# Summary

This lab consists of:

- Example model ( `train.py` , `validate.py` )
- MLFlow/project configuration ( `MLproject` , `python_env.yaml` , `requirements.*.txt` )
- Service to deply model as a REST API ( `run.py` )
- Model configuration ( `params.env` )
- Docker configuration for model deployment ( `Dockerfile` )
- Docker configuration for MLFlow server deployment ( `mlflow/Dockerfile` )
- Pipeline configuration for infrastructure and model deployment ( `.github/workflows/` )
- Kubernetes configuration for model deployment and autoscaling ( `k8s/*` )
- Terraform infrastructure configuration ( `terraform/` )

# Setup

Setup environment and install dependencies for mlflow:
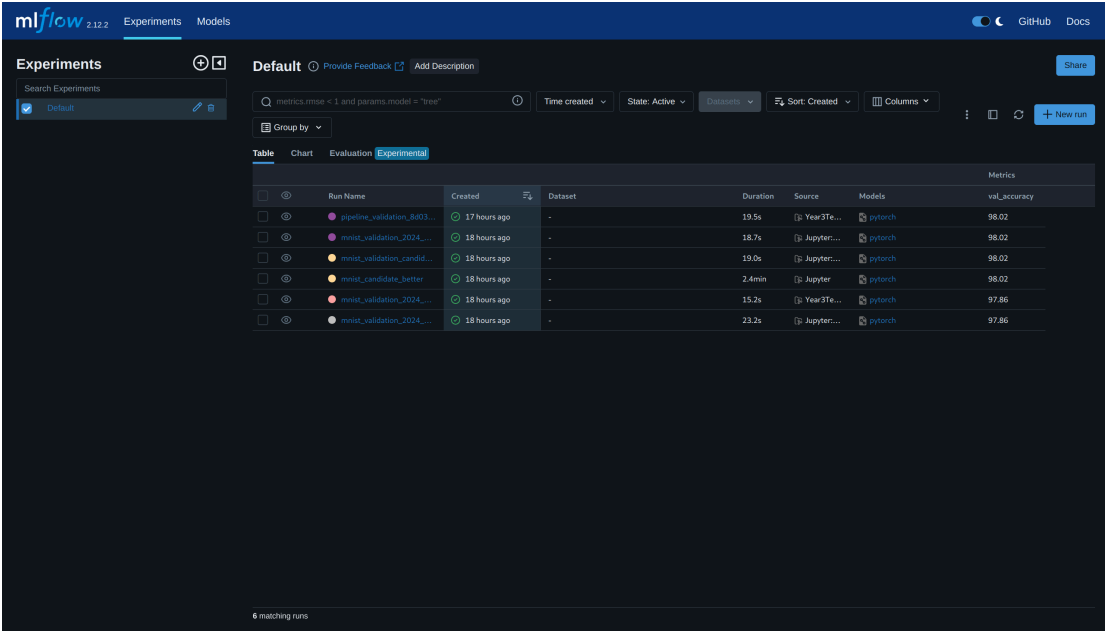
```
In [ ]:   ! python3.8 -m venv .venv
          ! source .venv/bin/activate
          ! pip install -r requirements.mlflow.txt
```
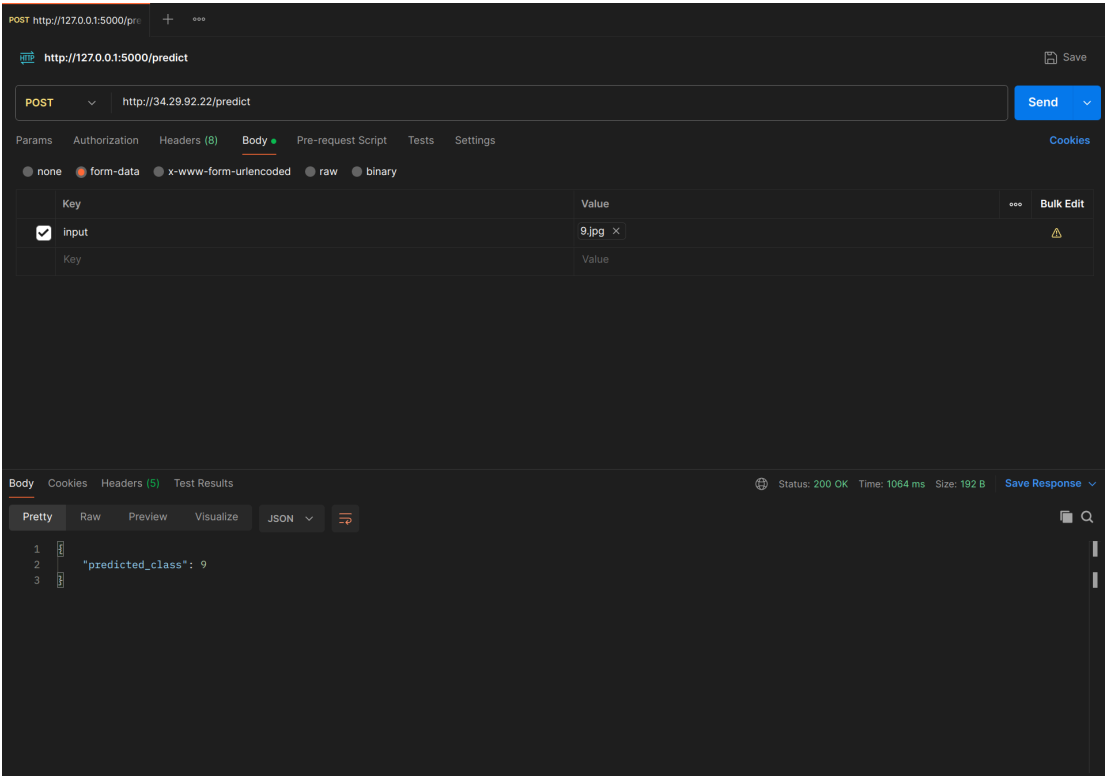
# Try it

