

# Data Version Control (DVC): Tutorial 2: Iris Demo Project

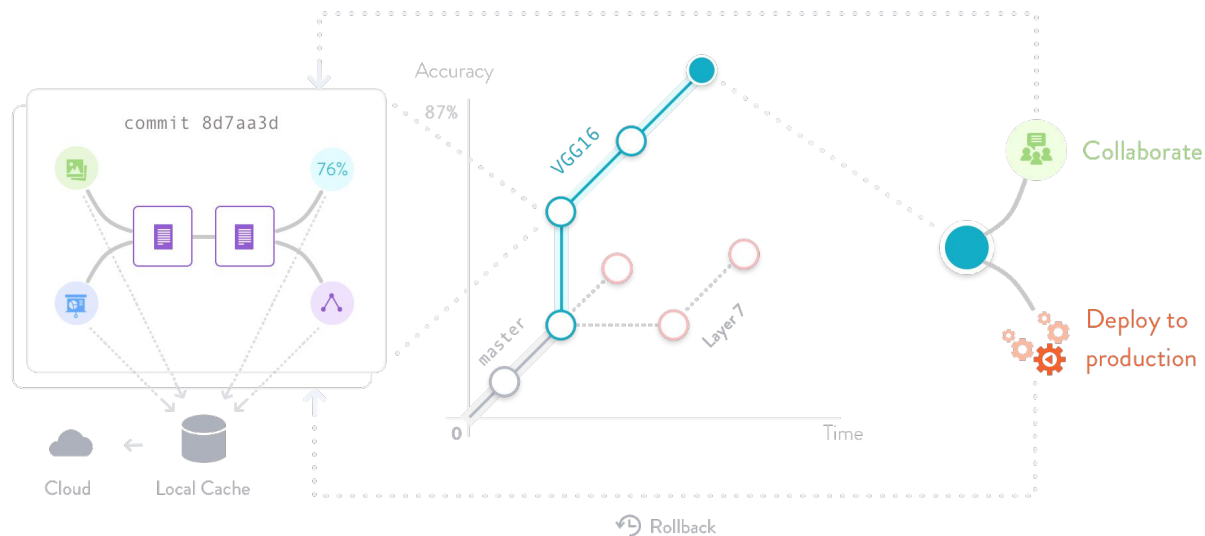


Mikhail Rozhkov

2019

# DVC core concepts

- Experiment
- State (experiment state)
- Reproducibility
- Pipeline
- Workflow
- Data files
- Data cache
- Cloud storage

