# Data Engineering Basics

Lacey Conrad
MSDS 610
Regis University

May, 2020

## 1    Introduction

The rise of Big Data has resulted in the development of methods to deal with the 5Vs associated with it: volume, velocity, variety, veracity, and value (Yang et al., 2016). Data Science, on one hand, deals with modeling and analyzing Big Data, whereas Data Engineering primarily focuses on developing and maintaining the data infrastructure (*state of data engineering* — Stitch benchmark report, 2019). Given the massive datasets data engineers are beginning to work with, they began to look for new ways to warehouse, process, and maintain the architecture of their data. Data engineers have started turning to cloud services for most of their needs, which has eliminated the need to maintain local servers infrastructure and software (Castrounis, 2017; Hufford, 2018).

The written portion of this assignment sought to provide foundational knowledge on the topic of data engineering. The technical portion of this exercise involved setting up a Google Cloud Platform account, including updating it, and then practicing Linux commands. After I had refreshed my Linux command knowledge, I followed the lab supplement and installed Docker on the Virtual Machine I had created in the previous step. Lastly, I performed a word count analysis on one of the works of Shakespeare, which is considered the "Hello World" analog for Hadoop.

## 2    Written Portion - Questions from week 1

1. What is Linux, and why is it important for data engineers and data scientists?

   Linux is an open source operating system, which means it relays instructions between programs and the hardware of a computer ("What is Linux?," 2019). Linux-based operating systems include a Linux kernel (which manages hardware resources) and software packages. Additional components are frequently included to aid in resource allocation, security, software installation, and more ("Understanding Linux," 2020).

   Linux has several characteristics that make it one of the most popular operating systems:

   (a) It is open source and free. One of its most popular features is its establishment as an open-source operating system which allows anyone to view the code and edit it. Also, this allows Linux to have superb support, since many people are knowledgeable of the operating system and have access to the code ("What is Linux?," 2019).

   (b) It is highly customizable. Linux has many distributions, allowing it to be highly customizable for the user and their project(s). Common distributions include Elementary, Mint, Fedora and Ubuntu ("What is Linux?," 2019).

   (c) It is secure. This is partly due to its handling of permissions. Without proper permissions, you cannot execute a .exe file. Also, since it is open source, there are many people capable of detecting places where the code could be more resilient to viruses ("What is Linux?," 2019).

   (d) Similarly, it is reliable. Server reboots, for example, only occur when a kernel is updated, which in some cases may be once every year or more ("What is Linux?," 2019).

Linux is vitally important to data scientists and engineers due to its computing power. The data handled by data scientists is so large that it becomes difficult to handle, and the computing speed offered by Linux makes it an ideal platform to use compared to other options. Additionally, the use of Docker makes Linux an ideal operating system for a data scientist. Docker, which enables programs to run simultaneously without interfering with each other, runs at its fastest speeds on a GPU. In many cases, the GPU based Docker containers will only run on Linux operating systems. Several of the characteristics listed above hugely benefit the data scientist when using Linux, but of central importance is the ability of Linux to handle the large amounts of data to be processed or analyzed by a data scientists and engineers. This increase in ability to handle Big Data is in part due to Linux's resource management (Misal, 2019).

2. Who are the main cloud services providers, and what sort of things do they provide?

Cloud computing enables users access to software whenever and wherever they are, assuming they have an internet connection. This is a great improvement over the past need to set-up and maintain physical servers and manage software. Cloud computing services are provided by a 3rd-party which hosts a network of remote servers that store and process data. In addition to these services, cloud service providers also manage the maintenance, performance and security of their servers (Hufford, 2018). The three main models of cloud services are as follows (although there are several other cloud services, they are not as frequently encountered):

(a) Software as a service (SaaS): This is a licensing and delivery model where software (or an application) is made available to a user via a web browser. The software lives and runs on servers within the cloud and does not need to be installed locally (Hufford, 2018).

