

Informe:

Práctica 4: Redes y Conectividad Avanzada y Contenedores

Jordi Blasco Lozano
Infraestructuras y Servicios Cloud
Universidad de Alicante
29 de octubre de 2025

Resumen

En esta práctica se ha diseñado y desplegado una arquitectura cloud para un sistema de recomendaciones, empleando segmentación de red en una VPC, balanceo de carga, despliegue canario y transición a contenedores. Se han creado subredes públicas y privadas distribuidas en dos zonas de disponibilidad, configurado reglas de control de acceso (SGs y NACLs), y se implementó un Application Load Balancer para distribuir tráfico entre versiones V1 y V2 del modelo.

Índice

| | |
|--|-----------|
| 1. Configuración de la VPC | 2 |
| 1.1. Segmentación y subredes | 2 |
| 1.2. Creación simplificada de la VPC | 2 |
| 1.3. NAT Gateway y Endpoint S3 automatizados | 2 |
| 2. Configuración del S3 | 3 |
| 2.1. Política del bucket | 3 |
| 3. Configuración de la seguridad | 4 |
| 3.1. NACL pública | 4 |
| 3.2. NACL privada | 4 |
| 3.3. Security groups | 5 |
| 4. Instancias | 6 |
| 4.1. Creación de la AMI | 6 |
| 4.2. Comprobación User Data | 7 |
| 5. Configuración del Balanceador de Carga (ALB) | 8 |
| 5.1. Target groups | 8 |
| 5.2. ALB y pruebas | 9 |
| 6. Transición a contenedores | 11 |
| 6.1. Nueva regla SG Backend | 11 |
| 6.2. Editar listener | 12 |
| 6.3. Pruebas endpoint del ALB | 12 |
| 7. Preguntas de reflexión | 13 |
| 7.1. Seguridad y NACLs | 13 |
| 7.2. WAF vs. NACL | 13 |
| 7.3. VPC Endpoints | 13 |
| 7.4. Redes de Contenedores | 13 |
| 7.5. Rollback Rápido | 13 |

1 Configuración de la VPC

1.1 Segmentación y subredes

Para el despliegue del sistema de recomendaciones se requiere una estructura específica en la VPC. La segmentación se realiza mediante bloques CIDR dedicados para cada subred, con máscara /24. El enunciado exige utilizar 10.0.10.0/24 para la subred pública principal y 10.0.20.0/24 para la privada principal. Al emplear dos zonas de disponibilidad, se crean cuatro subredes:

- Pública 1: 10.0.10.0/24
- Pública 2: 10.0.30.0/24
- Privada 1: 10.0.20.0/24
- Privada 2: 10.0.40.0/24

Cada subred permite hasta 256 direcciones IP individuales gracias a la máscara /24.

Number of Availability Zones (AZs) [Info](#)
Choose the number of AZs in which to provision subnets. We recommend at least two AZs for high availability.
1 | 2 | 3

► Customize AZs

Number of public subnets [Info](#)
The number of public subnets to add to your VPC. Use public subnets for web applications that need to be publicly accessible over the internet.
0 | 2

Number of private subnets [Info](#)
The number of private subnets to add to your VPC. Use private subnets to secure backend resources that don't need public access.
0 | 2 | 4

▼ Customize subnets CIDR blocks

Public subnet CIDR block in us-east-1a
10.0.10.0/24 256 IPs

Public subnet CIDR block in us-east-1b
10.0.30.0/24 256 IPs

Private subnet CIDR block in us-east-1a
10.0.20.0/24 256 IPs

Private subnet CIDR block in us-east-1b
10.0.40.0/24 256 IPs

Figura 1: CIDR subredes

1.2 Creación simplificada de la VPC

Para facilitar el proceso y evitar configuraciones manuales, AWS proporciona el botón de "VPC and more". Seleccionando esta opción durante la creación, es posible definir directamente los bloques CIDR personalizados para cada subred desde el menú desplegable, ajustando el tamaño y las direcciones antes de crear la VPC.

1.3 NAT Gateway y Endpoint S3 automatizados

Adicionalmente, desde el mismo "VPC and more" se permite añadir el NAT Gateway para las subredes privadas y el endpoint de S3. Esto implica que tanto el NAT Gateway como el endpoint de S3 quedarán automáticamente conectados a las tablas de rutas necesarias de cada subred, sin necesidad de modificar manualmente las tablas de ruta después de la creación de la VPC. Cuando salgamos de la pestaña de configuración obtendremos el siguiente esquema de nuestra VPC, listo y funcional.

