**Train In Data**

# Streamlining Model Deployment with Open-Source

# Reproducible ML Pipelines

# Machine Learning Pipeline

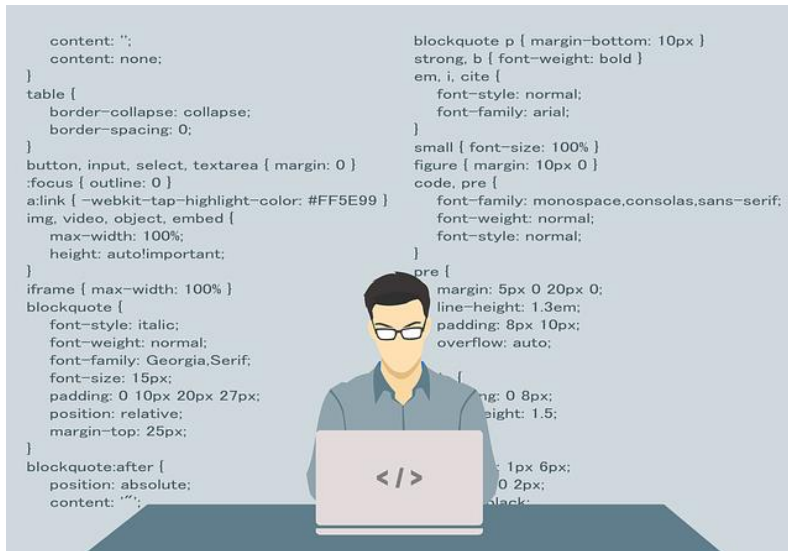| Data | Feature Engineering | Model Training / Scoring | Predictions |
|------|---------------------|--------------------------|-------------|

**Feature Engineering**
- Missing data imputation
- Categorical encoding
- Discretisation
- Transformation
- Feature extraction
  - Text
  - Datetime
  - Images
  - Time series
- Feature combination

- **Lots to code!**

**Model Training / Scoring**
- Linear Models
- Decision trees
- KNNs
- SVMs
- Neural Networks

- Stacking

- **Lots of parameters ➔ Reproducibility**

Train In Data

# Challenges

- A lot to code

- Repetitive

- Learn and store parameters

- Reproducibility

# Lots to code



- Time consuming

- Different versions across team

# Repetitive

- Multiple copies of same code

- Different versions of same code

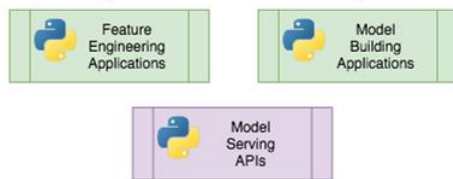- Difficult to keep track

# Learn and store parameters

- Multiple intermediate files with parameters

- Config or params file
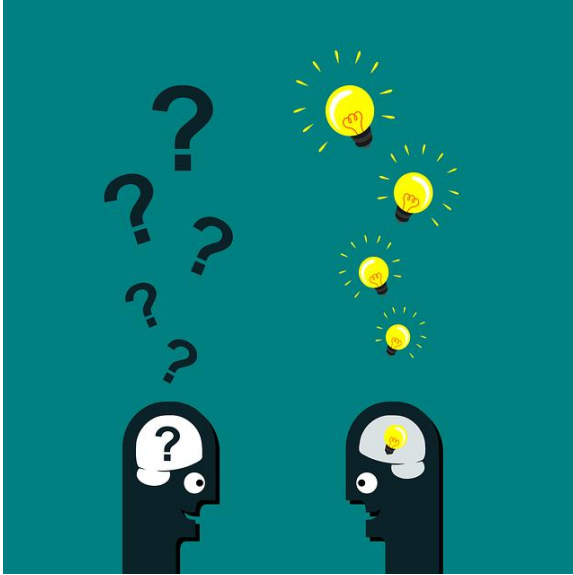
# Reproducibility in Deployment



Research Environment

Production Environment

- Re-write code

- Include tests

- Reproducibility

# Team Performance



- Decreased Performance

- Frustration

- Lack of reproducibility

- Increased deployment times

# Open-source



- Increase Performance

- Prevent Frustration

- Maximise reproducibility

- Minimise deployment times

# Open-source


open source

- No more coding

- Version tracking for reproducibility

- Classes and functions include tests – no need to recode for production

# Reproducible ML Pipelines

# Reproducible ML Pipelines