# DataScienceSeed Online #10

Data Science, Machine Learning, Artificial Intelligence Meetup a Genova

# MLOps, quando si smette di giocare

Simone Merello - Head of Deep AI, Perceptolab

**Simone Merello -** Head of Deep AI @ Perceptolab

Academic background:
- Computer scientist with focus on data science
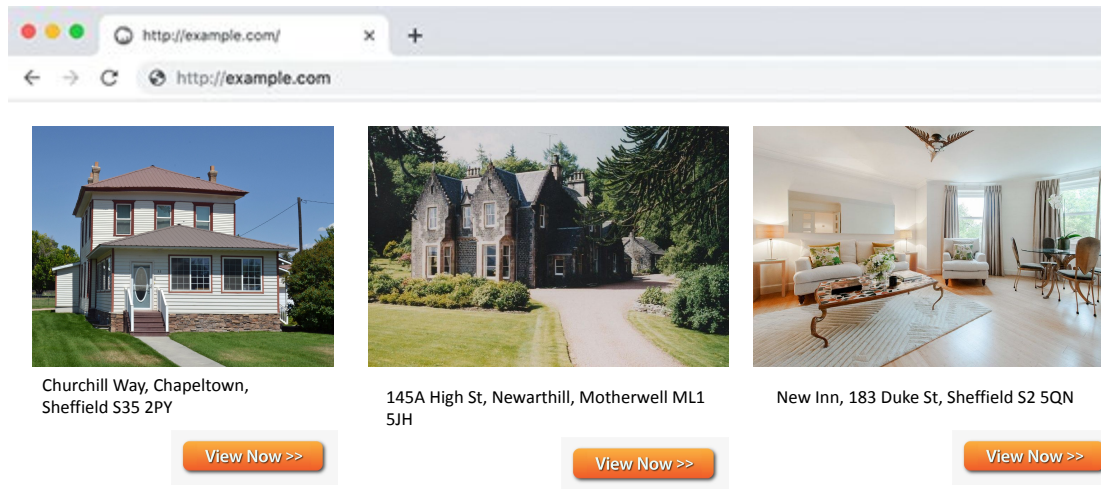- Researcher @ NTU university of Singapore

Now:
- Still ML Researcher with focus on Computer Vision: TensorFlow, Pytorch
- MLOps engineering to help team collaboration and automation of ML pipelines.

Leisure:
- Love traveling and love water sports!