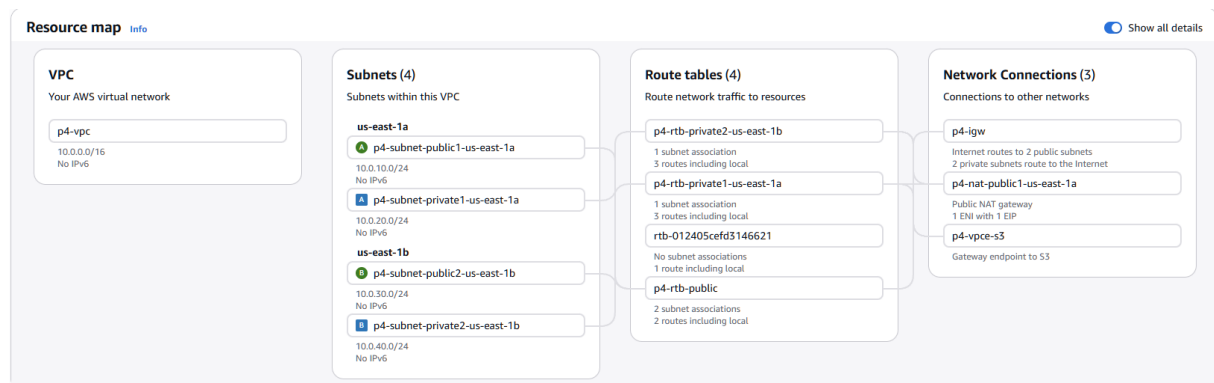


Figura 2: Esquema VPC

2 Configuración del S3

Para descargar los archivos y los modelos nos tendremos que conectar a un S3 que almacene estos datos. Esto lo haremos mediante el S3 endpoint que conectamos en el paso anterior.

Subiremos los archivos siguiendo la estructura del user data proporcionado, de forma que tengamos estos tres archivos (el modelo.pkl simplemente servirá como archivo de prueba sin valor, no se ejecutará, pero sí descargará):

- /backend/model_v1_prod.py
- /backend/model_v2_canary.py
- /backend/modelo.pkl

Summary
Destination
[s3://p4-buquet/backend/](#)

Files and folders | **Configuration**

Files and folders (3 total, 2.1 KB)

| Name | Folder |
|----------------------------------|--------|
| modelv1prod.py | - |
| modelv2canary.py | - |
| modelo.pkl | - |

2.1 Política del bucket

Como en la práctica anterior, en esta también tenemos que permitir el acceso al bucket usando la política, de forma que permita accesos de descarga únicamente si utilizamos nuestro endpoint específico. Para esto definiremos el siguiente json:

Bloque 1: política-s3.json

```
1 {
2   "Version": "2012-10-17",
3   "Statement":
4   [
5     {
6       "Sid": "AllowBackendFolderViaVPCE",
7       "Effect": "Allow",
8       "Principal": "*",
9       "Action": ["s3:GetObject"],
10      "Resource": "arn:aws:s3:::p4-buquet/backend/*",
11      "Condition":
12      {
13        "StringEquals":
14        {
15          "aws:SourceVpce": "vpce-0aa42de173d26c739"
16        }
17      }
18    }
19  ]
20 }
21 }
```

3 Configuración de la seguridad

La seguridad de esta arquitectura está estructurada en varias capas y cada componente (subred pública, subred privada, bastión, backend, ALB) tiene reglas específicas tanto en NACLs como en Security Groups, diseñadas para minimizar la superficie de ataque y garantizar que solo el tráfico necesario fluya entre los recursos.

3.1 NACL pública

Para las subredes públicas (10.0.10.0/24 y 10.0.30.0/24), las NACLs permiten la entrada de peticiones HTTP (puerto 80) desde cualquier origen para el ALB, y entrada SSH (puerto 22) únicamente desde mi IP pública, facilitando la gestión y el acceso seguro. Además, se permite la entrada de puertos efímeros provenientes de las subredes privadas, para soportar las respuestas a las conexiones originadas en las instancias backend. En cuanto a la salida, la NACL pública autoriza puertos efímeros hacia cualquier destino y la salida HTTP (puerto 80) hacia Internet para el ALB. Todo tráfico no contemplado se deniega, permitiendo sólo los flujos estrictamente necesarios.

