

ML Challenge – ACID LAB

Technical Manual

General Description

Project Source: <https://github.com/jossteeven/acid-labs-devops-challenge>

In this project you will find the following file structure:

- **.github/workflows/**: Github Actions Workflow pipelines
 - [main.yaml -> for the app implementation process]
 - [infrastructure_pipeline.yaml -> only for the infrastructure implementation, in this case Heroku]
- **app/**: Directory with all the files necessary to build the REST API app. Python app and ML model can be found here.
- **datasets/**: Group of datasets to train and test the model in use.
- **infrastructure/**: Terraform files to handle the Heroku app implementation.
- **.gitignore**: Git .gitignore to exclude not needed temporal files.
- **Dockerfile**: Dockerfile needed to build the app image to deploy to Heroku app.
- **model.ipynb**: Model Jupyter Notebook

Used Tools

Heroku: Cloud used to create a web app instance to expose REST API.

Docker: To test and build the image generated to run into the Heroku app instance.

Flask: Framework to load the .pkl model as a REST API.

Terraform: To handle the Heroku infrastructure implementation.

GitHub/GitHub Actions: To host the project and create CI/CD pipeline to automate the REST API deployment.

Project Implementation

REST API

To expose the pickle_model.pkl model as a REST API Flask was used as framework. This project can be found in the /app folder of this project.

You can refer to the code annotation marked with ## to get a little more of detail of the code flow.

This REST API is expose in a Heroku web app as shown below:

This Web app redirects the form inputs to the route **/predict** that is seclusively used for the web UI to predict the flight delay.

For a quick test of the REST API model the **/predict_api** endpoint was expose so the usage of any CLI will be great to have the prediction values. An example usage of this endpoint is shown below.

```

1 {
2   "data": [1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1]
3 }

```

```

1 {
2   "data": [1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1]
3 }

```

To pass the data array related to each of the 37 parameters to get a prediction you must pass as body a json as this:

```

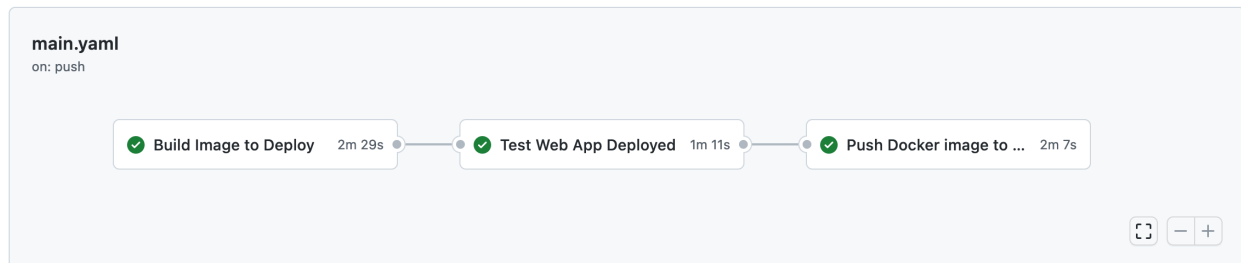
{
  "data": [1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1]
}

```

Pipeline

To be able to execute the automatic deployment from the project hosted in GitHub, a pipeline in GitHub Actions was built.

This pipeline has three stages as follows:



- **Build Image to Deploy.** - Stage responsible of generating the Docker image based on the Dockerfile you will find on the `/app` folder. This Dockerfile helps prepare the container environment with the necessary packages and resource to raise the REST API.
- **Test Web App Deployed.** – Stage responsible to generate the HTTP Benchmark as required to the challenge.

The screenshot shows the logs for the `Test Web App Deployed` job, which succeeded 2 hours ago in 1m 11s. The logs are displayed in a dark theme with a search bar at the top right.

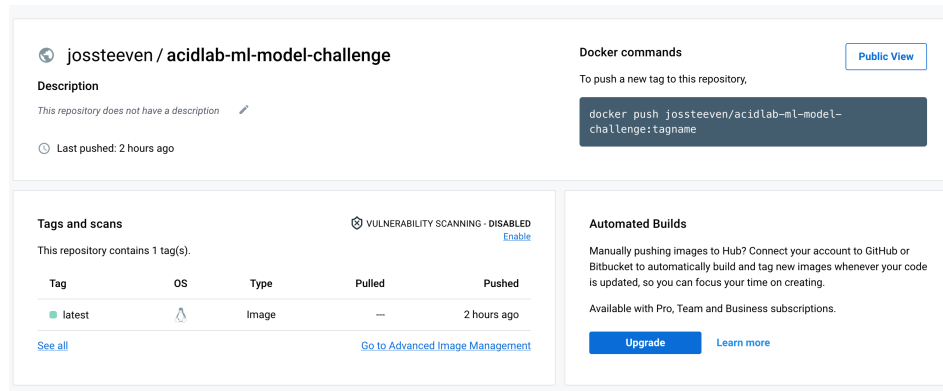
The job steps are:

- Set up job (0s)
- Initialize containers (2s)
- Latency Test - HTTP Benchmark (1m 7s)
- Stop containers (0s)
- Complete job (0s)

The `Latency Test - HTTP Benchmark` step contains the following log output:

```
1 Run (wrk -t12 -c50000 -d45s https://acid-lab-app.herokuapp.com/)
4 Running 45s test @ https://acid-lab-app.herokuapp.com/
5 12 threads and 50000 connections
6 Thread Stats Avg Stdev Max +/- Stdev
7 Latency 37.37ms 27.95ms 100.18ms 63.33%
8 Req/Sec 184.00 29.92 215.00 50.00%
9 90 requests in 1.11m, 69.47KB read
10 Socket errors: connect 21830, read 71, write 0, timeout 0
11 Non-2xx or 3xx responses: 90
12 Requests/sec: 1.35
13 Transfer/sec: 1.04KB
```

- **Push Docker Image to Heroku.** – Stage to generate and store the Docker image used to deploy to Heroku app. This Image can be pull from DockerHub as `"docker pull jossteeven/acidlab-ml-model-challenge:latest"`



```
wrk.method = "POST"  
wrk.body = "{\"data\": [1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1]}"  
wrk.headers["Content-Type"] = "application/json"
```

This configuration allows us to change the requests configuration and include a body with this array to generate the 50000 requests for test.

The Dockerfile used to generate this image can we found at test/Dockerfile.

(Spanish) Extra Questions:

1. ¿Cómo podrías mejorar el performance de las pruebas anteriores?

Se podría controlar dinámicamente la ejecución de cada hilo asincrónicamente y asociarlos en base a afinidad de cpu con el comando taskset y permitir así la ejecución paralela de miles de request.

Referencia: <https://github.com/wg/wrk/issues/261#issuecomment-284027586>

2. ¿Cuáles serían los mecanismos ideales para que sólo sistemas autorizados puedan acceder a esta API? (NO es necesario implementarlo).

Existen varios mecanismos que podrían limitar el acceso o el consumo de dicho API, entre ellos agregar un header de autenticación cuyo valor solo tengan las personas autorizadas y lo implementen al hacer llamados al API endpoint. La validación se realizaría en un API Proxy que podría estar ubicado en frente del API REST propio del modelo expuesto. Así también se podría enmascarar el endpoint de esta API REST desarrollada a través de API Proxies SaS como Apigee, configurando un proxy con Access Control para Ips específicas y agregar ciertas validaciones adicionales como el header adicional ya mencionado o algún tipo de autenticación basado en user/password.

- c. ¿Este mecanismo agrega algún nivel de latencia al consumidor? ¿Por qué?

El utilizar un Api Proxy por delante de una REST API agrega latencia de alguna forma, pudiendo llegar a afectar la calidad de experiencia de usuario. Esta latencia, aunque mínima, puede ser crítica para ciertas aplicaciones y es agregada porque la arquitectura presentaría ahora un nuevo filtro por el que los requests deben pasar. Para acceder al API REST (backend) deberán ahora pasar por las validaciones del API Proxy, agregando así latencia en el response.

- d. ¿Cuáles serían los SLIs y SLOs que definirías y por qué?

Ya que los SLO se relacionan al equipo de trabajo, aquellos que plantearía son los siguientes:

- Mejorar el rendimiento de respuesta del modelo ante la cantidad de peticiones distribuidas globalmente al API REST.
- Disminuir el tiempo de indisponibilidad del API REST controlando los recursos de las instancias que conforman su arquitectura.
- Incrementar la eficiencia y precisión del modelo de predicción realizando nuevos entrenamientos y por ende nuevos despliegues de código.

Y los SLI asociados para cada SLO respectivamente serian:

- Response Time
- Web App Container Health Metric over time
- Parámetros de confianza del modelo predictive y de las predicciones que está arrojando en producción.
-

Definiría estas métricas ya que al tratarse de un modelo predictivo y que posiblemente sea usado en una página web distribuida y de consumo masivo, el mantener el API REST activo 24/7 y con una tolerancia a grandes cantidades de tráfico es indispensable por el tipo de plataformas que consumirían esta API. Mencionando ejemplos, estas plataformas podrían ser sistemas de tráfico aéreo comercial o relacionadas. Así también, considero importante que el mantener el modelo lo más preciso posible permitirá tener un sistema aún más exacto, por lo que su constante mejora es necesaria a través de nuevos entrenamientos y actualizaciones de código.