

Engineering For Data Science with Docker

Joshua Cook

May 4, 2018

Contents

1	Configuration	2
1.1	The Modeling Problem and the Engineering Problem	2
1.2	Getting Started	2
1.2.1	Prerequisites	3
1.2.2	Useful tools	3
1.2.3	SSH Key Pairs	3
1.3	Amazon Web Services	5
1.3.1	Configure your AWS Account	6
1.4	Exercises	8
1.5	Summary	8
2	Elastic Compute Cloud	10
2.0.1	Launch a New AWS EC2 Instance	10
2.1	Configure Git & Github	13
2.1.1	Git and Github	15
2.1.2	Configuring Github	15
2.2	Docker & Docker Compose	19
2.2.1	Docker Compose	19
2.2.2	Installing and Configuring Docker	21
2.2.3	Installing and Configuring Docker	21
2.3	Summary	22
2.4	Exercises	22

Chapter 1

Configuration

1.1 The Modeling Problem and the Engineering Problem

As a practicing data scientist, it is likely that you spend the bulk of your time working toward the development of a model for a particular inference or prediction application. It is less likely that you spend time thinking of the equally complex problems stemming from your system infrastructure. We might trivially think of these two often orthogonal concerns as the **modeling problem** and the **engineering problem**. The typical data scientist is trained to solve the former, often in an extremely rigorous manner, but can often wind up developing a series of ad hoc solutions to the latter.

Since its introduction in 2013, Docker has quickly become a fundamental tool in the design and deployment of robust engineering infrastructure for many applications. From the smallest tech shops to Google, Docker is being used to

- modernize traditional software
- leverage cloud resources for application architecture
- streamline continuous integration and deployment pipelines
- build out microservices

These will certainly seem downright esoteric to the Data Scientist who is typically concerned with feature importances or how many epochs to run to train a certain neural network.

That said, developing a robust engineering practice with Docker at its core can only make for better data science. This includes immediate concerns such as environment configuration and replicability and presentation of results, but in learning how to use Docker properly the data scientist can ensure that their work is deployed correctly as part and product of their team's software.

Here, I discuss Docker as a tool for the data scientist, in particular in conjunction with the popular interactive programming platform Jupyter and the cloud computing platform Amazon Web Services (AWS). Using Docker, Jupyter and AWS, the data scientist can take control of their environment configuration, prototype scalable data architectures, and trivially clone their work toward replicability and communication.

1.2 Getting Started

In this first chapter, you will configure your local system and create an AWS instance in order to do data science work. To do this work, you will use the package management system `conda`, the containerization technology Docker, the version control software `git` and its counterpart cloud-based backup service Github.com, and the system design tool Docker Compose.

There is something of a bootstrap moment in this first chapter. I'm writing this book using Jupyter notebooks. The goal is that you are able to read it and execute the very same code from Jupyter notebooks while you're reading it. That said, you may not even have a Jupyter running on your system. The steps here I designed to take you through step-by-step to install and configure all of the tools you will need to do the work in this book.

1.2.1 Prerequisites

We will assume a basic knowledge of working in bash. This should include things like knowing that `~` is an alias for your home directory, that `pwd` shows your current location, `cd` changes directories and `ls` can be used to list files and in a directory.

1.2.2 Useful tools

Bash

If you are using a Mac OS X system or a Linux system, you will already have Bash available to you in an application called Terminal. If you are using a Windows system, you can use the Anaconda Prompt that will be installed next. If you are on a Mac, you may want to install iTerm (<https://iterm2.com>).

Conda

It is recommended that you install [Conda](#) on your local system. Conda is a package and environment management system for Linux, Mac, and Windows and is perfect for managing our local Python packages. Detailed instructions for installing Conda on your system can be found here: <https://conda.io/docs/user-guide/install/index.html>

Atom

A simple but extensible text editor such as Atom (<https://atom.io>) can be an invaluable tool. Atom is available for any modern operating system.

`vim`

Nearly all of the work we will be doing will be on remote systems. It can be useful to be able to edit text files in place on these remote systems. We have essentially three options to do this:

- Vim
- nano
- Emacs

I will emphasize Vim. Vim is a text editor, like Atom. It is extensible, if not exactly simple. It is available on nearly every system by default, but it can be very challenging to learn. With a little guidance, we will be able to do but we need to do. It should be noted that we will only be using Vim on remote systems. If you would like to spend some time acquainting yourself with Vim, you can type `vimtutor` at a Bash prompt.

1.2.3 SSH Key Pairs

All of the work that you will be doing will take place remotely. As such, there is very little configuration to be done for the local system. The one thing that you will need to do is configure a set of SSH Keys to enable secure connection to the remote system you bring online. We will be using the Secure Shell protocol (SSH) to do the vast majority of our command line work. It is considered a best practice to use an SSH key pair for authentication when using SSH. An ssh key pair is a set of two long character strings: a private key and a public key. Though they are both called keys, I prefer to think of them as a key and lock.

An SSH Key is a password-less method of authenticating to a remote system using public-key cryptography. Authentication is done using a key pair consisting of a public key (`id_rsa.pub`), which can be shared publicly, and a private key (`id_rsa`), which is known only to the user (See [fig 1.1](#)). One might think of the public key as the lock on your front door, accessible to anyone, and the private key as the key in your pocket so that only you are able to open your door and gain access to your home.

You will generate this key pair on our local system and then provide the public key to AWS so that it can be added to any system you wish to launch. You will keep the private key on our local system and use it whenever you wish to gain access.

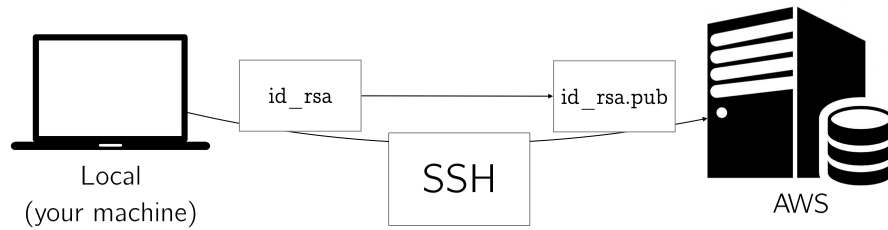


Figure 1.1: Connecting with SSH Keys

Check to See if SSH Key Pair Exists

You will use the Bash tool `ssh-keygen` to create a new key pair. To begin open a new terminal session¹, where you will examine whether or not you already have a key pair. Launching a new Bash session will put us in our home directory. The canonical location for storing SSH Keys is in a folder called `~/.ssh` in our home directory. Note that this directory begins with a `.` which makes it a hidden directory. In Bash Command 1, you use `ls -la` to display all of the contents of your home directory (`~`) in a list.

