

Uncovering Project and Community Insights Using Data Driven Methods

Oindrilla Chatterjee
Senior Data Scientist
Red Hat

SCHEDULE

- 1 Metrics
- 2 Project AI4CI
- 3 Operate First Community Cloud
- 4 AI tools for open source projects
- 5 Time to Merge Model

Let's talk metrics

Why?

- Better allocate resources
- Evaluate a project's success
- Advocate for the project
- Analyze the development community
- Growth and journey of community

How?

- Sources like repositories, communication channels, website, service endpoints
- Open end-to-end workflows, from collection to deployment

ML to aid project development?

- Predicting time-to-merge to get an idea of "estimated effort"
- Predicting time to respond on issues
- Optimal stopping point for long running tests



Not here to tell you

Which metrics to track

Why use metrics to support your open source community



Here to tell you

How AI driven metrics can help

How to use open source ML tools to achieve this

How this can be done on an open community cloud

Thursday, September 15 • 11:00 - 11:40

☐ Community Data: What to Measure and Why? - Cali Dolfi, Red Hat

Tuesday, September 13 • 09:55 - 10:35

☐ Dev Team Metrics that Matter - Avishag Sahar, LinearB

AI4CI: Open Source AIOps toolkit

<https://github.com/aicoe-aiops/ocp-ci-analysis>

Problem

- Need for **AIOps** - Automated monitoring, analysis, alerting with Ops (CI/CD, development processes)
- **Open Source data** originating from real world production systems is a rarity for public datasets.
- Lack of AI driven metrics for open source community health.

Opportunity

- **Open operations data** made available by running open source software and applications in production.
- Data includes CI/CD data, code, telemetry, logs, operational dashboards.
- Eg: Kubernetes testing infrastructure, Fedora make their testing data available open source.

Solution

- Collection of intelligent and open source **data science tools** to collect and analyze the CI/CD data.
- **AIOps models** like Github time-to-merge service, optimal stopping time prediction, build log classifier
- KPI and Metric dashboards
- Goal is to foster an open source AIOps community with open ops data, AI tools and services.

AI4CI supports CI/CD and software dev processes

What is AI4CI?

Collection of **Open Source AIOps tools** including scripts, notebooks, pipelines, dashboards and data sources.



Data collection

Collection of open operations data from Kubernetes testing platforms eg: Testgrid, Github, and Prow.



Metrics

Collects metrics and **KPIs** and visualization dashboards.



ML Services

ML services which can support CI/CD processes.



Open source AIOps template

Resource for open source AIOps communities (notebooks, scripts, automated ML pipelines, dashboards, services tools)


The Operate First Community Cloud

- Operate First makes **operations open source**.
- An initiative centered around learning and developing code and practices in an open **production community cloud**.
- Deploy and maintain apps in an open environment leading to **open operations data** which include logs, issues, metrics.




www.operate-first.cloud/

Open Source data science and engineering tools on Operate First

 **OPEN DATA HUB**
All Platforms powered by Open Source


OPERATE FIRST >>

Show All ▾




JupyterHub

A multi-user version of the notebook designed for companies, classrooms and research labs




Argo

Kubernetes native workflows, events, CI and CD




Superset

A modern, enterprise-ready business intelligence web application




Prometheus

Systems monitoring and alerting toolkit




Grafana

Visualization and analytics software




Spark

Unified analytics engine for large-scale data processing




Seldon

Platform for rapidly deploying machine learning models on Kubernetes.




