

Machine Learning Operations (MLOps)

Usage of Pipelines in the ML Lifecycle with
Tensor Flow Extended (TFX) and Kubeflow

Prof. Dr. Jan Kirenz
HdM Stuttgart

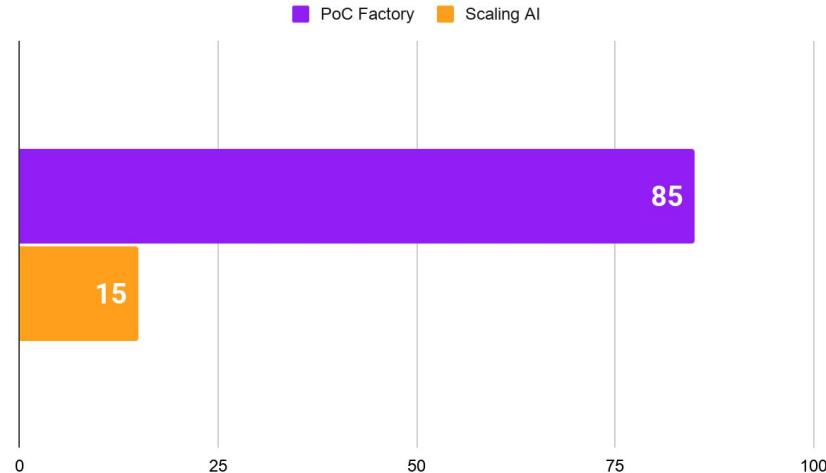
The Proof of Concept Factory

80-85% PoC Factory

Most companies...

- ... conduct AI experiments and pilots but achieve a low scaling success rate
- ... have significant under investment, yielding low returns

accenture



Gartner Top 10 Data and Analytics Trends, 2021



gartner.com/SmarterWithGartner

Source: Gartner
© 2021 Gartner, Inc. All rights reserved. CTMKT_1164473

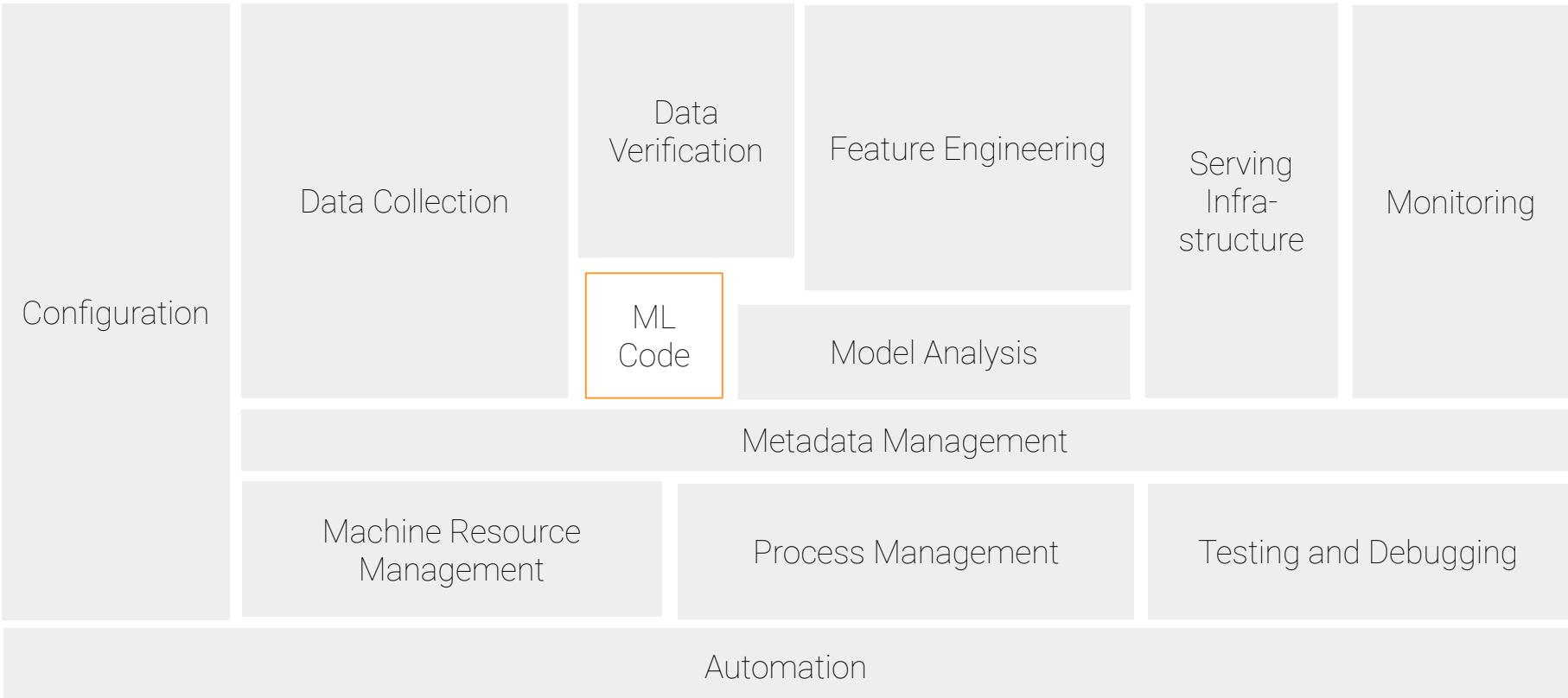
Gartner[®]

The problem with scaling AI

ML code is only a fraction of a production-ready ML project code

ML Project Code

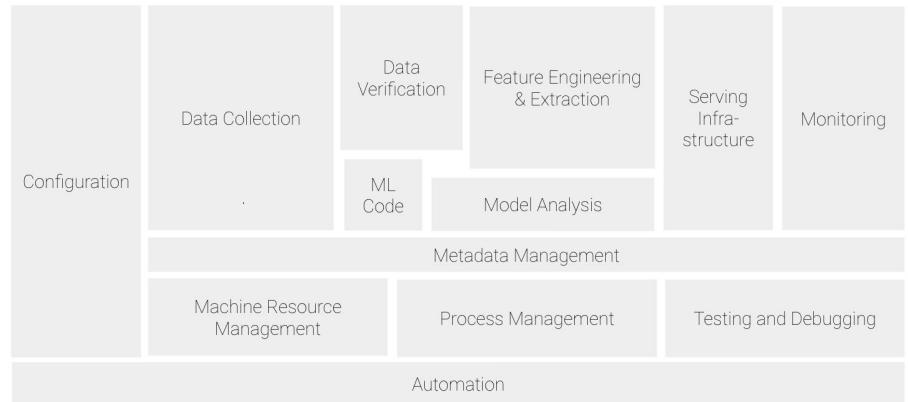




Hidden technical debt in machine learning systems

Machine learning operations (MLOps)

- ML Engineering culture and practice that aims at **unifying** ML System **development** (Dev) and ML system **operations** (Ops)
- Tools and principles to support workflow **standardization** and **automation** through the ML system lifecycle (e.g. with *pipelines*)

