# Deploy ML models in Production using AWS

Whenever we need to deploy our application on AWS, we need to change the host and port in which the following flask application is running. By default, the host address will be at '127.0.0.1' in other words localhost and port number by default is 5000. Now we need to change the default port number to 8080 and host address to '0.0.0.0'. For this, the following code needs to be added in place of app.run().
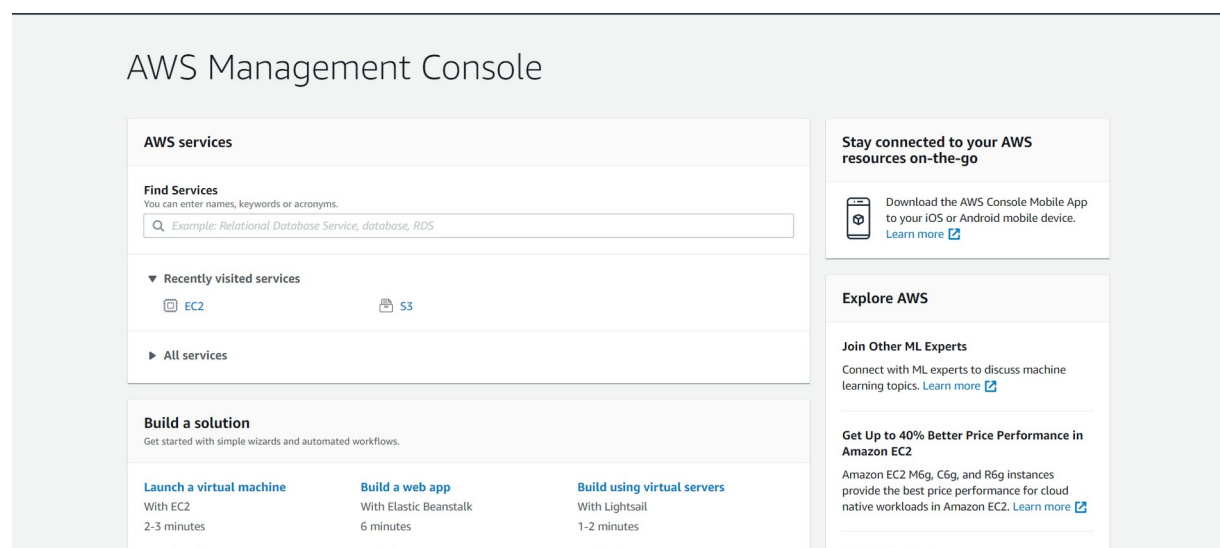
```
1  if __name__=='__main__':
2      app.run(host='0.0.0.0',port=8080)
```

Now if you run the flask application it will run but if you go to http://0.0.0.0:8080/ it will give you the error. The reason for this is address 0.0. 0.0 is a non-routable meta-address used to designate an invalid, unknown, or non-applicable target. Don't care about this error just leave about this and go on to the next step.

## Getting Started with AWS:-

Create free student account at: https://aws.amazon.com/cn/education/awseducate/

Now it's all set, the remaining thing to reach our goal is to deploy flask application to AWS Cloud. For this, we need to have an AWS account. If you are new to AWS, You need to create an AWS account by providing your credit card details to avail of free tier services free for one year. After creating the account you need to go to AWS Management Console and need to sign in to AWS by using the account that you have just created. Once you sign in the webpage looks similar to the below image.



**Creating an AWS EC2 Instance:**

Now we need to go to the search bar in the AWS Management Console and search for EC2 and need to click on running instances. Click on Launch Instance to create the new EC2 instance. Now we need to follow 4 steps to create an instance.

## I. Choose an Amazon Machine Image (AMI):-

We need to choose AMI which means choosing an OS for the server. In this case, we can choose the UBUNTU 18.04 which is a free tier and proceed next.



## II. Choose Instance Type:-

Now we need to choose instance type for our server. In this case, we can go with a free tier instance. If yours is a company you can choose the paid ones. Here is the snapshot of the choose instance type.



