

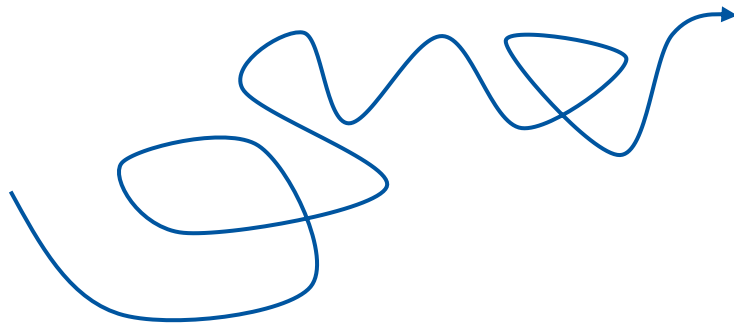
Machine Learning Operations - MLOps

Getting from Good to Great

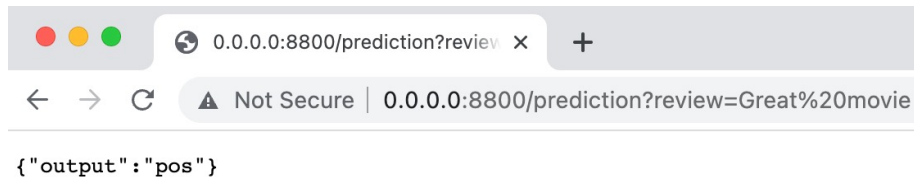
Slide credit: Michal Maciejewski, PhD



Acknowledgements: Dejan Golubovic, Ricardo Rocha, Christoph Obermair, Marek Grzenkowicz

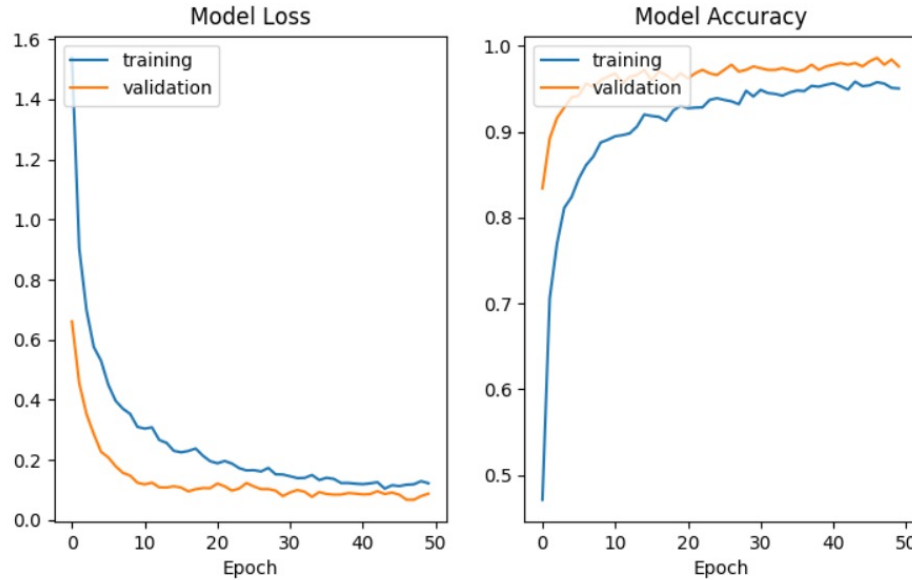


$Y = f(X)$



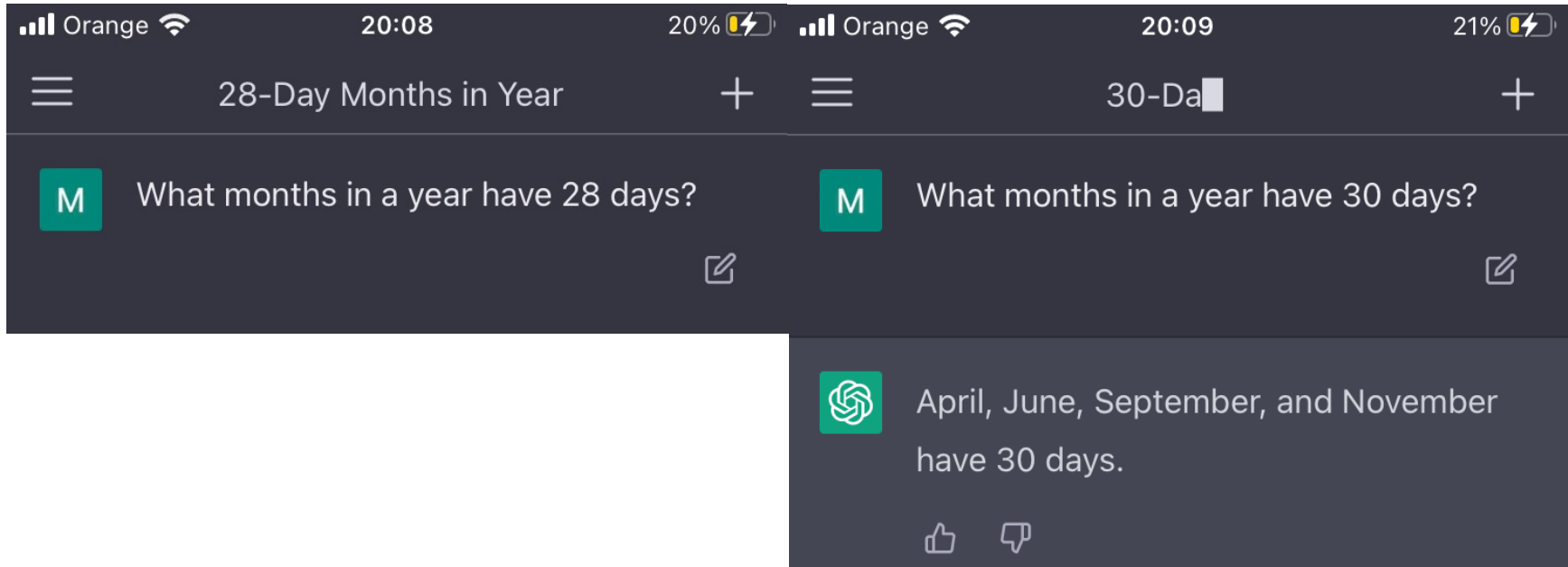
Let's share our model with users *aka* let's put it into production!

What Has to Go Right?



What is needed for an ML model to perform well in production?

What Can Go Wrong?



Concept and data drifts are one of the main challenges of production ML systems!