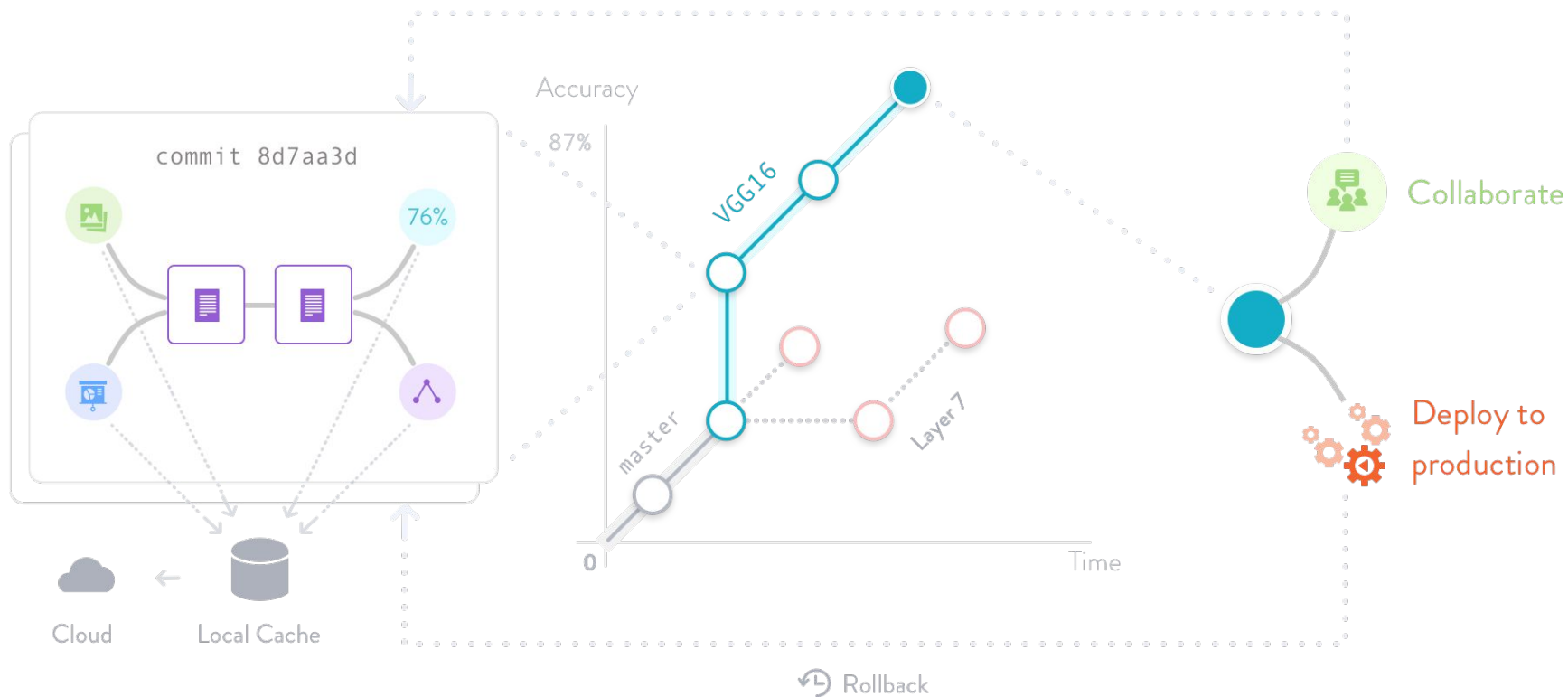


# Experimentation methodology

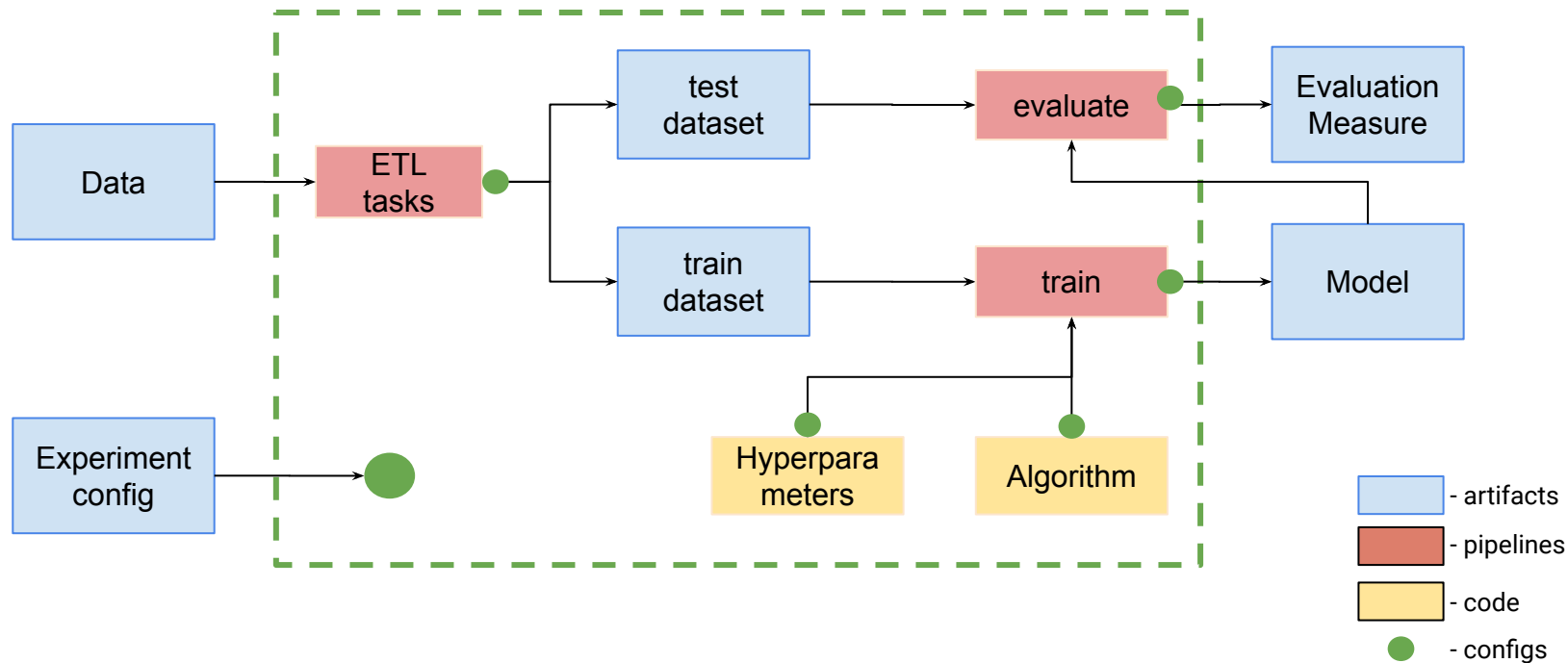
## DVC core concepts

- Experiment
- State (experiment state)
- Reproducibility
- Pipeline
- Workflow
- Data files
- Data cache
- Cloud storage

# DVC tracks ML models and data sets



# Experiment: pipelines, configs and artifacts



## Use Case:

## Iris Flowers Classification

- Task: classify Iris flowers
- Dataset: Iris dataset
- Metrics: F1



### References:

- [https://en.wikipedia.org/wiki/Iris\\_flower\\_data\\_set](https://en.wikipedia.org/wiki/Iris_flower_data_set)