Kafka

Distributed event streaming platform



Airflow

Platform to programmatically author, schedule, and monitor workflows

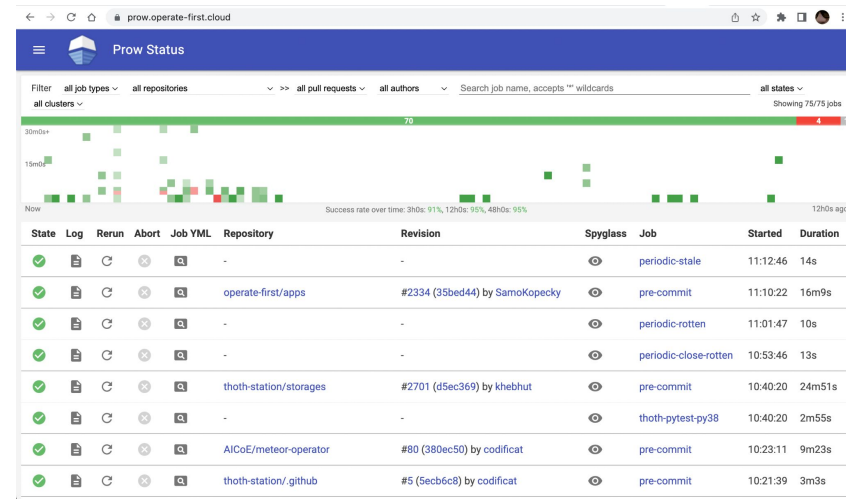


Hue

Data exploration platform for Hive and S3 storage

Open Source data science and engineering tools on Operate First

- Operations Data like: Logs, Metrics, Github issues, PRs, architectural decisions, blueprints.
- These open operations datasets can help **enable AI tools** to assist with cloud operations.