(b) Platform as a service (PaaS): In PaaS, a cloud service provider builds and maintains software development infrastructure so enterprises can create and manage applications without the need of interacting with the underlying infrastructure (Violino, 2019).

(c) Infrastructure as a service (IaaS): Here, the components that typically comprise a physical data center are hosted and managed by a cloud provider. IaaS services can include servers, storage, and network hardware in addition to the services needed to utilize, manage, and protect these resources (Rouse, 2018).

According to Hufford (2018), the following three cloud service providers are likely to account for roughly 80% of all cloud revenue in 2020:

(a) Amazon Web Services (AWS): AWS has a substantial tool kit with a focus on public cloud computing.

(b) Microsoft Azure: Microsoft Azure is the go-to cloud service for enterprises. It uses a hybrid cloud, allowing Azure to work with enterprise data centers. Microsoft Azure is particularly noted for their SaaS products such as Microsoft 365.

(c) Google Cloud: Google Cloud is the best option for technical cloud services pertaining to deep learning, artificial intelligence, machine learning, and data analytics (Harvey & Patrizio, 2020).

3. What can we do with the cloud providers related to data engineering?

Data engineering is moving at a very rapid rate to the cloud. Frequently, the datasets Data Engineers manage are simply too big to be housed on a local machine, and similarly, many these huge datasets require too much processing power to perform tasks in a reasonable amount of time (Castrounis, 2017). Cloud services can be used to warehouse huge datasets cheaply and can keep them in any format or size they choose. The data can be organized, merged (or separated), and analyzed in the cloud without any pressure on the local development machine. Cloud providers will also provide the support to manage the hardware, maintenance, server performance, and so on (Hufford, 2018).

Many cloud providers also promise elasticity and flexibility to data processing infrastructure. While this shows how important it is for data engineers to keep up-to-date with data engineering technology, it also shows how the cloud providers are promising to the data engineers that they (the cloud providers) will provide cutting edge data processing technology ("Data engineering and science in the cloud," 2019).

4. What is Docker, how is it useful, and why is it important to know?

In a way, Docker can be thought of as a lightweight virtual machine (Husain, 2018). Docker is a tool that uses containers to allow easier creation, deployment, and execution of applications. Docker sought to reduce the reliance on system resources previously encountered in virtual machine hypervisors by utilizing containers, which are much more efficient and secure (Vaughan-Nichols, 2018). Containers are portable units of software that can be used reliably across different computing environments. These lightweight units are packages of code, its dependencies, run-time, tools, libraries and settings. Containers include everything needed to execute an application, and in essence they have their own file system, storage, CPU, RAM, etc ("What is a container?," 2020; Vaughan-Nichols, 2018). Additionally, since containers use a shared operating system they place a lot less stress on system resources, reducing application size and increasing performance ("What is docker?," 2019; Vaughan-Nichols, 2018). This allows for four to six times the number of application instances running at once. Docker does require containers to use the same operating system and kernel, which is not the case with other options like Microsoft Azure (Vaughan-Nichols, 2018).

Docker provides Data Scientists with multiple benefits:

(a) Makes the work of data scientists and engineers easily reproducible. Using Docker containers allows you to wrap up your computing environment, making it easier for others to reproduce your work.

(b) Allows the data scientist to have a more portable computing environment. This allows data scientists to access the computing environment they need without worrying about having to re-create a local environment every time they change environments.

(c) Makes is easy to deploy applications, which increases accessibility to a data scientists work. This also allows for reproducibility and quick replication of research (Husain, 2018).

# 3    Methods and Code

To begin this lab, I created an account with Google Cloud Platform (GCP), and then set up a budget. After I set up my account, I was able to spin up my first GCP VM instance using the settings provided in the lab write-up:

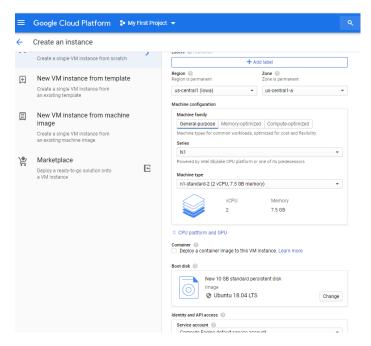## 3.1   Linux practice and Setting up Google Cloud Platform



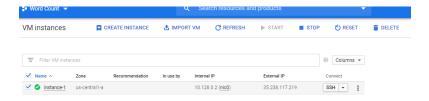Figure 1: Setting up a VM instance in Google Cloud Platform.



