

Building a data science solution for an NGO when you don't know what infrastructure it will run on

A case study predicting tutor supply and demand mismatch



DataKinduk



- Who am I?
- Who are DataKind UK & The Brilliant Club?
- What did we set out to do?
- How we did it!
- The final result

Adam Hill

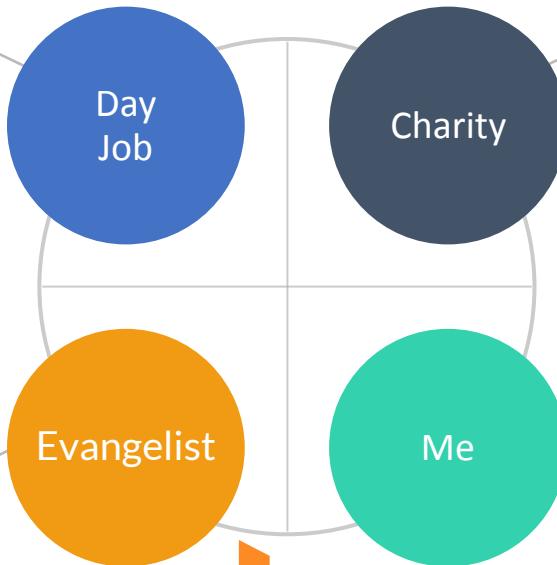
tldr;

**COMPLY
ADVANTAGE**

Senior
Data Scientist
Financial Crime Fighter

Former Royal Society
Entrepreneur in Residence
Recovering astrophysicist

UNIVERSITY OF
Southampton THE ROYAL
SOCIETY



DataKind volunteer
DataDive volunteer; Data Ambassador; Committee member



DataKind

Find me at

@astroadamh



www.linkedin.com/in/adambenhill/
Horsewithapointyhat.com

A circular inset photograph in the bottom-left corner shows four people laughing together in an office setting. One man in a red shirt is smiling, another in a black shirt is laughing heartily, and two others are partially visible, also smiling. The background shows office windows and a modern interior.

DataKind UK

We're a charity that
builds data science
capacity in charities,
social enterprises, and
the public sector





The Brilliant Club

The Brilliant Club supports less advantaged students to access and succeed in the UK's most competitive universities.

They do this by mobilising the PhD community to support students in schools via their courses and tutoring programme.

Team Effort!

Six dedicated DataKind
volunteers

- Meg
- Kate
- Leo
- Chris
- Chris





What was the
problem we
needed to solve?

The Challenge

The Brilliant Club's 21-26 strategy focuses on two goals: student access and student success. TBC wants more pupils from less advantaged backgrounds to access university

To achieve this they need:

- Scalable systems and processes
- The right number of PhD researcher volunteers in the right areas to meet school demand

How can we estimate the number of new tutors TBC needs to recruit over a coming school year?



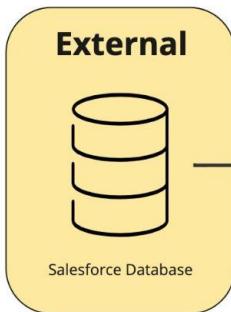
Current Solution

Spreadsheets Galore!

- Multiple teams, copy and paste snapshot data from database reports
- Manually build reports and charts, but it's easy to lose track of hard-coded numbers versus calculated figures