MLOps is about maintaining the trained model performance in production.*

*The performance may degrade due to factors outside of our control
so we ought to monitor the performance and if needed, roll out a new model to users.*

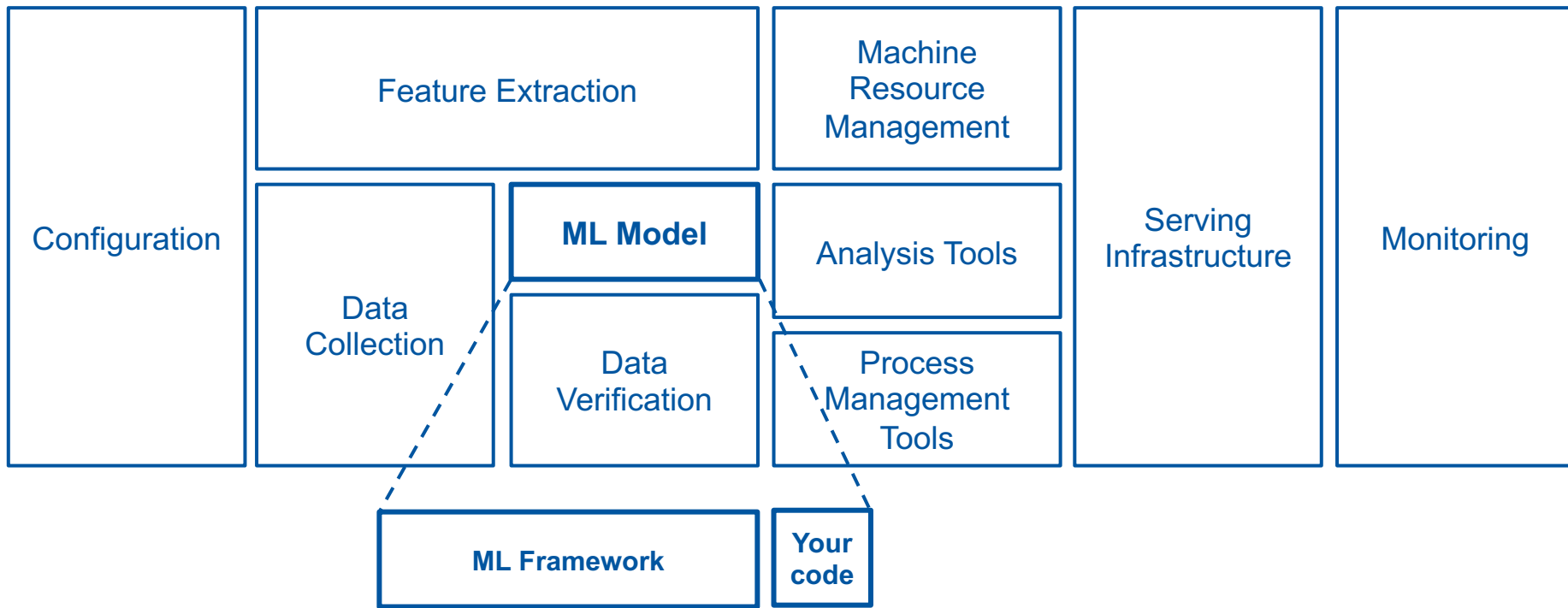
ML Model = Data + Code

MLOps = ML Model + Software

- + Algorithm
- + Weights
- + Hyperparameters

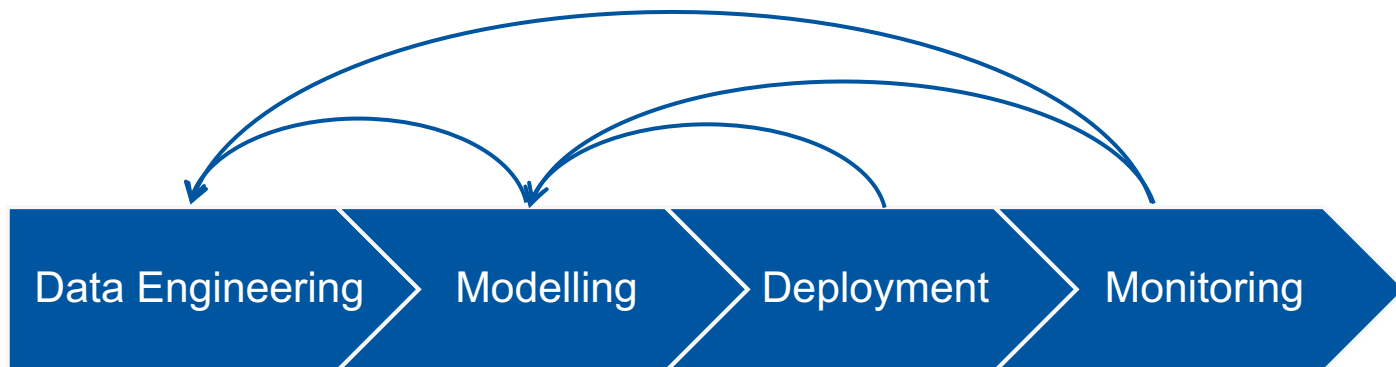
- + Scripts
- + Libraries
- + Infrastructure
- + DevOps

MLOps = ML Model + Software



Good news: most of these components come as ready-to-use frameworks

MLOps Pipeline



MLOps is a multi-stage, iterative process.

Data Engineering

Reproducibility

Traceability

Data-driven ML

$$f(\text{trash}) = \text{trash}$$

Exploratory Data Analysis

For structured data:

- schema as required tables, columns and datatypes

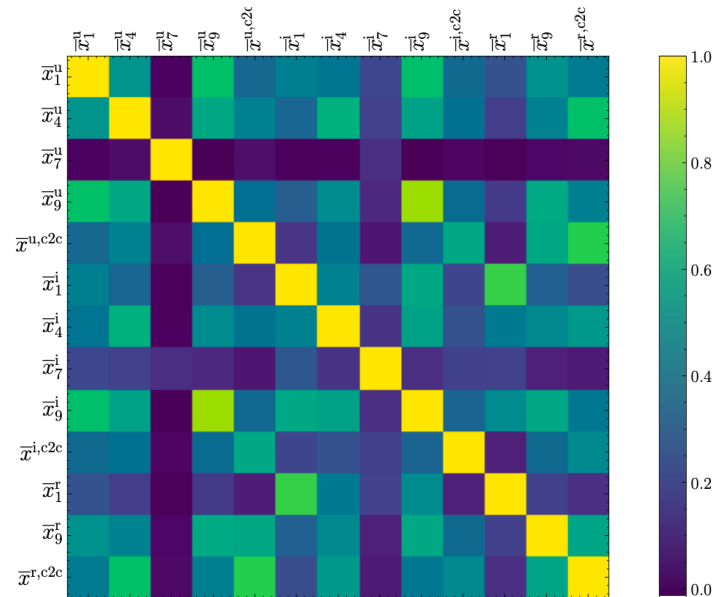
For unstructured data:

- resolution, image extension
- frequency, duration, audio codec

```
DataFrame.corr(method='pearson', min_periods=1, numeric_only=_NoDefault.no_default)
```

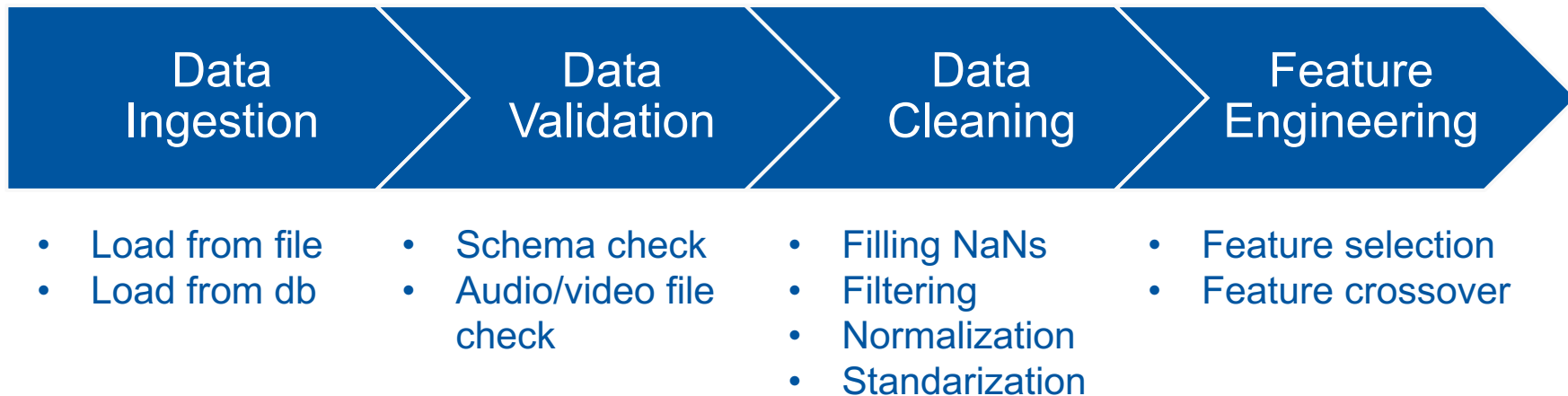
Compute pairwise correlation of columns, excluding NA/null values.

[\[source\]](#)



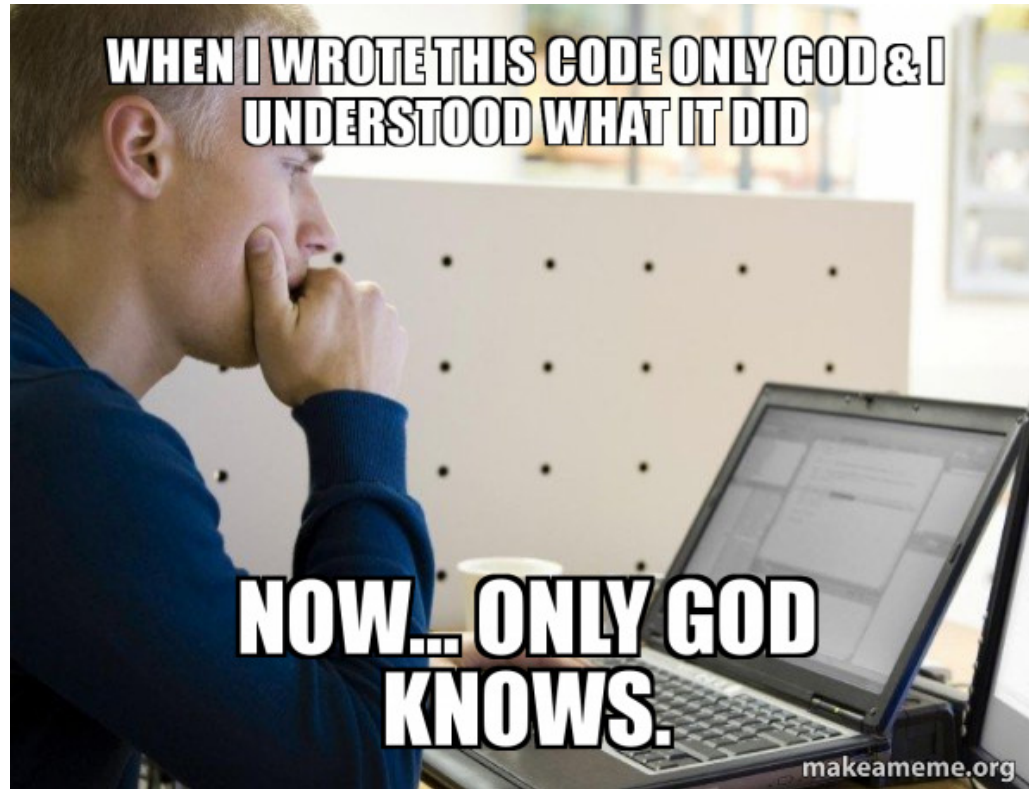
Initial exploration allows identifying requirements for input data in production.

Data Processing Pipeline



We need to reproduce some of those steps (e.g. subtracting mean) in production!

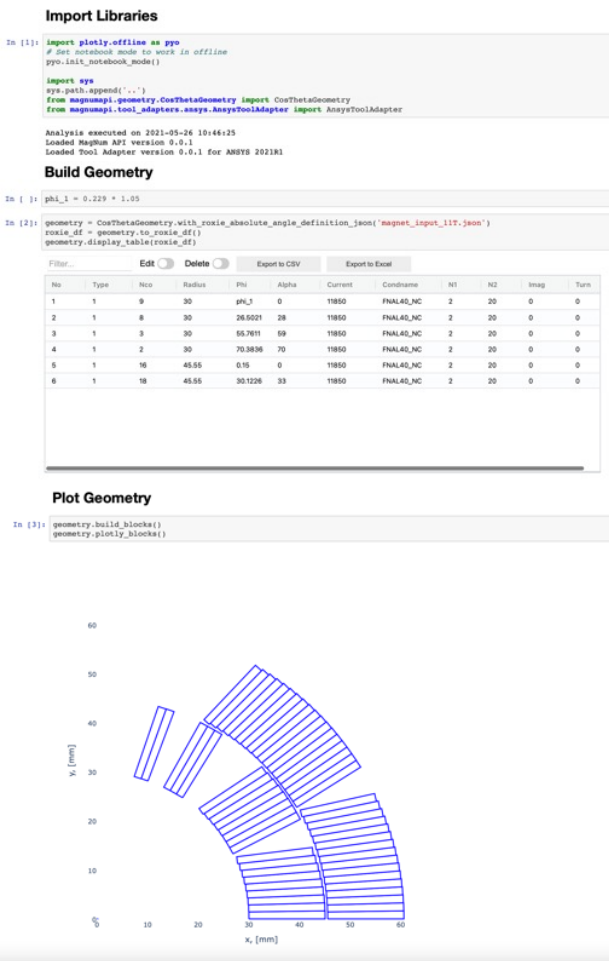
Reproducibility



Keeping Track of Data Processing

- Version Input Data – DVC framework
- Version Processing Script - GitLab
- Version Computing Environment - Docker

Data Provenance – where does data come from?
Data Lineage – how data is manipulated?