- [https://scikit-learn.org/stable/tutorial/statistical\\_inference/supervised\\_learning.html](https://scikit-learn.org/stable/tutorial/statistical_inference/supervised_learning.html)

Image source:

<https://medium.com/@jebaseelanravi96/machine-learning-iris-classification-33aa18a4a983>

Use Case:

# Iris Flowers Classification

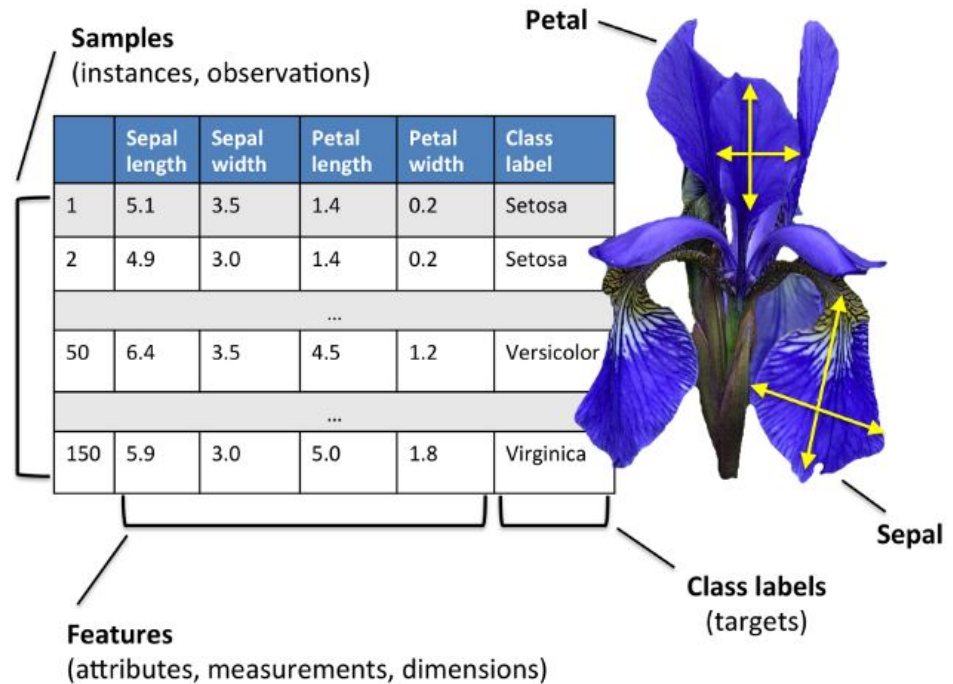


Image source:

<https://medium.com/@jebaseelanravi96/machine-learning-iris-classification-33aa18a4a983>