We will first check to make sure that you do not already have an SSH key pair.

Bash Command 1 List the contents of home directory

```

LOCAL:  ls -la ~

...
drwxr-xr-x  21 joshuacook  staff    714 Jul 31  2017 .pylint.d
drwx-----  11 joshuacook  staff    374 Jan 28 18:49 .ssh
drwxr-xr-x   6 joshuacook  staff    204 Feb  2 22:01 .vim
-rw-----   1 joshuacook  staff 20788 Feb 10 08:54 .viminfo
-rw-r--r--   1 joshuacook  staff  1263 Jul 26  2017 .vimrc
drwx-----@  5 joshuacook  staff    170 Aug 26 09:12 Applications
drwx-----+ 19 joshuacook  staff    646 Feb 11 09:32 Desktop
drwx-----+  6 joshuacook  staff    204 Feb  4 12:18 Documents
...
  
```

My local system is running Mac OS X and has the `.ssh` folder already. As the directory already exists, in Bash Command 2, I list the contents of my `.ssh` directory.

As can be seen, I already have an SSH Keypair named `id_rsa` and `id_rsa.pub`. If this is true for you, as well, you should skip the next step and not create new SSH Keys. If you do not have these keys, proceed to Bash Command 3.

Create a new SSH Key Pair

If when listing the home directory, you do not see a folder called `.ssh` or when displaying the contents of `.ssh` you do not see an SSH Keypair named `id_rsa` and `id_rsa.pub`, a new SSH Keypair will need to be created. In Bash Command 3, you create a new SSH Keypair using the `ssh-keygen` command line utility.

¹On a Mac or Linux system, simply open the Terminal application. On Windows, open Anaconda Prompt.

Bash Command 2 Display contents of the .ssh directory

LOCAL: `ls -la ~/.ssh`

```
total 24
drwxr-xr-x  8 jovyan users 272 Apr 24 23:07 .
drwxr-xr-x 21 jovyan users 714 Apr 25 00:39 ..
-rw-----  1 jovyan users 1679 Apr 24 23:06 id_rsa
-rw-r--r--  1 jovyan users  418 Apr 24 23:06 id_rsa.pub
```

During the creation of the SSH Keypair, you will be prompted three times. The first asks where you should save the SSH Keypair, defaulting to the `.ssh/id_rsa` in our home directory. In Bash Command 3, you see that this is being done at `/Users/joshuacook/.ssh/id_rsa` on my local system where my username is `joshuacook`. The second and third prompts will ask for a passphrase to be added to the key. For our purposes, leaving this passphrase empty will be fine. In other words, the default options are preferable and you may simply hit `<ENTER>` three times.

Bash Command 3 Create a new SSH Key Pair

LOCAL: `ssh-keygen`

```
Generating public/private rsa key pair.
Enter file in which to save the key (/Users/joshuacook/.ssh/id_rsa):
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in id_rsa.
Your public key has been saved in id_rsa.pub.
The key fingerprint is:
SHA256:p5KeEomPt6izFC5gaFphfx3zw8aAB+D8RiUA1/nEsUc joshuacook@LOCAL
The key's randomart image is:
+---[RSA 2048]----+
|  ..+ooo.E      |
|   +  o=oo      |
|  o o oo* .     |
|.. o o o..0     |
|o+...+ S B      |
|*.o oo . + .    |
|oo o .o .       |
|+ ..+. o        |
|o+...oo         |
+----[SHA256]-----+
```

You can verify the SSH Keypair you just created by displaying the Public Key in our shell (See Bash Command 4). Here, you use the `cat` command, which concatenates the contents of `id_rsa.pub` to the shell output.

That is the sum of the local configuration you will need to do in order to get started.

1.3 Amazon Web Services

If you have not already done so, set up an AWS account². You will be using AWS to manage the hardware upon which your data science platform will run. We will leave the details of what exactly “hardware”

²As of 2017/12/19, detailed instructions for doing this can be obtained here: <https://aws.amazon.com/premiumsupport/knowledge-center/create-and-activate-aws-account/>

Bash Command 4 Display Public SSH Key

LOCAL: `cat ~/.ssh/id_rsa.pub`

```
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDAQdnHPEiq1a40sDDY+g9luWQS8pCjBmR
64MmsrQ9MaIaE5shIcFB1Kg3pGwJypipZjoSh9pS55S9LckNsBfn8Ff42ALLj
R8y+WlJKVv/ODvDXgGVcCc0t/uTvVx0bRruYxLW167J89UnxnJuRZDLey9fD
OfIzSR5eglhCWVqi0zB+0sLqR1W04Xz1oStID78UiY5msW+EFg25Hg1wepYMC
JG/Zr43By0YPGseUrbCqFBS1KlQnzfWRfEKHZbtEe6HbWwz1UDL2NrdFXxZAI
XYYoCVt14WXd/WjDwSjbMmtf3BqenVKZcP2DQ9/W+geIGGjv0TfUdsCHennYI
EUfEEP joshuacook@LOCAL
```

means to AWS. This is to say that AWS may be allocating resources as a virtual machine, but for your purposes, the experience will be as if you are using a physical system across the room from you.

The most popular service offered by Amazon Web Services is the Elastic Compute Cloud (EC2), “a web service that provides secure, resizable compute capacity in the cloud”³. For our purposes, compute capacity means a cloud-based computer you will use to run your platform. What we learn should generalize to other cloud providers such as DigitalOcean or Google cloud platform.

If you are new to AWS you will be able to work through this text using the AWS Free Tier⁴. For the first 12 months following sign up, new users receive 750 Hours per month of EC2 time. This amounts to 31.25 days of availability and, provided that readers keep only one server running at a time, ensures that readers can work through this text at no cost.

1.3.1 Configure your AWS Account

The next thing you will need to do is configure your AWS Account.

This will involve:

1. Configure a Key Pair
2. Creating a Security Group

The AWS Key Pair is slightly misnamed as it is not in fact a pair, but rather is simply the public portion of the SSH Key Pair you have on our local system. You will simply add the Public Key from the SSH Key Pair you just created.

To begin, log in to your AWS control panel and navigate to the EC2 control panel (fig. 1.2). First, access “Services” (fig. 1.2, #1) then access “EC2” (fig. 1.2, #2). The Services link can be accessed from any page in the AWS website.