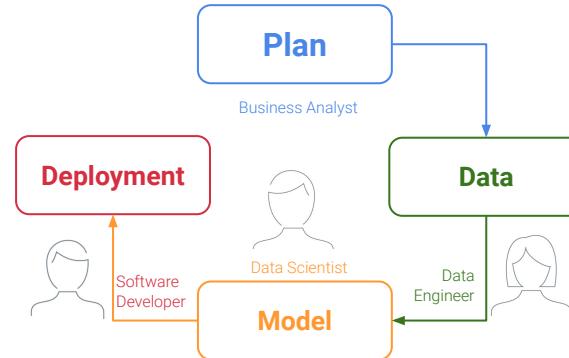


Machine learning **lifecycle**

Lifecycle

of an ML System

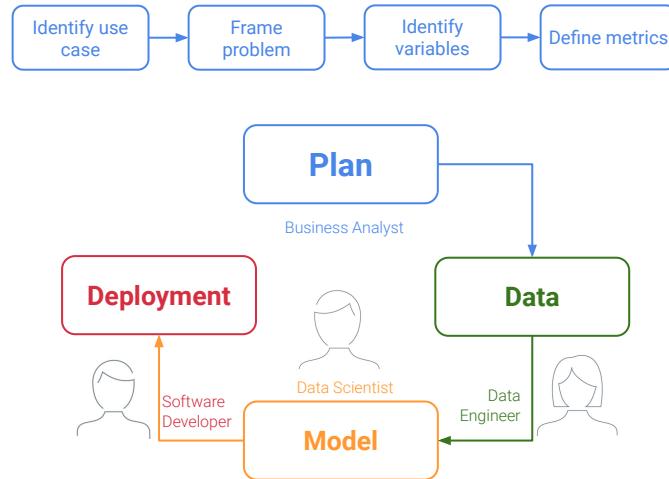
Plan | Data | Model | Deployment



Lifecycle

of an ML System

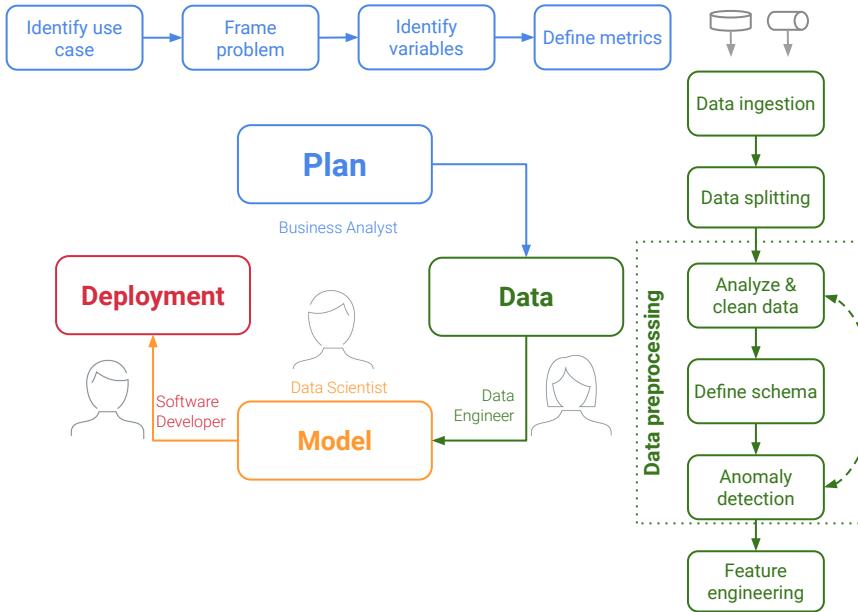
Plan | Data | Model | Deployment



Lifecycle

of an ML System

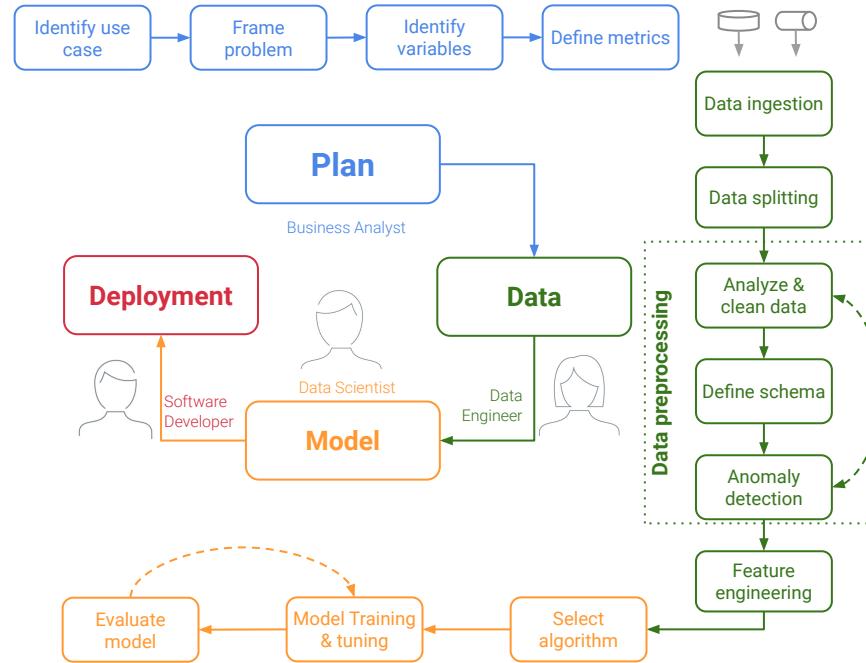
Plan | Data | Model | Deployment



Lifecycle

of an ML System

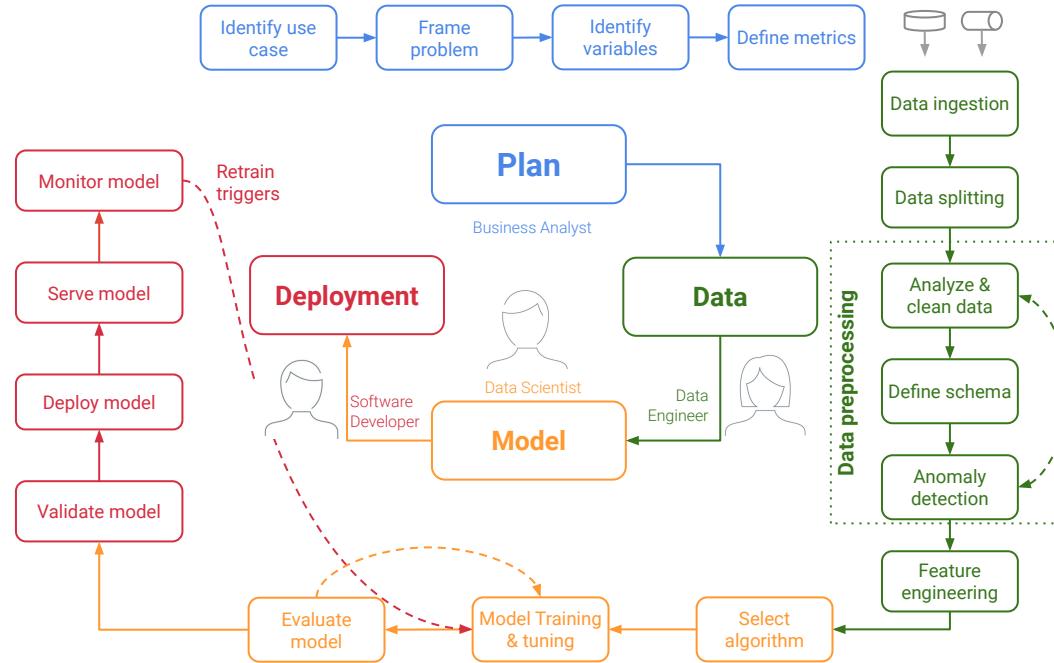
Plan | Data | **Model** | Deployment



Lifecycle

of an ML System

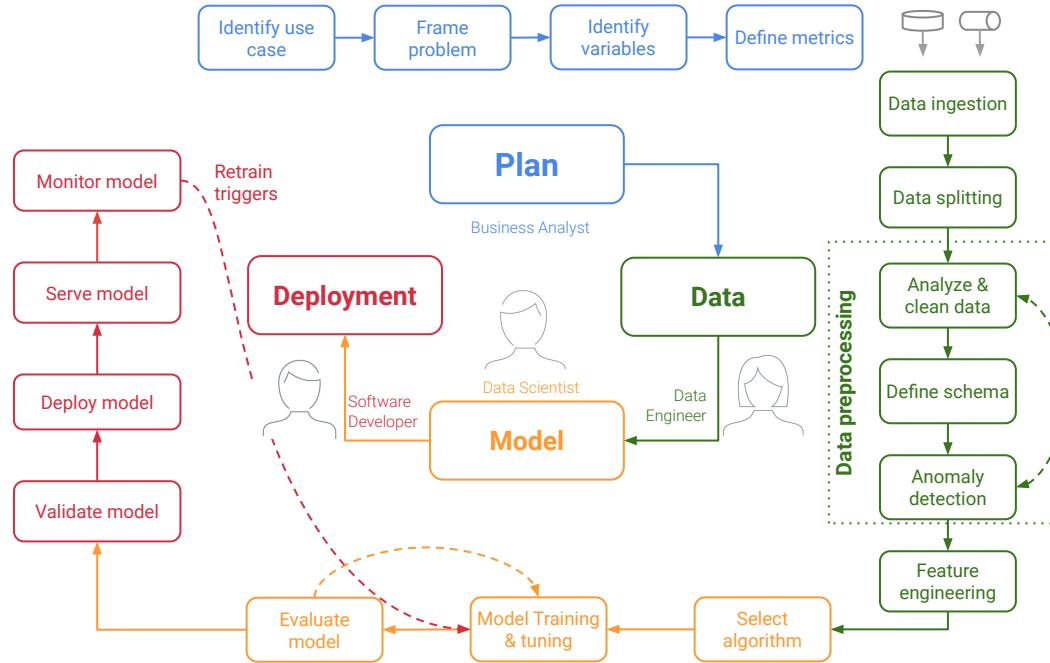
Plan | Data | Model | **Deployment**



Lifecycle

of an ML System

Plan | Data | Model | Deployment