Tooling for capturing metrics



Code Repository



Github API,
Thoth MI-Scheduler,
CHAOSS-Augur



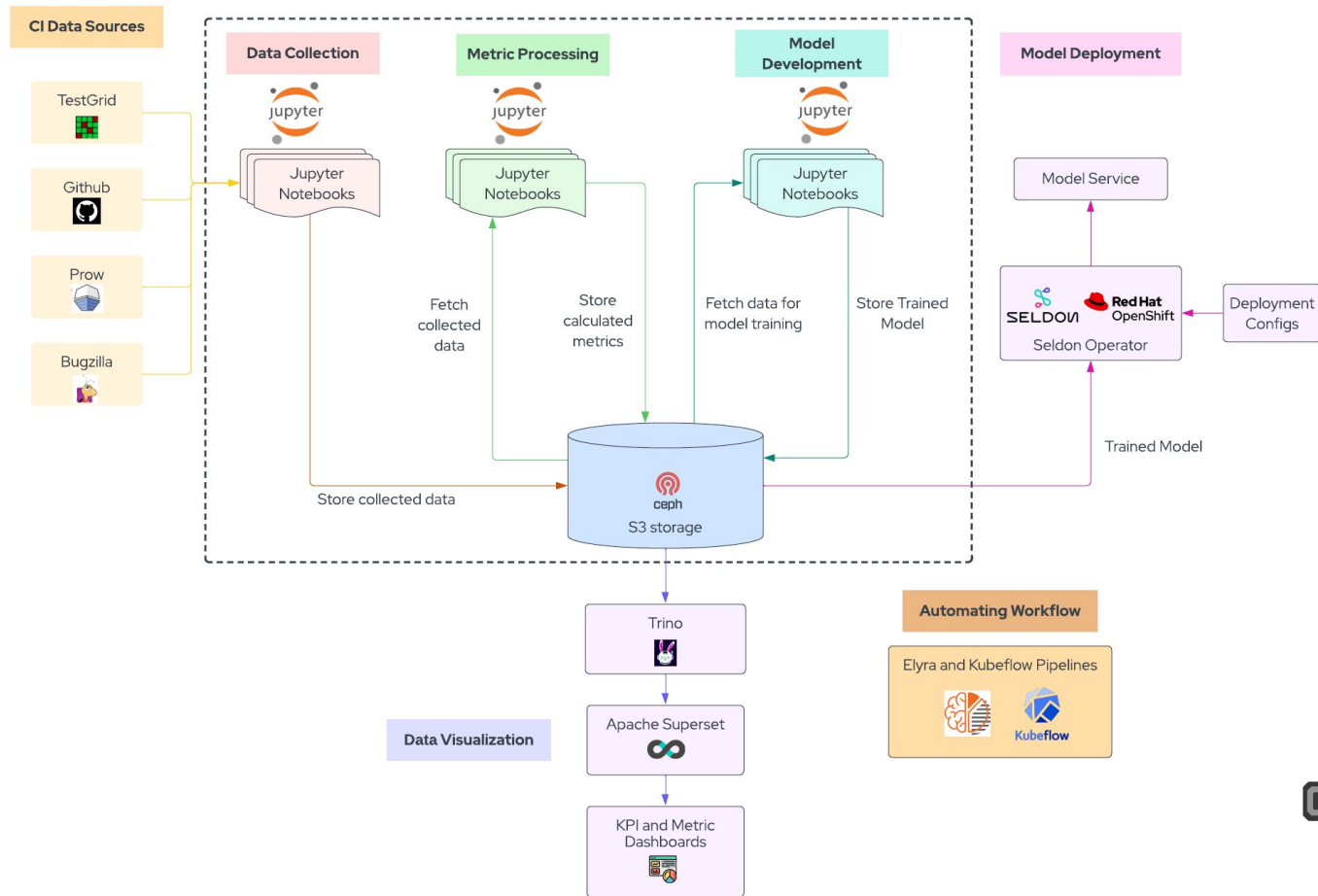
Feature exploration
& Engineering



Dashboard with key
metrics



AI4CI Architecture



Running on



OPEN DATA HUB
AI Platform powered by Open Source

OPERATE
FIRST>>

ML Service: Optimal Stopping Point Prediction



Sometimes tests/builds take longer than expected to run

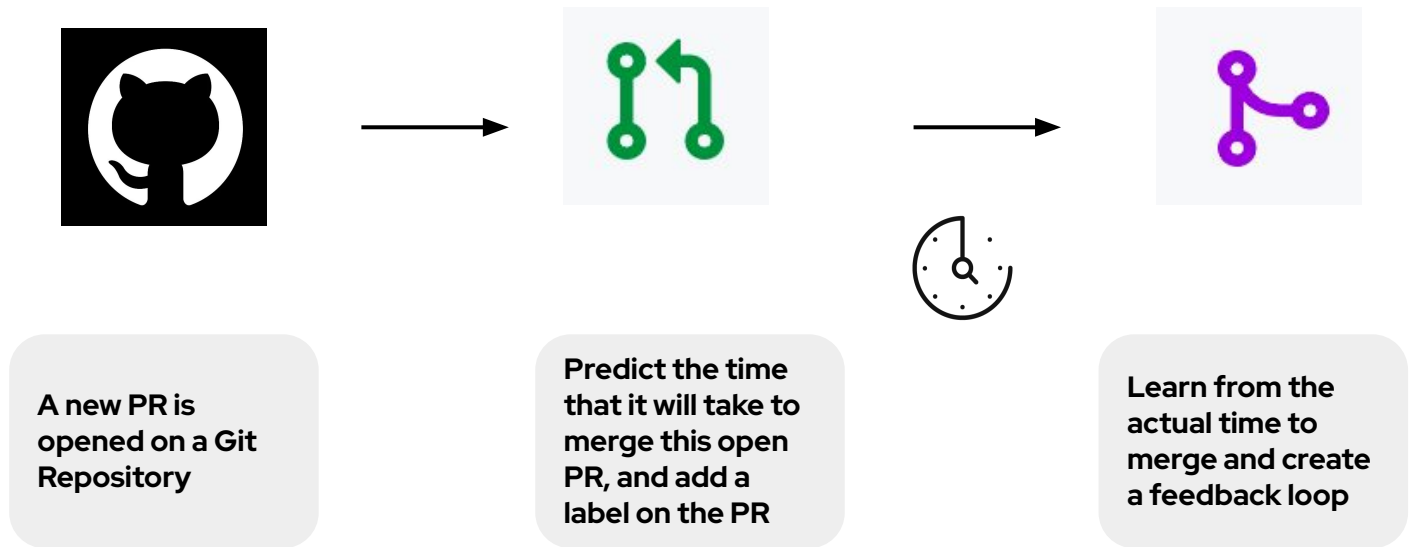


Find an Optimal Stopping Point after which the test will fail.



We can better allocate and save resources.