Figure 2: All created VM instances in GCP showing the instance I just created.

Clicking on the SSH button under Connect allowed me to open a terminal window. Here, I practiced several Linux commands (although I have a decent amount of experience in Linux, so I didn't spend much time here). Below I took a screen capture of the open terminal in addition to one of the Linux commands:

Figure 3: Left: After clicking on SSH in the figure above, a terminal window opens up. Right: A screen capture of Linux command practice, in this case `mad pwd`.

I then installed Python on my virtual machine, and checked to see if tmux is installed:



Figure 4: Installation of Python and Tmux, showing that Tmux is already installed.

I used `update` to update the current software, and `upgrade` to update the software to its most recent version:

Figure 5: Left: Updating the software in Ubuntu. Right: Upgrading the current software in Ubuntu.
Once the software was upgraded, I practiced creating a tmux session and attaching and detaching from it:



Figure 6: A tmux session window.

## 3.2   Installing and Running docker on Ubuntu

Since Docker will be an important component of this lab and class, it needed to be installed on my virtual
machine. I used `sudo apt-get` to install the packages outlined in the lab document. I have lightly
commented the code below to reflect its purpose whenever possible.

```
1   # The following is the code required to install Docker in Ubuntu.  It consists
2   # of the necessary packages to set up a Docker repository:
3
4   # Upgrading software to its current version:
5   $ sudo apt-get update
6   $ sudo apt-get upgrade
7
8   # Setting up the Docker repository:
9   $ sudo apt-get install apt-transport-https
10  $ sudo apt-get install ca-certificates
11  $ sudo apt-get install curl
12  $ sudo apt-get install gnupg-agent
13  $ sudo apt-get install software-properties-common
14
15  # Adding Docker's GPG key:
16  $ curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo apt-key add -
17
18  # Verifying the key's fingerprint:
19  $ sudo apt-key fingerprint 0EBFCD88
20
21  # Setting up the stable repository:
22  $ sudo add-apt-repository \
23          "deb [arch=amd64] https://download.docker.com/linux/ubuntu \
24          $(lsb_release -cs) \
25          stable"
26
```

```
27  # Updating the newly installed software:
28  $ sudo apt-get update
29
30  # Installing several Docker containers:
31  $ sudo apt-get install docker-ce docker-ce-cli containerd.io
```



Figure 7: Installing Docker.

Here, I am verifying the fingerprint of the key:



Figure 8: Verification of key with fingerprint `9DC8 5822 9FC7 DD38 854A E2D8 8D81 803C 0EBF CD88`.

The following shows my establishment of the Docker repository. I also updated the software/packages I had installed, and lastly, installed the Docker containers:



Figure 9: Setting up the stable Docker repository and updating the software.

7

My ID needed to be attached to Docker, which is shown in the following code. Also, I used a Docker `pull` to download the container for cloudera quickstart, and then ran the docker container:

```
1  # Added my user to the Docker group:
2  $ sudo usermod -a -G docker $USER
3
4  # Pulling/downloading the Docker container for Cloudera quickstart:
5  $ docker pull ngeorge/ubuntu-hadoop-quickstart
6
7  # Running the Docker container:
8  $ docker run -it ngeorge/ubuntu-hadoop-quickstart
```



Figure 10: Downloading the Docker container for cloudera quickstart.