| Nº regla | Tráfico | Protocolo | Puerto | Origen | Acción |
|----------|---------|-----------|------------|-----------------------------|--------|
| 100 | HTTP | TCP | 80 | 0.0.0.0/0 | Allow |
| 110 | SSH | TCP | 22 | Mi IP pública | Allow |
| 120 | efímero | TCP | 1024-65535 | 10.0.20.0/24 y 10.0.40.0/24 | Allow |
| * | All | All | All | All | Deny |

Cuadro 1: NACL Pública - Reglas de entrada (pública-1 y pública-2)

| Nº regla | Tráfico | Protocolo | Puerto | Destino | Acción |
|----------|------------|-----------|------------|-----------|--------|
| 100 | efímero | TCP | 1024-65535 | 0.0.0.0/0 | Allow |
| 110 | HTTP Resp. | TCP | 80 | 0.0.0.0/0 | Allow |
| * | All | All | All | All | Deny |

Cuadro 2: NACL Pública - Reglas de salida (pública-1 y pública-2)

3.2 NACL privada

Las NACLs de las subredes privadas (10.0.20.0/24 y 10.0.40.0/24) son aún más restrictivas. Permiten la entrada solamente desde las subredes públicas por el puerto 8000 (tráfico backend API), SSH desde la IP privada del bastión para gestión, y entrada a través del endpoint de S3 para descargas. En cuanto a la salida, estas NACLs permiten puertos efímeros únicamente hacia las subredes públicas, para que las respuestas de las aplicaciones backend puedan volver correctamente al ALB, y salida por el puerto 443 hacia el endpoint S3 para acceso seguro a los datos almacenados. Todo lo demás queda explícitamente bloqueado.

| Nº regla | Tráfico | Protocolo | Puerto | Origen | Acción |
|----------|-------------|-----------|--------|-----------------------------|--------|
| 100 | Backend In | TCP | 8000 | 10.0.10.0/24 y 10.0.30.0/24 | Allow |
| 110 | SSH | TCP | 22 | IP privada Bastión | Allow |
| 120 | S3 Endpoint | TCP | 443 | 0.0.0.0/0 | Allow |
| * | All | All | All | All | Deny |

Cuadro 3: NACL Privada - Reglas de entrada (privada-1 y privada-2)

| Nº regla | Tráfico | Protocolo | Puerto | Destino | Acción |
|----------|-------------|-----------|------------|-----------------------------|--------|
| 100 | efímero | TCP | 1024-65535 | 10.0.10.0/24 y 10.0.30.0/24 | Allow |
| 110 | S3 Endpoint | TCP | 443 | 0.0.0.0/0 | Allow |
| * | All | All | All | All | Deny |

Cuadro 4: NACL Privada - Reglas de salida (privada-1 y privada-2)

3.3 Security groups

A nivel de Security Group, cada recurso se configura con reglas mínimas y específicas. El ALB solamente acepta conexiones HTTP desde cualquier origen y, voy a poner también que acepte SSH desde mi IP por si necesitara conectarme. Las instancias backend restringen su entrada al puerto 8000 proveniente del SG del ALB y SSH desde el SG del bastión, evitando cualquier acceso no autorizado. Finalmente, las reglas de salida permiten el flujo de datos hacia el ALB, Internet (para descargas), y hacia el endpoint S3. El bastión está configurado para aceptar conexiones SSH únicamente desde mi IP, y permite que sus conexiones salientes alcancen las instancias backend.

| Tráfico | Protocolo | Puerto | Origen | Acción |
|---------|-----------|--------|---------------|--------|
| HTTP | TCP | 80 | 0.0.0.0/0 | Allow |
| SSH | TCP | 22 | Mi IP pública | Allow |

Cuadro 5: SG ALB - Reglas de entrada

| Tráfico | Protocolo | Puerto | Destino | Acción |
|--------------|-----------|------------|------------|--------|
| HTTP/Backend | TCP | 8000 | SG-Backend | Allow |
| efímero | TCP | 1024-65535 | 0.0.0.0/0 | Allow |

Cuadro 6: SG ALB - Reglas de salida

| Tráfico | Protocolo | Puerto | Origen | Acción |
|-------------|-----------|--------|------------|--------|
| API/Backend | TCP | 8000 | SG-ALB | Allow |
| SSH | TCP | 22 | SG-Bastión | Allow |

Cuadro 7: SG Backend - Reglas de entrada

| Tráfico | Protocolo | Puerto | Destino | Acción |
|-------------|-----------|------------|-----------|--------|
| efímero | TCP | 1024-65535 | 0.0.0.0/0 | Allow |
| S3 Endpoint | TCP | 443 | 0.0.0.0/0 | Allow |

