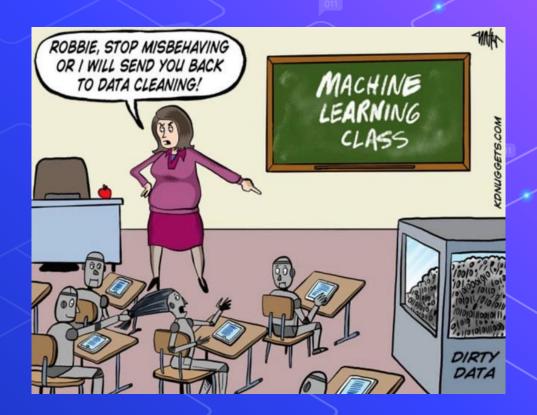






A wild DATASET appeared!

What data scientists spend 80% of their time on?



# Who remember that?

Harvard Business Review



ATA

## Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

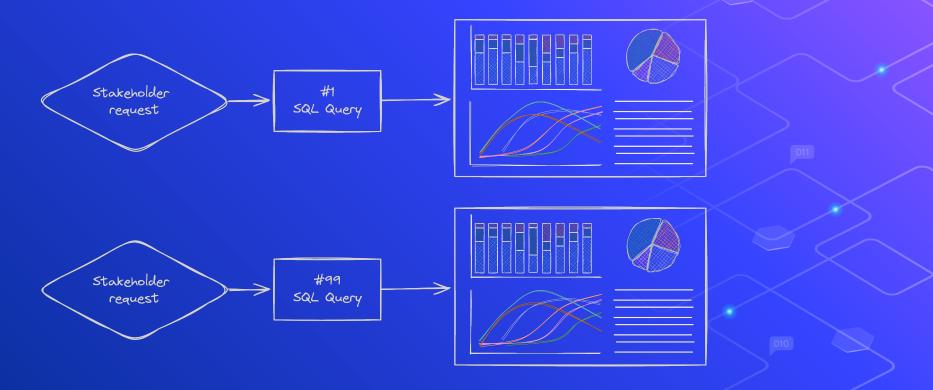
hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social experience. As one LinkedIn manager put it, "It was like arriving at a conference reception and realizing you don't know

Per Forrester Research, 60% of the data and analytics decision-makers surveyed said they are not very confident in their analytics insights. Only 10% responded that their organizations sufficiently manage the quality of data and analytics. Just 16% believe they perform well in producing accurate models.



https://xkcd.com/1838/

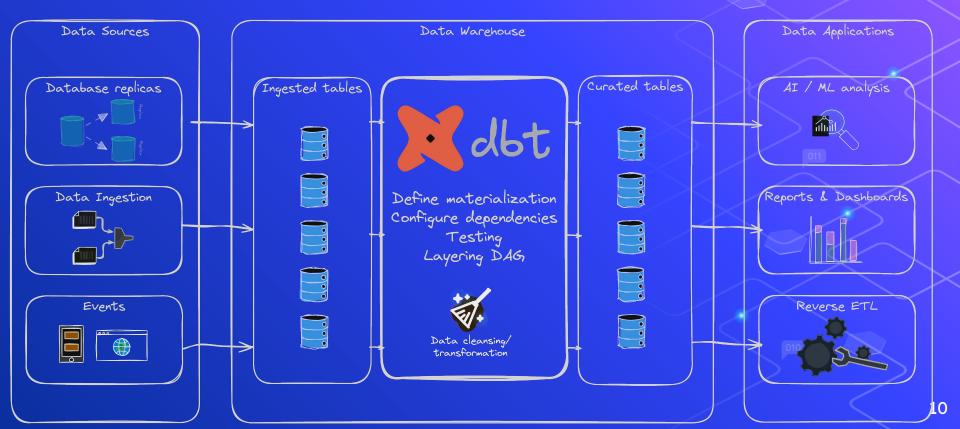
### This is not DRY



### The schedule

- 1 Best SoftEng practices in data?
- 2 'dbt run' for free?
- 3 Blue-Green-Deployment vs. DWH
- 4 Differences in CI/CD process between Snowflake & BigQuery
- 5 Landcape of dbt-enabled products

### **EL[T]:** T means transformation



### dbt for DataOps?

- Testing both code and data
- Versioning
- Proper management of branches and PRs
- Automation of the above
- Multiple environments
- Containerization
- Parameterization of the process

# Transformation flow: Types of data models

transform src backend.vml version: 2 ... stg\_customer.sql sources: {{ - name: warehouse  $\bullet \bullet \bullet$ dim\_customer.sql config( description: Data from materialized = "table" with customers as ( process. tables: select \* }} from {{ ref('stg\_customers') }} - name: customers columns: with source as ( - name: customer state as ( select \* tests: select \* from {{ source('warehouse) - not\_null from {{ ref('stq\_state') }} - name: orders select c.customer\_id, renamed as ( columns: c.zipcode, select customer\_id, - name: order\_id c.city, tests: zipcode, c.state\_code, - not\_null city, - unique s.state\_name, state\_code, - name: cust\_id c.datetime\_created, datetime created::TIM tests: c.datetime\_updated. datetime\_updated::TIM - relationshi c.dbt\_valid\_from::TIMESTAMP as valid\_from, dbt\_valid\_from, to: source dbt\_valid\_to CASE field: cus from source WHEN c.dbt\_valid\_to IS NULL THEN '9999-12-31'::TIMESTAMP ELSE c.dbt\_valid\_to::TIMESTAMP - name: state END as valid\_to select \* from renamed from customers c join state s on c.state\_code = s.state\_code

mart

stage

### Open source means free, right?









### Blue-Green Deployment

- o develop
- build on a copy of prod
- swap (schemas)