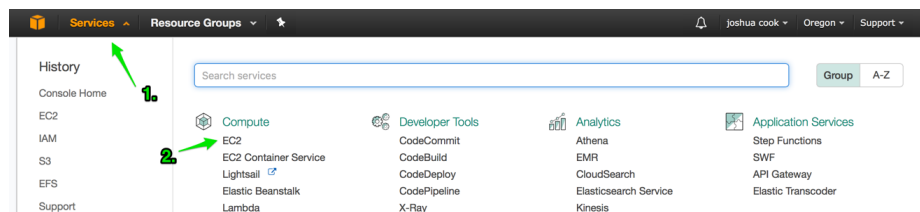


Figure 1.2: Access EC2 Dashboard

Configure a Key Pair

Once at the EC2 control panel, access the Key Pairs pane using either link (fig. 1.3).

³<https://aws.amazon.com/ec2/>

⁴<https://aws.amazon.com/free/>

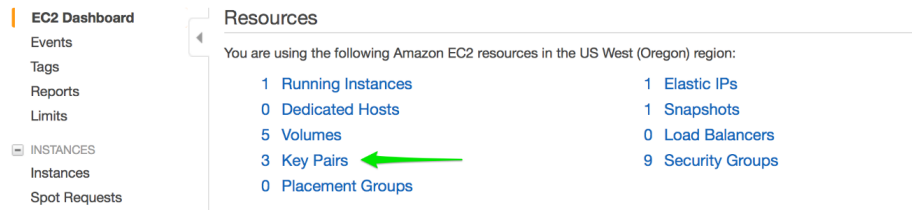


Figure 1.3: Access Key Pairs from the EC2 control panel

From the Key Pairs pane, choose “Import Key Pair.” This will activate a modal that you can use to create a new key pair associated with a region on your AWS account. Make sure to give the key pair a computer-friendly name, like `from-MacBook-2018`. Paste the contents of your public key (`id_rsa.pub`) into the public key contents. Prior to clicking Import, your key should appear as in fig. 1.4. Click Import to create the new key.

Import Key Pair

Click Browse and navigate to your public key. You may change the name of your key if necessary. Alternatively, you can copy and paste the contents of your public key into the dialog.

Load public key from file
no file selected

Key pair name

Public key contents

```
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDdnHPEiq1a4OsDDY+g9luWQS8pCjBmR64MmsrQ9
MalaE5shlcFB1Kg3pGwJpypiZjoSh9pS55S9LckNsBfn8Ff42ALLjR8y+WlJKVv0DvDXgGVcCc0t
/uTvXVx0bRruYxLW167J89UnxnJuRZDLeY9fDOftzSR5eglhCWVqiOzB+OsLqR1W04Xz1oStlD7
8UiY5msW+EFg25Hg1wepYMCJG/Zr43ByOYPGseUrbCqFBS1KIQnzfWRfEKHZbtEe6HbWwz1
UDL2NrdFXxZAlXYyOCVtl4WXd/WjDwSjbMmtf3BqenVKZcP2DQ9/W+gelGGjvOTfUdsCHennYl
EUfEEP ubuntu@ip-172-31-43-19
```

Figure 1.4: Import a New Public Key

You have just created a key pair between AWS and your local system. When you create a new instance, you will instruct AWS to provision the instance with this public key and thus you will be able to access the cloud-based system from your local system using your private key.

Ports & Security Groups

A security group is a set of ports public facing Internet. This may be a new concept for you and we will not dig very deeply into it. Suffice it to say that we need a short list of ports to be available to us for accessing the different services we might configure. A port is a number appended to an IP address or domain name with a colon like this

`192.168.99.100:3000`

Here, `192.168.99.100` is the IP address and `3000` is the port. A port can be thought of as a channel over which a service will listen for requests. You have already been using a few ports without knowing that you are. All internet traffic is routed using ports `80` and `443` by default. What this means is that visiting `http://google.com:80` is equivalent to visiting `http://google.com` and visiting `http://google.com:443` is equivalent to visiting `https://google.com`.

As we are largely concerned with learning, we are not concerned with the specifics of high-availability, and for us networking best practices will consist of making sure that things work. If you intend on putting work that you do here into production, you should certainly consult with your local Site Reliability Engineer.

For security purposes, by default Amazon closes all ports to outside traffic. This is why we will need to create a security group to open the ports that we need. These include the following:

port	service
22	ssh
80	http
443	https
5000	miscellaneous
5432	PostgreSQL
6379	Redis
8888	Jupyter
27017	Mongo

Create a New Security Group

From the EC2 Control panel, access Security Groups (See fig. 1.5).

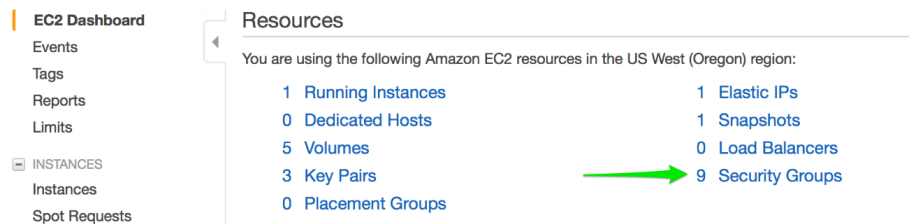


Figure 1.5: Access Security Groups from the EC2 control panel

From the Security Group pane, click “Create Security Group.” Give the security group a computer friendly group name like `ds_engineering`. Give the security group a description like “Open access to important Data Science related services”. Use the default VPC.

Access the Inbound tab (See fig. 1.6, #1) and configure the security rules listed above. **Make sure to set Source to “Anywhere”.**

1.4 Exercises

1.5 Summary

1.

Create Security Group

Inbound

Outbound

Type ⓘ	Protocol ⓘ	Port Range ⓘ	Source ⓘ	Description ⓘ	
SSH	TCP	22	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
HTTP	TCP	80	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
HTTPS	TCP	443	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
Custom TCP Ru	TCP	5000	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
PostgreSQL	TCP	5432	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
Custom TCP Ru	TCP	6379	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
Custom TCP Ru	TCP	8888	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕
Custom TCP Ru	TCP	27017	Anywhere ⓘ 0.0.0.0/0, ::/0	e.g. SSH for Admin Desktop	✕

Add Rule

Cancel

Create

Figure 1.6: Create a New Security Group

Chapter 2

Elastic Compute Cloud

2.0.1 Launch a New AWS EC2 Instance

AWS EC2 `t2.micro`

AWS EC2 virtual machines are available to meet a host of different applications and purposes. The `m` series and the `t` series are considered in general purpose and are adequate for our needs. Additionally, the `t2.micro` from the `t` series is considered a machine “free-tier” and can be run for free under certain circumstances. We will use this type of machine as often as possible.