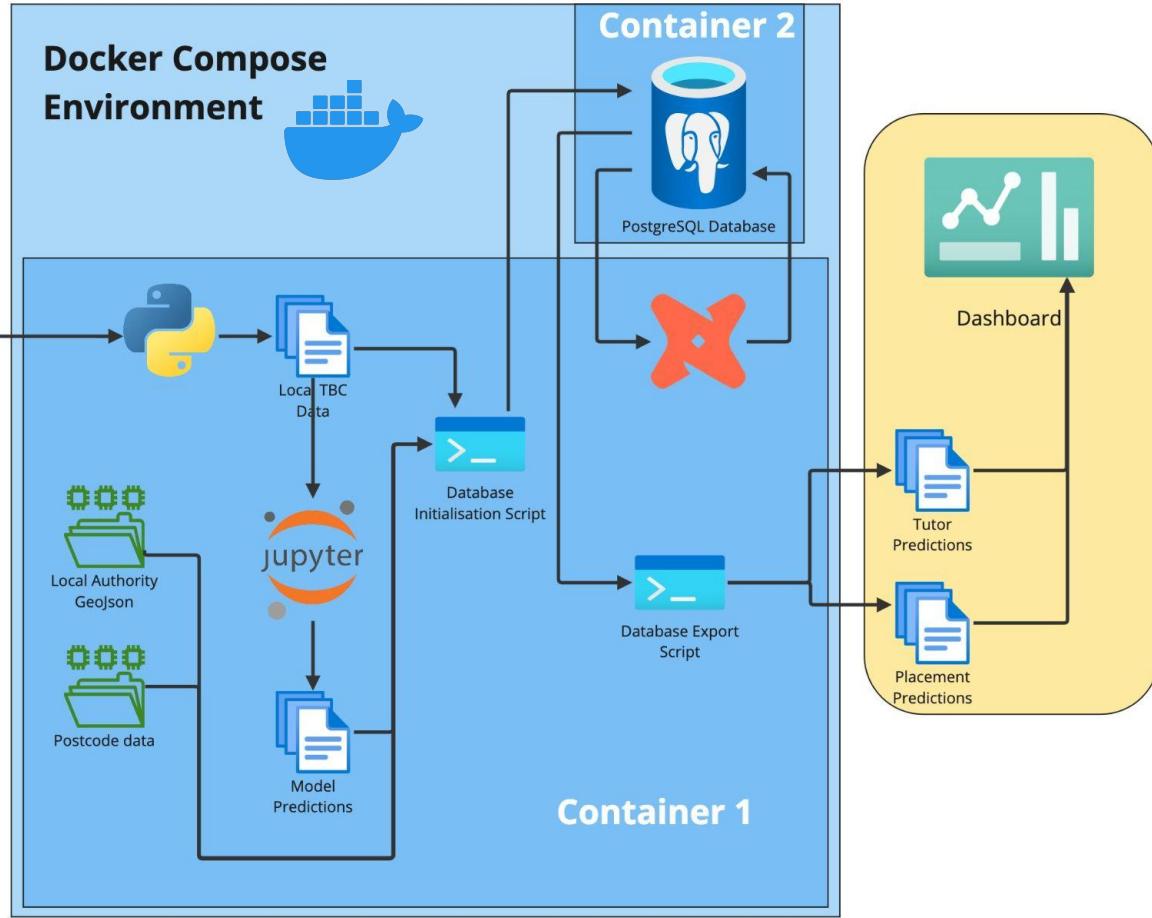


tldr; Our Solution



Easy-to-use data pipeline

- One click data build process
- Generates fresh data directly into an established dashboard



Collaboratively working as a volunteer team

Coordination and support

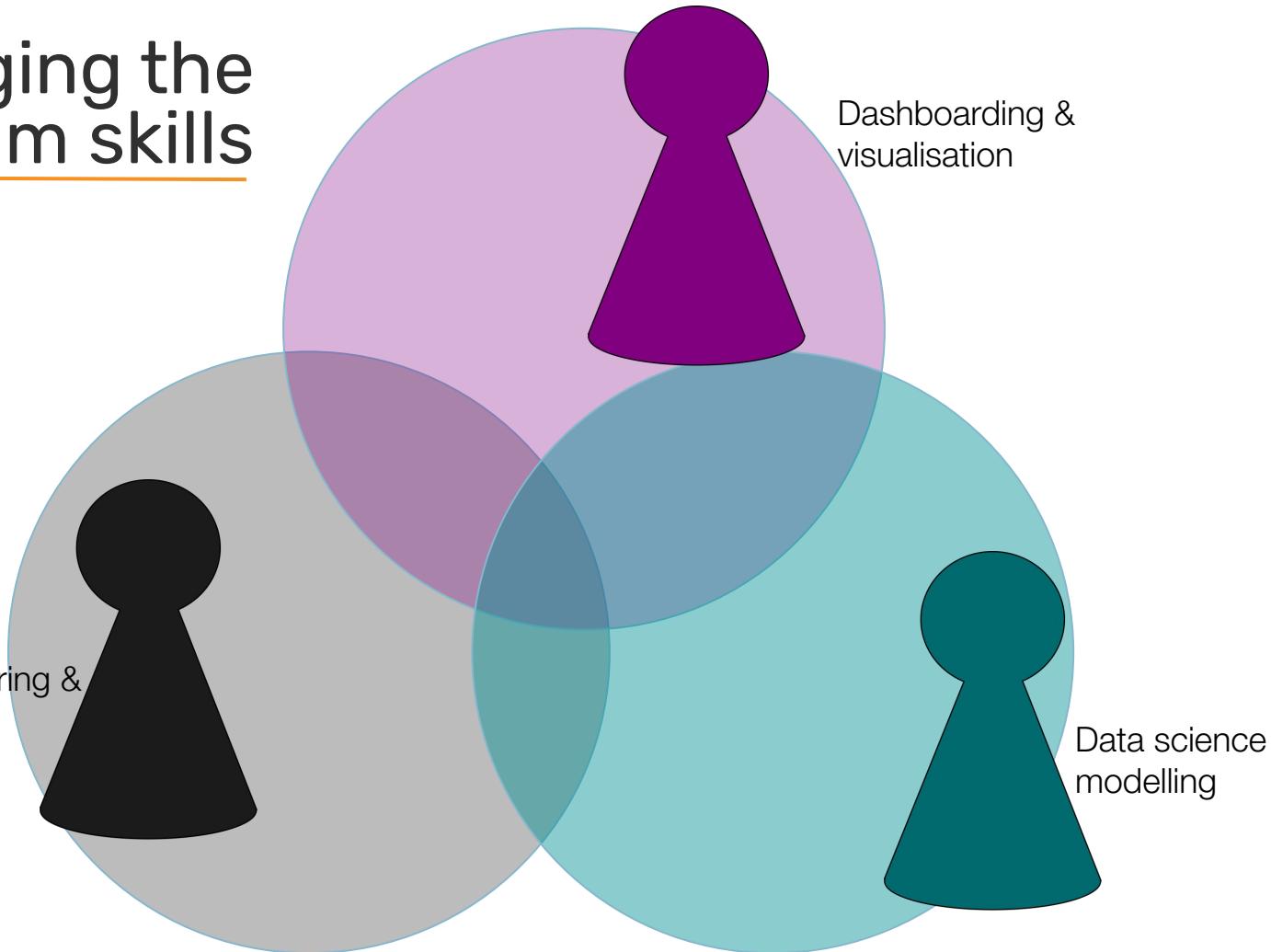
As volunteers working in our spare time we have a challenge to make sure we can move forward coherently but with flexibility to account for everyone's lives.