ML Service: Time to Merge Prediction

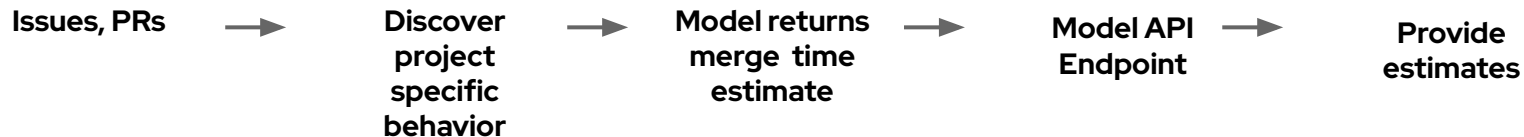
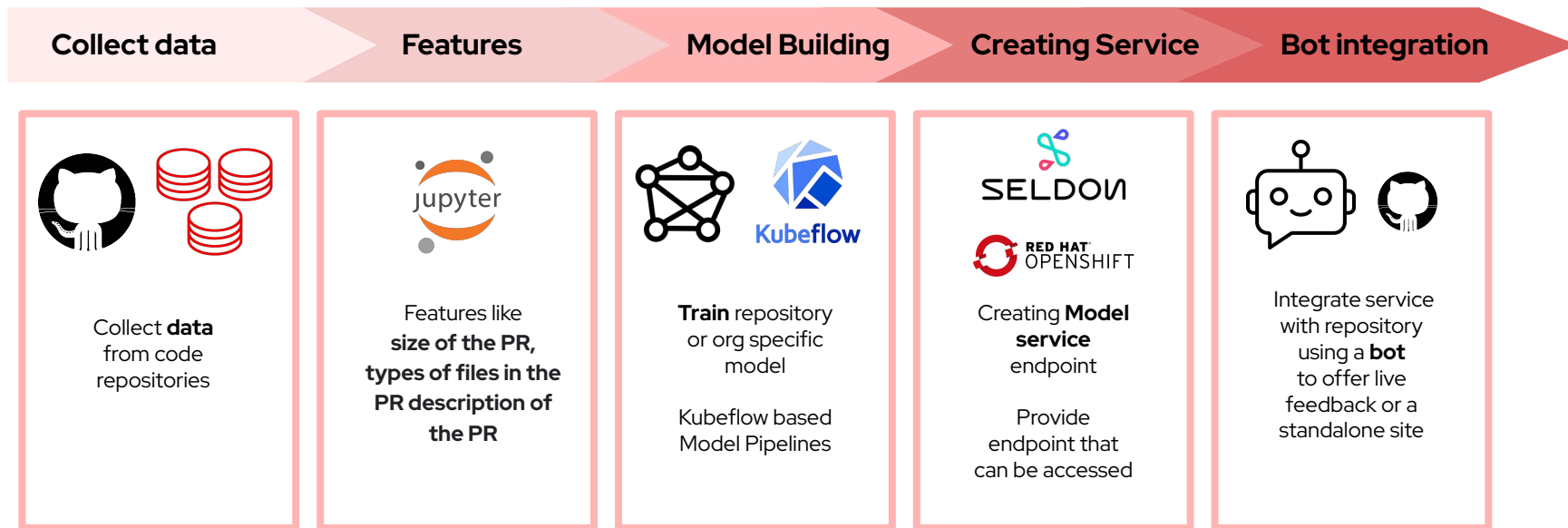


Time to merge prediction service for community health



- Identify **bottlenecks** in development process
- Leverage the rich **historical data** of consisting of Issues, Commits, PRs
- Give **new contributors** of an estimate of when their PR will be reacted upon
- Most importantly - develop an **AI driven mindset** for community health

Current workflow: Github time to merge prediction service



ML Service: Github Time to Merge Model

COLLECT DATA

 [openshift / origin](#)