Launch Instance

To create a new instance, start from the EC2 control panel and click the Launch Instance button (fig. 2.1).

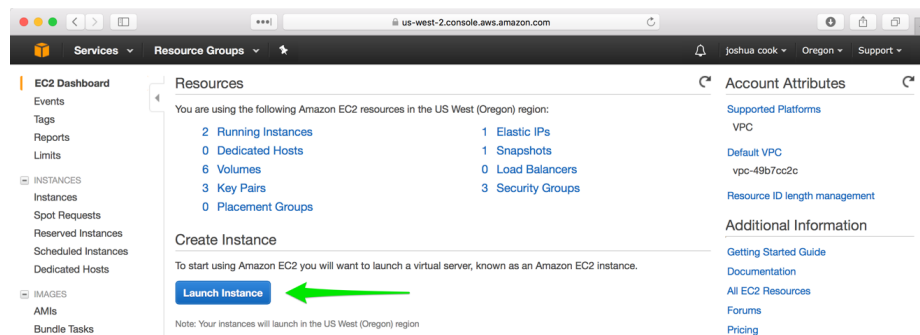


Figure 2.1: Begin the launch process for a new instance

Step 1: Choose an Amazon Machine Image (AMI)

The launching of a new instance is a multi-step process that walks the user through all configurations necessary. The **first tab** is “Choose AMI.” An AMI is an Amazon Machine Image¹, and contains the software you will need to run your sandbox machine. I recommend choosing the latest stable Ubuntu Server release that is free-tier eligible. At the time of writing, this was `ami-efd0428f`, Ubuntu Server 16.04 LTS (HVM), SSD Volume Type (fig. 2.2).

¹<http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AMIs.html>

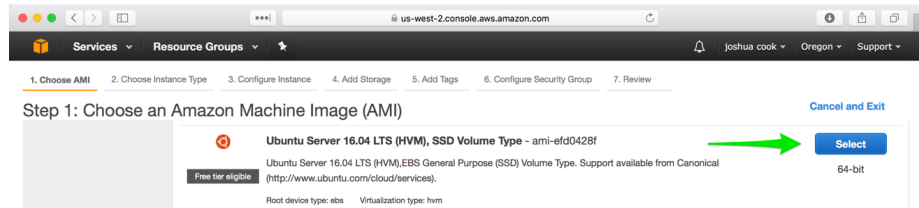


Figure 2.2: Choose the latest stable Ubuntu Server release as AMI

Step 2: Choose Instance Type

The **second tab** is “Choose Instance Type.” In practice, I have found that the free tier, `t2.micro` (fig. 2.3), is sufficient for many applications. Furthermore, the instance type may always be changed later should the need present itself.

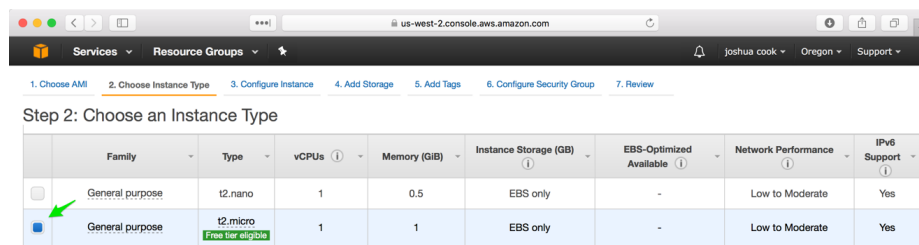


Figure 2.3: Choose `t2.micro` for Instance Type

Step 3: Configure Instance Details

The **third tab**, “Configure Instance,” can be safely ignored.

Step 4: Add Storage

The **fourth tab** is “Add Storage.” This option is also specific to intended usage. It should be noted that Jupyter Docker images can take up more than 5GB of disk space in the local image cache. For this reason, it is recommended to raise the value from the default 8GB to 30GB. Furthermore, as noted on this tab:

Free tier eligible customers can get up to 30 GB of EBS General Purpose (SSD) or Magnetic storage.

Step 5: Add Tags

The fifth tab, “Add Tags,” can be safely ignored.

Step 6: Configure Security Group

The sixth tab, “Configure Security Group,” is critical for the proper functioning of your systems. By default this tab will be set up to “Create a **new** security group”. This will not work for us! Ultimately, we will be accessing our system via a web browser which we require at a minimum that port 80 is open. We recommend simply using the default group which will open our system on all ports. If greater security is required for your specific application a more restrictive security group may be defined and used.

Select the “default” security group (fig. 2.4).

Note: You may receive a Warning stating, “Rules with source of 0.0.0.0/0 allow all IP addresses to access your instance. We recommend setting security group rules to allow access from known IP addresses only.” This is expected and is okay.

Step 6: Configure Security Group

A security group is a set of firewall rules that control the traffic for your instance. On this page, you can add rules to allow specific traffic to reach your instance. For example, if you want to set up a web server and allow Internet traffic to reach your instance, add rules that allow unrestricted access to the HTTP and HTTPS ports. You can create a new security group or select from an existing one below. [Learn more](#) about Amazon EC2 security groups.

Assign a security group: ☐ Create a new security group
☒ Select an existing security group

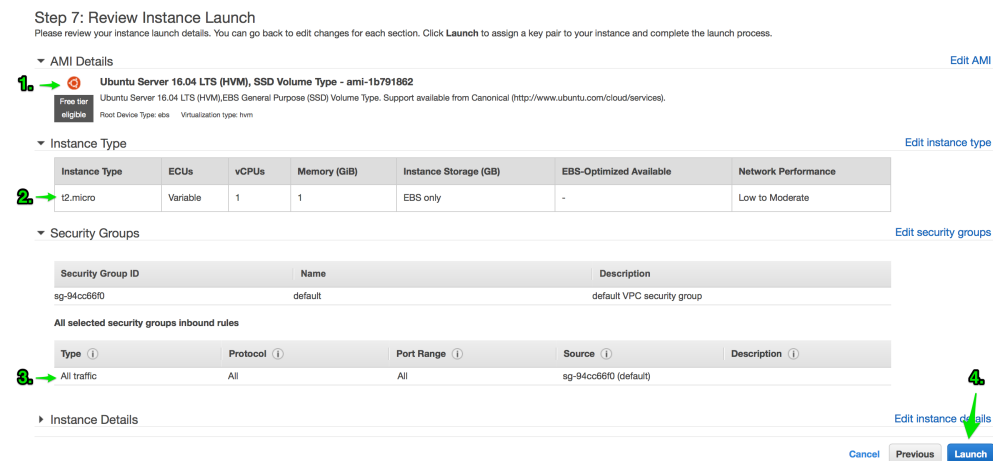


Security Group ID	Name	Description	Actions
sg-94cc66f0	default	default VPC security group	Copy to new

Figure 2.4: Choose the latest stable Ubuntu Server release as AMI

Step 7: Review Instance Launch

Finally, click “Review and Launch.” Here, you see the specific configuration of the EC2 instance you will be creating. Verify that you are creating a `t2.micro` (fig. 2.5, #2) running the latest free tier-eligible version of Ubuntu Server (fig. 2.5, #1) and that it is available to all traffic (fig. 2.5, #3), and then click the Launch button (fig. 2.5, #4).