# A practical example



Churchill Way, Chapeltown, Sheffield S35 2PY

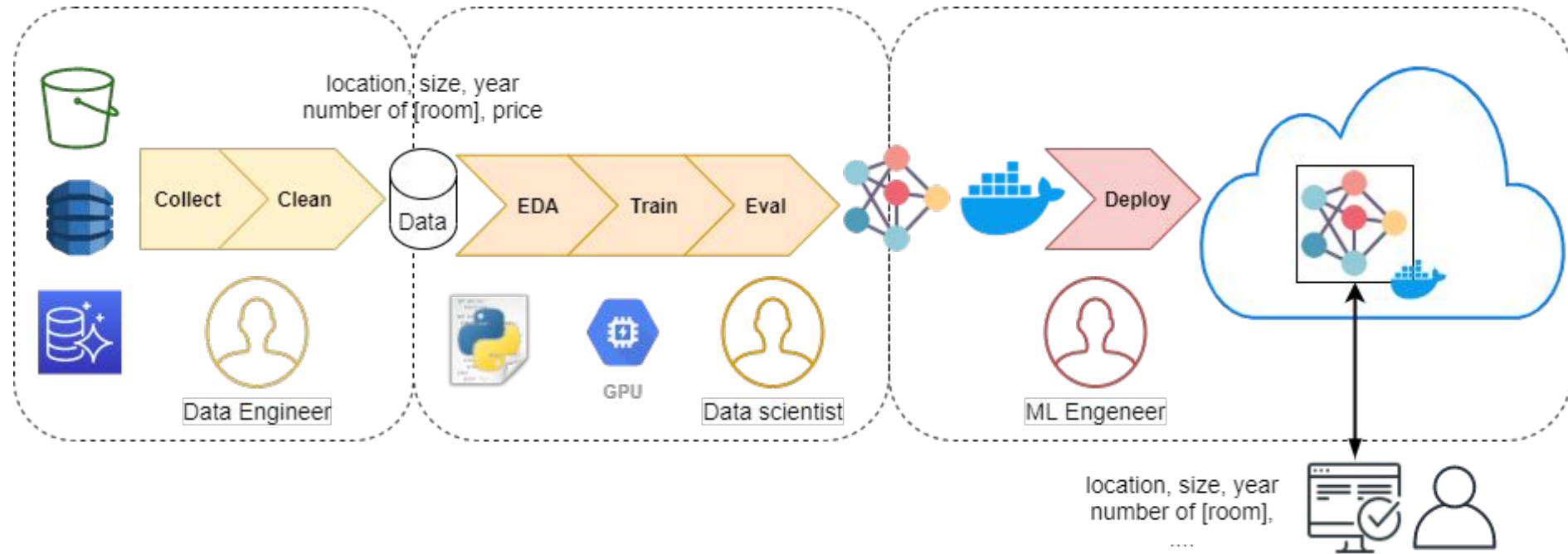**View Now >>**

145A High St, Newarthill, Motherwell ML1 5JH

**View Now >>**

New Inn, 183 Duke St, Sheffield S2 5QN

**View Now >>**

→ New feature: Real Estate Appraisal

# ML Development cycle



location, size, year
number of [room], price

Collect — Clean

Data

EDA — Train — Eval

Deploy

GPU

Data Engineer

Data scientist

ML Engeneer

location, size, year
number of [room],
....

# But then something happens… [examples]

- New feature doesn't work as expected  → **MONITORING**

- More data comes in → **FAST ITERATION**

- New team members → **NEED FOR COOPERATION**

- Multiple projects / data / models / experiments → **TRACKING**

- Performance degradation with time → **AUTOMATIC RE-TRAINING**

# ML Systems require organization!

## [1] Some Issues of ML Systems

1. **Entanglement**: "Changing Anything Changes Everything"

2. **Tracking dependencies:** data, code, env, input models