Our manner of working needs to allow

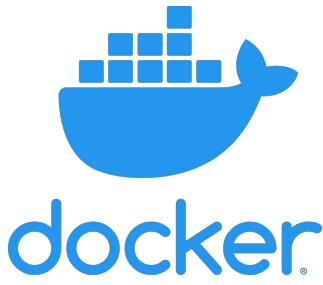
- Asynchronous work patterns
- Minimise dependency between individuals
- Leverage everyone's skills maximally



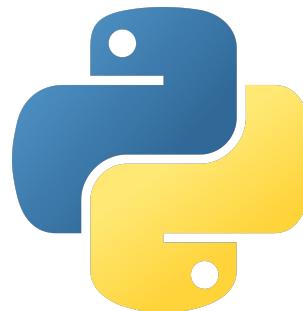
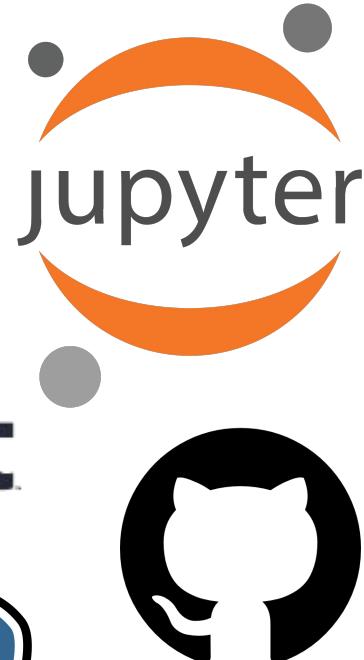
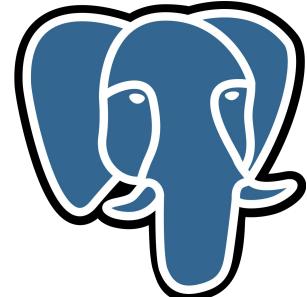
Leveraging the team skills



Open Source free* Tools



docker dbt



*I know Tableau isn't free but it's what the charity already uses

Control the Environment

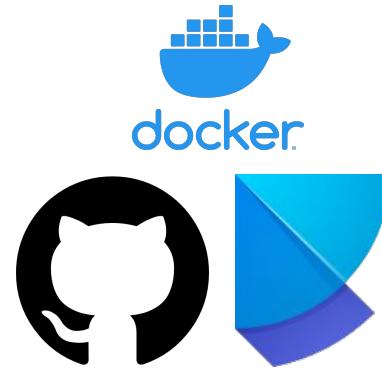
Github

Makes it easy to share the code between us - and with the charity.