Notebook Good Practices

- Linear flow of execution
- Little amount of code
- Extract reusable code into a package
- Pre-commit for cleaning notebook before committing to a repository
- Set parameters on top so that notebook can be treated as a function (papermill and scrapbook packages)

From Model-driven to Data-driven ML

	Model-driven ML	Data-driven ML
Fixed component	Dataset	Model Architecture
Variable component	Model Architecture	Dataset
Objective	High accuracy	Fairness, low bias
Explainability	Limited	Possible

Modelling

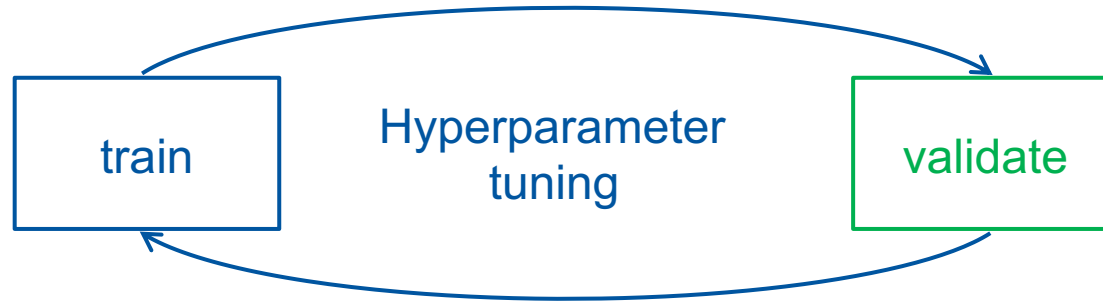
Training challenges

Rare events

Analyzing results

Selecting Data for Training

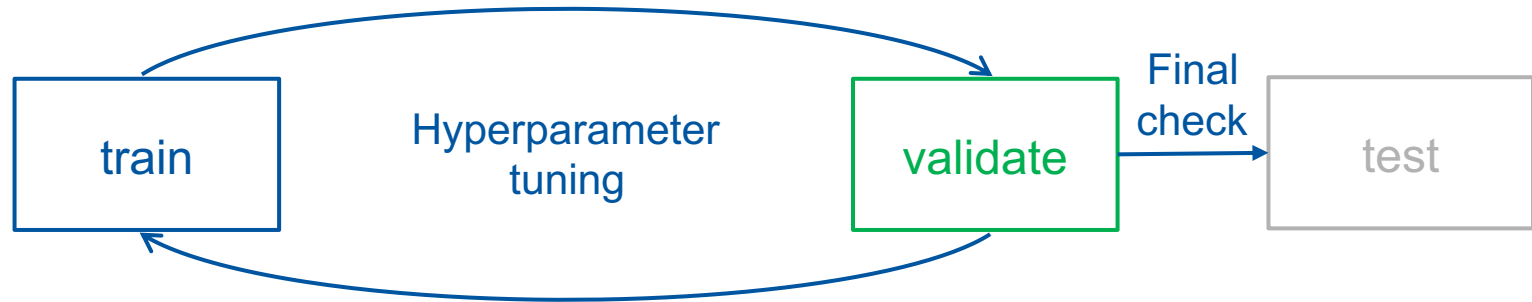
Dataset



With this approach, the model eventually sees the entire dataset.

Selecting Data for Training

Dataset



Splitting dataset in three allows to perform a final check with unseen data.

Balancing Datasets

Consider a binary classification problem with a dataset composed of 200 entries. There are 160 negative examples (no failure) and 40 positive ones (failure).

Expected:

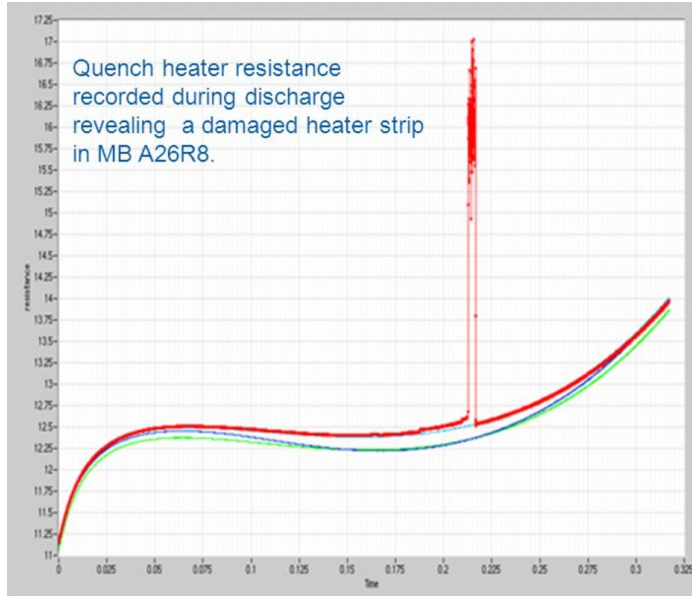
Training	Validation	Test
75%	15%	10%
(120 + 30)	(24+6)	(16+4)

Random:

Training	Validation	Test
75%	15%	10%
(131 + 19)	(19+11)	(10+10)

For continuous values it is important to preserve statistical distribution. Although for big datasets it is not an issue, it is still a low-hanging-fruit.