3. **Cascading**: the output of a model A might affect input of an [undeclared] model B

4. **Feedback Loops**: models influencing each other if they update over time

5. **Staleness**: if the input changes during time, the model has to adapt

## [2] ML Systems Best practices

1. **Data management:** Ensuring availability, accessibility, quality and versioning of data.

2. **Pipelines:** supporting data preprocessing, train, test and deployment

3. **Automation** of training and deployment pipelines allows fewer deployment issues

## [3] ML Systems readiness

1. **Features and data:** assert expectations, cost/benefit tradeoff, fast addition of new features, tested features creation

2. **Model development:** versioning, evaluate {metrics = KPI, staleness, fairness}.

3. **Infrastructure:** reproducibility, integration tests, canary testing, quick rollback

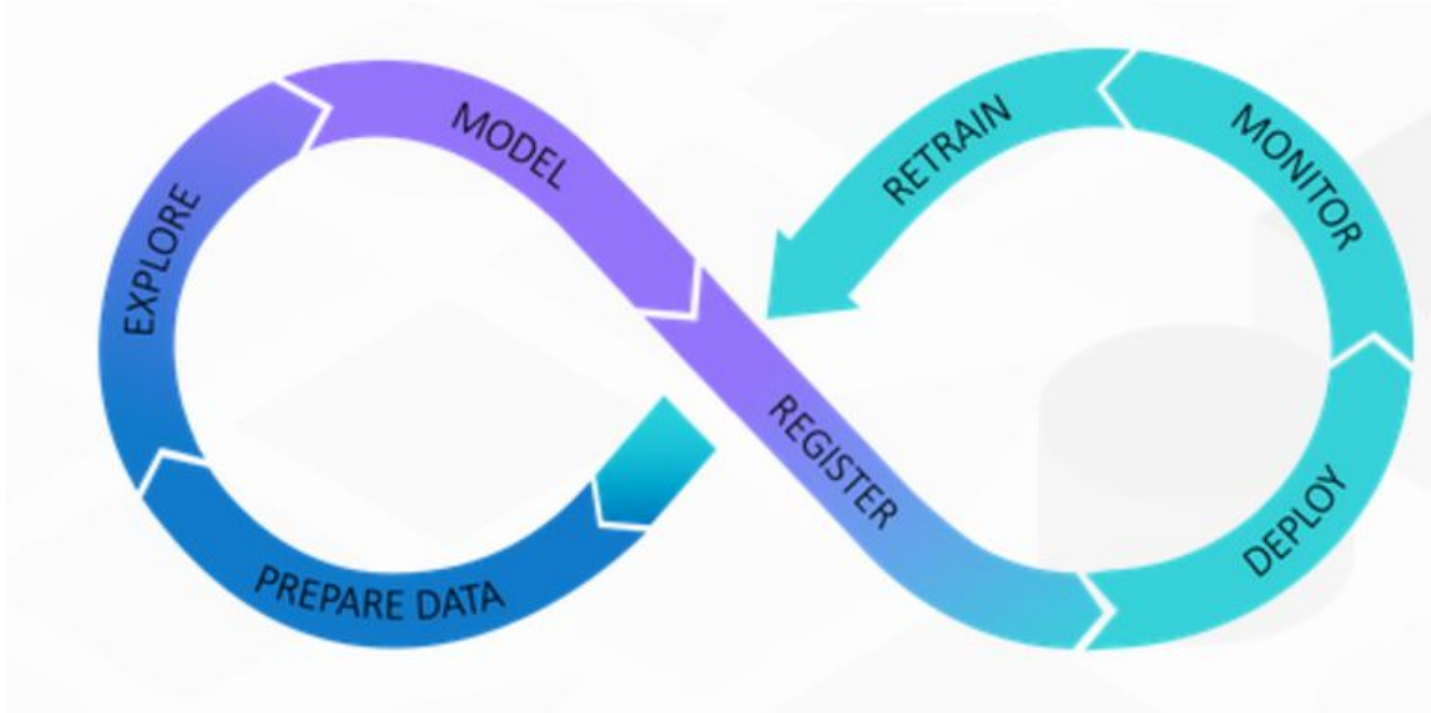4. **Monitoring:** monitor {changes in dependencies, input expectations, staleness}

