

Set up a virtual server on AWS (in anticipation of hosting Shiny apps)

Ronald (Ryy) Glenn Thomas

2023-06-20

Table of contents

1	Introduction	2
2	Hosting	2
2.1	Select a hosting service	3
2.2	AWS Working Environment	4
2.2.1	Ssh key pair	5
2.2.2	Firewall	6
2.2.3	Static IP address	6
2.2.4	Domain Name	6
2.3	Select and launch instance	7
3	Access server	8
4	Appendix: Tip 1	8



Photo by Nathan Waters on
Unsplash

1 Introduction

i Motivation for this post:

Ok! I've got my Shiny app running just the way I want it. Works great! Now, how do I get it up on the web and shared with my collaborators?

Assuming we have a working shiny app, we often next need to also address the task of how to host the app on the web to share with our collaborators. How to accomplish this? Below we describe a recommended process we've found to be straightforward. In this post, we'll describe how to 'spin up' a server on Amazon Web Service EC2 using the EC2 console. In a future post we'll show how, through the application of Docker, R, Shiny, and Caddy (webserver) functionality we can have a fully functional and secure app available on the web in just a few steps.

2 Hosting

No matter what procedure we use to host a shiny app (named, for example, `power1_shiny`) online we'll need to complete the following set of tasks one way or another:

Pre-launch tasks:

1. obtain a static IP address
2. obtain a domain name (e.g. `rgtlab.org`) and associate it with the IP address
3. configure a virtual server. Configuration entails selecting number of CPUs, amount of memory, OS, etc.
4. define a security model, aka a firewall for the server
5. associate the IP, domain name and firewall with the server

Post-launch tasks:

1. install and configure a webserver
2. obtain and install an SSL certificate (to allow encrypted communication)
3. setup an authentication method (password protection for app access)
4. configure a reverse proxy method i.e. translate https (port 443) requests to Shiny (port 3838). This avoids the need to explicitly name the port Shiny is listening on i.e. requiring URLs like `https://rgtlab.org:3838/power1_shiny`

i Not to worry:

At first glance these requirements can appear daunting, but on closer inspection all can be met with relative ease through the use of right tools.

i Actually ...:

Technically, if the goal is simply to get the app up on the web, its not required to have a static IP, or a domain name, or a firewall, or an authentication method or an SSL certificate, or even a reverse proxy. but if these elements of the process are skipped the server will only be able to communicate via the unencrypted HTTP protocol and the site URL will be something like `111.222.333.444:3838/power1_shiny`, and anyone with the URL will be able to reach the site. Also the IP address will change everytime the server is rebooted.

2.1 Select a hosting service

There are a number of cloud based server hosting options to choose from: for example Microsoft Azure, Oracle, Google Cloud, Amazon AWS EC2, Digital Ocean or Hetzner to name a few. Each has their own approach to setting up a custom virtual server. Several have free or low-cost service tiers available.

In this post we'll provide a step-by-step description of a process using Amazon Web Services Elastic Compute Cloud (AWS EC2) infrastructure.

AWS is, in our view, a reasonable choice for setting up a small custom server. Its not the cheapest option, but the system is well documented and, in our experience, reliable.

The first step is to get set up with AWS. To start, open the EC2 console by visiting the URL:

```
https://aws.amazon.com/console
```

(see margin figure)

In the console window choose regional service. For us its “N. California”.

Next create an account, or sign in, and once you're logged in navigate to the EC2 dashboard. Its through this dashboard (aka console) that we'll define the parameters for the type of server to launch and the mechanisms for communicating with it. We'll refer to these as “Pre-Launch” tasks.

2.2 AWS Working Environment

Along with selecting a server we need to set up a working environment. We recommend setting up the working environment before launching the server, as it saves some back and forth with the console, but the order is not critical. The working environment consists of four main components:

1. A secure shell (ssh) key-pair to allow remote and secure login to the virtual server once its launched.
2. A firewall or security model which will restrict server access to only secure connections. The firewall closes off all incoming traffic except through those ports specifically named.

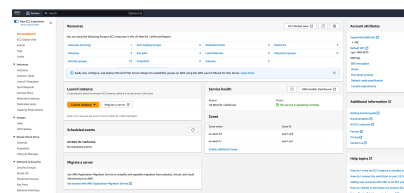


Figure 1: AWS console

3. A static IP address. A static IP is required for maintaining the link between the domain name and the server when rebooting. (The default is for the instance/server to be assigned a new IP address each time its rebooted). and
4. A domain name, say **rgtlab.org**. A domain name is not required but will facilitate collaborator access by not needing to use the IP address directly.

These working environment components are not directly tied to any specific server. In fact, you can define multiple instances of each component. The only requirement is that you pick one of each to associate with each server.

