# DevOps Project Report

-Aditya Ankush

-Shiloni Pipalia

-Suraj Singh

## 1. Introduction

### A. Why DevOps?

Adopting DevOps strategies leads to benefits such as:

- Shorter development cycles, Faster innovation
- Reduced deployment failures, rollbacks, and recovery time
- Improvement in communication leading to more collaboration
- Increase in efficiency
- Reduced cost and manpower

### B. About the Application

The Application is an ML based project for handwriting recognition of numeral characters using the MNIST dataset. The application is a web-based portal where the user can choose to train or test the model, upon training the trained model is saved on the user's system and is utilized for testing until overwritten by a new train call. The output is finally saved in the form of a csv file stored on the user's system.

Choose whether to train or test(IMG 1):



Upload Train/Test file(IMG 2):

# 2. Software Development Lifecycle

## A. Scope of the Project

The scope of the project is to build a ML based handwriting recognition tool and use DevOps tools for deploy, integrate, test and monitor. Our focus was on using DevOps tools. In order to achieve our goal, we have successfully used DevOps tools like GitHub, Jenkins, docker, Rundeck, Filebeat and ELK.

## B. Project architecture, workflows

We tried to build a completely automated pipeline. The stages included in the pipeline are:

- Clone/pull from GitHub.
- Build the docker image
- Push to Docker hub.
- Trigger Rundeck for deployment (More information present in the Rundeck section).
- ELK for monitoring

## C. SCM

The source control tool used as mentioned above is GitHub. We linked the Jenkins pipeline from the Source Control Management option and linked it to the git repository. Whenever a change is pushed to git, the SCM polling in Jenkins keeps checking for updates and will find one when an update is made. Some of the advantages of Git are as follows:

- Git is one of the most widely-used version control systems in use today. Git is Open Source and helps us implement access control.
- The project can be worked on simultaneously by multiple team members using branches. It is easy to understand and merge the branches.
- Git is faster for performing network operations such as download and upload of project files to the file server.
- Git supports projects up to 1 Gb. Since our application wouldn't cross a GB, we preferred to use this instead of other SCM tools like GitLab, BitBucket etc.

## D. Artifact

An artifact is one of many by-products produced during the software development.

The artifact of our project was one docker image of the application itself which was pushed to DockerHub.

Image pushed to DockerHub(IMG 3):



Advantages:

- Simple and faster configuration.
- Docker enables you to build a container image and use that same image across every step of the deployment process.
- Rapid deployment
- Docker ensures your applications and resources are isolated and segregated.

## E. Deploy

For continuous deployment we have used Rundeck. In Rundeck we create jobs by defining a single step or a workflow that can execute any set of commands, scripts, or tools on any number of local or remote nodes. These jobs can be triggered by a scheduler or on demand via the web interface or API. Advantages:

- Ability to provide inputs to jobs at runtime.
- Nodes can be physical (servers, network devices, etc.) or logical (VM, service accounts, containers, etc.) endpoints.
- Node data can be pulled from anywhere, including Chef, Puppet, Amazon EC2, OpenStack, Docker, VMware, in-house CMDBs, or even monitoring tools.

- Multiple user-friendly ways to watch your jobs execute in real time. There is a summary view, node-oriented view, workflow view, and collated log view.

## F. Monitor

ELK stack stands for Elasticsearch, Logstash and Kibana. We use this tool for continuous monitoring. Elastic search is used for deep search and data analytics. It's an open source distributed NoSQL database, built in Java and based on Apache Lucene. Lucene takes care of storing disk data, indexing, and document scanning while Elastic Search keeps document updates, APIs, and document distribution between Elastic Search instances in the same cluster.

Logstash is used for centralized logging, log enrichment, and parsing. It's an ETL (Extract, Transfer, Load) tool that transforms and stores logs within Elastic Search.

Kibana is a web-based tool through which Elastic Search databases visualize and analyze stored data.