Each dev creates their own branch when developing. Allows independent experiment but we require **eventual consistency**.

Poetry

Manages equivalent coding environments for each user.



```
[tool.poetry]
name = "brilliantclub"
version = "0.1.0"
description = "DataCorps tutor forecasting project for the Brilliant Club"
authors = ["TBC DataCorps <brilliantclubproject@datakind.org.uk>"]
license = "CC0 1.0 Universal"

[tool.poetry.dependencies]
python = "3.9"
pandas = "1.4.2"
jupyter = "1.0.0"
jupyterlab = "3.4.7"
simple-salesforce = "1.12.1"
PyYAML = "6.0"
openpyxl = "3.0.10"
plotly = "5.9.0"
absl-py = "1.2.0"
tqdm = "4.64.0"
dbt-postgres = "1.3.0"
matplotlib = "3.5.3"
matplotlib-venv = "0.11.7"
xlsxwriter = "3.0.3"
sqlalchemy = "1.4.42"
psycopg2-binary = "2.9.5"
icecream = "2.1.3"
setuptools = "65.6.0"
scikit-learn = "1.2.1"
joblib = "1.2.0"
prefect = "2.8.0"
prefect-shell = "0.1.5"
prefect-dbt = "0.3.1"
prefect-jupyter = "0.2.0"
```

Control the Environment

Dockerfile > ...

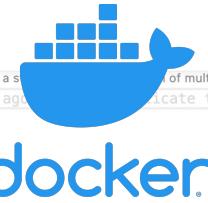
```
You, 3 months ago | 2 authors (Chris Watkins and others)
1 # syntax=docker/dockerfile:1
2 FROM python:3.9-slim-buster
3
4 # Install build-essentials to facilitate installation of Python dependencies
5 RUN apt-get update --fix-missing && apt-get install -y build-essential
6 RUN apt-get install -y git postgresql-client
7
8 # Skip writing .pyc files
9 ENV PYTHONDONTWRITEBYTECODE=1 \
    PYTHONUNBUFFERED=1 \
    # Pin a specific Poetry version for maximum compatibility
12 POETRY_VERSION=1.2.2
13
14 WORKDIR /code/
15
16 RUN pip install "poetry==$POETRY_VERSION"
17 RUN poetry config virtualenvs.create false
18
19 COPY poetry.lock pyproject.toml /code/
20 RUN poetry install --no-interaction
21
22 # Expose default jupyter notebook port.
23 EXPOSE 8888
24
25 # Copy Jupyter launch script
26 COPY launch.py /code/
27
28 ENTRYPOINT [ "/bin/bash" ]
```



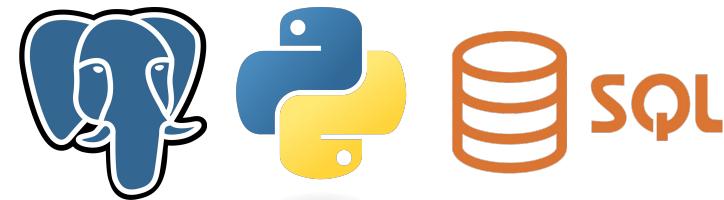
docker-compose.yaml > [version](#)

docker-compose.yml - The Compose specification establishes a standard way to declare multi-container platform-agnostic app configurations. It's designed to be portable between different environments, and to indicate tty in docker compose ...

```
1 version: "3.9"          LeoTheGriff, 7 months ago
2 services:
3   jupyterlab:
4     build: .
5     tty: true
6     env_file:
7       - .env
8     environment:
9       POSTGRES_HOST: postgres
10      REFRESH_DATA:
11      ports:
12        - "8888:8888"
13        - "8889:8889"
14        - "8000:8080" # For dbt docs serve
15      working_dir: /app
16      volumes:
17        - ./src:/code/src
18        - ./notebooks:/code/notebooks
19        - ./data:/data
20        - ./exported_data:/exported_data
21        - ./app
22      depends_on:
23        - postgres
24      entrypoint: "python launch.py"
25      restart: always
26
27
28 postgres:
29   image: postgres:15-alpine
30   environment:
31     POSTGRES_USER: REDACTED
32     POSTGRES_PASSWORD: REDACTED
33     POSTGRES_DB: REDACTED
34   volumes:
35     - ./data:/data
36     - ./exported_data:/exported_data
37     - ./init.sql:/init.sql
38     - ./export.sql:/export.sql
39   ports:
40     # NOTE: "5432:5432" makes the Postgres server accessible to both the host
41     # developer machine *and* the "app" container in Docker. If you don't want
42     # it available on the host machine, change this to simply "5432".
43     - 5432:5432
44   restart: always
```