Step 7: Review Instance Launch

Please review your instance launch details. You can go back to edit changes for each section. Click **Launch** to assign a key pair to your instance and complete the launch process.

AMI Details [Edit AMI](#)

1. **Free tier eligible** Ubuntu Server 16.04 LTS (HVM), SSD Volume Type - ami-1b791862
Ubuntu Server 16.04 LTS (HVM), EBS General Purpose (SSD) Volume Type. Support available from Canonical (<http://www.ubuntu.com/cloud/services>).
Root Device Type: ebs Visualization type: hvm

Instance Type [Edit instance type](#)

2. **t2.micro**

Instance Type	ECUs	vCPUs	Memory (GiB)	Instance Storage (GB)	EBS-Optimized Available	Network Performance
t2.micro	Variable	1	1	EBS only	-	Low to Moderate

Security Groups [Edit security groups](#)

Security Group ID	Name	Description
sg-94cc66f0	default	default VPC security group

All selected security groups inbound rules

Type	Protocol	Port Range	Source	Description
3. All traffic	All	All	sg-94cc66f0 (default)	

Instance Details [Edit instance details](#)

[Cancel](#) [Previous](#) [Launch](#) 4.

Figure 2.5: Review and launch the new instance

Add an SSH Key

In a final confirmation step, you will see a modal titled “Select an existing key pair or create a new key pair.” Select the key pair you previously created. Check the box acknowledging access to that key pair and launch the instance (fig. 2.6).

□ **Note:** If this step is not done correctly, that is, if the correct key pair is not added to the launching instance, the instance will need to be terminated and a new instance will need to be launched. There is now way to add a key pair to a running instance.

You should see a notification that the instance is now running. Click the View Instances tab in the lower right corner to be taken to the EC2 Dashboard Instances pane, where you should see your new instance running.

Examining the Newly Launched Instance

Make note of the IP address of the new instance (fig. 2.7).

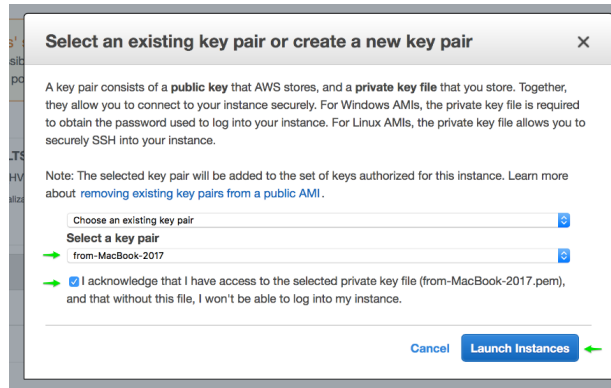


Figure 2.6: Add a key pair to the instance

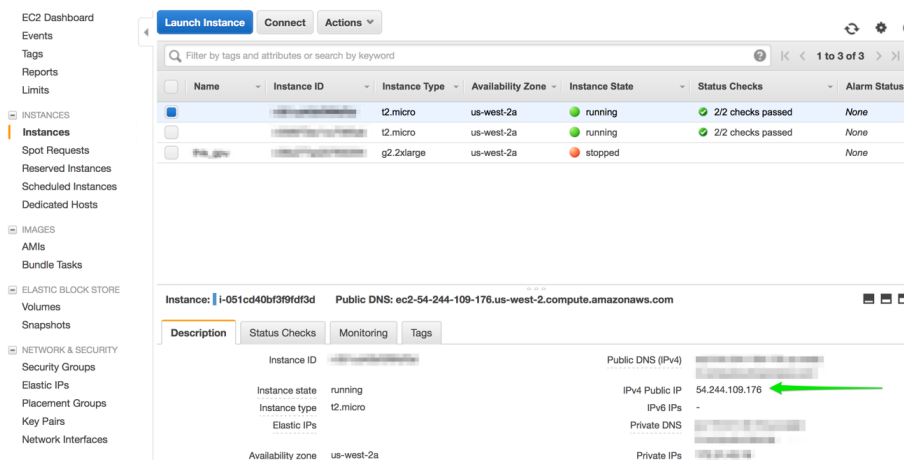


Figure 2.7: Add a key pair to the instance

2.1 Configure Git & Github

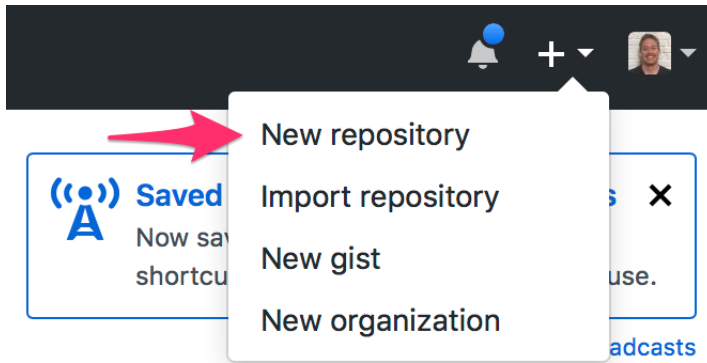
Armed with our new ssh key, we will add the public key to Github so that we can use the key to access our github account.

1. copy the public key
2. choose

The first thing we will do is create a new local git repository and a remote Github repository. The local repository is where we will do all of our work. We will use the remote Github repository to track all of the work that we will do. It is a common beginning misconception to think of git and Github as the same thing. git is a command line tool we will use to track changes to our project. Github is a cloud-based service designed to work with git to help us to make sure our work is always backed up on an additional system.

Create the new repository on Github

At github.com, in the upper-right hand corner, click the plus icon and select “New Repository”.



Give your new repository a name. It does not need to match `engineering-for-data-science`, but it is generally considered a best practice to make sure the repository on Github and the local directory containing your files have the same name. You do not need to provide a description nor create a README file.

