

Running Algorithms Locally and Remotely via Docker

Your computing committee

Contents

1	Introduction	2
2	Running docker locally in your computer	3
2.1	Install docker	3
2.2	Configure a docker image	3
2.3	Building a docker container	3
2.4	Running an algorithm in the docker container	4
3	Remotely running one instance of the algorithm	5
3.1	Configure the server	5
3.2	Copying your files to the remote server	5
3.3	Remote execution and monitoring	5
4	Remotely running multiple instances of your algorithm	6
4.1	Configure the servers	6
4.2	Remote execution and monitoring	7
4.3	Resetting everything	7

1 Introduction

This document explains the current software used to run locally and remotely algorithms inside Docker. The focus is on explaining the salient points of different scripts so that you can modify them to suit your algorithm. We won't explain the technicalities involved in using `docker`, `tmux` and `ssh`. For a basic introduction to those tools, the reader is referred to the wiki maintained by the statistics department at stat.cmu.edu/computewiki.

The example project is a simplified version of Sparse Information Filter for Fast Gaussian Process Regression (Kania et al. 2021), which proposes a method for training sparse Gaussian processes via mini-batches rather than full-batches.

The structure of the folder that contains this file is as follows

- `Dockerfile` configures a docker image. At a basic level, it is a declaration of the environment where the algorithm will run.
- `DockerfileGPU` configures a docker image with GPU support.
- `.dockerignore` configures exceptions when creating the docker image
- `rundocker.sh` runs the algorithm inside the container.
- `src` contains the source code of the algorithm
- `data` contains the data used for the experiments
- `results` is the folder where the results of the experiments are saved
- `utilities` contains useful scripts for remotely running algorithms
- `README` contains the files needed to generate this document
- `paper.pdf` The paper whose source code is used in this project. The algorithms implemented inside the `src` folder are explained in it. `supp.pdf` contains additional information.

Many of the code snippets shown in this document have been substantially simplified to reflect the essence of the scripts. However, the scripts in the root folder differ since they target a specific algorithm rather than a generic one.

2 Running docker locally in your computer

2.1 Install docker

Get the latest version of docker desktop from docs.docker.com/desktop. Once the installation finishes, please start docker desktop.

2.2 Configure a docker image

`Dockerfile` configures the docker image, that is, the environment containing all the dependencies necessary for running your algorithm.

We start by declaring a base image on which we will install further dependencies. For example, the following code provides a `linux` system with `python 3.8` support.

```
1 FROM python:3.8-slim-buster
```

Consult the `Dockerfile` for further instructions on how to use a base image with `Pytorch` and `Tensroflow` with GPU support.

On it, we might install further dependencies via `pip`. It is recommended to specify the version of the required dependencies. Otherwise, `pip` will install the latest version compatible with our other packages making the docker image non-deterministic.

```
1 # We install additional packages required for this particular project
2 RUN pip install pandas==2.0.3
```

Then, we copy the code required to run the algorithm and experiments

```
1 # We copy the folder containing the source code of our algorithm
2 # into the folder /program/src in the container
3 COPY src /program/src
```

We do not copy the data or results folders into the container. They will be mounted to the Docker container during runtime.

2.3 Building a docker container

Our next step is to start a Docker container based on our image. The script `rundocker.sh` simplifies the building procedure. To wit, it

- creates a docker container based on the image defined in `Dockerfile`
- mounts the data folder instead of copying it into the container. This avoids making images unnecessarily large due to large datasets.
- mounts the results folder instead of copying it into the container. Therefore, if the computer crashes, the saved data will remain after reboot.

- runs an algorithm inside the container.

The most important part of the script is indicating the location of your algorithm's executable. Everything else in the file might be left untouched if you follow the conventions adopted by this tutorial.

```
1 #####
2 # Location of the executable inside the container
3 #####
4 readonly executable=/program/experiments/runner.py
```

2.4 Running an algorithm in the docker container

To run your container, execute the script `rundocker.sh`. You will immediately see the output of our algorithm.

```
1 bash rundocker.sh
```

If you wish to pass arguments to your executable, you can do so by running

```
1 bash rundocker.sh your_arguments
```

Then, you can then parse `your_arguments` inside your script.