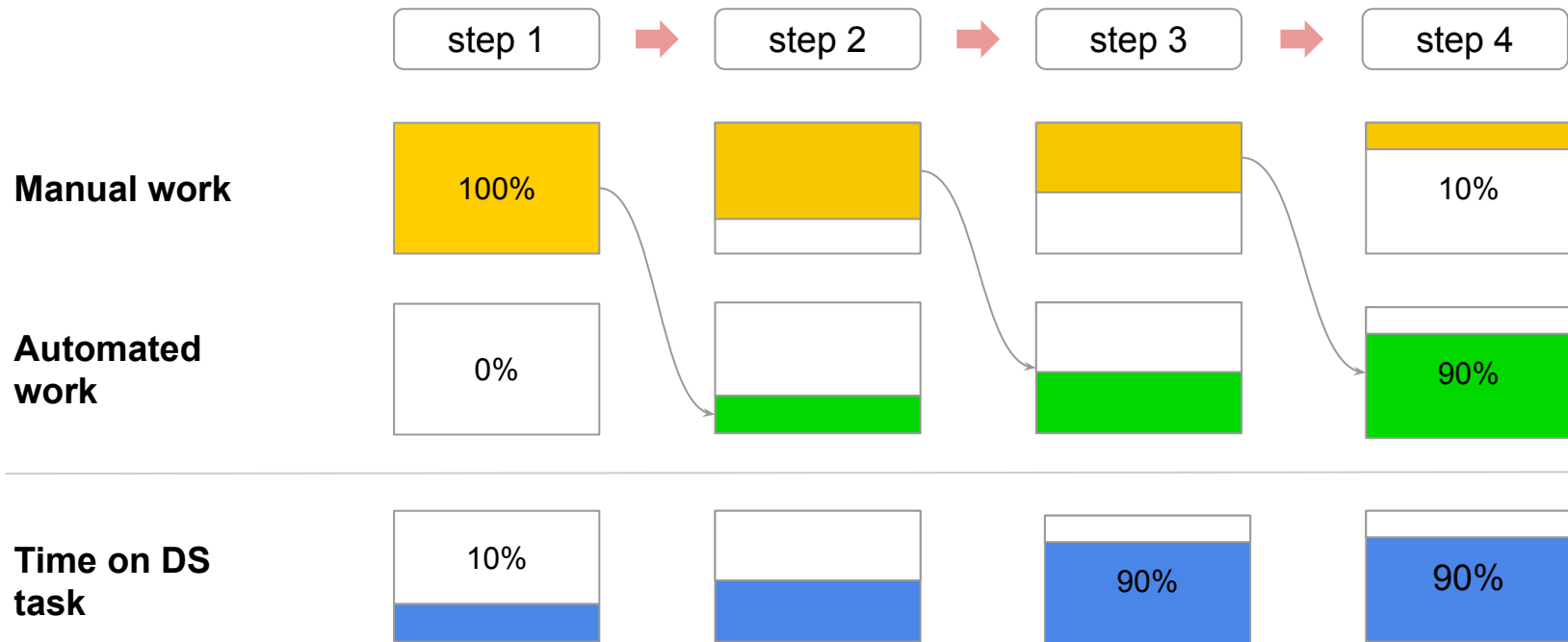
# ML Reproducibility checklist

1. Automated pipelines
2. Control run params
3. Control execution DAG
4. Code version control
5. Artifacts version control (models, datasets, etc.)
6. Use shared/cloud storage for artifacts
7. Environment dependencies control