[1] Sculley, David, et al. "Hidden technical debt in machine learning systems" .2015
[2] Amershi, Saleema, et al. "Software engineering for machine learning: A case study." *2019*
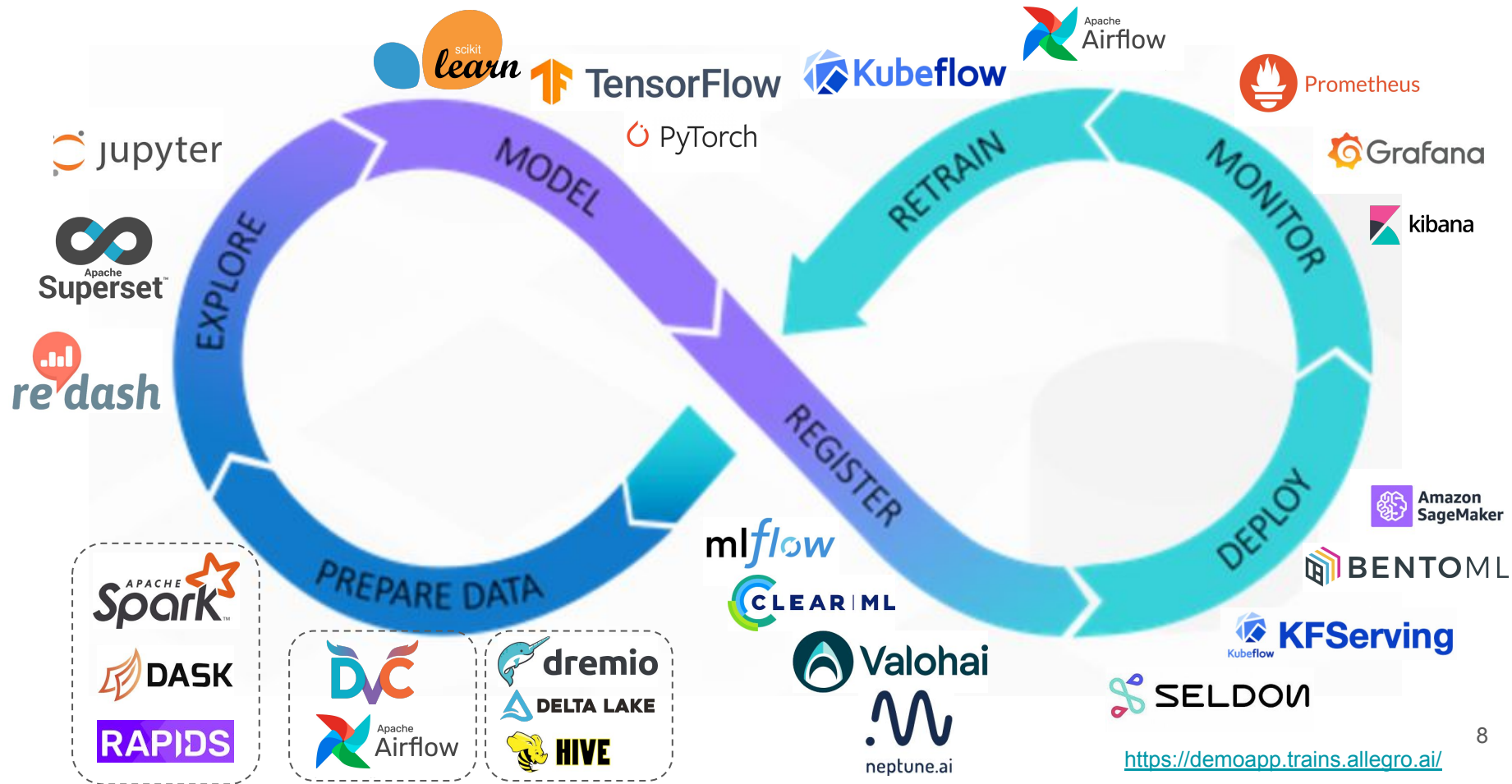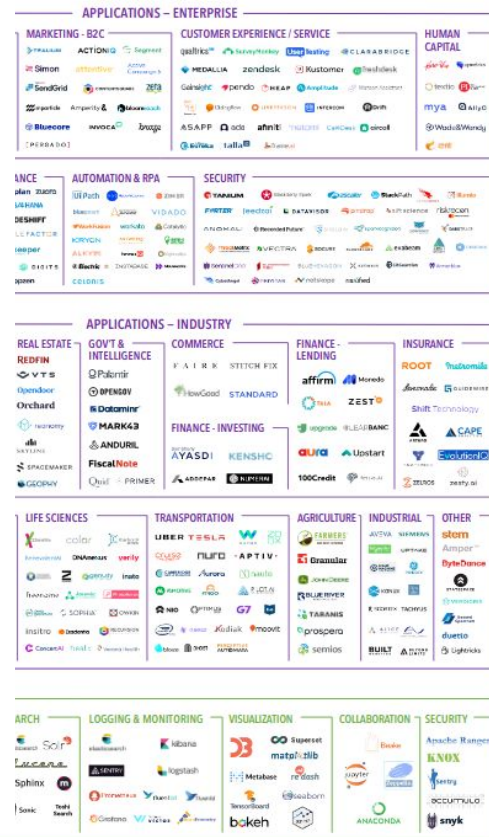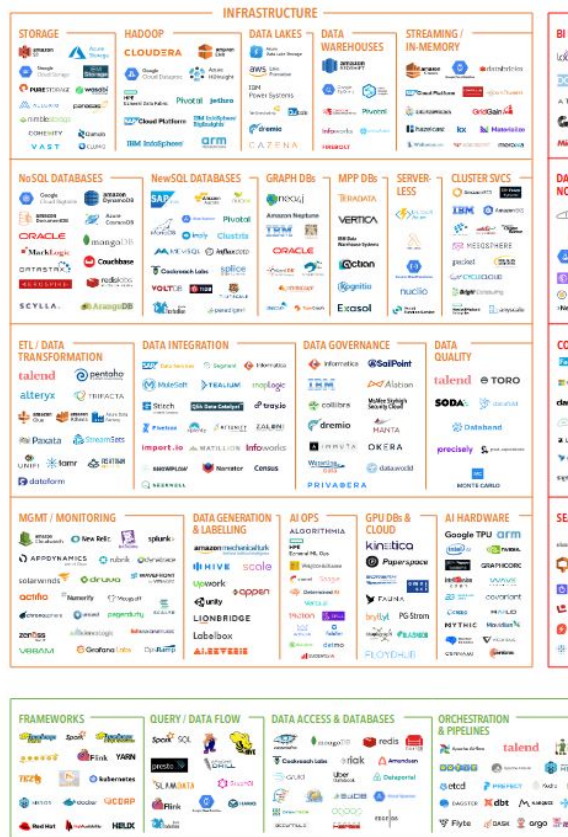[3] Breck, Eric, et al. "The ml test score: A rubric for ml production readiness and technical debt reduction." *2017*