## 3.3 Word count example

Now that I have a running Docker container, I can practice using Hadoop. I will do this by running MapReduce on a text file as described in the lab document. First off, I needed to obtain the code for the MapReduce from git hub by using `git clone`. Then, I downloaded a text of Shakespeare's work.

```
1  # Cloning the git hub repo where the MapReduce code is located:
2  $ git clone https://github.com/Regis-University-Data-Science/
3         simple_Hadoop_MapReduce_example.git
4
5  # Downloading the Shakespeare text:
6  $ wget http://norvig.com/ngrams/shakespeare.txt
```



Figure 11: Cloning the github repo where the MapReduce code is stored. Then, downloading the text to be analyzed

I then needed to create directories in hdfs for the raw and processed data. The hdfs directories are separate from local directories, so we will also need to copy our text document to the hdfs file system.

```
1  $ hdfs dfs -mkdir /shakespeare
2  $ hdfs dfs -mkdir /shakespeare/input
3  $ hdfs dfs -copyFromLocal shakespeare.txt /shakespeare/input
4  $ hdfs dfs -ls /shakespeare/input
```

Figure 12: Creating directories in preparation for the MapReduce output.

Now, it's time to run the MapReduce. The MapReduce implementation we are using reads in a text file and then returns how many times a word is encountered in the file. It reads in the file line by line, and as the line it being read in, it separates the line into a list of words which were separated by white space. The mapper program will then parse through each word, and simply assign a 1 to it, creating key-value pairs. The output of the mapper is sorted, in our case alphabetically, and fed into the reducer program. The reducer takes all the occurrences of a word and adds up how many times that words was encountered, thus producing a word count of all words in the document.

```
1  # The commands to run the MapReduce programs.  Here, our Shakespeare file
2  # is ran through the mapper, and then the reducer, and the output
3  # is saved in output6 (yes, I had a few issues that ended up being
4  # really silly):
5  $ mapred streaming \
6  -mapper mapper.py \
7  -reducer reducer.py \
8  -input /shakespeare/input \
9  -output /shakespeare/output6
```



Figure 13: Output of running the MapReduce streaming command.

Finally, I viewed the results of the MapReduce analysis, and copied the output to my local file system:

```
1  $ hdfs dfs -ls /shakespeare/output6
2  $ hdfs dfs -cat /shakespeare/output6/part-00000
3  $ hdfs dfs -copyToLocal /shakespeare/output6/part-00000 result
4  $ sort -gr -k 2 result | head
```

## 4   Results and Output

While most of this lab involved setting up our google cloud platform, and then Docker containers, there were some data produced as a result of the MapReduce ran on the Shakespeare text. In the first image

9

below we see the unsorted results of the MapReduce in the format of (word, number of instances). The MapReduce program reads in files line by line while separating each line at its white spaces. This is why we see punctuation in our output, and some odd breaks in words.
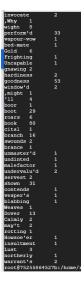


Figure 14: Results of running MapReduce on Shakespeare text.

Here is the same data after running the sort command, which sorts the output by most frequently encountered word.



Figure 15: Word count produced by MapReduce, sorted.

# 5 Analysis

The MapReduce ran on the Shakespeare text produced the results as we would expect; without any sorting the results are simply a list of words and their counts. The sorted results, which look to be sorted, in descending order, according to the most frequently seen "word," show results that if we were actually analyzing a text file would likely not be of much use. The most counted word is a comma, and other punctuation marks follow, along with the most frequently used English words such as "the", "and", "I", "to". The text would be better analyzed after proper cleaning where we could remove unwanted characters such as the punctuation, make all letters the same case (so island and Island aren't 2 different words), and removing unwanted text that tends to be found at the beginning and end of a text document.

# 6 Conclusion

I have heard of Google Cloud Platform, but have never utilized it for anything. In fact, I barely knew what it did until this lab. Same can be said about cloud computing and Docker. After doing the written portion of the lab, I began filling in some of the holes in my knowledge, while also learning about technologies I knew little about. Then, seeing many of these technologies in action really helped to drive home their uses.