Rare Events



There were 3130 healthy signals ($Y=False$) and 112 faulty ones ($Y=True$)

Rare Events

```
1  import pandas as pd
2
3
4  def run_prediction(signal: pd.DataFrame) → bool:
5      return False
6
```

Rare Events

		Ground truth	
		Y = True	Y = False
Model	Y = True	0 <i>true positive</i>	0 <i>false positive</i>
	Y = False	112 <i>false negative</i>	3130 <i>true negative</i>

$$\text{Avg accuracy} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 97\%$$

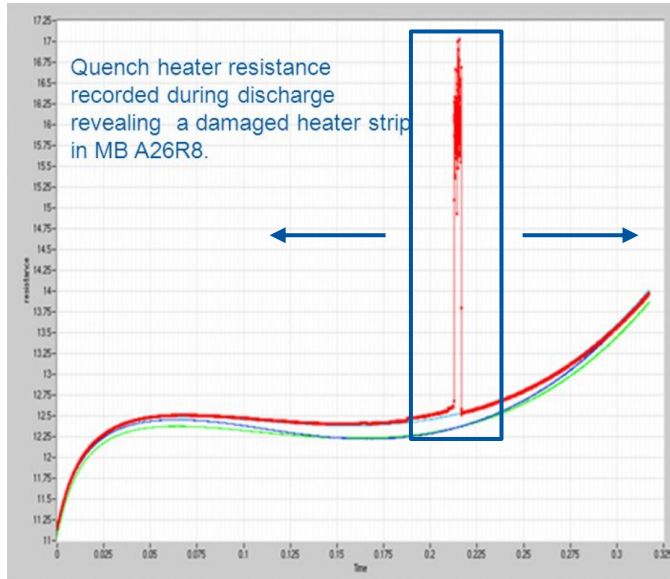
$$\text{Precision} = \frac{TP}{TP + FP} = \frac{0}{0}$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{0}{0 + 112} = 0$$

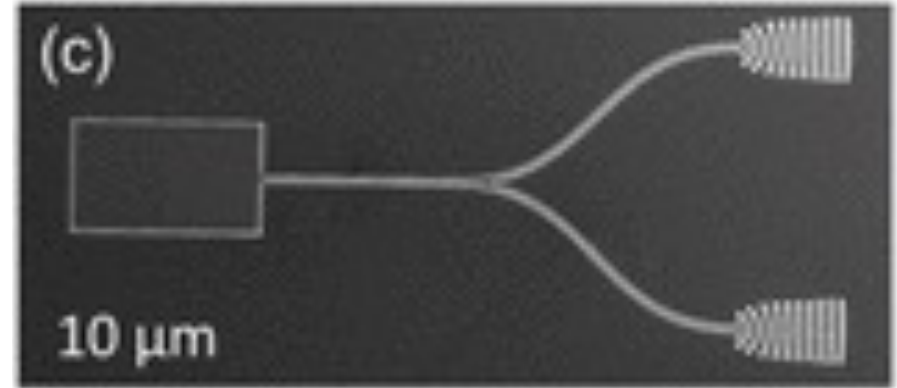
$$\text{F1}_{\text{score}} = \frac{2}{1/\text{Precision} + 1/\text{Recall}}$$

It is a valuable conversation to decide if precision or recall (or both) is more important.

Data Augmentation



New examples obtained by shifting the region left and right

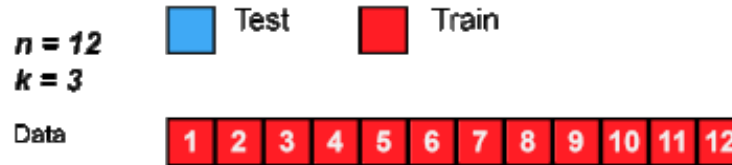


New examples obtained by rotating/shifting/hiding

What else can we do?

When one of the values of Y is rare in the population, considerable resources in data collection can be saved by randomly selecting within categories of Y . [...]

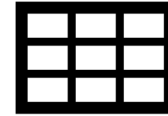
The strategy is to select on Y by collecting observations (randomly or all those available) for which $Y = 1$ (the "cases") and a random selection of observations for which $Y = 0$ (the "controls").



We can also collect more data of particular class (if even possible).

Training Tracking

1. Pen & Paper
2. Spreadsheet
3. Dedicated framework
 - Weights and Biases
 - Neptune.ai
 - Tensorflow
 - ...



Error Analysis

#	Signal	Noise	Gap in signal	Bias	Wrong sampling
1	Magnet 1	x	x		
2	Magnet 2			x	x
3	Magnet 3	x	x		

Such analysis may reveal issues with labelling or rare classes in data.

For unstructured data, a cockpit could help in analysis.

Useful in monitoring of certain classes of inputs.



Andrej Karpathy ✓

@karpathy



When you sort your dataset descending by loss you are guaranteed to find something unexpected, strange and helpful.

6:23 am · 2 Oct 2020



230 Retweets **41** Quote Tweets **2,173** Likes



Deployment

Degrees of automation

Modes of deployment

Reproducible environments

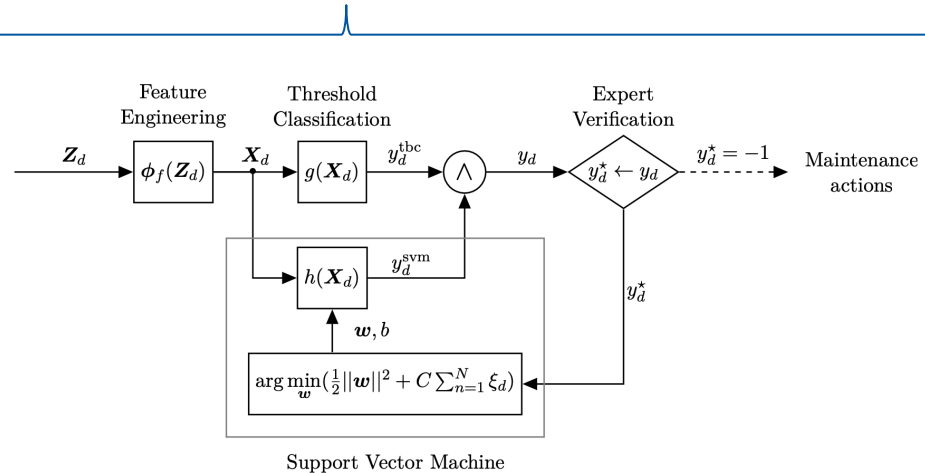
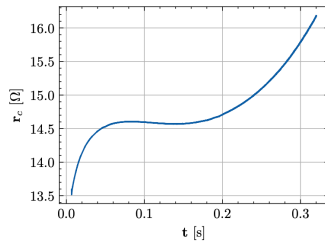
Degrees of Automation