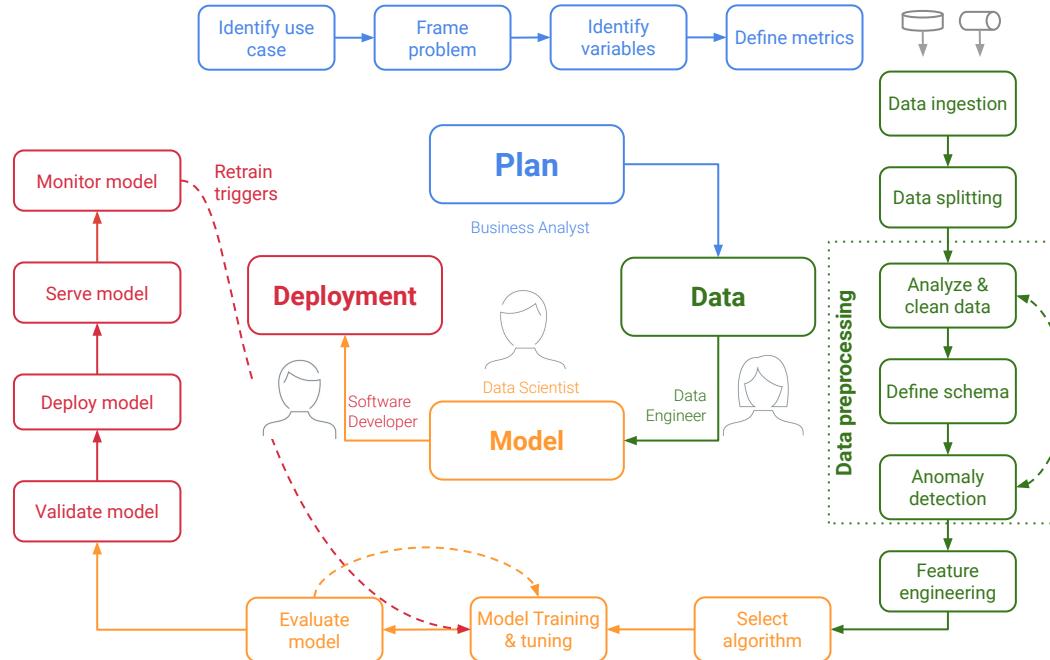
Common issues which lead to a PoC to production gap

- Lack of reuse and duplication
- Inconsistency (data, code, models)
- Manual and slow transition from PoC to production

Lifecycle

of an ML System

Plan | Data | Model | Deployment

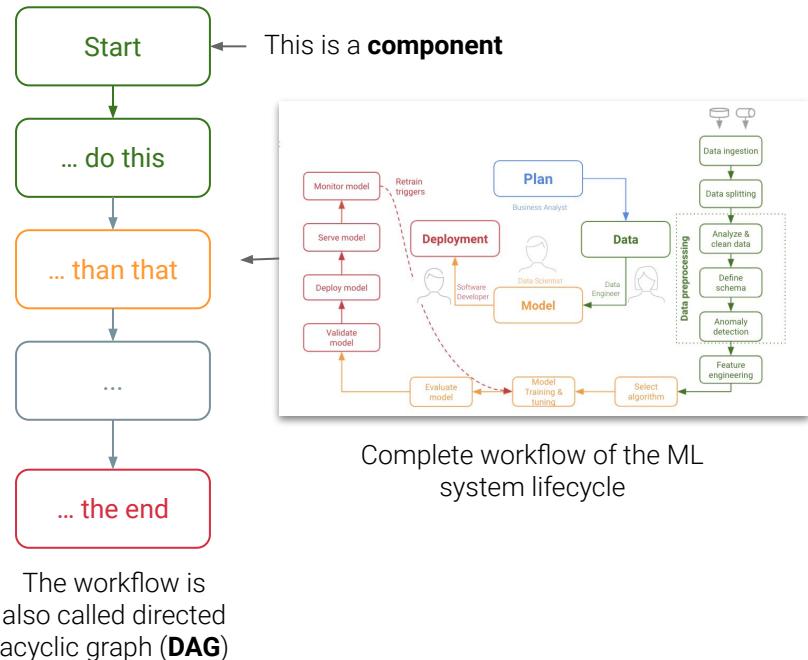


MLOps components



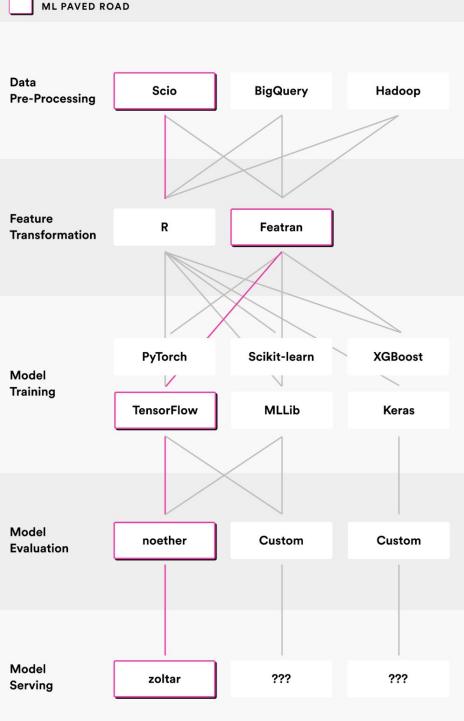
What is a pipeline?

- Description of an ML **workflow**
- A pipeline **component** is a self-contained set of user code that performs one step in the pipeline
- Includes the definition of the **configuration** and **inputs** required to run the pipeline (e.g. model hyperparameters)