I also found myself going back after completing the written questions, and reading more about cloud computing, GCP, Docker, and Hadoop. In particular, I wanted to understand Docker better, as I was

struggling to understand what it was and why we used it. I believe I better understand it now, and I am even working through some tutorials at `training.play-with-docker.com` to beef up my Docker knowledge, and to feel more comfortable using it. Now that I have a better idea of what it does and why it's important to Data Science, I can understand why it is in such high demand right now. I can also see why it hasn't and likely wont replace virtual machines anytime soon.

As far as does data science or data engineering interest me more, I am still figuring that out. After discussing this topic on the forums, I do feel that the skills and requirements of each of these careers are more clear to me, yet I still want to learn more about data engineering before I can fully commit to either one.

# 7 References

1. *5 top cloud service providers companies in the world.* (2019, January 30). DataFlair. `https://data-flair.training/blogs/cloud-service-providers-companies/`

2. Castrounis, A. (2017, September 20). *Cloud computing and architecture for data scientists.* DataCamp Community. `https://www.datacamp.com/community/blog/data-science-cloud`

3. *Data engineering and science in the cloud.* (2019, August 29). Cloudera. `https://www.cloudera.com/products/data-science-and-engineering/data-engineering.html`

4. Harvey, C., & Patrizio, A. (2020, March 17). *AWS vs. Azure vs. Google: Cloud comparison.* Datamation: Emerging Enterprise Tech Analysis and Products. `https://www.datamation.com/cloud-computing/aws-vs-azure-vs-google-cloud-comparison.html`

5. Hufford, J. (18, February 6). *Cloud vs SaaS: What's the difference?* — nChannel blog. eCommerce, ERP, POS & 3PL System Integration — nChannel. `https://www.nchannel.com/blog/cloud-vs-saas/`

6. Husain, H. (2018, January 17). *How docker can help you become a more effective data scientist.* Towards Data Science. `https://towardsdatascience.com/how-docker-can-help-you-become-a-more-effective-da`

7. *LinkedIn's 2017 U.S. Emerging Jobs Report.* 2017. LinkedIn.

   `news.linkedin.com/2017/12/introducing-linkedins-2017-u-s--emerging-jobs-report.`

8. Misal, D. (2019, October 27). *Linux vs Windows: Which is the best OS for data scientists?* Analytics India Magazine.

   `https://analyticsindiamag.com/linux-vs-windows-which-is-the-best-os-for-data-scientists/`

9. Rouse, M. (2018, September 2). *What is infrastructure as a service (IaaS)? - Definition from WhatIs.com.* SearchCloudComputing. `https://searchcloudcomputing.techtarget.com/definition/Infrastructure-as-a-Se`

10. Saltz, J. S., Yilmazel, S., & Yilmazel, O. (2016). Not all software engineers can become good data engineers. textit2016 IEEE International Conference on Big Data (Big Data).

11. *The state of data engineering — Stitch benchmark report.* (2019). Stitch. `https://www.stitchdata.com/resources/the-state-of-data-engineering/?thanks=true`

12. Steven J. Vaughan-Nichols. (2018, March 21). *What is docker and why is it so darn popular?* ZDNet. `https://www.zdnet.com/article/what-is-docker-and-why-is-it-so-darn-popular/`

13. *Understanding Linux.* (2020). Red Hat - We make open source technologies for the enterprise. `https://www.redhat.com/en/topics/linux`

14. Violino, B. (n.d.). *What is PaaS? platform-as-a-service explained.* InfoWorld. `https://www.infoworld.com/article/3223434/what-is-paas-software-development-in-the-cloud.html`

15. *What is a container?* (2020). Docker. `https://www.docker.com/resources/what-container`

16. *What is docker?* (2019). Opensource.com. urlhttps://opensource.com/resources/what-docker

17. *What is Linux?* (2019). Opensource.com. `https://opensource.com/resources/linux`

18. Yang, C., Huang, Q., Li, Z., Liu, K., & Hu, F. (2016). Big data and cloud computing: Innovation opportunities and challenges. *International Journal of Digital Earth*, 10(1), 13-53. `https://doi.org/10.1080/17538947.2016.1239771`