Cuadro 8: SG Backend - Reglas de salida

| Tráfico | Protocolo | Puerto | Origen | Acción |
|---------|-----------|--------|---------------|--------|
| SSH | TCP | 22 | Mi IP pública | Allow |

Cuadro 9: SG Bastion - Reglas de entrada

| Tráfico | Protocolo | Puerto | Destino | Acción |
|---------|-----------|------------|------------|--------|
| SSH | TCP | 22 | SG-Backend | Allow |
| efímero | TCP | 1024-65535 | 0.0.0.0/0 | Allow |

Cuadro 10: SG Bastion - Reglas de salida

4 Instancias

Después de haber configurado todo el entorno debemos de lanzar las instancias. Para ahorrar tiempo instalaré en una instancia de la subred publica todas las dependencias, tales como: python awscli, etc. Esto lo hare para poder crear una AMI de esta instancia con todas las dependencias instaladas. Posteriormente creare una instancia privada a partir de la AMI, a esta instancia me conectaré mediante el bastión y probare que todos los comandos del user data me funcionen correctamente. Una vez que me asegure de que el user data funciona podremos configurar el balanceador de carga.

4.1 Creación de la AMI

Una vez tengamos todo instalado creamos la AMI. He tenido que cambiar un poco los comandos, sobre todo el de awscli porque no me dejaba instalarlo de la forma en la que estaba en la práctica (no encontraba el paquete). He instalado únicamente lo siguiente en la AMI:

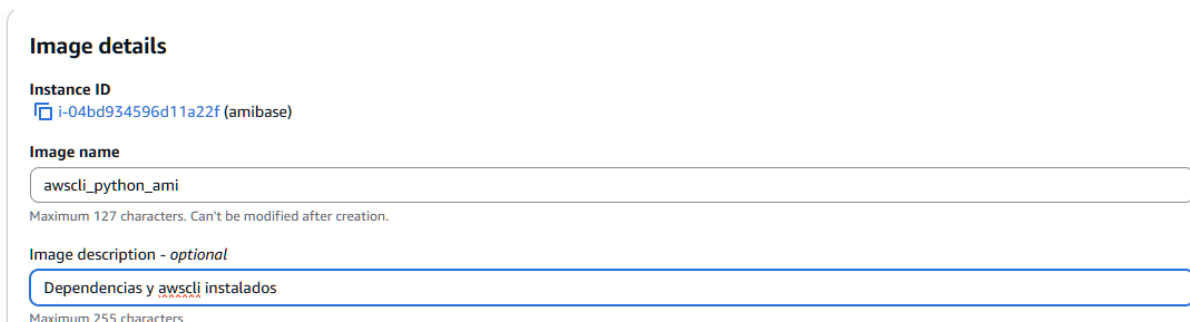
Bloque 2: dependencias.sh

```
1 sudo apt-get update -y
2 sudo apt-get install -y python3-pip python3-venv git nodejs npm unzip curl
3 sudo apt install python3-flask
4 pip3 install joblib
5 cd /tmp
6 curl "https://awscli.amazonaws.com/awscli-exe-linux-x86_64.zip" -o "awscliv2.zip"
7 unzip awscliv2.zip
8 sudo ./aws/install
9 cd -
```

Bloque 3: salida terminal instancia base

```
1 ubuntu@ip-10-0-10-133:/tmp$ sudo ./aws/install
2 You can now run: /usr/local/bin/aws --version
3 ubuntu@ip-10-0-10-133:/tmp$ cd -
4 /home/ubuntu
5 ubuntu@ip-10-0-10-133:~$ aws --version
6 aws-cli/2.31.23 Python/3.13.7 Linux/6.14.0-1011-aws exe/x86_64.ubuntu.24
7 ubuntu@ip-10-0-10-133:~$
```

Después instalar las dependencias cerramos la instancia y creamos la AMI a partir de la instancia anterior.



The screenshot shows the 'Image details' section of the AWS console. It includes the following fields:

- Instance ID:** i-04bd934596d11a22f (amibase)
- Image name:** awscli_python_ami
- Image description - optional:** Dependencias y awscli instalados

Below the 'Image name' field, there is a note: 'Maximum 127 characters. Can't be modified after creation.' Below the 'Image description' field, there is a note: 'Maximum 255 characters'.