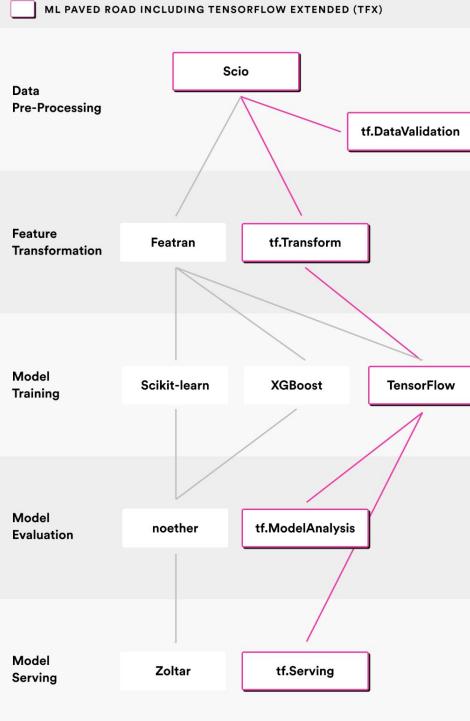




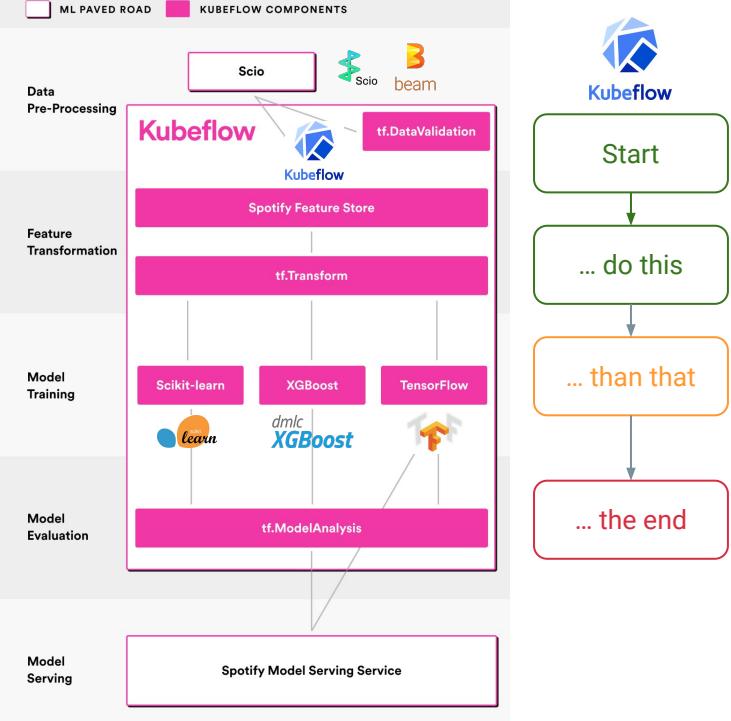
ML Platform 2018



ML Platform 2019



ML Platform 2020



TensorFlow Extended (TFX)

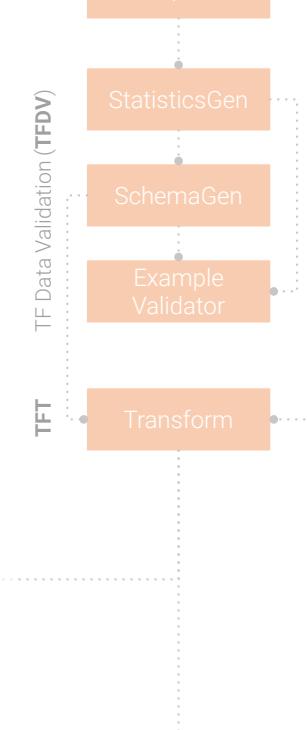
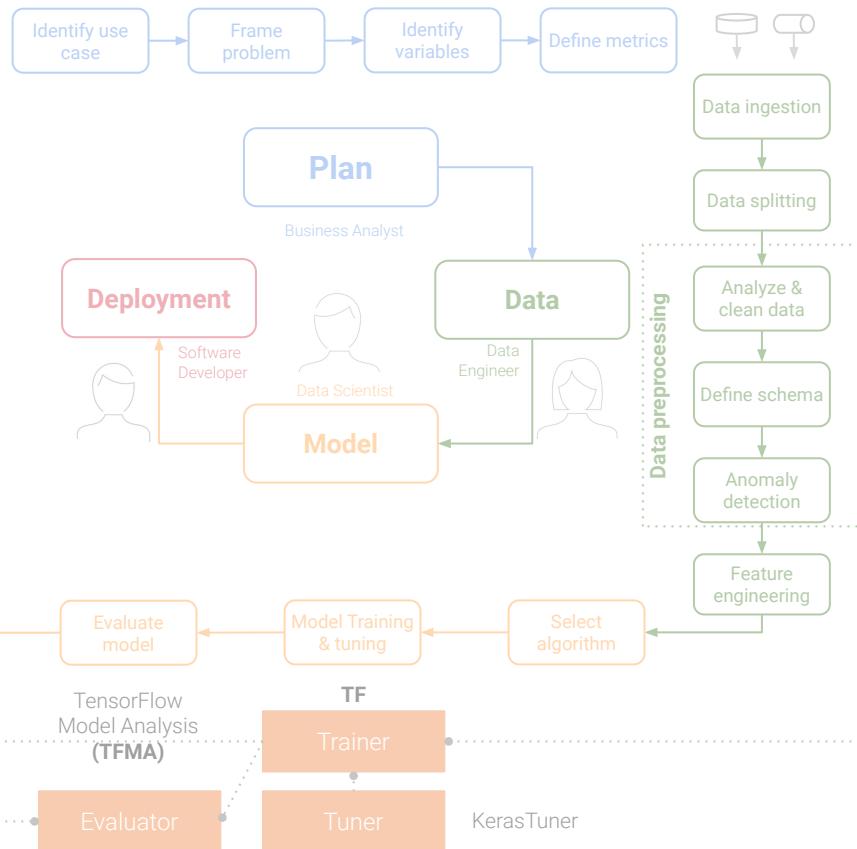
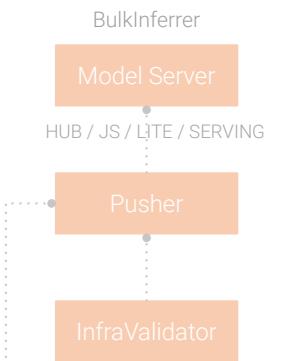
- Google-production-scale machine learning (ML) platform based on TensorFlow
- Portable to multiple environments (Azure, AWS, Google Cloud, IBM, ...)
- Python based toolkit; can be used with notebooks
- Helps you orchestrate your ML process: Apache Airflow, Apache Beam or Kubeflow pipelines



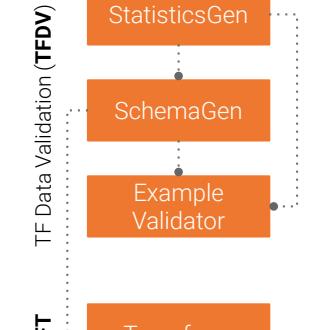
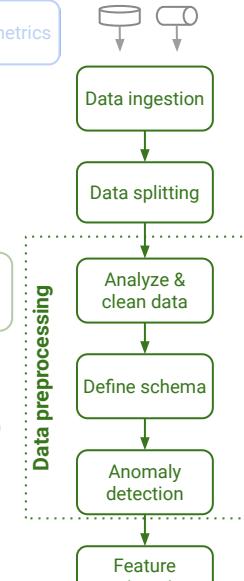
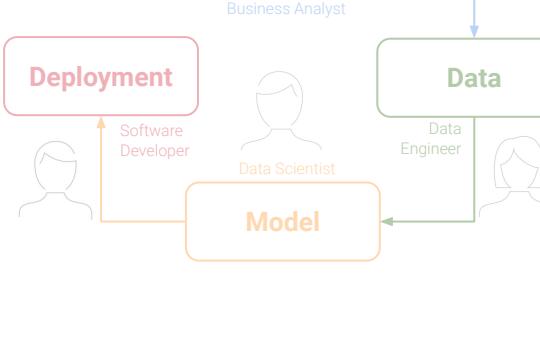
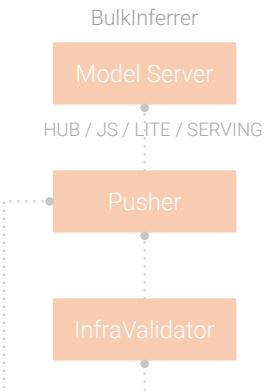
TFX 1.0 (19.05.21)

- Enterprise-grade support
- Security patches and select bug fixes for up to three years
- Guaranteed API & Artifact backward compatibility