# MLOps

*Management of ML Systems: operationalize the steps than produce, serve and improve ML models*

# Some tools [divided by their native usage]

https://demoapp.trains.allegro.ai/

# Tools landscape

https://mattturck.com/data2020/

# ML tools

https://about.mlreef.com/images/blog/ml_landscape.png          https://docs.google.com/spreadsheets/

# [SOME] Features you might look for

| ML PERSPECTIVE FOR BATCH LEARNING TASKS | | |
|---|---|---|
| DATA | MODELLING | DEPLOY [1] |
| Versioning | Track experiments | Automated CI/CD |
| Availability | Reproducibility (dependencies) | Track Deployment |
| Training - Serving consistency | Track outputs (models / performances) | Canary / Shadow testing |
| Schedule jobs | Compare experiments | Automatic Retraining |
| Exploration | Hyperparams Optimization | Monitoring data (outliers /  dist. shift) |
| Data quality checks | Infrastructure handling | Monitoring model performances |
| Cataloging | Peer reviewing | Explaining predictions |
| Labelling | | Automatic Scalability |
| Handle Real time data | | |
| Infrastructure | | |
| Storage scalability | | |

[1] Klaise, Janis, et al. "Monitoring and explainability of models in production." arXiv preprint arXiv:2007.06299 (2020).

# Revised Steps



**Automatic Retrain**

Apache Airflow

Kubeflow

**Train**:
- Reproducibility
- Versioning
- Experiment Tracking
- Model Registry

Data scientist

MODEL

RETRAIN

MONITOR

EXPLORE

REGISTER

**CI/CD**:
- Evaluation
- Deployment
- Track deployment

ML Engeneer

DEPLOY

PREPARE DATA

Data

**Data pipelines**:
- Solve data deps.
- Scheduling
- Backfilling

Data Engineer

Apache Airflow

Pipeline tools

Kubeflow

Apache Airflow

DVC

KubeFlow: **Pipelines**, **Notebook Servers**, **Katib** (hyperparameter tuning)**, Artifact Store**, **KFServing**

**Airflow: DAG** (chained Operators)**,
Scheduler, Executor**

## DAGs

| | DAG | Owner | Runs | Schedule | Last Run | Recent Tasks | Actions | Links |
|---|---|---|---|---|---|---|---|---|
| | **All 26**  Active 10  Paused 16 | | | Filter DAGs by tag | | | Search DAGs | |
| ● | example_bash_operator  example example2 | airflow | 2 | 0 0 * * * | 2020-10-26, 21:08:11 | 6 | ▶ C 🗑 | ⋯ |
| ● | example_branch_dop_operator_v3  example | airflow | | */1 * * * * | | | ▶ C 🗑 | ⋯ |
| ○ | example_branch_operator  example example2 | airflow | 1 | @daily | 2020-10-23, 14:09:17 | 11 | ▶ C 🗑 | ⋯ |
| ● | example_complex  example example2 example3 | airflow | 1 1 | None | 2020-10-26, 21:08:04 | 37 37 | ▶ C 🗑 | ⋯ |
| ● | example_external_task_marker_child | airflow | 1 | None | 2020-10-26, 21:07:33 | 2 | ▶ C 🗑 | ⋯ |
| ● | example_external_task_marker_parent | airflow | 1 | None | 2020-10-26, 21:08:34 | 1 | ▶ C 🗑 | ⋯ |
| ● | example_kubernetes_executor  example example2 | airflow | | None | | | ▶ C 🗑 | ⋯ |
| ● | example_kubernetes_executor_config  example3 | airflow | 1 | None | 2020-10-26, 21:07:40 | 5 | ▶ C 🗑 | ⋯ |
| ● | example_nested_branch_dag  example | airflow | 1 | @daily | 2020-10-26, 21:07:37 | 9 | ▶ C 🗑 | ⋯ |
| ○ | example_passing_params_via_test_command  example | airflow | | */1 * * * * | | | ▶ C 🗑 | ⋯ |

14

# DVC

Why?
1. Easy to setup and use
   `$ pip install dvc`
2. Can be used for many MLOps steps
3. Is it the best one? NO (It depends on your needs)

# Experiment & Data versioning



train.py
evaluate.py
data.zip
model.h5

```
$ git commit -am ".."
$ git push
```

```
$ dvc add model.h5 data.zip
```

```
$ dvc push
```

commit: 2d01b8e

train.py
evaluate.py
data.zip.dvc
model.h5.dvc

md5: 2bc...

md5: 21e...