As you develop the model, you can use MLFlow to store metrics and artifacts, using command:

```
In [ ]:   ! MLFLOW_TRACKING_URI='https://mlflow-y3omahtamq-uc.a.run.app' mlflow run
```

Then visit the link in `MLFLOW_TRACKING_URI` , to see the effect:

The deployed model can also be accessed, current address is
`http://34.29.92.22/` . `POST` requests to `http://34.29.92.22/predict`
with request body parameter `input` as an image of a number will yield a
prediction. Example:

# How it works

## MLFlow commands

The project has several options for commands: `main`, `validate` and `deploy`. `main` trains and validates the model, `validate` performs validation only, and `deploy` builds a docker image for the model to run on the cloud as a REST service. Use `main`, or just `mlflow run .`, to develop the model, `validate` to test the model, and `deploy` to ship model to production.

## Pipeline

Used in this project is GitHub Actions. The following steps are run:

- Check code formatting
- Test/validate model
- Push of metrics and artifacts (model, ROC graphs for each class) into the cloud-hosted MLFlow server
- Build model as a REST API, and package it as a docker container
- Push docker container into the artifact registry
- Apply updated Tearrform configuration, which includes kubernetes cluster, to update it to use newly built docker image

Nuances of implementation:

- Variables like mysql user and password and github credentials ares stored using GitHub secret manager.

# Terraform layout

- Cloud storage bucket for syncing terraform state
- Managed kubernetes cluster for model deployment with horizontal autoscaling and load balancing.
- MySQL instance to store MLFlow data
- Cloud storage to store MLFlow artifacts
- IAM configuration for GitHub actions, MLFlow, Deployed instance and Kubernetes cluster
- Cloud run deployment for hosting MLFlow server docker image
- Artifact registry configuration to store REST model services
- Firewall configuration for access to MySQL database

# Artifact tracking

During train runs, model learning curve of loss, and its parameters are tracked. During validation, models validation accuracy and ROC curves for each class are tracked. For all runs, model is stored to artifact registry.