Once you selected the instance we need to click on "Review and Launch".

## III. Edit Security Groups:-

Once the above steps are done, you need to click on the "Edit Security Groups" in the "Review and Launch" page. Once if you go "Edit Security Groups" Page you need to create a security group.

What is the reason behind creating Security Groups?

It helps us to edit which kind of traffic that should be allowed while accessing the instance and it helps us to decide which IPs can access our application which allows us to restrict the usage.

In this case, we are allowing all traffic and source to anywhere so that every network from anywhere in the world can be able to access our application. Once it is done we need to click on "Review and Launch". Here is the Snapshot of editing security groups.



## IV. Creating a Private Key Pair:-

Once you followed with the above steps you will be able to click on the launch button in "Review Instance Launch". You will be getting a pop-up for choosing a private key pair. Let's click on create a new private key pair and choose a keypair name. In this case, I am choosing it as the iris. Now we need to click on the download key pair button to download it.

An Important Note is mentioned here below:-



Now store this file in your project folder and click on launch instances. Here is the snapshot of how to create a private key pair.

In this way, we can be able to create an EC2 Instance. Now we need to go to view instances where you could be able to find the instance you have created.

**Downloading Required Softwares:**

The Required Softwares to be downloaded which helps in connecting the local machine with the AWS EC2 instance. The list of Softwares to do be downloaded are

1. WinSCP:- It is a free and open-source SFTP, FTP, WebDAV, Amazon S3, and SCP client for Microsoft Windows. Its main function is secure file transfer between a local and a remote computer.
2. Putty:- It is a free and open-source terminal emulator, serial console, and network file transfer application. It supports several network protocols, including SCP, SSH, Telnet, rlogin, and raw socket connection. It can also connect to a serial port.

**Converting private key from PPM to PPK Format using PuttyGen:**

You need to convert a private key from PPM format which you have downloaded previously to PPK format using PuttyGen.It will be very much useful in the future for getting access to transfer files from local machine to Ubuntu Server. To convert firstly you need to open PuttyGen Software which is already installed default when you installed Putty. The Snapshot of PuttyGen looks like shown below.

Now you need to click on load and add the PPM file you recently downloaded. Then you need to click on "Save private key" and store it. Now close the PuttyGen window and continue to follow the next step.

**Transferring Project files using WinSCP:**

Now you need to transfer project files from local machine to ubuntu server using WinSCP Software. Firstly you need to go to the AWS Management Console and right-click on the instance where you will find many options. Now click on connect and you will be able to get the similar page shown below:-

Now you need to copy the Public DNS which is shown in point no 4.Then you need to open WinSCP and need to paste this DNS address in the host input value. The window looks similar to below shown one.



Now click on advanced options and go to "Authentication" and insert the created private key in PPK format. The window looks similar to below shown one.

Now click on OK and go to log in. Now you can be able to connect with the ubuntu server. Now select the project files you wanted to add to the server and drag and drop to the right side section in the window. The file transfer starts and after a while files will be transferred. The window looks similar to below shown one.

**Connect to the Instance using Git Bash:**

Now required file transfer is done and the remaining thing is we need to get an access terminal to compute the files. For that open git bash and navigate to the project source. Once you have done with that you need to go to instances in AWS Management Console click on connect then you will be able to find a similar page as shown below.



Now you need to go to the example section and copy the code given there. Now paste the command on git bash and click enter and you will be getting a question to accept the security. Enter "yes" and click on enter. Now you are connected to the instance from your local machine. The git bash looks similar to this when you successfully connected to the server instance.

```
  ubuntu@ip-172-31-42-158: ~                                    —    □    ✕
$ ssh -i "iris.pem" ubuntu@ec2-3-20-236-91.us-east-2.compute.amazonaws.com
Welcome to Ubuntu 18.04.4 LTS (GNU/Linux 5.3.0-1023-aws x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

  System information as of Sun Aug  9 14:58:32 UTC 2020

  System load:  0.0                Processes:           91
  Usage of /:   14.0% of 7.69GB    Users logged in:     0
  Memory usage: 15%                IP address for eth0: 172.31.42.158
  Swap usage:   0%


0 packages can be updated.
0 updates are security updates.


Last login: Sun Aug  9 14:50:17 2020 from 175.101.68.181
To run a command as administrator (user "root"), use "sudo <command>".
See "man sudo_root" for details.

ubuntu@ip-172-31-42-158:~$ |
```