Finally, click “Create Repository”.

Create a new repository
A repository contains all the files for your project, including the revision history.

Owner: joshuacook /

Enter Your Project's Name Here.

Great repository names are short and memorable. Need inspiration? How about `super-duper-octo-happiness`.

Description (optional)

☒ **Public**
Anyone can see this repository. You choose who can commit.

☐ **Private**
You choose who can see and commit to this repository.

☐ **Initialize this repository with a README**
This will let you immediately clone the repository to your computer. Skip this step if you're importing an existing repository.

Add .gitignore: **None** | Add a license: **None** ⓘ

Create repository

Make a new directory

We use the `mkdir` command to create a new directory to hold our work.

Bash Command 5 []

```
LOCAL: mkdir -p engineering-for-data-science
```

Change directories to the new directory

Next, use the `cd` (change directory) command to change directories to the new directory we just created.

Bash Command 6 []

```
LOCAL: cd engineering-for-data-science
```

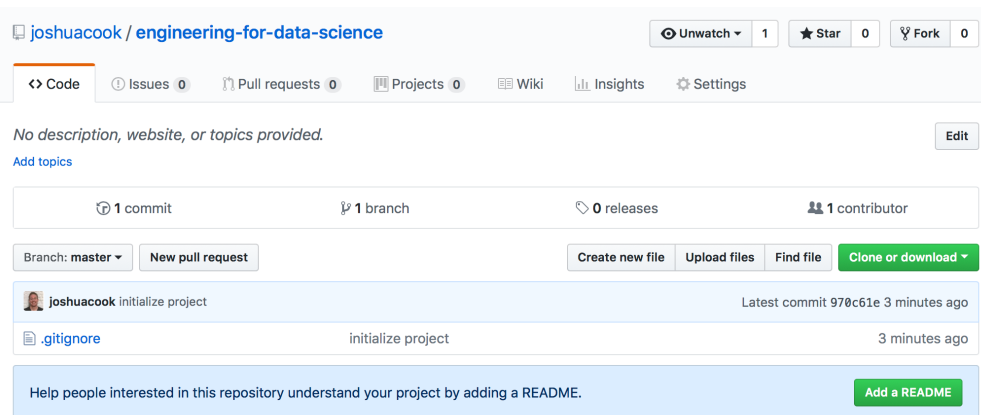
Bash Command 7 []

```
LOCAL: git init
```

```
Initialized empty Git repository in /home/jovyan/engineering-for-data-science/.git/
```

Bash Command 8 []

```
LOCAL: mkdir chapter-01-introduction
```



2.1.1 Git and Github

As you work through this text, you will be developing a series of data science projects. Tracking software development work is typically done using version control software. One of the most popular version control tools is `git`. Additionally, it can be useful to use a version control hosting service as a remote backup for work being tracked using `git`. The remote service we will use is Github.com. In my experience, learners who are new to version control often confuse `git` and Github, so it bears repeating – we will use `git` to track changes we make to our code and Github as a remote backup for these changes.

2.1.2 Configuring Github

We will assume that you have a Github account. Once this has been done, you will need to configure an SSH connection between AWS and Github. This next part may potentially create a confusion. We are actually going to need a new SSH Keypair, this one associated with our AWS instance. This is because it is our AWS instance that will be connecting to Github, not our local machine (See fig. 2.8).

Bash Command 9 []

```
LOCAL: echo ".ssh" > .gitignore
```

Bash Command 10 []

```
LOCAL: ls

chapter-01-introduction
```

Create a New Key Pair

In [at:create_new_ssh_key_remote], you create a new key pair on your remote AWS instance. In [at:ssh_into_new_instance], you connect to your new AWS instance. To do this we will use the IP address we made note of in fig. 2.7. We use SSH to connect to our remote AWS instance. Note that we use the username, ubuntu, the default username for the Ubuntu 16 AMI provided by AWS.

Listing: Create a new SSH Keypair

```
$ ssh ubuntu@54.244.109.176
Welcome to Ubuntu 16.04.2 LTS (GNU/Linux 4.4.0-64-generic x86_64)
```

□ **Note:** The first time you access your EC2 instance, you should see the following message: The authenticity of host '54.244.109.176 (54.244.109.176)' can't be established ... Are you sure you want to continue connecting (yes/no)? This is expected. You should hit <ENTER> to accept or type yes and hit <ENTER>.

In [at:create_new_ssh_key_remote], you create a new SSH Keypair on our remote AWS instance. Again, during the creation of the SSH Keypair, you will be prompted three times. The first asks where you should save the SSH Keypair, defaulting to the .ssh/id_rsa in our home directory. In [at:create_new_ssh_key_remote], you see that this is being done at /home/ubuntu/.ssh/id_rsa². The second and third prompts will ask for a passphrase to be added to the key. For our purposes, leaving this passphrase empty will be fine. In other words, the default options are preferable and you may simply hit <ENTER> three times.

Listing: Create a new SSH Keypair

```
$ ssh-keygen
Generating public/private rsa key pair.
Enter file in which to save the key (/home/ubuntu/.ssh/id_rsa):
Created directory '/home/ubuntu/.ssh'.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /home/ubuntu/.ssh/id_rsa.
Your public key has been saved in /home/ubuntu/.ssh/id_rsa.pub.
The key fingerprint is:
SHA256:ZSpFpGSRgRqlQom8yVBG2dZo1tgkPQdrmUGgMXGDtRY
ubuntu@ip-172-31-43-19
The key's randomart image is:
+---[RSA 2048]----+
|o=XBE/*..o      |
|==+=0**0.      |
|o++o *o. o      |
```

²This should be the same for everyone now, as you should be working on an AWS t2.micro running an Ubuntu system where the user's name is ubuntu

Bash Command 11 []

LOCAL: `ls -la`

```
total 4
drwxr-xr-x  5 jovyan users 170 Apr 24 23:53 .
drwxr-xr-x 20 jovyan users 680 Apr 24 23:51 ..
drwxr-xr-x  2 jovyan users  68 Apr 24 23:52 chapter-01-introduction
drwxr-xr-x 10 jovyan users 340 Apr 24 23:51 .git
-rw-r--r--  1 jovyan users   5 Apr 24 23:53 .gitignore
```

Bash Command 12 []

LOCAL: `git add .gitignore && git commit -m 'initialize project'`

```
*** Please tell me who you are.
```

```
Run
```

```
git config --global user.email "you@example.com"
git config --global user.name "Your Name"
```

```
to set your account's default identity.
Omit --global to set the identity only in this repository.
```

```
fatal: empty ident name (for <jovyan@eb5eb401ee58.(none)>) not allowed
```

```
|o+ . . . +
|      . S
|      .
|
|
|
+-----[SHA256]-----+
```

As before, you can verify the SSH Keypair you just created by displaying the Public Key in your shell ([@lst:cat_pub_key_remote]). Again, you use the `cat` command, which concatenates the contents of `id_rsa.pub` to the shell output.