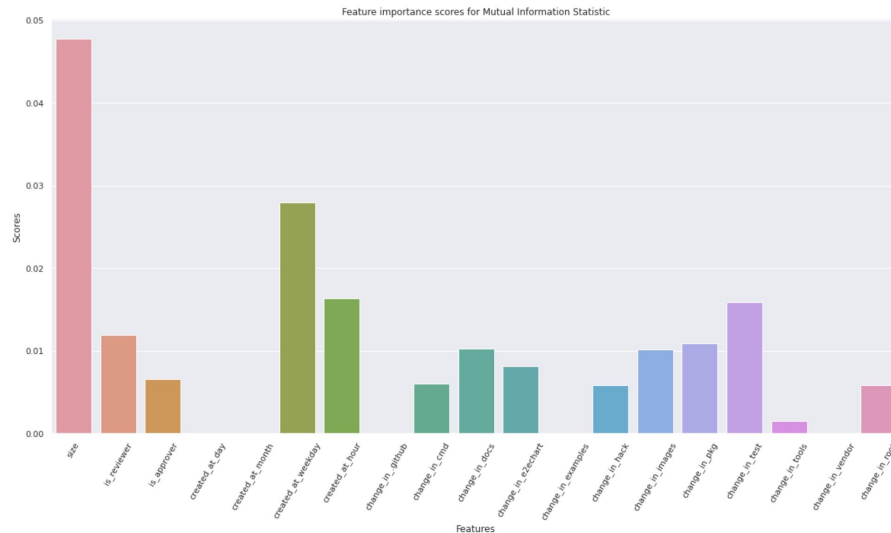
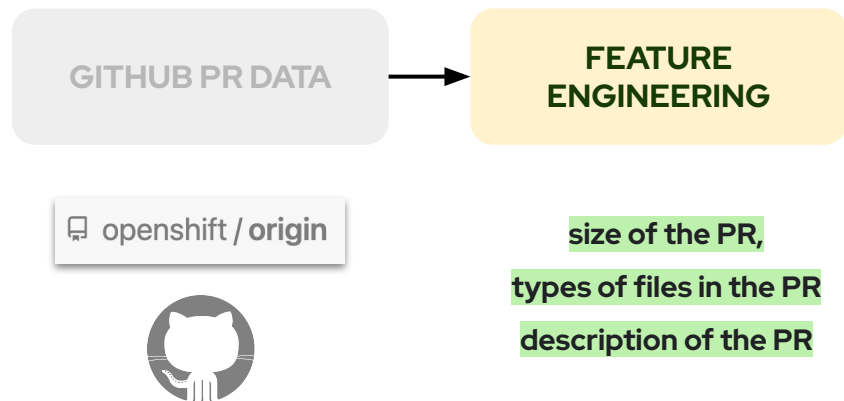