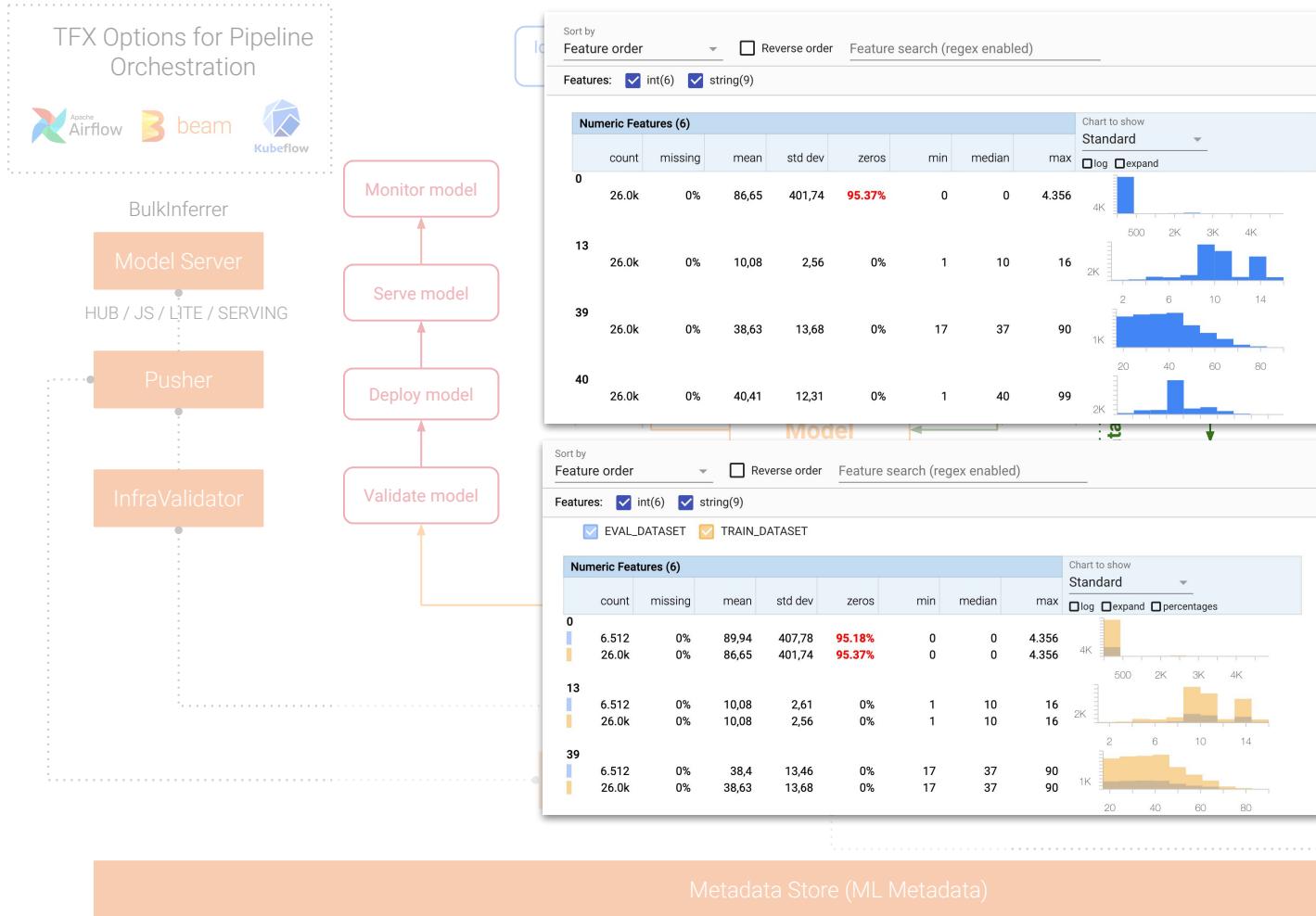
TFX Options for Pipeline Orchestration



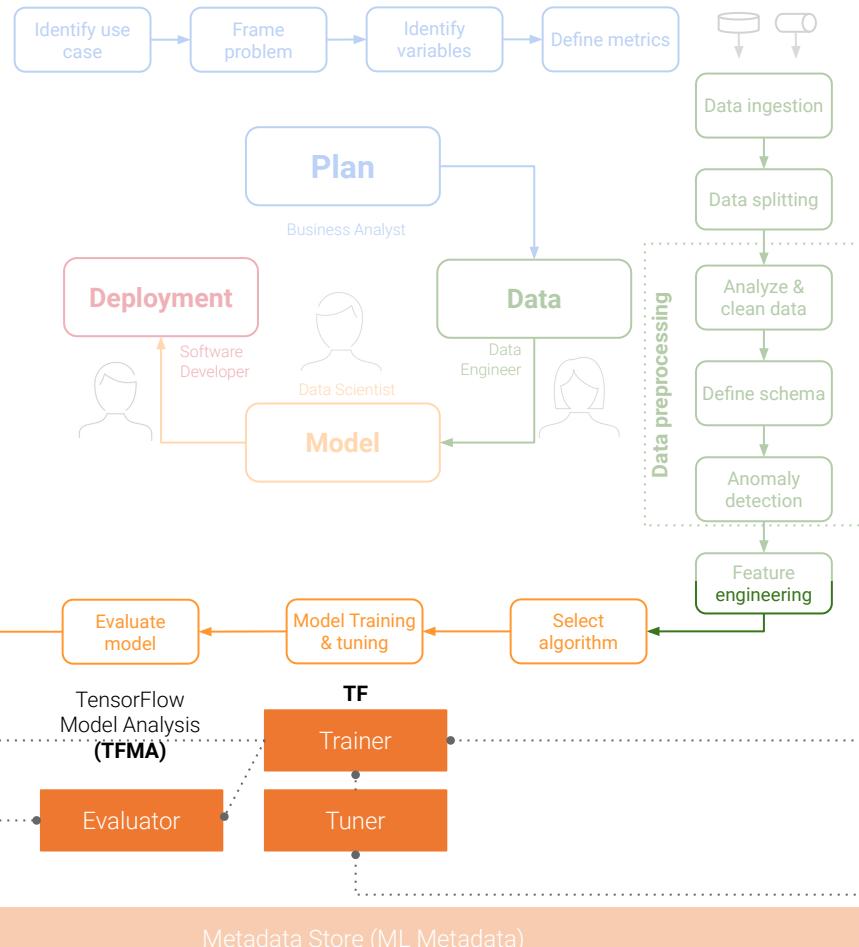
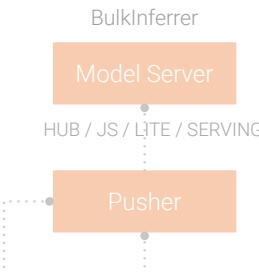
Metadata Store (ML Metadata)



TensorFlow Extended



TFX Options for Pipeline Orchestration



TensorBoard.dev

SCALARS

- Show data download links
 Ignore outliers in chart scaling

Tooltip sorting method: default

Smoothing

Horizontal Axis

STEP RELATIVE WALL

Runs

Write a regex to filter runs

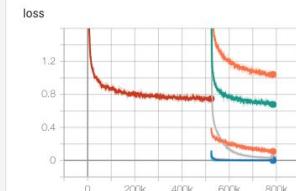
- cnn_dailymail_v002
- glue_v002_proportional
- pretrain
- squad_v010_allanswers
- super_glue_v102_proportional
- wmt15_enfr_v003
- wmt16_enro_v003
- wmt_t2t_ende_v003

TOGGLE ALL RUNS

experiment EvNO346lT0lYbmeaWmoNCQ

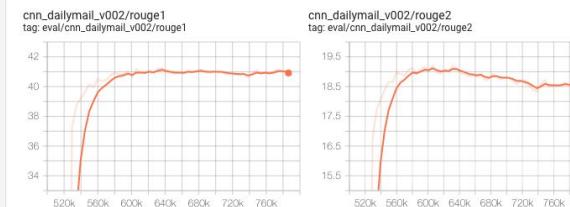
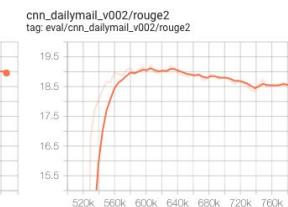
loss

Tags matching /loss/



eval

PREVIOUS PAGE

cnn_dailymail_v002/rouge1
tag: eval/cnn_dailymail_v002/rouge1cnn_dailymail_v002/rouge2
tag: eval/cnn_dailymail_v002/rouge2TensorFlow
Model Analysis
(TFMA)

Evaluator

Trainer

Tuner

Metadata Store (ML Metadata)

Model Card for Census Income Classifier

Model Details

Overview

This is a wide and deep Keras model which aims to classify whether or not an individual has an income of over \$50,000 based on various demographic features. The model is trained on the UCI Census Income Dataset. This is not a production model, and this dataset has traditionally only been used for research purposes. In this Model Card, we review quantitative components of the model's performance and data, as well as information about the model's intended uses, limitations, and ethical considerations.