Advantages:

• Elastic search provides flexibility and scalability. It operates in real time.

• Logstash can collect data from multiple systems into a central system.

• Kibana offers powerful and easy to use features such as histograms, line graphs, pie charts, heat maps, etc.

As shown in the figure, the pie chart shows the different ids (partitions) when the website is visited. We can also visualize the pie chart with the different IP addresses that reflect that the website is visited from different devices.

## 3. CI/CD Pipeline

In the pipeline, we have used the SCM git option and provided it with the GitHub project repository URL ([https://github.com/Shiloni/classification-of-handwritten-digits-using-MNIST-dataset-on-kaggle](https://github.com/Shiloni/classification-of-handwritten-digits-using-MNIST-dataset-on-kaggle)). We have provided it with the credentials of GitHub which were already saved in Jenkins. In the script path, we have mentioned it to use Jenkinsfile.

GitHub Repo linking(IMG 4):

Once the SCM triggered, the build will start. Jenkins will search for the file 'Jenkinsfile' because we have mentioned it in the pipeline. All the instructions which Jenkins need to provide are written in this file. Along with this, the GitHub repository along with branch name is mentioned and the credentials are also specified.
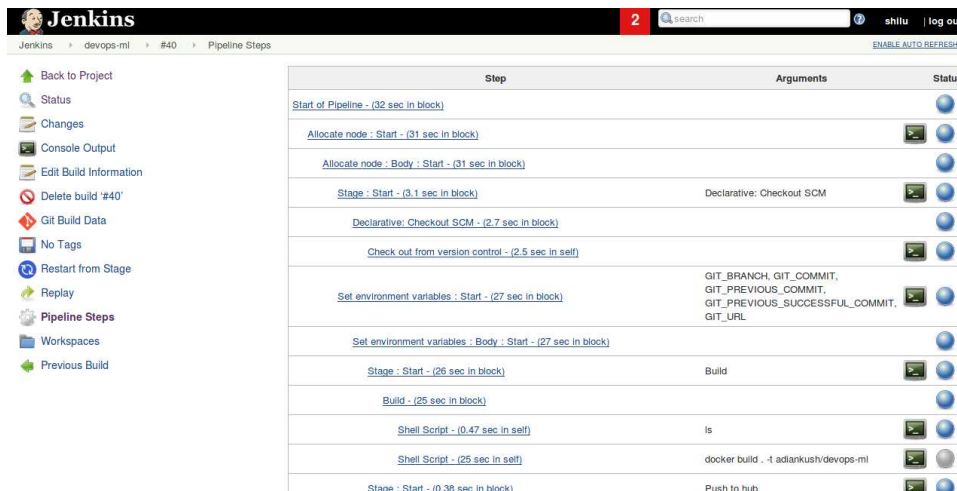
Providing Jenkinsfile path(IMG 5)



For source code management we use GitHub which contains all the files relevant to the code including the Dockerfile and Jenkinsfile, We have used jenkins for continuous integration in which we created a pipeline project (named devops-ml) linked to the Github repository consisting of three stages:

1. SCM polling to pull the files from the repository every time it is updated.
2. Building the Docker image utilizing the pulled files.
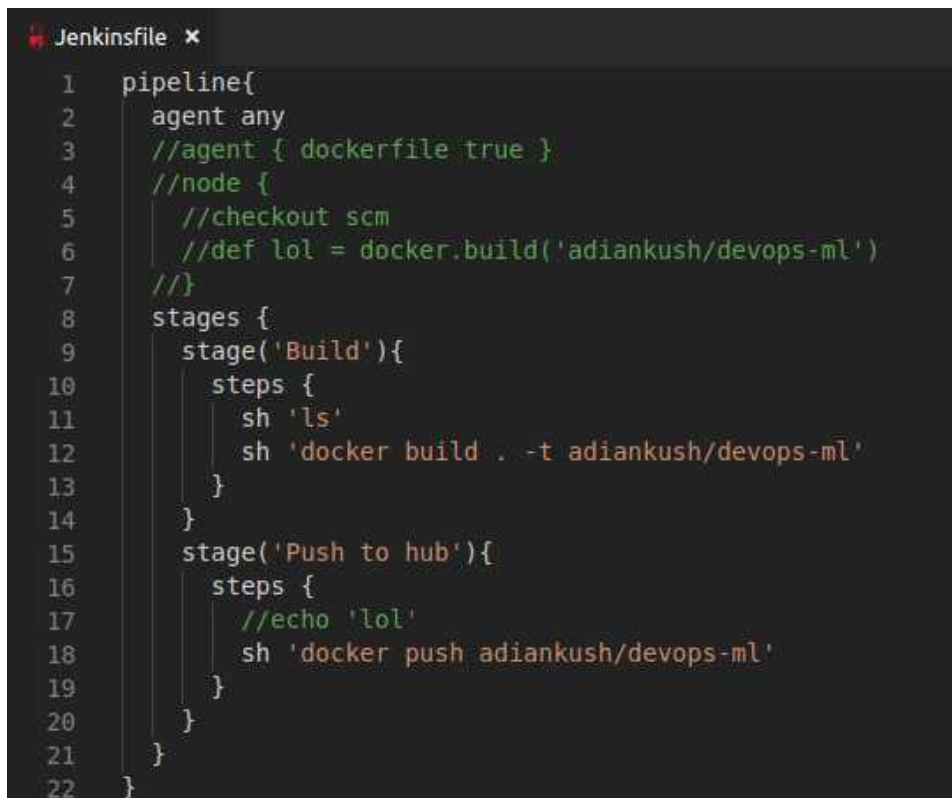3. Pushing the Docker image to DockerHub.

Jenkins Pipeline(IMG 6):



Jenkins file(IMG 7):

```
 1    pipeline{
 2      agent any
 3      //agent { dockerfile true }
 4      //node {
 5        //checkout scm
 6        //def lol = docker.build('adiankush/devops-ml')
 7      //}
 8      stages {
 9        stage('Build'){
10          steps {
11            sh 'ls'
12            sh 'docker build . -t adiankush/devops-ml'
13          }
14        }
15        stage('Push to hub'){
16          steps {
17            //echo 'lol'
18            sh 'docker push adiankush/devops-ml'
19          }
20        }
21      }
22    }
```

Pipeline GUI(IMG 8):

Docker file(IMG 9):



```
Dockerfile ×
1   # Use an official Python runtime as a parent image
2   FROM python:3.6.5-slim
3
4   # Set the working directory to /app
5   WORKDIR /train
6
7   # Copy the current directory contents into the container at /app
8   COPY . /train
9
10  # Install any needed packages specified in requirements.txt
11  RUN pip install --trusted-host pypi.python.org -r requirements.txt
12
13  # Make port 80 available to the world outside this container
14  EXPOSE 5000
15
16  # Define environment variable
17  ENV NAME World
18
19  # Run app.py when the container launches
20  CMD ["python", "Hello.py"]
21
```

On DockerHub we have created a project name adiankush/devops-ml where the latest artifact is pushed every time, Concluding the Continuous integration step.

DockerHub Artifact(IMG 10):

Upon completion of the above pipeline a free-style project triggers as the post-build step(named rundeck-ml), This project is linked to the RunDeck profile containing a predefined job whose uid is provided while creating the project. This then calls the RunDeck job which then executes the job steps on the nodes linked to the RunDeck account, The job steps are as follows:

- Docker pull (To pull the latest image from DockerHub)
- Docker create (To create the container from pulled image on the remote nodes)
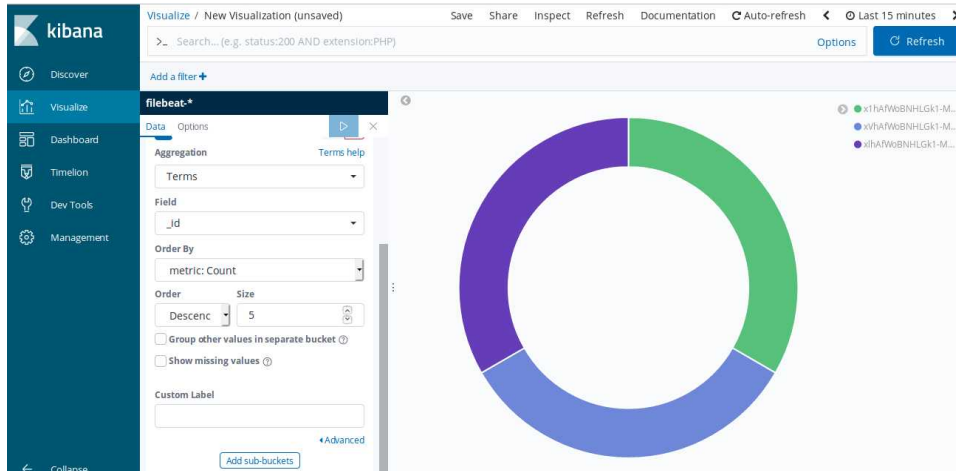
RunDeck Job(IMG 11):



```
! job3.yaml  ×
1    defaultTab: summary
2    description: ''
3    executionEnabled: true
4    id: '12345'
5    loglevel: INFO
6    name: job3
7    nodeFilterEditable: false
8    scheduleEnabled: true
9    sequence:
10     commands:
11     - description: pulling the image from docker hub
12       exec: docker pull adiankush/devops-ml
13     - description: creating container
14       exec: docker create -i -t -p 5000:5000 -v ~/log_fold:/train/logs --name final1 adiankush/devops-ml /bin/bash
15     keepgoing: false
16     strategy: node-first
17   uuid: '12345'
18
```

We have used ELK stack and Filebeat for continuous monitoring to which we provide the Log File generated by the flask python library (used to create website) within the code. The logfile is then generated within the container which we then create a copy of on the local machine using the –v(volume) extension on the Docker create step.

We provide the log file's(copy) path to Filebeat as the config file that gets sent to Elastic search for the creation of a pattern index which can then be visualized through Kibana.

Kibana Visualization(IMG 12):



## 4. Results and Discussions

We were successfully completed the application and deployed it using the aforementioned tools. It can be observed in IMG 9 that all of the pipeline stages have been built without any error. The memory required for the project is 252MB. The time taken for total pipeline to complete is (Varies with internet speed): 5min 41 secs. Since we have employed the DevOps approach our application is scalable. This is because when we add/update a feature, all we need to do is, push it to GitHub and our pipeline will handle the rest. Theoretically, this should work but we cannot comment on it with confidence as our application is not public. The application performed as expected in our production environment, but this might not be true in the actual deployment environment where we can have millions of users.

## 5. Future Work
- We would like to make the website more user-friendly.
- The project can also be improved by having multiple visualizations based on different parameters.
- We could also allow one user to train and store multiple models simultaneously

## 6. Conclusion

We were successful in developing a simple ML based handwriting recognition application and deploying it using DevOps tools. The DevOps tools used were GitHub, Jenkins, Docker, Rundeck and ELK. A pipeline architecture was made, and the above-mentioned tools were integrated, and the build was automatically triggered. By using the DevOps approach, we could configure, automate, monitor and maintain the project easily.

## 7. Acknowledgement

We would like thank Prof. Thangaraju and all the TA's (Shayaan and Vineeth) for their endless and selfless help throughout the project. Thanks to help us grow.

## 8. References

https://www.rundeck.com/what-is-rundeck

https://github.com/rundeck-plugins/openssh-node-execution

https://github.com/deviantony/docker-elk

https://logz.io/learn/complete-guide-elk-stack/