srcopsmetrics 2.11.1

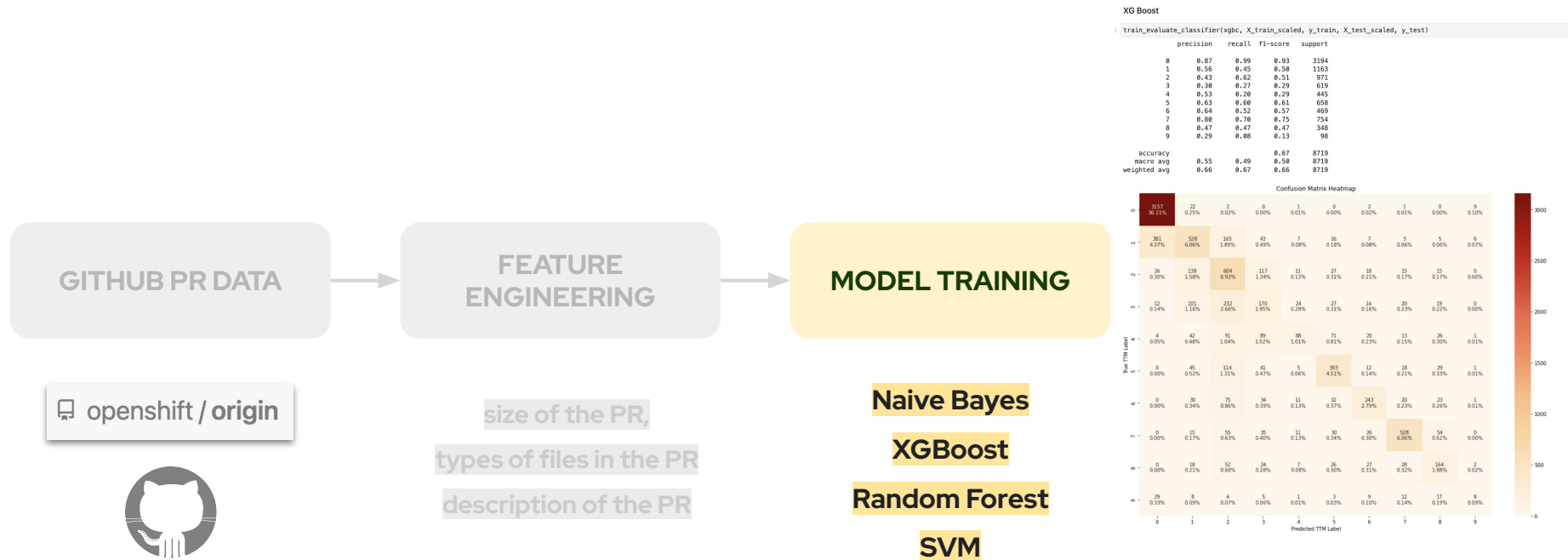
```
pip install srcopsmetrics
```



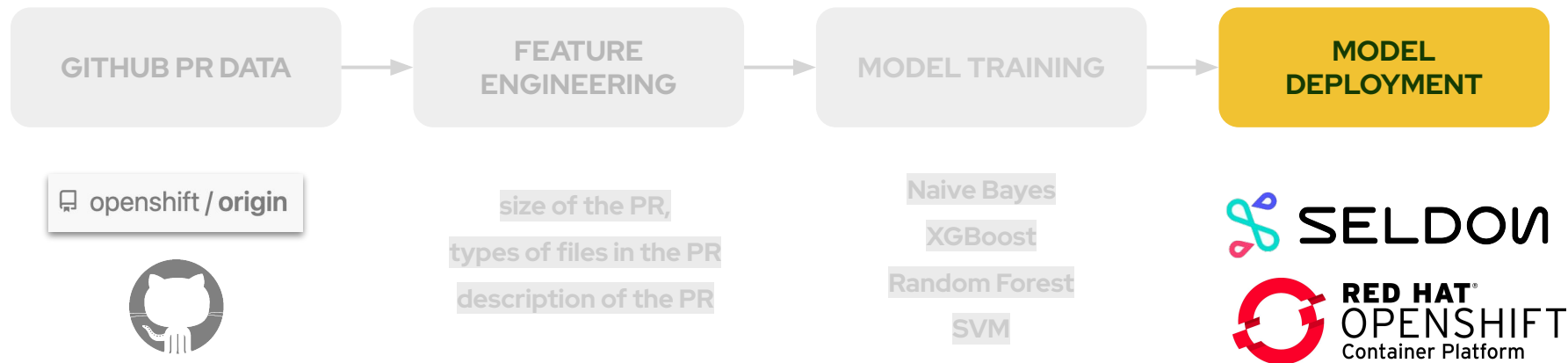
ML Service: Github Time to Merge Model



ML Service: Github Time to Merge Model





ML Service: Github Time to Merge Model



ML Service: Github Time to Merge Model


Services > Service details

 **github-pr-ttm-seldon-github-pr-ttm-predictor**
Managed by  [github-pr-ttm-seldon](#)

[Details](#) [YAML](#) [Pods](#)


Service details

Name
github-pr-ttm-seldon-github-pr-ttm-predictor

Namespace
 [ds-ml-workflows-ws](#)

Labels [Edit](#)

- app.kubernetes.io/managed-by=seldon-core
- seldon-app=github-pr-ttm-seldon-github-pr-ttm-predictor
- seldon-deployment-id=github-pr-ttm-seldon

Pod selector
 seldon-app=github-pr-ttm-seldon-github-pr-ttm-predictor

Annotations
[1 annotation](#)

**MODEL
DEPLOYMENT**



NEXT STEPS

- Service as a bot on Github PRs
- Live feedback from bot / service
- Iterate on time-to-merge models for better performance
- Toolified API for using the training service.

Resources

<https://github.com/aicoe-aiops/ocp-ci-analysis/blob/master/docs/get-started.md>



Open Data
Sources



Notebooks



Dashboards



Model Endpoints



Automated
Workflows




Video Playlist



<https://tinyurl.com/aiforci>

Thank you!

@oindrilla_chat 

www.linkedin.com/in/oindrilla-chatterjee/ 