Automate the data pull



Connect to Salesforce

Pull the latest TBC data

- Registered tutors
- Tutor availability reports
- Confirmed school placements
- Candidate school placements

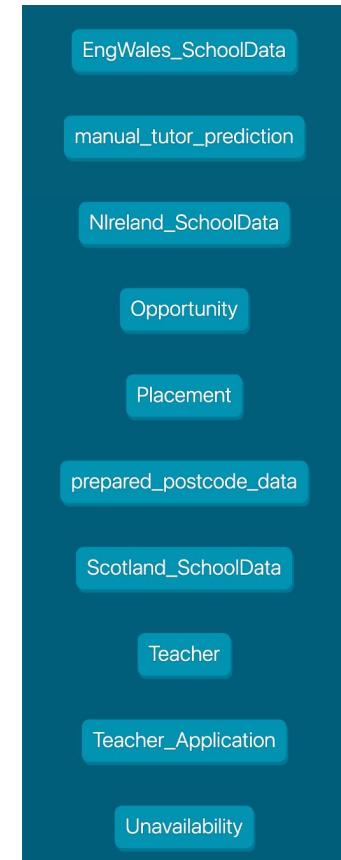
Data Enrichment

- Local Authority geo boundaries
- Postcode geo data

PII Data Security

SHAhash256 any personal data

SQL scripts to seed
the Postgres DB



EDA, model training & prediction



Jupyter Lab EDA

- Explored initial dataset which enabled us to identify missing or “dirty” data with TBC
- Started to experiment with potential features sets
- Identified the typical return rate(s) when tutors were asked to submit their upcoming availability

Model Building

- Scikit-learn Random Forest Classifier
- From the trained model extract the individual tutor availability probabilities to facilitate better aggregate population statistics
- Convolve the availability probability with the population probability that an availability response will be submitted

Primary Feature Generation

Predicting an individual tutors availability

- Time to PhD Completion
- Previous number of placements
- Mean number of placements per term
- Number of terms since they registered
- Whether they did a placement last term.

Final output is a per-tutor availability probability for 1,2 & 3 terms in the future

Reusability of the modelling



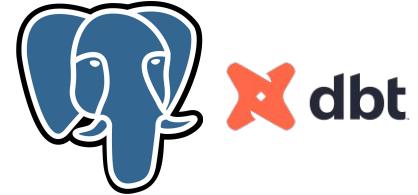
Turning notebooks into scripts

- Refactored the code into local python files that can be imported into notebooks or scripts
- Make use of dedicated dataclasses for tutors, school terms etc
- De-couple the code to make it easier to substitute into different models or add new features

```
Running - FutureTerm_1
Predict: 472 tutors in SchoolTerm(year=2023, name='Summer')
Running - FutureTerm_2
Predict: 436 tutors in SchoolTerm(year=2023, name='Autumn')
Running - FutureTerm_3
Predict: 489 tutors in SchoolTerm(year=2024, name='Spring')
```

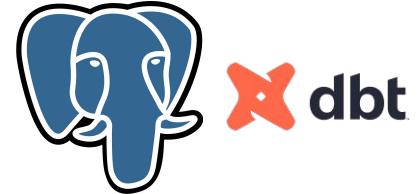
Jupyter Lab Training & Prediction

- Two notebooks
 - a. Train a new model and report accuracy metrics
 - b. Load an appropriate model and make new predictions
- Keep the key details visible to users



Transform and load the data





Transform and load the data





Orchestrate the pipeline

```
1 import os
2 from pathlib import Path
3
4 from prefect import flow, task
5 from prefect_shell import ShellOperation
6 from prefect_dbt.cli.commands import DbtCoreOperation
7 from prefect_jupyter import notebook
8
9 from icecream import ic
10
11 ic.configureOutput(prefix='Pipeline | ')
12
13 # Setup env variables
14 ## Postgres
15 PGUSER = os.environ.get("PGUSER")
16 PGHOST = os.environ.get("PGHOST")
17 PGPORT = os.environ.get("PGPORT")
18 PGDB = os.environ.get("PGDB")
19
20 # Jupyter Modelling
21 NOTEBOOKS_PATH = Path("/app/notebooks/pipeline")
22 MODEL_NOTEBOOK_NAME = "02.Tutor-Modelling"
23 PREDICT_NOTEBOOK_NAME = "03.Tutor-predictions"
24 MODEL_YEAR_CUTOFF = os.environ.get("MODEL_YEAR_CUTOFF", 2019)
25 MODEL_PREDICT_START_YEAR = os.environ.get("MODEL_PREDICT_START_YEAR")
26 MODEL_PREDICT_START_TERM = os.environ.get("MODEL_PREDICT_START_TERM")
27
28 @task(name="Validate Model Variables",
29       description="Task to validate that the model variables are all set with reasonable values",
30       task_run_name="validate-model-variables",
31       retries=0)
32 > def env_validation_task(): ...
33
34 @task(name="Check location files",
35       description="Task that checks the pre-downloaded location files are where they are expected",
36       task_run_name="check-location-files-present",
37       retries=0)
38 > def check_presence_location_files(): ...
```