Version

name: 36dea2e860670aa74691b5695587afe7

Owners

• Model Cards Team, model-cards@google.com

References

• interactive-2020-07-28T20_17_47.911887

Considerations

Use Cases

- This dataset that this model was trained on was originally created to support the machine learning community in conducting empirical analysis of ML algorithms. The Adult Data Set can be used in fairness-related studies that compare inequalities across sex and race, based on people's annual incomes.

Limitations

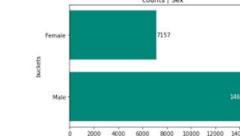
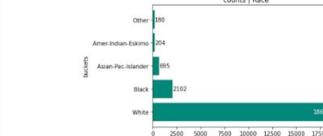
- This is a class-imbalanced dataset across a variety of protected classes. The ratio of male-to-female examples is about 2:1 and there are more examples with the "White" attribute than the "Other" attribute compared to the "Asian-Pac-Islander" or "Amer-Indian-Eskimo" groups. At a threshold of \$50,000 less than or equal to \$50,000, most elements is just over 5:1. Due to the imbalance across income levels, we can see that our true negative rate seems quite high, while our true positive rate seems quite low. This is true to an even greater degree when we only look at the "Female" subgroup. We can't make any claims about the fairness of the model's predictions without running a formal audit to verify these examples. To avoid this, we can try various remediation strategies in future iterations (e.g. undersampling, hyperparameter tuning, etc), but we may not be able to fix all of the fairness issues.

Ethical Considerations

- Risk: We risk expressing the viewpoint that the attributes in this dataset are the only ones that are predictive of someone's income, even though we know this is not the case.
- Mitigation Strategy: As mentioned, some interventions may need to be performed to address the class imbalances in the dataset.

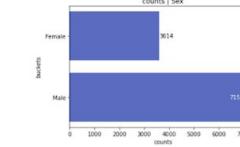
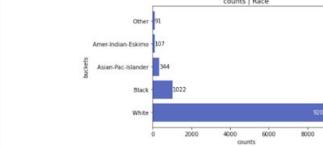
Train Set

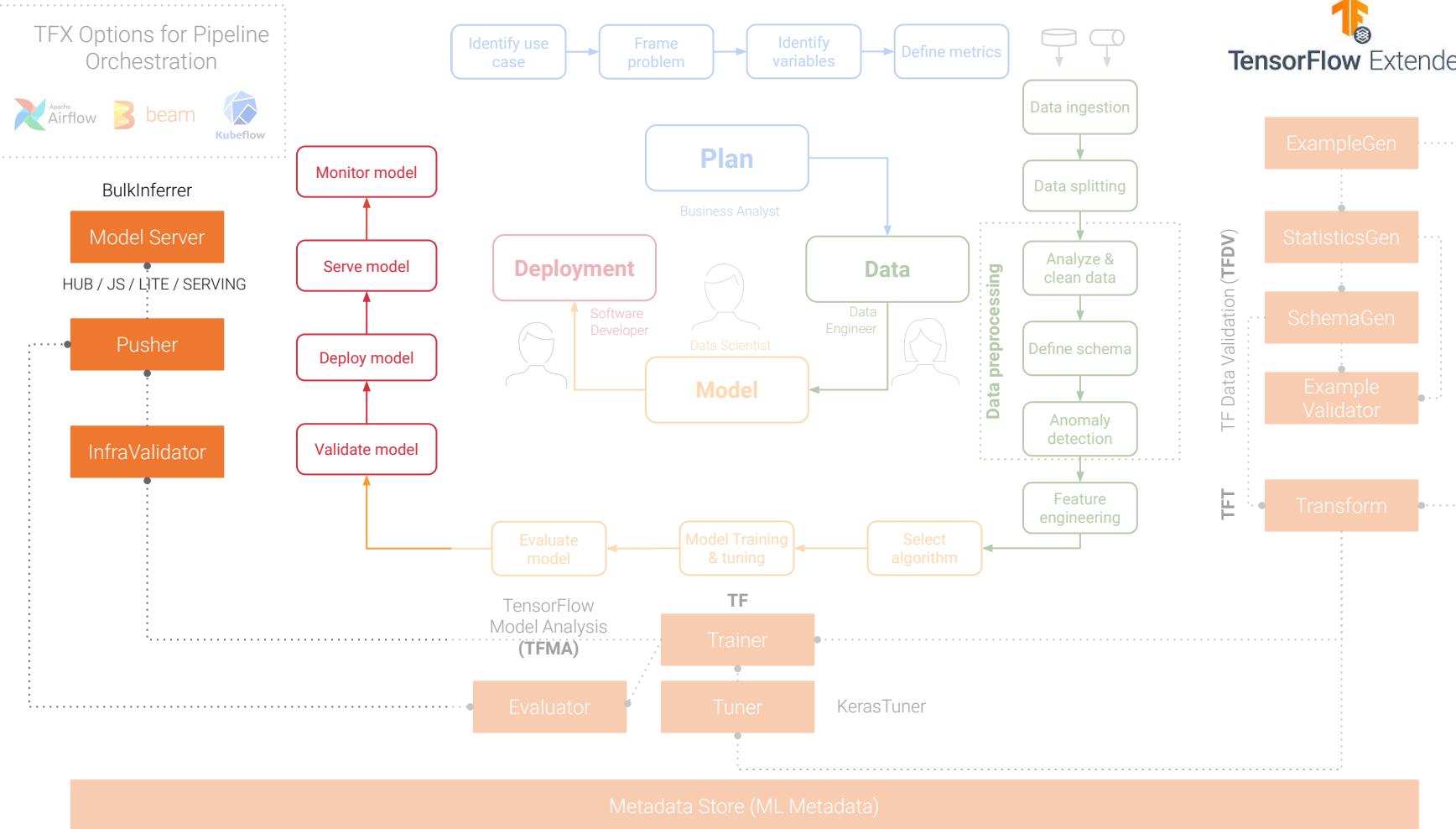
This section includes graphs displaying the class distribution for the "Race" and "Sex" attributes in our training dataset. We chose to show these graphs in particular because we felt it was important that users see the class imbalance.



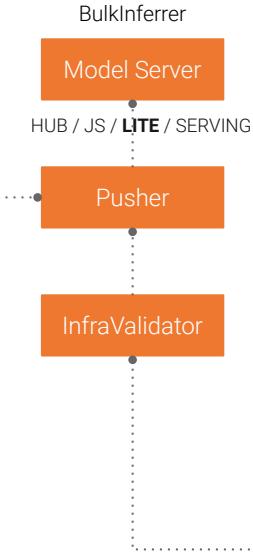
Eval Set

Like the training set, we provide graphs showing the class distribution of the data we used to evaluate our model's performance.

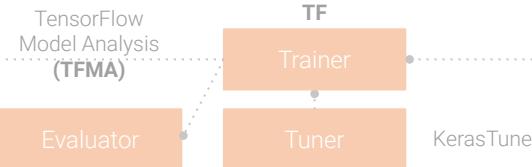




TFX Options for Pipeline Orchestration

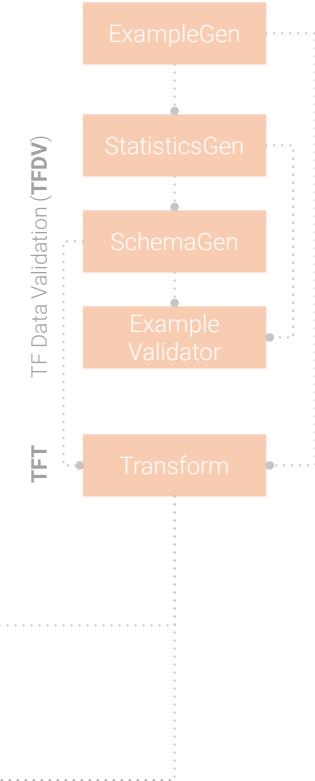


TensorFlow Lite is a set of tools that enables on-device machine learning by helping developers run their models on mobile, embedded, and IoT devices.