data.zip

model.h5

# Pipelines: Training & CI/CD

```
                                +----------+
                                | prepare  |
                                +----------+
                                     *
                                     *
                                     *
                                +----------+
                                | featurize |
                                +----------+
                               **          **
                             **              *
                            *                  **
                        +-------+                *
                        | train |                **
                        +-------+                  *
                             **          **
                               **      **
                                 *    *
                              +----------+
                              | evaluate |
                              +----------+
```

dvc.yaml

```yaml
stages:
 [...]
 featurize:
    cmd: python features.py data/ features
    deps:
    - data/
    - features.py
    outs:
    - features/
  train:
    cmd: python train.py features model.pkl
    deps:
    - features
    - train.py
    outs:
    - model.pkl
[...]
```

Rerun if something change

$ dvc repro

Git commit: 2d01b8e

```yaml
stages:
 featurize:
    cmd: python featurization.py data/ features/
    deps:
    - path: data/
      md5: 20b78...
    - path: featurization.py
      md5: 28946...
    outs:
    - path: features/
      md5: 52c1f...
  train:
    cmd: python train.py features/ model.pkl
    deps:
    - path: features/
      md5: 52c1f...
    - path: train.py
      md5: 3ffc5...
    outs:
    - path: model.pkl
      md5: b4c48...
```

# Experimentation: Pick up the best model

```
$ dvc exp run --queue -S train.min_split=8
Queued experiment 'd3f6d1e' for future execution.
$ dvc exp run --queue -S train.min_split=64
Queued experiment 'f1810e0' for future execution.
$ dvc exp run --queue -S train.min_split=2 -S train.n_est=100
Queued experiment '7323ea2' for future execution.
$ dvc exp run --queue -S train.min_split=8 -S train.n_est=100
Queued experiment 'c605382' for future execution.
$ dvc exp run --queue -S train.min_split=64 -S train.n_est=100
Queued experiment '0cdee86' for future execution.
$ dvc exp run --run-all --jobs 2
```

```
$ dvc exp show --no-timestamp \
            --include-params train.n_est,train.min_split
```

| Experiment | avg_prec | roc_auc | train.n_est | train.min_split |
|------------|----------|---------|-------------|-----------------|
| workspace  | 0.56191  | 0.93345 | 50          | 2               |
| master     | 0.55259  | 0.91536 | 50          | 2               |
| ├── exp-bfe64 | 0.57833 | 0.95555 | 50        | 8               |
| ├── exp-b8082 | 0.59806 | 0.95287 | 50        | 64              |
| ├── exp-c7250 | 0.58876 | 0.94524 | 100       | 2               |
| ├── exp-b9cd4 | 0.57953 | 0.95732 | 100       | 8               |
| ├── exp-98a96 | 0.60405 | 0.9608  | 100       | 64              |
| └── exp-ad5b1 | 0.56191 | 0.93345 | 50        | 2               |

```
$ dvc exp apply exp-98a96
```

# Deploy: CI/CD pipeline

```
stages:
 test_performances:
    cmd: python test_performances.py model.pkl
    deps:
    - test_performances.py
    - model.pkl
    outs:
    - test_result.md
  deploy:
    cmd: python deploy.py test_result.txt model.pkl
    deps:
    - deploy.py
    - test_result.md
    - model.pkl
```

```
name: train-my-model
on: [push]
jobs:
  run:
    runs-on: [ubuntu-latest]
    container: docker://dvcorg/cml-py3:latest
    steps:
        - uses: actions/checkout@v2
        - name: cml_run
          env:
             repo_token: ${{ secrets.GITHUB_TOKEN }}
          run: |
             dvc pull model.pkl
             dvc repro
             git config [...]
             git add dvc.lock test_results.txt
             git commit "CI/CD pipeline" --allow-empty
             git push -u origin HEAD"
```

```
deploy:
    deps:
    - path: model.pkl
      md5: 20b...
    - test_results.txt
      md5: 21e...
```

Git commit:
2d01b8e

19

# How to begin:

1. **Tools must be useful**: reduce troubles and takes less time from the team, not more

2. **Start manually, then automate**: difficult to choose what to automate without knowing what issues are there

3. **Consider lock-ins:** easier to adopt a new tool than to leave it

4. **Give some extra points to "mature" tools**

**Simone Merello**

Head of Deep AI,
Perceptolab

Really happy to discuss about these topics further!

simone.merello@smartlab.ws