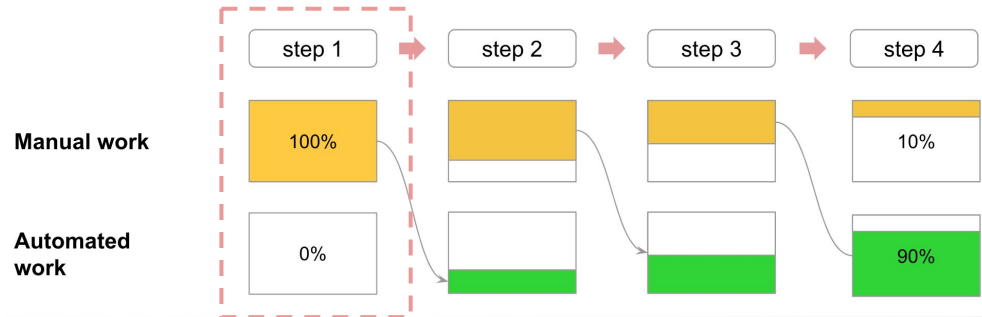


# How to start?



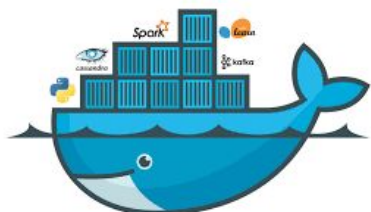
# Step 1: Jupyter Notebook

- code in Jupyter Notebook
- everything in Docker



# ML Reproducibility checklist

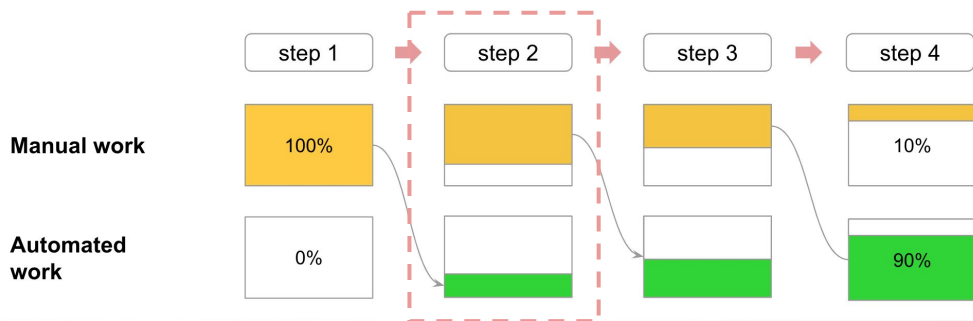
1. Automated pipelines
2. Control run params
3. Control execution DAG
4. Code version control
5. Artifacts version control (models, datasets, etc.)
6. Use shared/cloud storage for artifacts
7. Environment dependencies control
8. Experiments results tracking



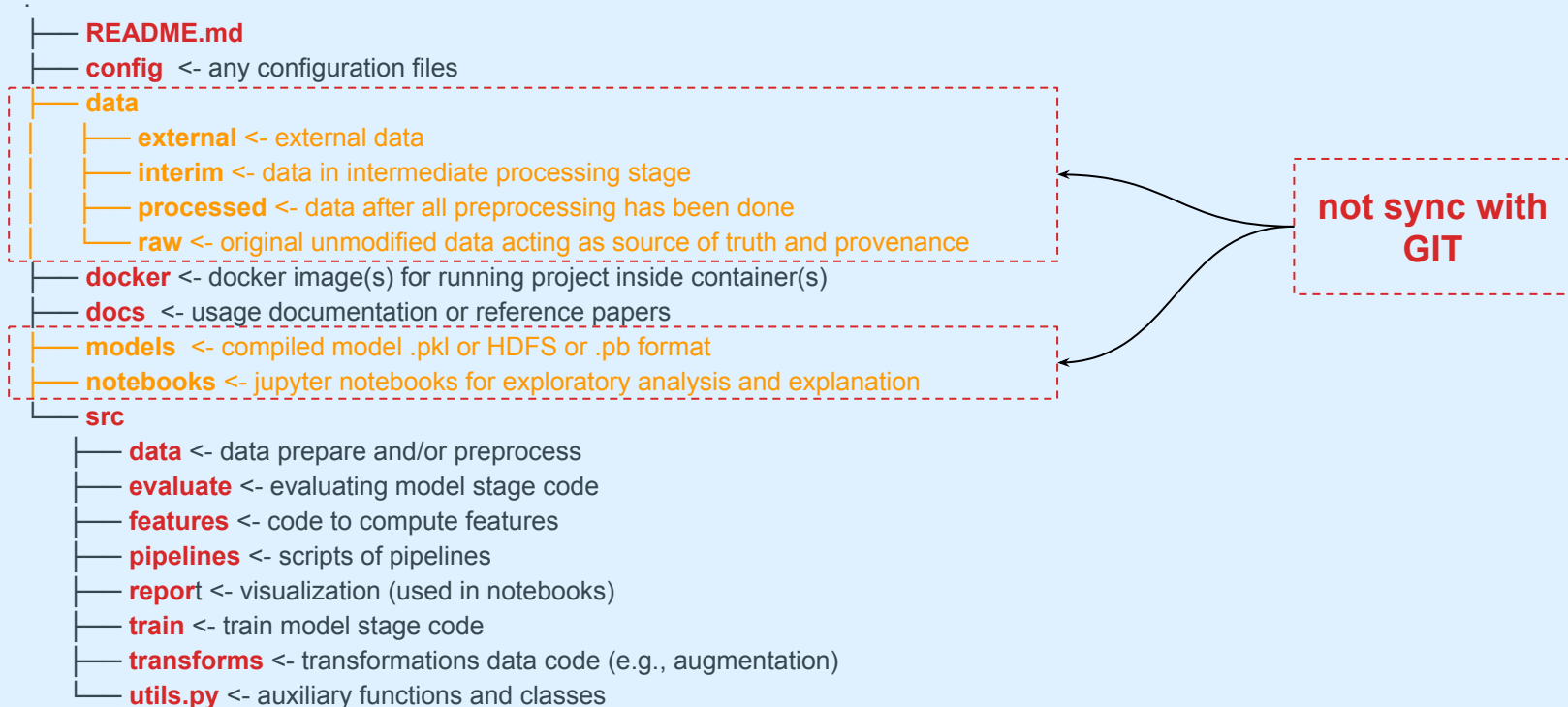
## Step 2:

### Build pipelines

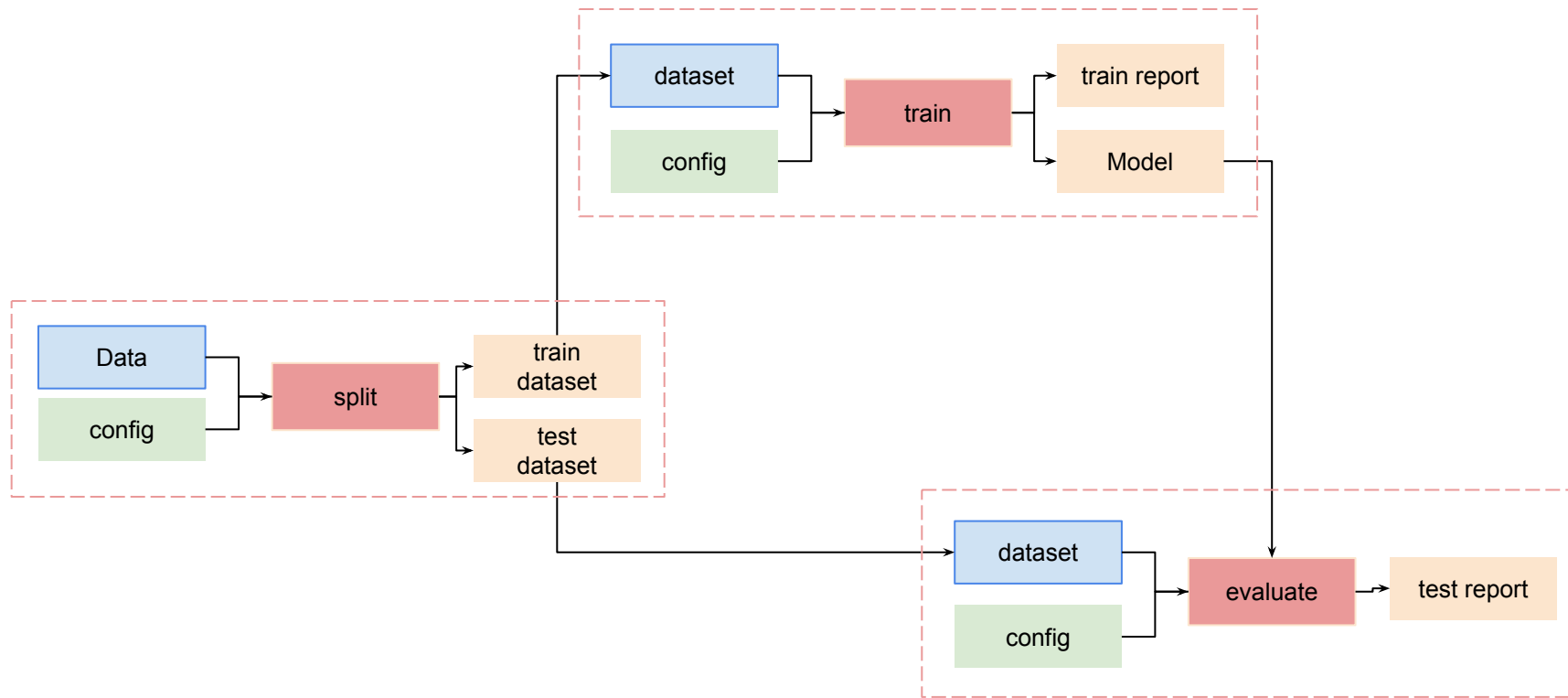
- move common code into .py modules
- build pipelines
- everything in Docker
- run experiments in terminal or Jupyter Notebook



