

BIRKBECK, UNIVERSITY OF LONDON

Pneumonia Detection from Chest X-Ray Images

Author:
Baran Buluttekın

Supervisor:
Dr. George Magoulas



*A project report submitted in fulfillment of the requirements
for the degree of MSc Data Science*

in the

Department of Computer Science

July 1, 2020

Declaration

I have read and understood the sections of plagiarism in the College Policy on assessment offences and confirm that the work is my own, with the work of others clearly acknowledged. I give my permission to submit my report to the plagiarism testing database that the College is using and test it using plagiarism detection software, search engines or meta- searching software.

The report may be freely copied and distributed provided the source is explicitly acknowledged.

Abstract

New test text. New words

Contents

Declaration	i
Abstract	ii
1 Introduction	1
1.1 Aims and Objectives	1
1.1.1 Objectives	1
1.2 CI/CD Pipeline	2
1.3 Engineering Related Challenges Specific to Machine Learning Systems	3
1.4 Project layout and design	3
1.5 Reproducibility Guidance	4
2 Neural Networks	5
2.1 Building Blocks of Neural Networks	5
2.2 Exploding and Vanishing Gradients	5
2.3 Optimization	5
2.4 Regularization and Over-fitting	5
3 Data	6
3.1 Data Augmentation	6
3.2 Limitations of the Dataset	6
3.3 Data Proccession	6
References	7

1 | Introduction

Medical diagnosis and specifically computer-aided diagnosis (CAD) is a hot topic in the field of technology. One of the main reasons for becoming a hot topic is the recent innovation and breakthroughs achieved by computer vision research. Combined with poor healthcare coverage around the globe, CAD systems offer a promising solution to mitigate the devastating impact of fatal diseases such as pneumonia. Controversial topics such as whether or not artificial intelligence will replace the radiologist in the future aside, these automated systems can offer answers for patient's questions in absence of medical help or to very least offer much needed second opinion in the face of unsatisfied diagnoses. Given all the mentioned possible benefits of the CAD systems, this project is focused on building classification CAD systems for diagnosing pneumonia from the chest X-ray images.

1.1 Aims and Objectives

The aim of this project is to build a fully functional chest X-ray image classification pipeline that implements CI/CD principals to experimentation and deployment.

1.1.1 Objectives

Project will be implemented with execution of following objectives:

- **Carrying out general data exploration:** This part involves general check on dataset.
- **Data pre-processing and augmentation:** Preparing the data for model ready state.
- **Building baseline model with well known neural network architectures:** This step involves setting additional benchmarks with out of the box models from section 4.
- **Using pre-trained network to increase model performance:** Using pre-trained networks to help training and accuracy of the model.
- **Visualizing neural network to ensure learning quality:** For making sure model learning as intended and focusing on correct parts of the image.

- **Model ensembling:** Using ensemble method with different neural network architectures.
- **Applying different deployment options:** Implementation of different deployment options. Based on their trade offs.

The intuition behind choosing these objectives instead of attempting to build novel architecture from scratch is the process of choosing architecture has a very large search space and requires a lot of iteration and experimentation. Due to the limited time frame of this project attempting new architecture would not be feasible.

1.2 CI/CD Pipeline

In this section, I will give a brief introduction to the CI/CD pipeline to explain what CI/CD is and why it is chosen as a preferred way to build this project.

Continuous integration (CI) is a workflow strategy that helps ensure everyone's changes will integrate with the current version of the project in the typical software engineering team. This lets members of the team catch bugs, reduce merge conflicts, and increase overall confidence that your software is working. While the details may vary depending on the development environment, most CI systems feature the same basic tools and processes. In most scenarios, a team will practice CI in conjunction with automated testing using a dedicated server or CI service. Whenever a developer adds new work to a branch, the server will automatically build and test the code to determine whether it works and can be integrated with the code on the main development branch. The CI server will produce output containing the results of the build and an indication of whether or not the branch passes all the requirements for integration into the main development branch. By exposing build and test information for every commit on every branch, CI paves the way for what's known as continuous delivery, or CD, as well as a related process, called continuous deployment. The difference between continuous delivery and continuous deployment is that CD is the practice of developing software in such a way that you could release it at any time. When coupled with CI, continuous delivery lets you develop features with modular code in more manageable increments. Continuous development is an extension of continuous delivery. It's a process that allows you to actually deploy newly developed features into production with confidence, and experience little if any, downtime. Even though the benefits of using CI/CD pipelines are more prominent in the software teams, integration automated testing will help even individual projects such as this by reducing time for debugging.