Human inspection

Shadow mode

Human in the loop

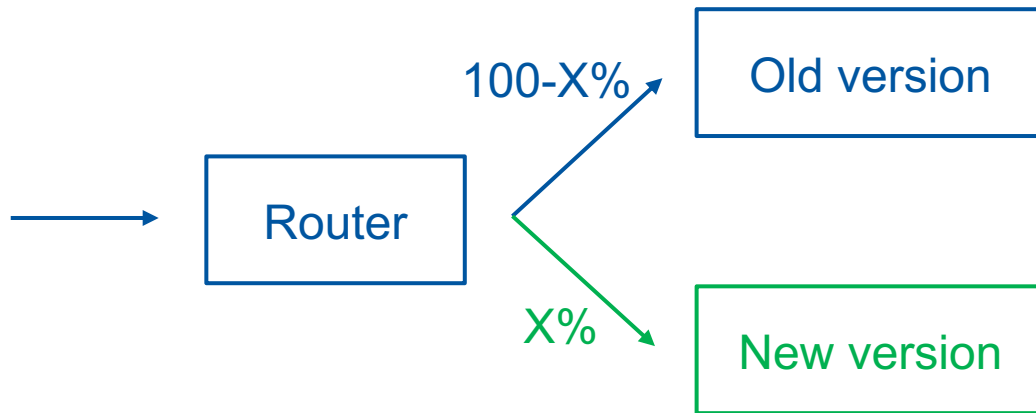
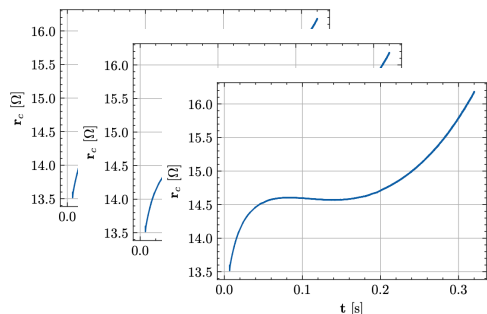
Full Automation



Starting from *Shadow mode* we can collect more training data in production!

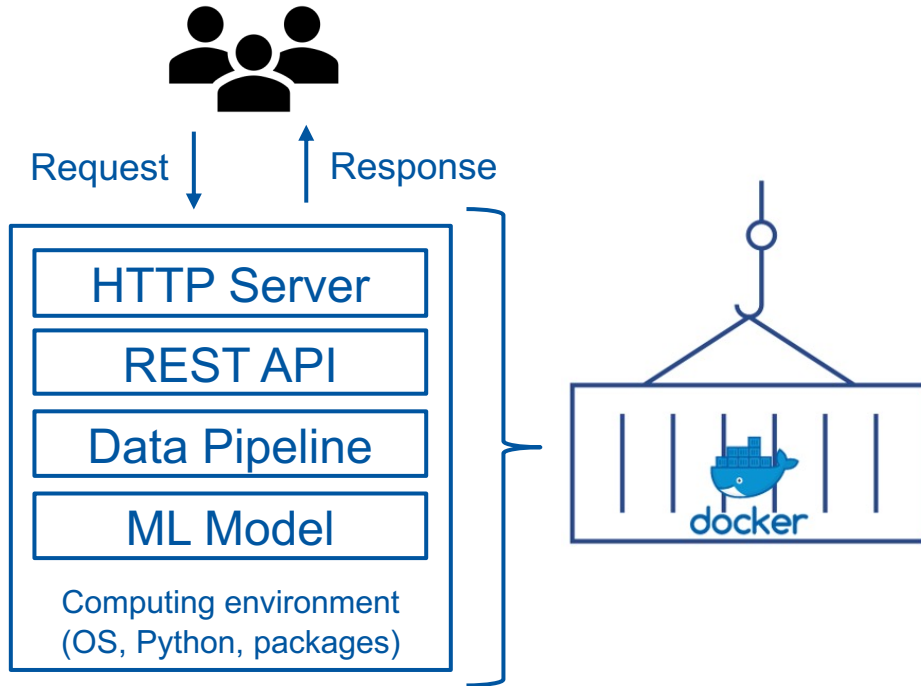


Modes of Deployment

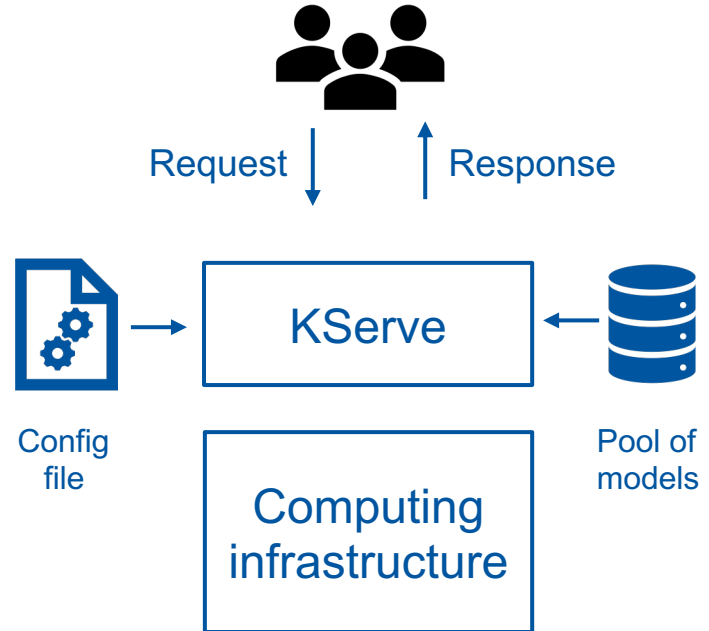


- In **Canary** deployment there is a gradual switch between versions
- In **Blue/green** deployment there is an on/off switch between versions

Reproducible Environments



Docker Containers



Serverless compute

We will play with those during the exercise sessions!