2.2.1 Ssh key pair

In order to securely communicate with the server we need to exchange an ssh key pair with AWS. The pair consists of a **private** key and a **public** key. We can identify an ssh key pair in one of two ways in EC2. Either, generate the pair locally, on our workstation and upload the public key to EC2, or have EC2 generate the key pair and download the private key.

For the first option we create a directory on our workstation to hold the keys and navigate to it, e.g. `~/.ssh`. In the `~/.ssh` directory generate the keys with the command

```
ssh-keygen -m PEM
```

“PEM” defines the key format. More information on public key authentication can be found [here](#). In the interactive dialog that follows name the key prefix something like **power1_app**. The dialog will ask for a passphrase. You can enter a phrase for an additional level of security, but its not required. The **ssh-keygen** program will generate two files: **power1_app.pem** and **power1_app.pub**

To complete the process return to the EC2 dashboard and select **Actions** and then **Import key pair** in the left panel. Enter the name **power1_app** and select the **Browse** button. Navigate

to the file `power1_app.pub` in the directory `~/.ssh` and select the **Import key pair** button at the bottom of the page.

For the second approach select **Create key pair** button in the upper right of the console page. A form will appear and ask for a name. Enter something like `power1_app`. Select **RSA** for key pair type and `.pem` for key file format. The keys will be created and the private key `power1_app.pem` will be downloaded to our local machine and should be placed in the `~/.ssh` directory. Lastly, change the access permissions for the private key with the following command:

```
sudo chmod 600 power1_app.pem
```

2.2.2 Firewall

To create a firewall click on **Network and Security** settings in the left hand panel. Choose **Create security group** and select **Allow SSH traffic** and **Allow HTTPS traffic**. This will create a firewall that leaves open only ports 22 and 443, for `ssh` and `https` incoming traffic respectively. Lastly, name the security group something like `power1_firewall`.

2.2.3 Static IP address

You can use the `elastic IP` service to get a static IP. Navigate to **Network and Security** again and select **Allocate Elastic IP**. An IP will be assigned from the EC2 pool of available IPv4 IP addresses e.g. 13.57.139.31.

2.2.4 Domain Name

To obtain a dedicated domain name leave the EC2 dashboard and go to Amazon route 53 service to select a domain name and associate it with our static IP.

Once a domain name is acquired, e.g. `rgtlab.org`, you can associate it with any IP address, static or dynamic. This can be done via the Route 53 service. For example, to associate

domain name `rgtlab.org` with the elastic IP 13.57.139.31 do the following in Route 53:

- click on **hosted zones** in the side panel
- click on `rgtlab.org` in center panel
- click on checkbox for `rgtlab.org type=A` line
- then click on edit record in right panel
- change IP address to the assigned static IP .

2.3 Select and launch instance

2. From **Quick Start** in the EC2 dashboard select an operating system for the server. We recommend the **Ubuntu** OS. Ubuntu is a mature Linux distribution based on Debian Linux. Click the **Ubuntu** button. (see margin figure)

- Name the server, say `power1`

3. Next choose an instance **type**, e.g. `t2-micro`. Different instance types are combinations of number of processors, memory, storage capacity, and network performance.

4. select **Configure Instance Details**

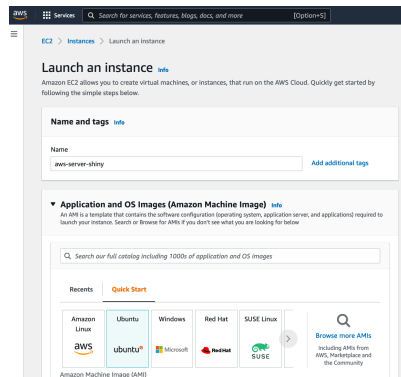
5. choose a Key pair e.g. select `power1_app` from your environment.

6. Add security group, e.g. use `power1_firewall` from your environment.

7. choose a storage amount. e.g. enter 30 GB of EBS **General Purpose (SSD)** or **Magnetic storage**. Thirty GBs is the maximum allowed in the ‘Free tier’ of servers on AWS. In our experience smaller disk sizes can lead to problems.

8. click **Launch Instance**

to launch the server.



3 Access server

On your laptop log into server with

```
ssh -i "~/.ssh/power1_app.pem" ubuntu@rgtlab.org
```

4 Appendix: Tip 1

Tip 1.

For convenience, construct a `config` file in `~/.ssh` as:

```
Host rgtlab.org #domain name
HostName 13.57.139.31 # static IP
User ubuntu # default user on ubuntu server
Port 22 # the default port ssh uses
IdentityFile ~/.ssh/power1_app.pem # private key
```

then we can ssh into the new server with the abbreviated command

```
sh> ssh rgtlab.org
```