PREFEC

```

0 > @flow(name="TBC-DKUK_End-to-End modelling pipeline", -
1 def end_to_end_pipeline() -> str:
2     """ Runs the complete TBC data modelling pipeline"""
3
4     # Validate the environment
5     _ = env_validation_task()
6
7     # Check for presence of location files
8     _ = check_presence_location_files()
9
10    # Download the data
11    with ShellOperation(commands = ["python -m src.utils.salesforce_to_csv --reports=false --write_to /data",
12                                working_dir="."]) as download_data:
13        download_data_process = download_data.trigger()
14        download_data_process.wait_for_completion()
15        output_lines = download_data_process.fetch_result()
16
17    msg = "Completed Data Download!"
18    ic(msg)
19
20    # Run the Jupyter training notebook
21    modelling_nb = notebook.execute_notebook(
22        body = notebook.export_notebook(modelling_nb)
23        output_path = NOTEBOOKS_PATH / f"{MODEL_NOTEBOOK_NAME}_latestResult.ipynb"
24    )
25    with open(output_path, "w") as outfile:
26        msg = "Completed Rebuild of Tutor Model"
27    ic(msg)
28
29    # Run the Jupyter prediction notebook
30    prediction_nb = notebook.execute_notebook(
31        body = notebook.export_notebook(prediction_nb)
32        output_path = NOTEBOOKS_PATH / f"{PREDICT_NOTEBOOK_NAME}_latestResult.ipynb"
33    )
34    with open(output_path, "w") as f:
35        msg = f"Completed prediction of future Tutor availability starting from {MODEL_PREDICT_START_YEAR}-{MODEL_PREDICT_START_TERM}"
36    ic(msg)
37
38    You, now + Uncommitted changes
39    # Initialise the database
40    with ShellOperation(commands = ["psql -U {PGUSER} -h {PGHOST} -p {PGPORT} -d {PGDB} -f 'init.sql'",],
41                        init_database_process = init_database.trigger()
42                        init_database_process.wait_for_completion()
43                        output_lines = init_database_process.fetch_result())
44
45    msg = "Completed PostgreSQL database initialisation ..."
46    ic(msg)
47
48    # DBT Commands
49    dbt_result = DbtCoreOperation(
50        msg = "Completed DBT data pipeline build ..."
51    )
52
53    # Export data files
54    with ShellOperation(commands = [f"psql -U {PGUSER} -h {PGHOST} -p {PGPORT} -d {PGDB} -f 'export.sql'",],
55                        export_database_process = export_database.trigger()
56                        export_database_process.wait_for_completion()
57                        output_lines = export_database_process.fetch_result())
58
59    msg = "Completed PostgreSQL prediction data export"
60    ic(msg)
61
62    return "Pipeline Run completed!"
63
64 if __name__ == "__main__":
65     # Run the complete pipeline
66     end_to_end_pipeline()

```



Present the data visually

Dashboarding

Some parts of The Brilliant Club business are already using Tableau so they chose that component. Alternative choices could have been:

- Streamlit
- Plotly
- PowerBI

The aim of the dashboard is to easily surface insight into the predicted and actual number of tutors/placements across the country and within different geographic boundaries