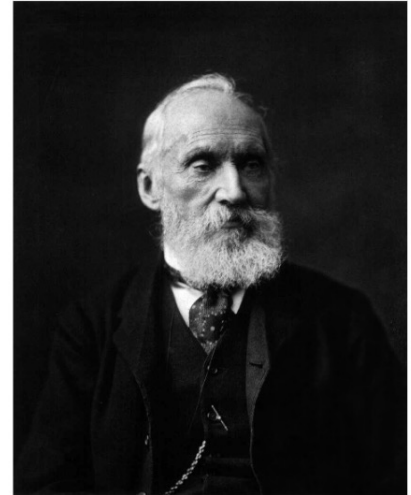
Monitoring

Useful metrics

Relevant frameworks

*If you can't measure it,
you can't improve it*

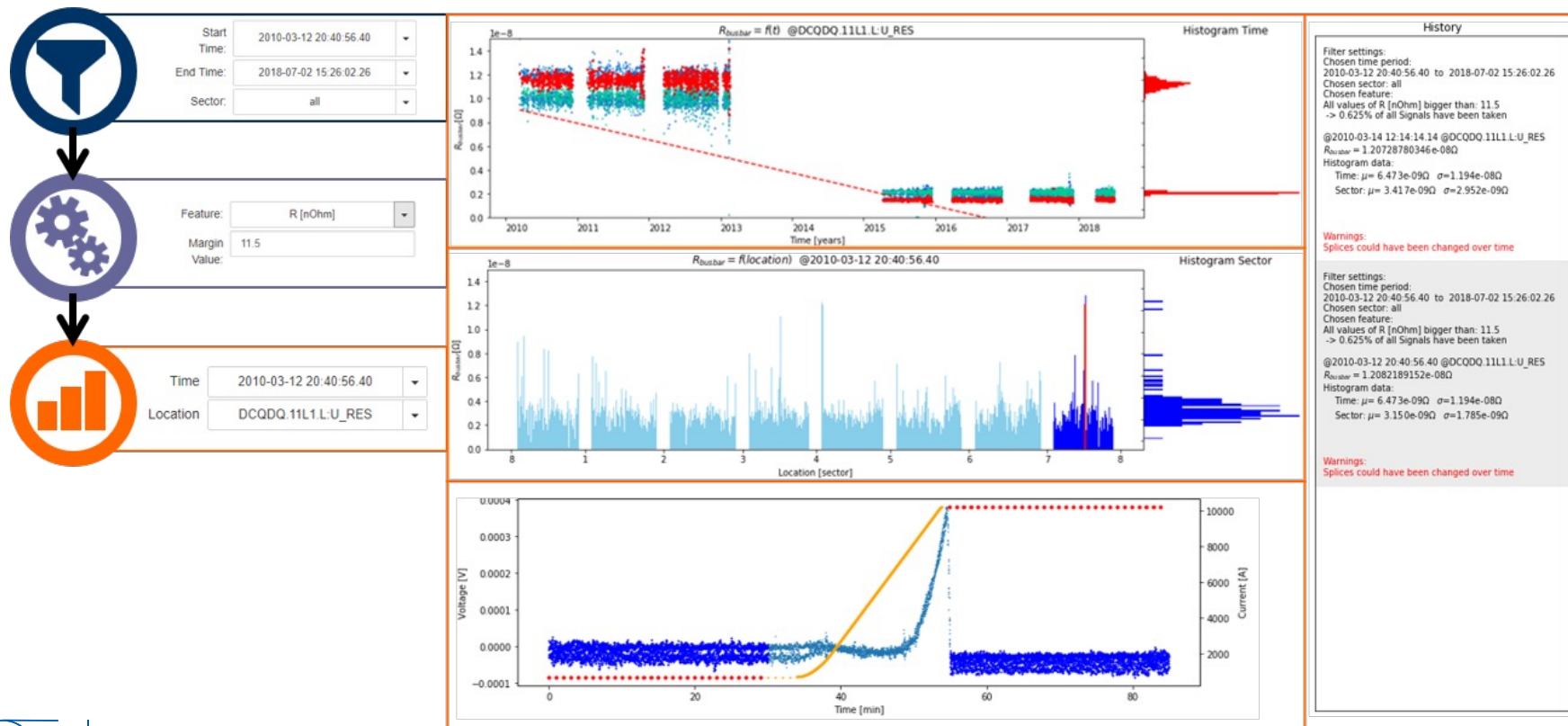
William Thomson, Lord Kelvin

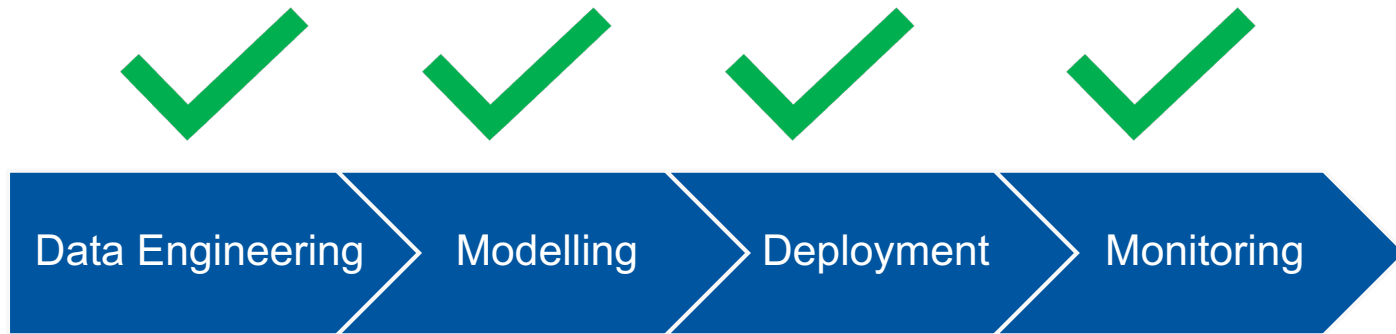


Relevant Metrics

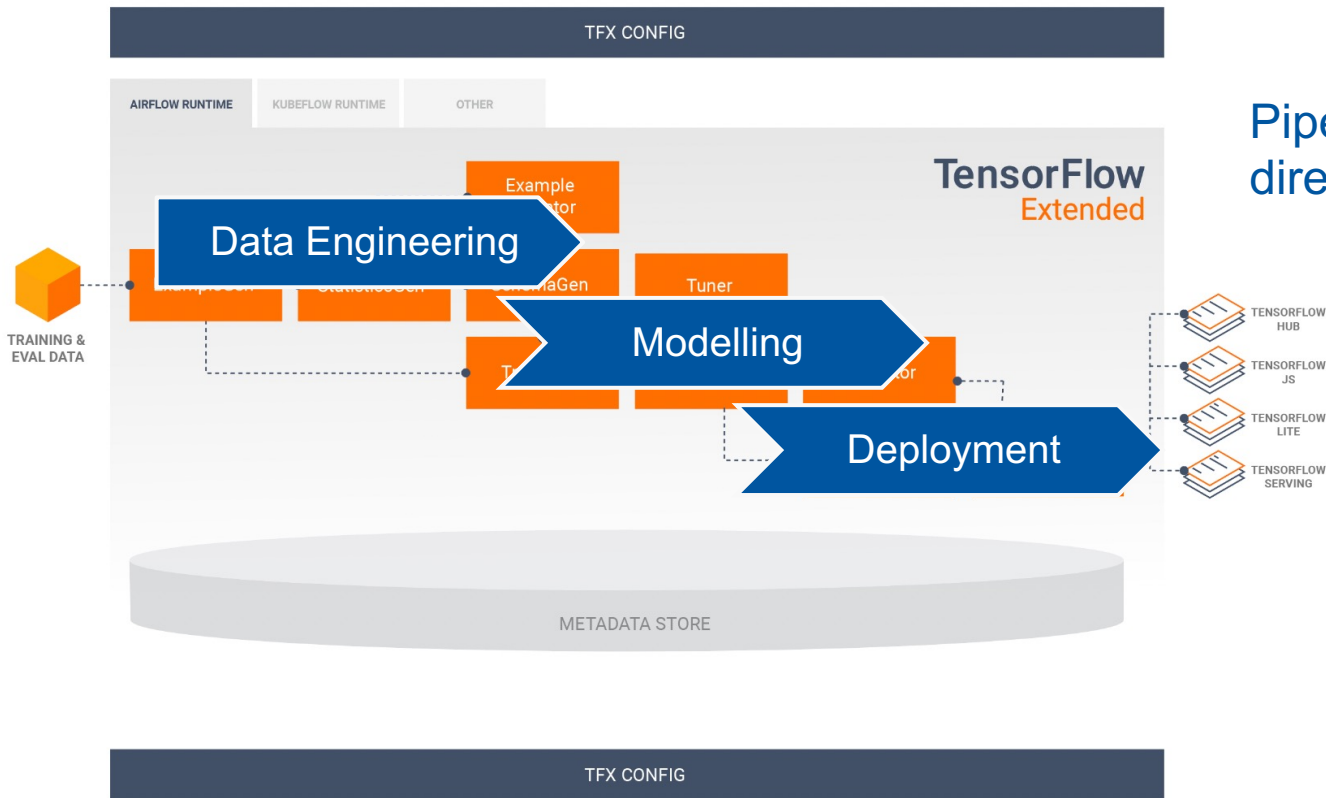
- Model metrics
 - Distribution of input features – data/concept drift
 - Missing/malformed values in the input
 - Average output accuracy/classification distribution – concept drift
- Infrastructure metrics
 - Logging errors
 - Memory, CPU resources utilization
 - Latency and jitter