Listing: Display Public SSH Key

```
$ cat ~/.ssh/id_rsa.pub
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQDQ896GUMgCMAIW79gwF3ojRjcUYCKUKc8b+q
iQlah2jtr7sOK4WRGjkt0y3lCCH0+1UK/GrzY1Y4VxCKoKJDH3G9N5UzyGhlxa/2Ah
kKxzHht1knyh/mkVGqYUhuHpXfxUQAstCFrIdp3G0MDPiko2qeJcBF7JSv1lLMbIuM
XuVU/Mzq6BU+tEogScYytmLckyEe1j8RJ+e5nBURwmkgj3UAN1DzmU/1VwL1ltEpmC
Dl0e14yEXAw8yBwM3GwjahfiBThvBHpsc43HxWrkM8Yi/kdDnvsDZYxU4zhXZPsPab
UY/LfxEod9c6Sui5W8GtAfdi6krnqbzxrKt81Mradh ubuntu@ip-172-31-43-19
```

Add the Public Key to Github

Previously, you added your local SSH public key to your AWS account. Now, you will add your AWS SSH public key to your Github account³. First, access the **Settings** for your account by clicking the profile photo

³<https://help.github.com/articles/adding-a-new-ssh-key-to-your-github-account/>

Bash Command 13 []

```
LOCAL: git config --global user.email "me@joshuacook.me"
```

Bash Command 14 []

```
LOCAL: git config --global user.name "Joshua Cook"
```

in the upper-right corner of any page on Github. Next, in the user settings sidebar, select **SSH and GPG keys**. On the SSH and GPG Keys page, click **New SSH key**. On the next page, give your key a descriptive title e.g. “AWS Feb 2018” and then paste your AWS public key in to the “Key” field . Finally, click **Add SSH key** and confirm your Github password , if prompted.

Learning to read the Bash Prompt

During your work you will no doubt notice that an idle SSH connection may become disconnected and/or unresponsive. Should this happen, simply close the terminal session, launch a new one, and reconnect to the remote instance.

The most important thing is that you are aware of which system your current shell session is connected to. Shell prompts are designed to relay this information to you immediately. If you are new to working with Bash, you may need to train yourself to being aware of the prompt when typing. [@lst:default_AWS_prompt] shows the default AWS Bash prompt. The information contained is the user-name, ubuntu, and the private IP address of the AWS instance. **This is not the public address you use to connect**. What is useful about this, is that we can immediately see that the user is ubuntu. This tells us we are connected to AWS.

Listing: The default AWS Bash prompt

```
ubuntu@ip-172-31-21-89:~$
```

Your local system will no doubt display something different (See [@lst:other_prompt]). Again, the important thing is to take note of what is displayed by the prompt and to learn to associate that prompt with the correct system. As you become a more advanced Bash user, you may wish to personalize your prompt, but for now it is imperative that you learn to read the prompt in order to always know to which system you are connected.

Listing: A local Bash prompt

```
joshuas-macbook-pro:~$
```

Test your SSH Connection to Github

Having added you AWS Public Key to your Github account, you should verify your SSH connection from your AWS instance. In [@lst:verify_github_ssh], we attempt to connect to Github via SSH. As before, we receive a message about the authenticity of the connection. Again, type yes, and continue. If successful, you will see a message telling you have successfully authenticated but that Github does not provide shell access.

Listing: Verify Github SSH Key

```
ubuntu@ip-172-31-21-89:~$ ssh -T git@github.com
The authenticity of host 'github.com (IP ADDRESS)' can't be
established.
RSA key fingerprint is
16:27:ac:a5:76:28:2d:36:63:1b:56:4d:eb:df:a6:48.
```

Bash Command 15 []

```
LOCAL: git add .gitignore && git commit -m 'initialize project'

[master (root-commit) 970c61e] initialize project
1 file changed, 1 insertion(+)
create mode 100644 .gitignore
```

Bash Command 16 []

```
LOCAL: git remote add origin git@github.com:joshuacook/engineering-for-data-science.git
```

```
Are you sure you want to continue connecting (yes/no)? yes
Hi username! You've successfully authenticated, but GitHub does
not provide shell access.
```

2.2 Docker & Docker Compose

Having configured our SSH connections and provisioned a new AWS EC2 instance, it is time to get to the business of building your data science platform. To do this you will use the containerization platform Docker and its Docker Compose tool. While Docker is very easy to use, it can be difficult to understand for the uninitiated. In an earlier work, *Docker for Data Science*⁴, I wrote:

[Using Docker] we add a layer of complexity to our software, but in doing so gain the advantage of ensuring that our local development environment will be identical to any possible environment into which we would deploy the application.

It may be simpler, however, to simply think about using Docker as a way to manage a running process. Your system will be running two processes: an IPython shell and a PostgreSQL server. Were you to not use Docker, you would need to ensure that the AWS instance had all of the libraries required to run both of those processes (and keep those libraries up to date).

Instead, you will let Docker manage the processes using a container for each process. Each respective container will be run using a predefined image built using best practices and ready to run their respective process. The exchange is this: you will take on the cognitive burden of *understanding* what Docker is doing and Docker (and the Docker community) will take over the burden of making sure that your processes run.

2.2.1 Docker Compose

Docker Compose is a tool built for managing an application consisting of multiple containers. Using Docker Compose, it is possible to completely define an application using a simple text file. To make this conversation less abstract, let's have a look at the `docker-compose.yml` file you will use to define your first application (See [lst:docker_compose]).