# Project structure

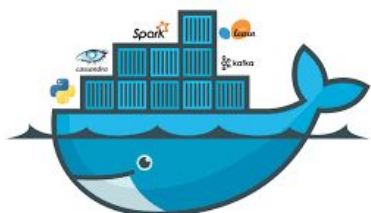


# Setup pipelines



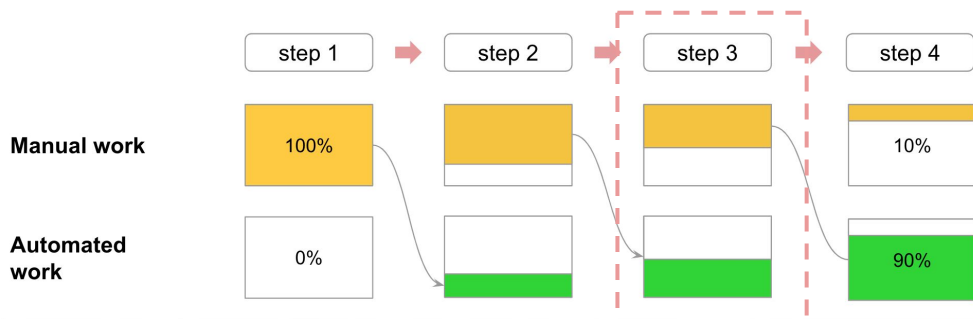
# ML Reproducibility checklist

- ✓ 1. Automated pipelines
- ✓ 2. Control run params
- 3. Control execution DAG
- ✓ 4. Code version control
- 5. Artifacts version control (models, datasets, etc.)
- 6. Use shared/cloud storage for artifacts
- ✓ 7. Environment dependencies control
- 8. Experiments results tracking



## Step 3: Pipelines automation

- add pipelines dependencies under DVC control
- add models/data/configs under DVC control
- models/data/configs under DVC control
- same code in .py modules
- same pipelines
- everything in Docker
- run experiments in terminal or Jupyter Notebook

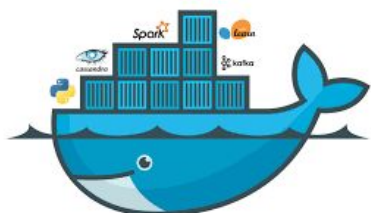




# ML Reproducibility checklist



1. Automated pipelines
2. Control run params
3. Control execution DAG
4. Code version control
5. Artifacts version control (models, datasets, etc.)
6. Use shared/cloud storage for artifacts
7. Environment dependencies control
8. Experiments results tracking

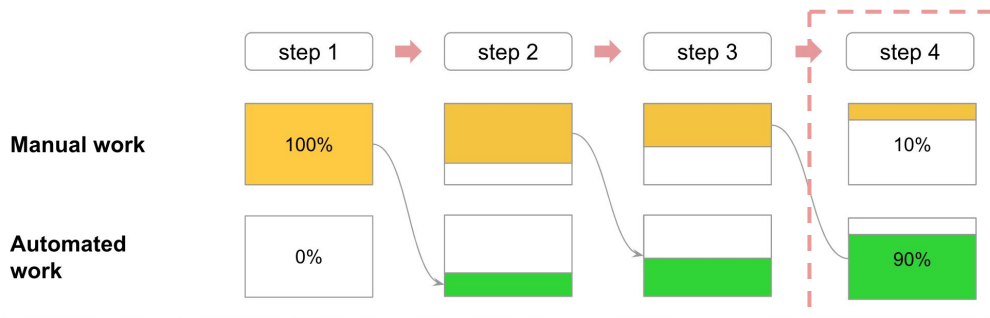


python



## Step 4: Experiments management

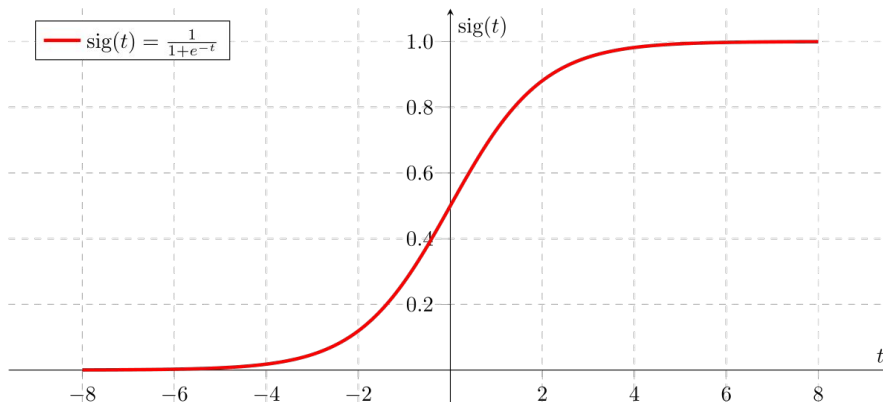
- add metrics tracking
- add experiments benchmark (DVC)
- pipelines dependencies under DVC control
- models/data/configs under DVC control
- same code in .py modules
- same pipelines
- everything in Docker
- run experiments in terminal or Jupyter Notebook