Monitoring Matters





MLOps Pipeline with Tensorflow



Pipeline represented as DAG
directed acyclic graph

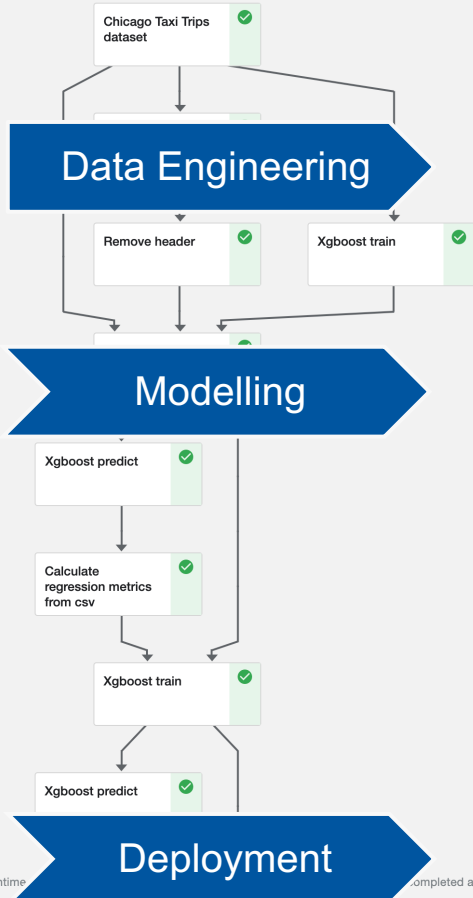
<https://www.tensorflow.org/tfx/guide>

MLOps Pipeline with Kubeflow

Graph

Run output

Config

☐ Simplify Graph

<https://ml.cern.ch>

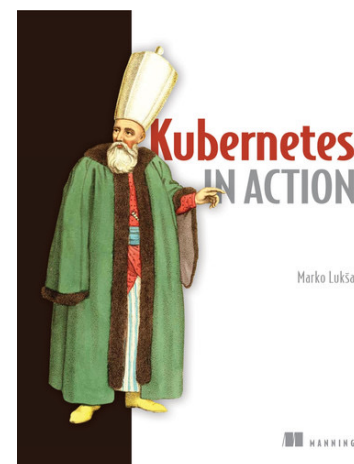
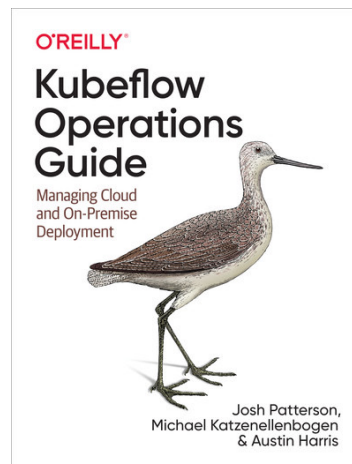
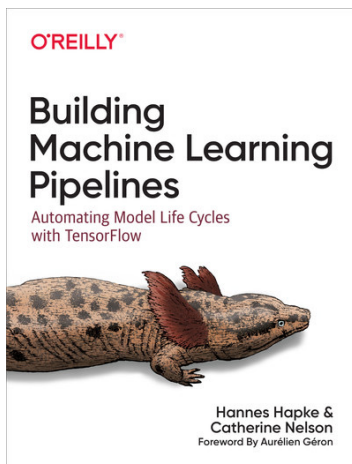
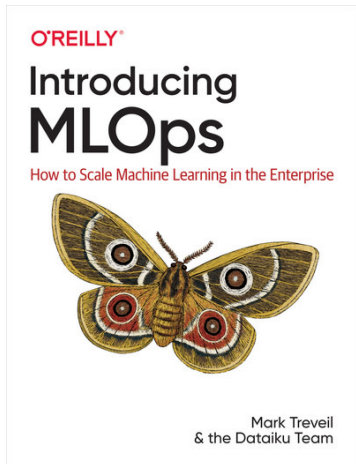
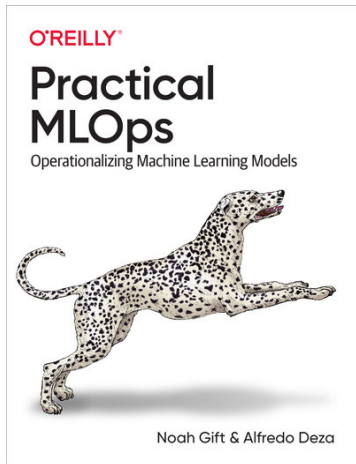
<https://www.kubeflow.org/docs/started/>

Conclusion

	Development ML	Production ML
Objective	High-accuracy model	Efficiency of the overall system
Dataset	Fixed	Evolving
Code quality	Secondary importance	Critical
Model training	Optimal tuning	Fast turn-arounds
Reproducibility	Secondary importance	Critical
Traceability	Secondary importance	Critical

I do hope the presented MLOps concepts will allow your models to transition
From Good to Great

Resources



Machine Learning Engineering for Production (MLOps) Specialization