Metadata Store (ML Metadata)

 **TensorFlow Extended**



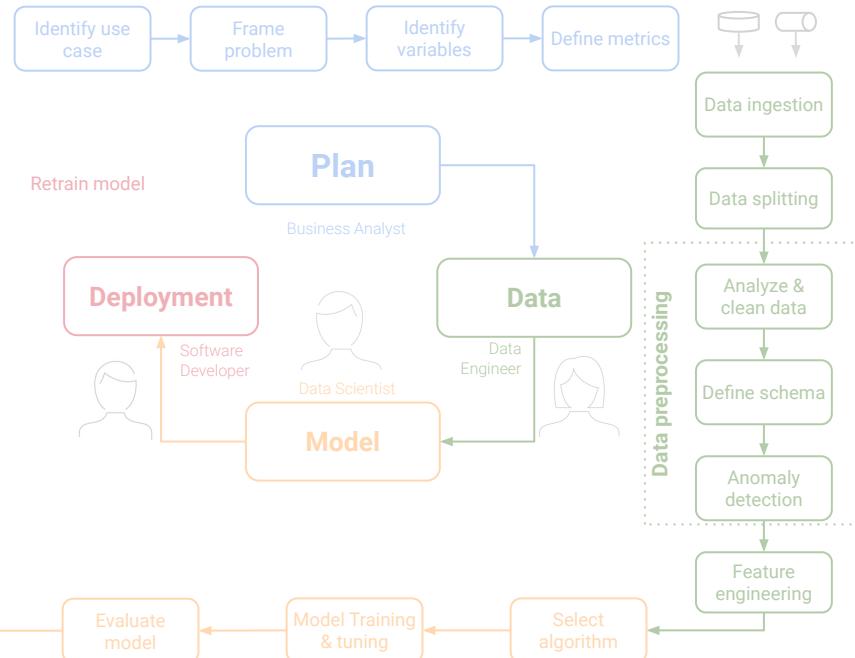


TFX Options for Pipeline Orchestration



Production phase:

automate the execution of the ML pipeline based on a schedule or certain triggering conditions.



Development phase: run the ML experiment, instead of manually executing each step.

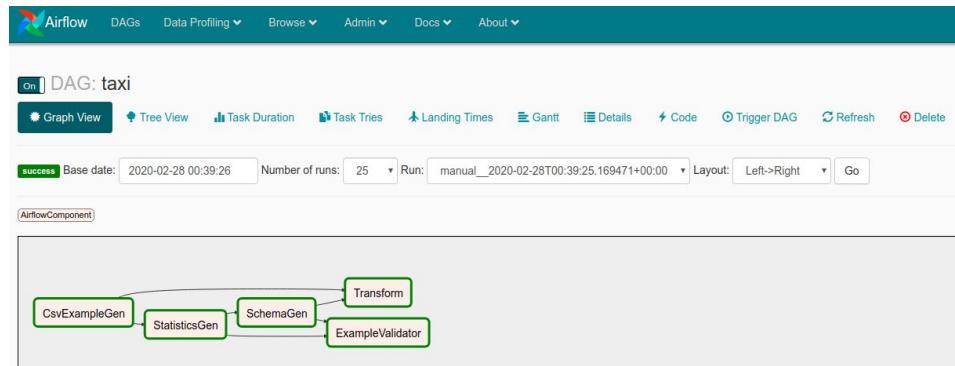


Pipeline orchestration

TFX & Apache Airflow



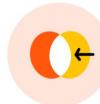
- Programmatically author, schedule and monitor workflows with **Python** code.
- **User interface** to visualize pipelines running in production, monitor progress, and troubleshoot issues.



TFX & Apache Beam



- Provides a framework for running **batch** and **streaming** data processing jobs that run on a variety of runners (Spark, Flink, ...).
- Beam provides an abstraction layer which enables TFX to run on any supported runner without code modifications
- TFX only uses the Beam **Python API**



Unified

Use a single programming model for both batch and streaming use cases.



Extensible

Write and share new SDKs, IO connectors, and transformation libraries.



Portable

Execute pipelines on multiple execution environments.



Open Source

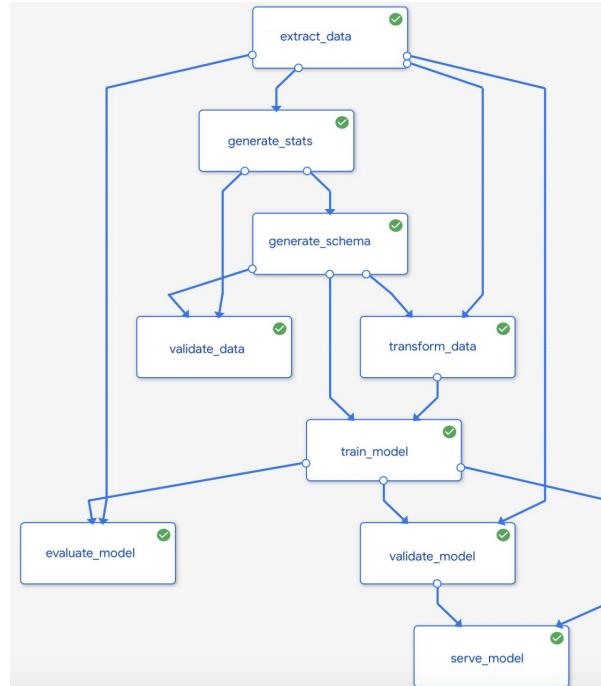
Community-based development and support to help evolve your application and use cases.

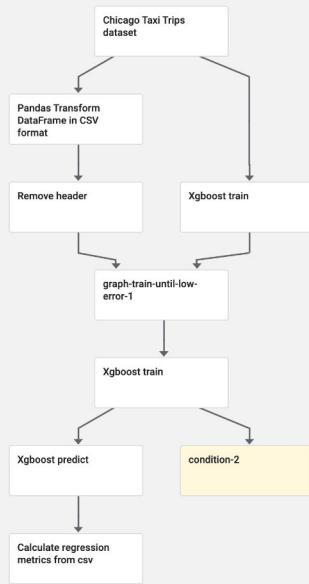
TFX & Kubeflow pipelines

The Kubeflow Pipelines platform consists of:

- An engine for scheduling multi-step ML workflows (using **Kubernetes**).
- **User interface** (UI) for managing and tracking experiments, jobs, and runs.
- **Python SDK** for defining and manipulating pipelines and components.
- **Notebooks** for interacting with the system using the SDK

Kubeflow Pipelines is available as a core component of Kubeflow or as a standalone installation.