# Experiment 1:

## Tune Logistic Regression

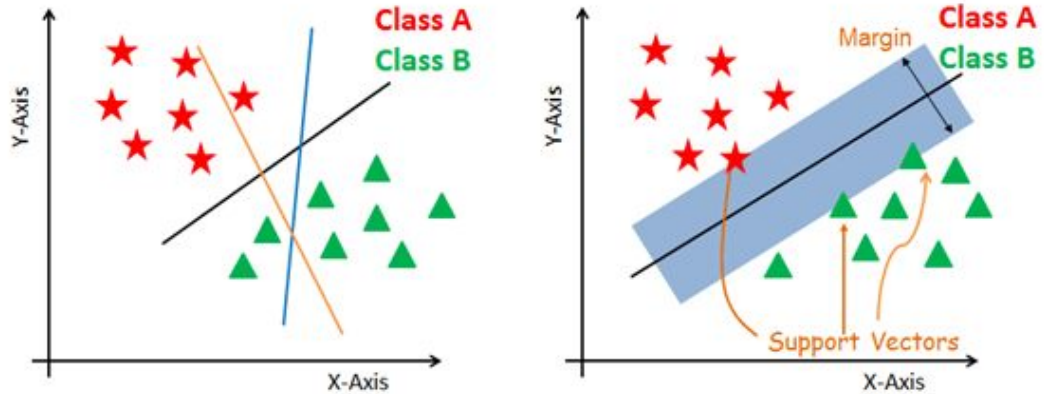
- hyperparameters tuning



## Experiment 2:

### Use SVM

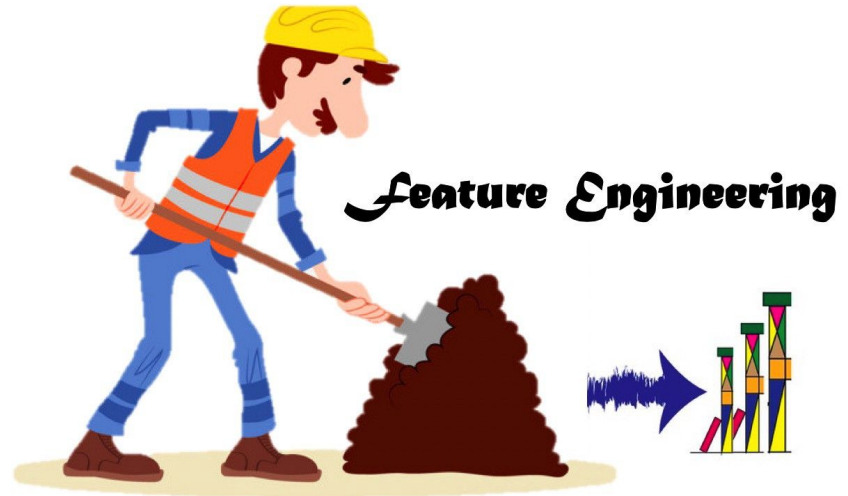
- use SVM estimator
- hyperparameters tuning



## Experiment 3:

### Add new features

- add squared features



## Compare experiments

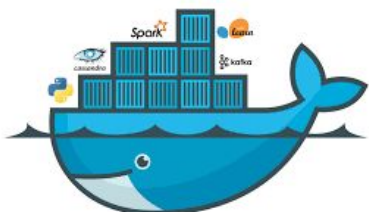
- **dvc metrics show**
- **dvc metrics show -a**
- **dvc metrics show -t json -x f1\_score -a**
- **dvc metrics show -T**



# ML Reproducibility checklist



1. Automated pipelines
2. Control run params
3. Control execution DAG
4. Code version control
5. Artifacts version control (models, datasets, etc.)
6. Use shared/cloud storage for artifacts
7. Environment dependencies control
8. Experiments results tracking



python



mlflow

# Conclusions

1. pipelines - not difficult
2. start where you detect a “copy-paste” pattern
3. artifacts version control - MUST
4. discipline in a team is important
5. more benefits for high complexity and large team projects



# Contact me



**Mikhail Rozhkov**

mail: [mnrozhkov@gmail.com](mailto:mnrozhkov@gmail.com)

ods: @Mikhail Rozhkov