In more granular detail, this system works with central version control services and in this project central version control service used is Github. GitHub uses a communication tool called *webhooks* to send messages to external systems about activities and events that occurred in the project. For each event type, subscribers will receive messages related to the event. Generally, events refer to action involving the software such as new commit push, pull (merge) request, or other software

related action. In this case, whenever a new commit is pushed to any branch of the project, a message from Github will be sent out to a third party system called *travis*.¹ Travis is a hosted CI service that allows build and test software hosted in version control services. When travis receives the webhook call it will fetch the most recent version of the project and run the tests associated with it. When the test runs completed with the latest version of the software, test results will be sent back to relevant commits as status information using GitHub API. This information can either be used by developers for making decisions such as whether to accept the pull request versus reject it or if applicable can be used by service to initiate the deployment process for the software. In all cases, CI/CD will work as automation for software quality assurance process to speed up the development and improve the overall reliability of the software.

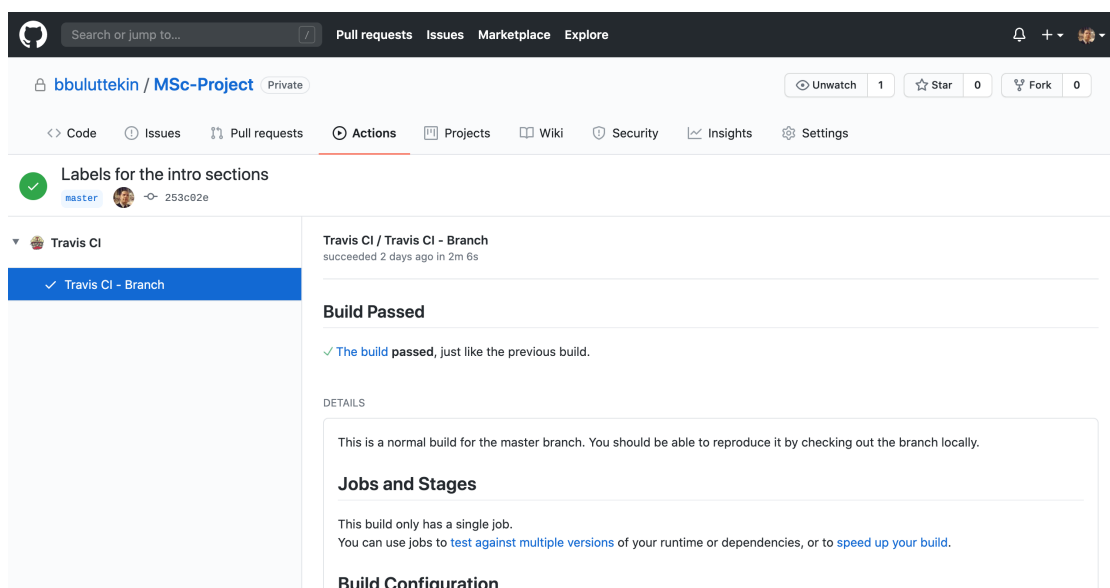


Figure 1.1: CI feedback received from Travis.

1.3 Engineering Related Challenges Specific to Machine Learning Systems

Role of hardware accelerators. Debugging related challenges of Neural Networks

1.4 Project layout and design

Overall project design such as folders and library interactions.

¹<https://travis-ci.org/>

1.5 Reproducibility Guidance

Reproducibility is a fundamental

As a scientific project, it is very important that anyone can reproduce the experiments and findings in this project to verify the conclusions reached are accurate. [1]

to reproduce the experiments and findings in this project to verify the conclusions reached are accurate.

Parts to consider before reproducing. Repo is not public Most runs carried out in colab notebooks should be positioned in same level as the src library running shell utility files requires kaggle api key

2 | Neural Networks

General history of artificial neural networks.

2.1 Building Blocks of Neural Networks

2.2 Exploding and Vanishing Gradients

2.3 Optimization

2.4 Regularization and Over-fitting

3 | Data

Reminder of the dataset. [2] Data set comprises of jpeg files stored in classification based folders.

3.1 Data Augmentation

Data augmentations applied and example results of the augmentation.

3.2 Limitations of the Dataset

Issues related to validation set size. Variance of the image resolution.

3.3 Data Proccession

Information about tf.data processing pipeline and advantages compare to other data feeds.

Talk about the decisions for image size and how it relates to information loss and preservation. Briefly mention the high variance of the image sizes and image resizing options and choice. Discussing different representation of the dataset will be covered.

References

- [1] Daniel Kermany and Kang Zhang. “*Labeled Optical Coherence Tomography (OCT) and Chest X-Ray Images for Classification*”. 2018. DOI: <http://dx.doi.org/10.17632/rscbjbr9sj.2>.
- [2] *Open Access Biomedical Image Search Engine*. <https://openi.nlm.nih.gov/>. Accessed: 2019-03-12.