Summary

Hide

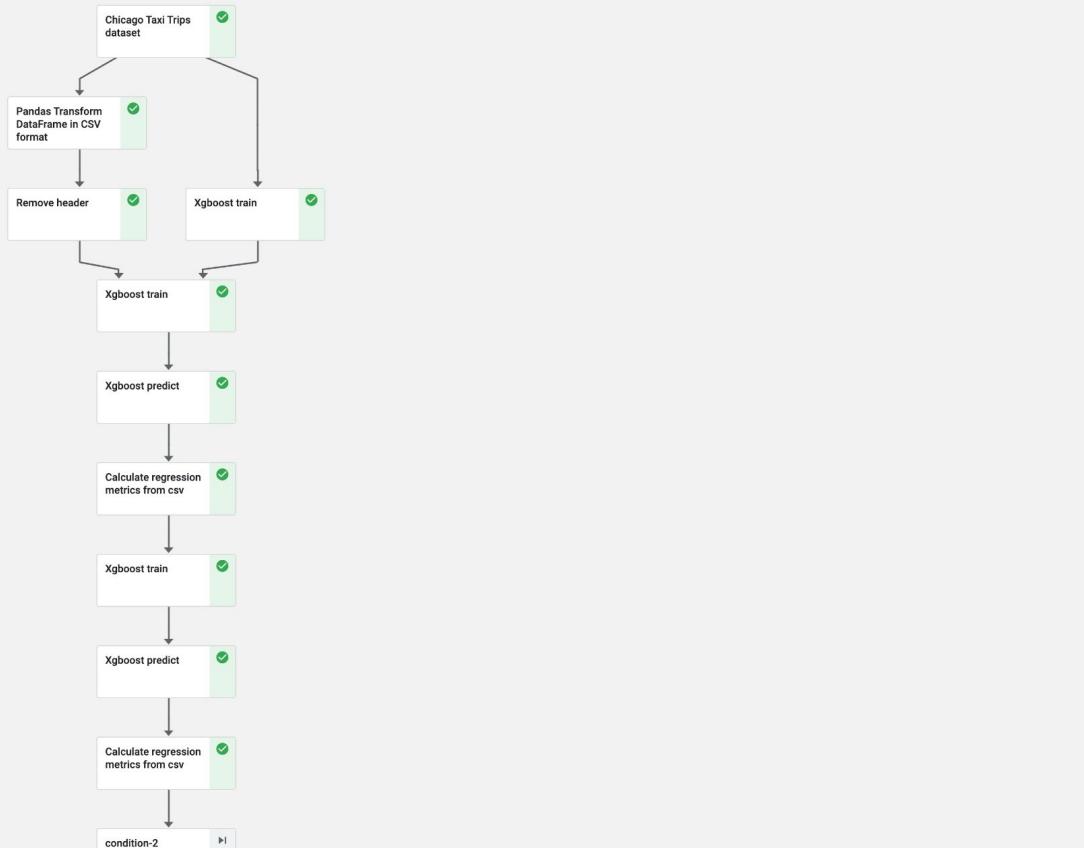
ID
2a56d15a-7680-4c01-a4fb-1b8a52a8de94Version
[Demo] XGBoost - Iterative model training ▾

Version source

Uploaded on
1.7.2021, 08:36:13Description
source code This sample demonstrates iterative training using a train-eval-check recursive loop. The main pipeline trains the initial model and then gradually trains the model some more until the model evaluation metrics are good enough.

Static pipeline graph

← ✓ Run of [Demo] XGBoost - Iterative model training (025a6)

[Graph](#) Run output Config[Simplify Graph](#)

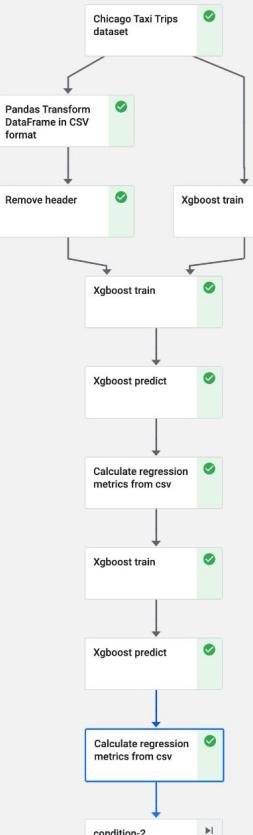
ⓘ Runtime execution graph. Only steps that are currently running or have already completed are shown.

- [Home](#)
- [Notebooks](#)
- [Tensorboards](#)
- [Models](#)
- [Schemas](#)
- [Volumes](#)
- [Experiments \(AutoML\)](#)
- [Experiments \(KFP\)](#)
- [Pipelines](#)
- [Runs](#)
- [Recurring Runs](#)
- [Artifacts](#)
- [Executions](#)
-
- [Manage Contributors](#)
-
- [GitHub](#)
-
- [Documentation](#)

Privacy + Usage Reporting

Graph
Run output
Config

Simplify Graph


Input/Output
Visualizations
ML Metadata
Details
Volumes
Logs
Pod
Events

train-until-good-pipeline-fv8rx-3435799838

xgboost-predict-predictions

```

minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-419816986/xgboo
st-predict-predictions.tgz
-3.267633914947566e-02
2.5738757301330564e-03
3.151332378387451172e+00
-5.836945722171020508e-02
8.7031834304809570e-01
-4.532337188720703125e-03
2.230248348920492e-02
1.3318153886962606e-02
1.53213143348693847e-02
5.172791481018066406e+00
-2...
...

```

remove-header-table

```

minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-1917203908/reme
ve-header-table.tgz
0.0
0.0
3.35
0.0
1.0
0.0
0.0
0.0
0.0
0.0
5.0
0.0
0.0
0.0
0.0
0.0
0.0
3.0

```

Output parameters

mean_squared_error 0.00828073701765554

Output artifacts

max_absolute_error
minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-3435799838/calculate-regression-metrics-from-csv-max_absolute_error.tgz
0.627315139705078

mean_absolute_error
minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-3435799838/calculate-regression-metrics-from-csv-mean_absolute_error.tgz
0.05721542470884323

mean_squared_error
minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-3435799838/calculate-regression-metrics-from-csv-mean_squared_error.tgz
0.00828073701765554

root_mean_squared_error
minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-3435799838/calculate-regression-metrics-from-csv-root_mean_squared_error.tgz
0.0099885503059123

main-logs
minio://mlpipeline/artifacts/train-until-good-pipeline-fv8rx/train-until-good-pipeline-fv8rx-3435799838/main.log
WARNING: Running pip as the 'root' user can result in broken permissions and conflicting



Kubeflow

KubeFlow



TensorFlow
Extended



[Kubeflow on AWS](#)
[Kubeflow on Azure](#)
[Kubeflow on GCP](#)
[Kubeflow on IBM Cloud](#)
[Kubeflow Operator](#)
[Kubeflow on OpenShift](#)

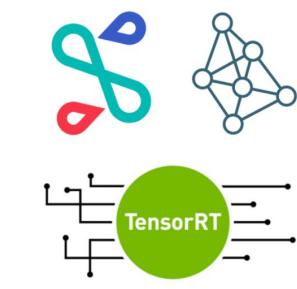
ML toolkit for Kubernetes



Notebooks



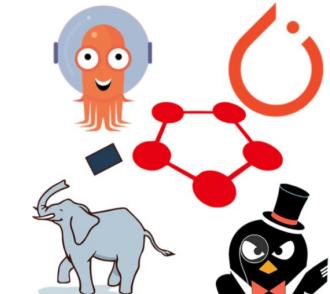
TensorFlow model training



Model serving



Pipelines



Multi-framework

Google's Vertex AI

Launched in May 2021



Google Cloud

Google Cloud Platform strikepose Search products and resources

Vertex AI

- Dashboard
- Datasets
- Features
- Labeling tasks
- Notebooks
- Pipelines
- Training
- Experiments
- Models
- Endpoints
- Batch predictions
- Metadata

Get started with Vertex AI

Vertex AI empowers machine learning developers, data scientists, and data engineers to take their projects from ideation to deployment, quickly and cost-effectively. [Learn more](#)

Region: us-central1 (Iowa)

Recent datasets:

- yogaposes (13 days ago)

+ CREATE DATASET

Recent models:

- yogaposes_202142173741 (12 days ago, Average precision: 0.9)

+ TRAIN NEW MODEL

Recent endpoints:

ONLINE TRAFFIC REQUESTS ERROR RATE

No data is available for the selected time frame.



ML Pipelines | wrap-up

By using a ML pipeline, you can:

- Automate your ML process, which lets you regularly retrain, evaluate, and deploy your model.
- Utilize distributed compute resources for processing large datasets and workloads.
- Increase the velocity of experimentation by running a pipeline with different sets of hyperparameters.

To learn more visit the following tutorials @:

<https://kirenz.github.io/>



MLOps tutorials on how to:

- Install TF and TFX
- Build your first TFX pipeline
- Install Kubeflow
- Build your first Kubeflow pipeline

Jan Kirenz

www.kirenz.com