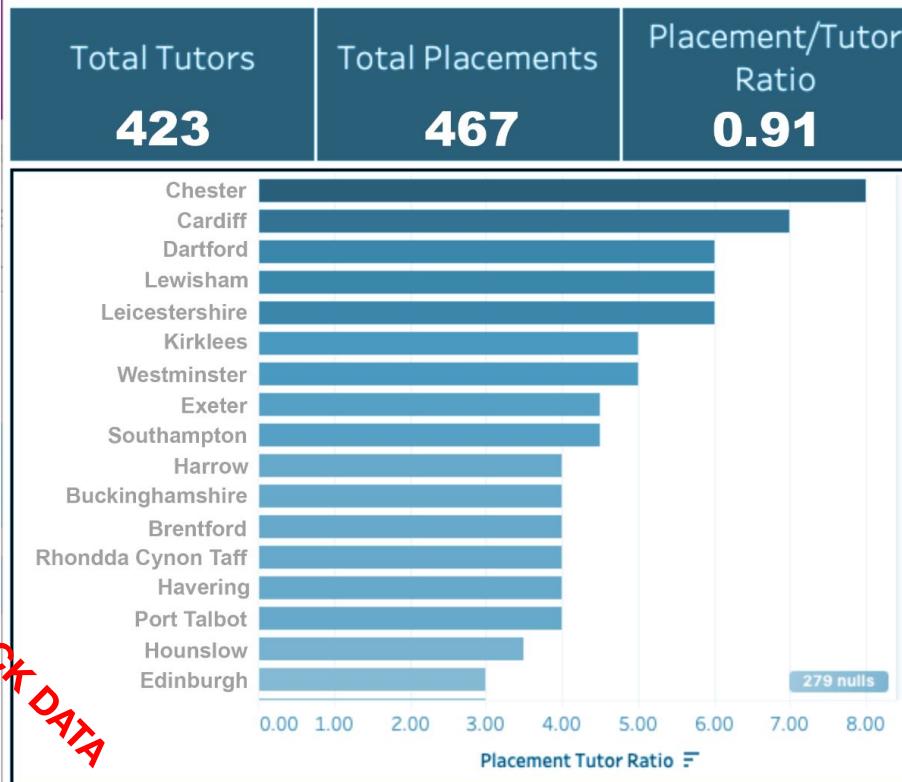
Overall Summary

Placement : Tutor Ratio

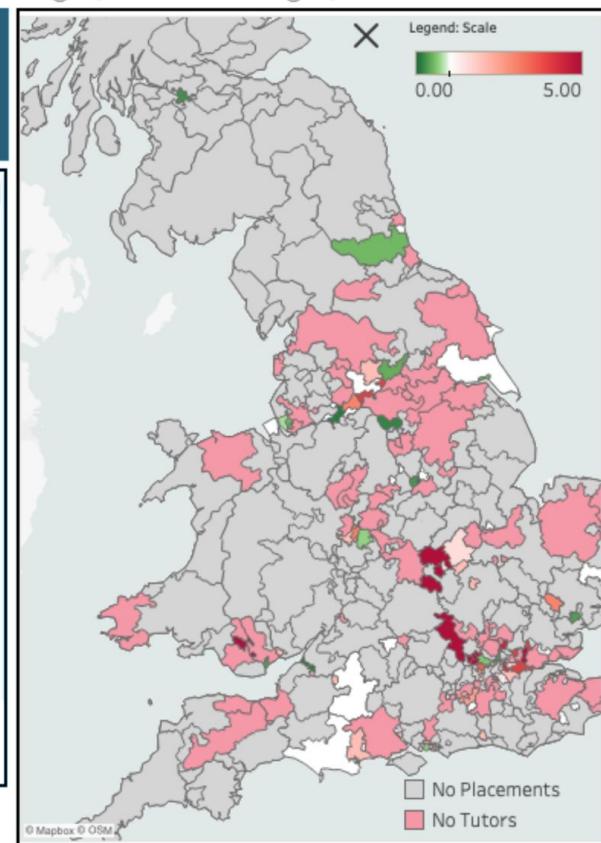


MOCK DATA

Term/Year	Select Metric	Select Geographic View	Region Name	Use placement prediction? Use tutor prediction?
Spring - 2021-22	Placement : Tutor Ratio	Local Authority	(All)	<input type="radio"/> Yes <input type="radio"/> Yes <input checked="" type="radio"/> No <input type="radio"/> No