# Characteristic of your DWH is crucial for right CI/CD setup



### ZCC

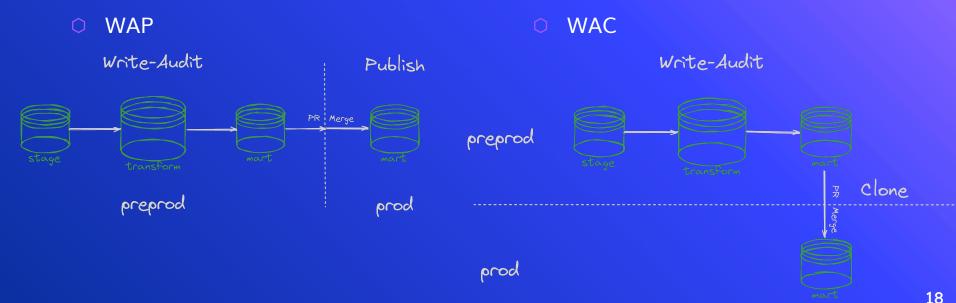
Source: macro by leremy Yeo



- No schema/dataset rename
- No database swapping

# Characteristic of your DWH is crucial for right CI/CD setup





# CI/CD with GH Actions Workflows

```
name: On PR
 pull_request:
   branches:
     - main
 on_pr:
   name: PR
   runs-on: ubuntu-latest
   permissions:
     contents: 'read'
     id-token: 'write'
   steps:
   - uses: actions/checkout@v4
   - uses: actions/setup-python@v4
       python-version: '3.9.x'
   - uses: 'google-github-actions/auth@v2'
       credentials ison: ${{ secrets.DBT GOOGLE BIGOUERY KEYFILE CI
       pip install dbt-bigguery
       dbt deps
   - run: |
       dbt run -t preprod
       dbt test -t preprod
```

```
push.yml
name: On Push
 push:
    branches:
     - main
iobs:
  deploy:
    name: Push to Prod
    runs-on: ubuntu-latest
    permissions:
     contents: 'read'
      id-token: 'write'
    steps:
    - uses: actions/checkout@v4
    - uses: actions/setup-python@v4
        python-version: '3.9.x'
    - uses: 'google-github-actions/auth@v2'
        credentials_json: ${{ secrets.DBT_GOOGLE_BIGQUERY_KEYFILE_PROD }}
        pip install dbt-bigquery
        dbt deps
    - run: l
        dbt run -t prod
    - run: |
        dbt test -t prod
```

```
schedule:
 - cron: '20 2 * * 1-5'
 inputs:
   on target:
     description: Choose the target to be used
     default: 'preprod'
     options:

    preprod

     required: true
     type: choice
 if: contains(fromJSON('["preprod", ""]'), github.event.inputs.on_target)
 runs-on: ubuntu-latest
 timeout-minutes: 30
 steps:

    uses: actions/checkout@v4

 - uses: actions/setup-python@v4
     python-version: '3.9.x'
 - uses: 'google-github-actions/auth@v2'
     credentials ison: ${{ secrets.DBT GOOGLE BIGOUERY KEYFILE CI }}
     pip install dbt-bigguery
     dbt run -t preprod
on_prod:
 if: contains(fromJSON('["prod", ""]'), github.event.inputs.on_target)
 runs-on: ubuntu-latest
 timeout-minutes: 30
 steps:

    uses: actions/checkout@v4

 - uses: actions/setup-python@v4
     python-version: '3.9.x'
 - uses: 'google-github-actions/auth@v2'
     credentials ison: ${{ secrets.DBT GOOGLE BIGOUERY KEYFILE CI }}
     pip install dbt-bigquery
     dbt dens
 - run: I
     dbt run -t prod -m mart+
```

# Built-in lineage and documentation

Source: github.com/bogumilo/boostrank

budget\_neighbors

nss\_neighbors

competitors\_area\_year

top\_neighbors

gsheets.ranking

stg\_ranking

top\_year\_avg



### What next?

#### dbt for Machine Learning

- dbt\_ml\_preprocessing package for light feature engineering
- Article on ELT-ML pipeline inside dbt and how ML team could benefit
- Meetup talk video by Sam Swift, VP of Data & Al at Bowery Farming
- Data Science use cases breakdown directly from the dbt Labs

#### dbt + Airflow

- Cosmos package for creating Airflow tasks/jobs from dbt project
- dbt+Airflow by Astronomer authors of the package

#### dbt + Dagster

• <u>Dagster + dbt</u> - official article introducing to dbt support

#### Write-Audit-Publish

- WAP: dbt + BigQuery intro video
- WAP: dbt + BigQuery PDF from the above presnetation
- <u>Conf talk referencing WAP strategy</u> by employees video intro by Michelle Ufford, Staff Engineer at Netflix

#### dbt-core technicalities

- analytics engineering handbook by Hifly Labs
- Talk on custom CI setup for dbt with automated reports a'la terraform plan in GitHub Actions
- Article on User Defined Functions managed inside dbt project and included in lineage

#### other

- DataOps Manifesto that you can sign on
- awesome-dbt collection on GitHub

#### Companies deploying and utilizing dbt

Company	Location	Technologies and Tools	Founding Year	Clients	Website
Rittman Analytics	London, UK	dbt, Looker, Snowflake, Google BigQuery, Segment,	2017	Rittman Analytics Clients	rittmananalytics.com

# Thanks!

LI: /in/pawelbogumilo

E: pawel@bogumilo.co



QR for feedback