3 Remotely running one instance of the algorithm

3.1 Configure the server

A ssh connection to the server can be configured at `/.ssh/config`. For example, using the following code, you can configure access to `hydra1.stat.cmu.edu`.

```
1 Host hydra1
2   Hostname hydra1.stat.cmu.edu
3   Port 22
4   user your_username
5   IdentityFile ~/.ssh/your_public_key_file
6   ServerAliveInterval 180
```

Thus, we execute

```
1 ssh hydra1
```

in the terminal, `ssh` will automatically connect to `hydra1.stat.cmu.edu` using the public key `/.ssh/your_public_key_file`.

3.2 Copying your files to the remote server

We can copy all the files in this folder to `hydra1` using `scp`

```
1 scp -r $PWD hydra1:~/remote/sparsegp
```

3.3 Remote execution and monitoring

After running the command in the previous section, `rundocker.sh` will be located at `/remote/rundocker.sh`. You can use the utility `distribute.sh` to connect to the remote server and run your `rundocker.sh` script.

```
1 bash distribute.sh --server hydra1 --session TEST --rundocker "~/remote/
  sparsegp"
```

`distribute.sh` connects to the `hydra1` server, and runs the `rundocker.sh` script inside a `tmux` session called `TEST`.

To observe the output of our algorithm, we can use the `process.sh` utility. Specify the `tmux` session name and the server as follows.

```
1 bash process.sh --session TEST --server hydra1
```

4 Remotely running multiple instances of your algorithm

You can run one instance per server of your algorithm with different parameters in each server. This is useful if you run the same experiments multiple times with different parameters.

4.1 Configure the servers

We configure four servers at `/.ssh/config`, two ghidorahs and two hydras.

```
1 Host ghidorah2
2   Hostname ghidorah2.stat.cmu.edu
3   Port 22
4   user your_username
5   IdentityFile ~/.ssh/your_public_key_file
6   ServerAliveInterval 180
7 Host ghidorah3
8   Hostname ghidorah3.stat.cmu.edu
9   Port 22
10  user your_username
11  IdentityFile ~/.ssh/your_public_key_file
12  ServerAliveInterval 180
13 Host hydra1
14   Hostname hydra1.stat.cmu.edu
15   Port 22
16   user your_username
17   IdentityFile ~/.ssh/your_public_key_file
18   ServerAliveInterval 180
19 Host hydra2
20   Hostname hydra2.stat.cmu.edu
21   Port 22
22   user your_username
23   IdentityFile ~/.ssh/your_public_key_file
24   ServerAliveInterval 180
```

Subsequently, we specify the servers where you want to run the instances of your algorithm in `servers.sh`

```
1 # Declare servers to be used for running your scripts
2 servers=(
3   ghidorah2
4   ghidorah3
5   hydra1
6   hydra2
7 )
```

4.2 Remote execution and monitoring

The process is analogous to running one instance, but we must iterate over the declared servers in `server.sh`. The utility `experiments.sh` gives one example of how to achieve this. In short, we load the server configuration by executing

```
1 #####
2 # load servers
3 #####
4 . servers.sh
```

Then, we declare an array with the parameters that we want to pass to the algorithm

```
1 #####
2 # Declare parameters that change from server to server
3 #####
4 parameters=(
5   param_for_ghidorah2
6   param_for_ghidorah3
7   param_for_hydra1
8   param_for_hydra2
9 )
```

and iterate over the servers

```
1 #####
2 # We iterate over the servers and run the script
3 # with different parameters in each one of them
4 #####
5 for i in "${!parameters[@]}"; do
6   ./distribute.sh --server "${servers[i]}" \
7   --session PARAMETERS \
8   --rundocker "~/remote/" \
9   "${parameters[i]}" \
10 done
```

We can run the script via bash

```
1 bash experiments.sh
```

To monitor the four instances, we can execute

```
1 bash process.sh --session TEST
```

Note that we specify the session but not the server. The utility `process.sh` will iterate over the servers declared in `servers.sh` and connect to each session.

4.3 Resetting everything

If you want to stop the container in all the servers specified in `servers.sh`, execute

```
1 bash clean.sh
```