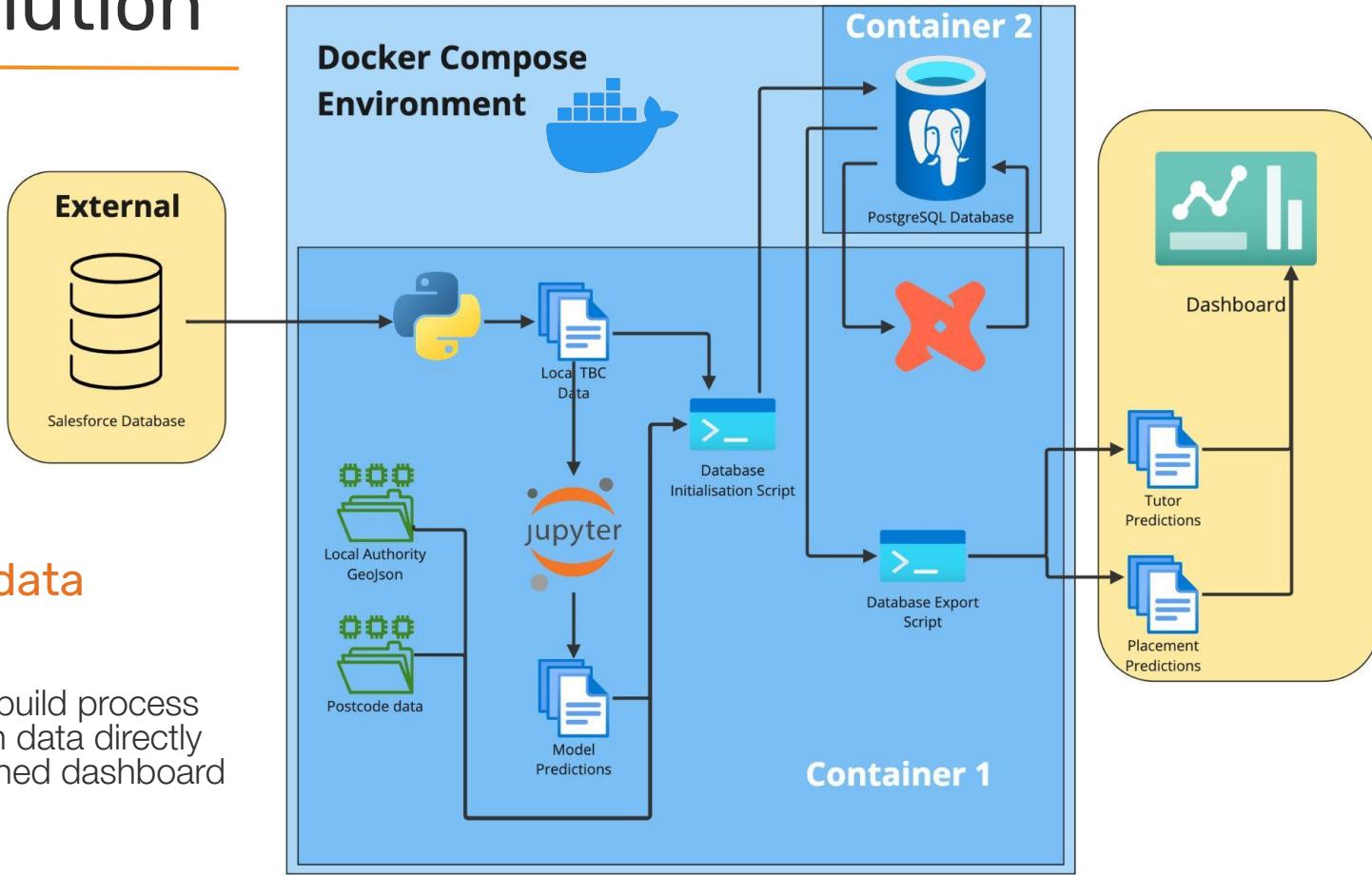
Note: "Use placement prediction" and "Use tutor prediction" will use placement/tutor predictions respectively, if selected. Placement predictions are based on confirmed placements, plus a percentage of likely Opportunities. Tutor predictions are based on past tutor behaviour.



The Solution

Easy-to-use data pipeline

- One click data build process
- Generates fresh data directly into an established dashboard



What impact will this have?

The goal

Use data on tutor availability and on sales placements to forecast and visualise the mismatch between tutor supply and demand, regionally to enable The Brilliant Club to decide which areas to prioritise for recruitment, so that sufficient tutors are available to meet the needs of partner schools.

Desired outcomes

1. School sales are easier to manage with less back and forth conversation with the programme team about how placing is going/whether we have tutors available in a given area
2. TBC more accurately meet recruitment targets regionally, securing tutors in the areas with highest school demand.
3. TBC need to reduce fewer school sales due to being unable to predict tutor need.
4. TBC is able to meet placement demands without undue burden on tutors who already have a placement or The Brilliant Club staff

An impact assessment will happen at a future date

Questions?



contact@datakind.org.uk

www.datakind.org.uk

@datakinduk

www.linkedin.com/in/adambenhill/

Horsewithapointyhat.com

@astroadamh