Listing: Your Data Science Application

```
version: "3"
services:
  ipython_shell:
    image: jupyter/scipy-notebook
  database:
```

⁴<https://www.apress.com/us/book/9781484230114>

Bash Command 17 []

LOCAL: `git push -u origin master`

```
Counting objects: 3, done.
Writing objects: 100% (3/3), 222 bytes | 0 bytes/s, done.
Total 3 (delta 0), reused 0 (delta 0)
To git@github.com:joshuacook/engineering-for-data-science.git
 * [new branch]      master -> master
Branch master set up to track remote branch master from origin.
```

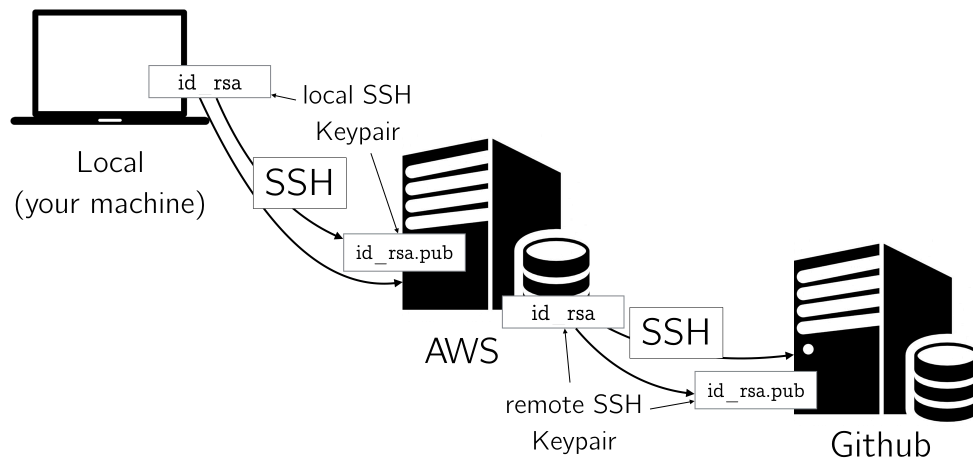


Figure 2.8: SSH Connections

```
image: postgres
volumes:
  - postgres_data:/var/lib/postgresql/data
volumes:
  postgres_data
```

That's it. This simple file completely defines a fully-functioning Data Science Application. In it, we define the two services we need: `ipython_shell` and `database`. These two services are defined using the `jupyter/scipy-notebook` and `postgres` images. When we launch the application, the images will be pulled from Docker Hub into our local memory and then launched. The one other thing we do is create a data volume `postgres_data`. We will use this as the data volume for our database server so that if for some reason we have to shut our system down, we do not lose our data. The data will exist on this volume independent of the services.

□ **Note:** Throughout this text, when discussing infrastructure, I may casually refer to containers, services, and processes. At the risk of annoying your local site reliability engineer, you may treat these as terms as synonymous. Care should be taken, however, not to confuse services/containers/processes and images. An image defines a service, but a service should be thought of as a living and active thing. You may loosely compare the service-image relationship to the object-class relationship in Object-Oriented Programming. A service is a running container defined by an image, just like an object is an instance of a class that exists in memory.

2.2.2 Installing and Configuring Docker

Installing Docker on your AWS instance is a downright trivial process. It consists of running an install script that can be obtained from Docker and then adding your user to the Docker group. In [ist:install_docker], we run these two commands. First, we download the install script from <https://get.docker.com>, then immediately pipe the script into the shell (`| sh`).

□ **Note:** It is generally considered to be a significant security vulnerability to execute arbitrary code obtained from an unknown, or untrusted source. For our purposes, the source (<https://get.docker.com>) is considered trustworthy, we are using SSL to perform the curl, and in practice this is the method I use to install Docker. Still, it may make the security minded more comfortable to curl the script, inspect, and then run it.

Listing: Install Docker via a Shell Script

```
$ curl -sSL https://get.docker.com/ | sh
# Executing docker install script, commit: 1d31602
+ sudo -E sh -c apt-get update -qq >/dev/null
...
```

Client:

```
Version: 18.02.0-ce
API version: 1.36
Go version: go1.9.3
Git commit: fc4de44
Built: Wed Feb 7 21:16:33 2018
OS/Arch: linux/amd64
Experimental: false
Orchestrator: swarm
```

Server:

```
Engine:
Version: 18.02.0-ce
API version: 1.36 (minimum version 1.12)
Go version: go1.9.3
Git commit: fc4de44
Built: Wed Feb 7 21:15:05 2018
OS/Arch: linux/amd64
Experimental: false
```

...

When the script completes there is one last thing to be done. In [ist:add_to_docker_group], you add the ubuntu user to the docker group. By default, the command line docker client will require sudo access in order to issue commands to the docker daemon. You can add the ubuntu user to the docker group in order to allow the ubuntu user to issue commands to docker without sudo.

Listing: Add the Ubuntu User to the Docker Group

```
$ sudo usermod -aG docker ubuntu
```

Finally, in order to force the changes to take effect, you should disconnect and reconnect to their remote system. You can achieve this by typing `exit` or `ctrl-d` and then reconnecting via `ssh` to your EC2 instance.

2.2.3 Installing and Configuring Docker

Recall that regardless of your local operating system, you are working on an AWS EC2 Instance running the Linux variant, Ubuntu. As such, `docker-compose` can be installed using the instructions provided here:

<https://github.com/docker/compose/releases>, which are written specifically for Linux machines. As of the writing of this book, this consists of two steps.

In [①:curl_docker_compose], you use `curl` to retrieve the `docker-compose` binary from Github. As of the writing of this book, the latest version of `docker-compose` was 1.19.0. You should retrieve the latest version from the above url.

Listing: Retrieve `docker-compose` binary from Github

```
$ sudo curl -L https://github.com/docker/compose/releases/download/1.19.0/docker-compose-`uname -s`-`un
```

In [①:chmod_docker_compose], we use the `chmod`⁵ utility to allow `docker-compose` to be executed (+x).

Listing: Enable Docker Compose to be Executed

```
$ sudo chmod +x /usr/local/bin/docker-compose
```

Finally, in [①:docker_compose_version], we check the version of `docker-compose` against what we expect to have installed.

```
$ docker-compose -v
docker-compose version 1.19.0, build 9e633ef
```

2.3 Summary

TODO

2.4 Exercises

1. What went wrong here? What should you do?

```
failed to register layer: Error processing tar file(exit status 1):
write /usr/bin/python2.7: no space left on device
```

2. What went wrong here? What should you do?

```
docker: Error response from daemon: driver failed programming external
connectivity on endpoint cocky_swartz (08124b75d2f031def6d36c6bc819549c009
391e3bd76f3fe3b4e06e11be6fbad): Bind for 0.0.0.0:80 failed: port is
already allocated.
```

3. What does this command do?

```
curl -sSL https://get.docker.com | sh
```

4. What are two ways to display the contents of a text file from the command line?
5. What are the two modes in vim?
6. How do you save the changes in a file in vim?
7. How would you find all files in your current directory that contain the string "bash"?
8. What are the steps to launching a Jupyter Notebook Server on a running AWS instance?

⁵The unix "change mode" utility. I pronounce it "shmod".