Install the Required Libraries and run the Flask Application:-

Now once the terminal is ready now we need to run some commands before running the flask app. Firstly we need to update the OS and update pip to install libraries in the future. The command for performing above process is

```
sudo apt-get update && sudo apt-get install python3-pip
```

Now you need to install all required libraries to run the flask application. The required commands for this application are as follows:-

```
pip3 install numpy
pip3 install flask
pip3 install sklearn
```
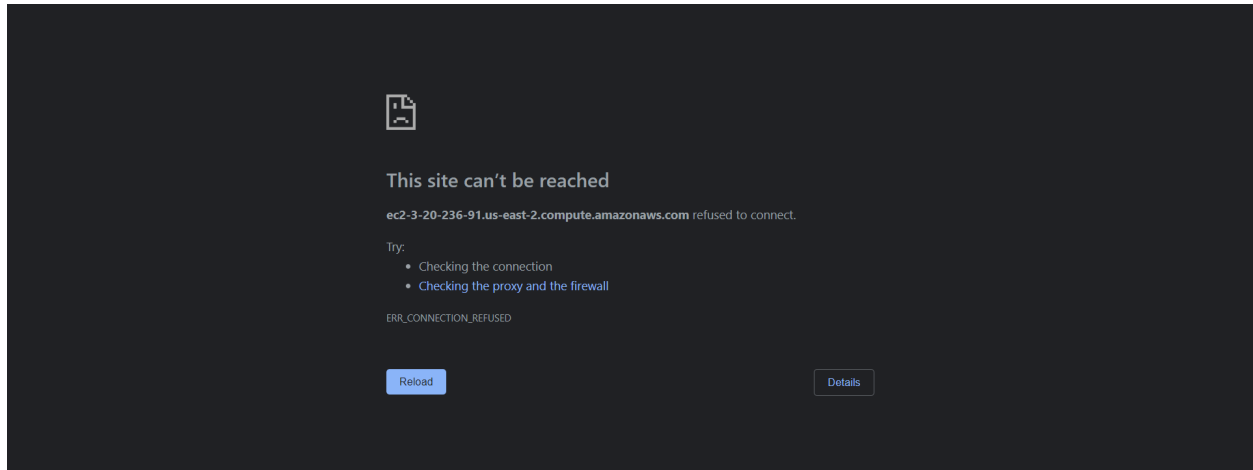
Now all the libraries are installed now we need to run the application by entering following command on git bash.

```
python3 app.py
```

Now the application started running. You can view the application at the URL as <your DNS>::<your port no>.In my case it is http://ec2-3-20-236-91.us-east-2.compute.amazonaws.com:8080/

# Important Note:-

Now your application is hosted but there is a problem you will face once you close your git bash. The problem is once if you close git bash and you reload the URL your screen looks similar to below.



**How to Solve this problem?**

To solve this issue we need to use Screen. Screen or GNU Screen is a terminal multiplexer. In other words, it means that you can start a screen session and then open any number of windows (virtual terminals) inside that session. Processes running in Screen will continue to run when their window is not visible even if you get disconnected.

For this, we need to run the following command on the git bash.

```
screen -R deploy python3 app.py
```

You can exit (detach) the screen by pressing "Ctrl-A" and "d" but the server will keep running.

After you detach the screen, let say you are disconnecting your **SSH** session and going home. In your home, you start to **SSH** again to your server and you want to see the progress of your download process. To do that, you need to restore the screen. You can run this command:

```
screen -r
```

And you will see that the process you left is still running.

When you have more than **1 screen** session, you need to type the screen session **ID**. Use screen **-ls** to see how many screen are available.

```
screen -ls
```