Figura 3: Creación AMI

4.2 Comprobación User Data

A partir de esta AMI podremos crear las instancias privadas y posteriormente descargar los archivos. Descargaremos los archivos de forma manual para comprobar los comandos del user data, así si nos da error, los solucionaremos antes de tener que volver a lanzar otra instancia. Para conectarnos nos bastará con conectarnos a una instancia intermedia, la cual llamamos bastión, a esta le pasaremos nuestra clave .pem mediante otra terminal usando:

```
scp -i ./v2.pem ./v2.pem ubuntu@ec2-23-20-27-162.compute-1.amazonaws.com:/home/ubuntu/. 
```

 Al hacer ls en nuestra terminal del bastión ya nos saldrá la clave y podremos conectarnos a nuestra instancia privada.

Bloque 4: salida terminal bastión

```
1 ubuntu@ip-10-0-10-22:~$ ls
2 v2.pem
3 ubuntu@ip-10-0-10-22:~$
```

Ahora desde el bastión nos conectamos a la instancia privada y vamos probando las credenciales y comandos de descarga del user data. Después de probar bastantes veces me di cuenta que los comandos que se indican en la práctica para ingresar las claves de awscli, ingresaban las claves en `/home/ubuntu/.aws/credentials`, y no en `/root/.aws/credentials`, por lo que cambié las credenciales de lugar y sí que me funcionó la descarga del s3. Dejo a continuación la prueba que hice:

Bloque 5: salida terminal instancia privada

```
1 ubuntu@ip-10-0-20-189:~/app$ sudo HOME=/home/ubuntu aws s3 cp s3://p4-buquet/backend
  /modelv1prod.py .
2 fatal error: Unable to locate credentials
3 ubuntu@ip-10-0-20-189:~/app$ cat /home/ubuntu/.aws/credentials
4 cat: /home/ubuntu/.aws/credentials: No such file or directory
5 ubuntu@ip-10-0-20-189:~/app$ sudo cat /root/.aws/credentials
6 [default]
7 aws_access_key_id = ASIA2HPTD7N2H2PSGE5C
8 aws_secret_access_key = 2Qc1gPNpKoHJYZegKopCGZMQx3Xer9ehBS7KvNaf
9 ubuntu@ip-10-0-20-189:~/app$ sudo cp /root/.aws/credentials /home/ubuntu/.aws/
  credentials
10 cp: cannot create regular file '/home/ubuntu/.aws/credentials': No such file or
  directory
11 ubuntu@ip-10-0-20-189:~/app$ sudo mkdir -p /home/ubuntu/.aws
12 ubuntu@ip-10-0-20-189:~/app$ sudo cp /root/.aws/credentials /home/ubuntu/.aws/
  credentials
13 ubuntu@ip-10-0-20-189:~/app$ sudo chown ubuntu:ubuntu /home/ubuntu/.aws/credentials
14 ubuntu@ip-10-0-20-189:~/app$ sudo chmod 600 /home/ubuntu/.aws/credentials
15 ubuntu@ip-10-0-20-189:~/app$ sudo HOME=/home/ubuntu aws s3 cp s3://p4-buquet/backend
  /modelv1prod.py .
16 download: s3://p4-buquet/backend/modelv1prod.py to ./modelv1prod.py
17 ubuntu@ip-10-0-20-189:~/app$ ls
18 modelv1prod.py
19 ubuntu@ip-10-0-20-189:~/app$
```

Después de comprobar lo que fallaba, simplemente introduje el traslado de las credenciales de lugar en el user data justo antes de las descargas del s3. Finalmente lancé una instancia para probar que todo iba bien con este user data actualizado. La prueba del user data salió como esperaba y ya podemos comenzar con el paso siguiente.

Bloque 6: logs instancia

```
1 [ 23.650346] cloud-init[915]: Completed 673 Bytes/673 Bytes (6.4 KiB/s) with 1
  file(s) remaining
2 download: s3://p4-buquet/backend/modelv1prod.py to ./modelv1prod.py
3 [ 24.701384] cloud-init[915]: Completed 896 Bytes/896 Bytes (7.7 KiB/s) with 1
  file(s) remaining
4 download: s3://p4-buquet/backend/modelv2canary.py to ./modelv2canary.py
5 [ 25.763120] cloud-init[915]: Completed 561 Bytes/561 Bytes (7.0 KiB/s) with 1
  file(s) remaining
6 download: s3://p4-buquet/backend/modelo.pkl to ./modelo.pkl
7 [ 25.949240] cloud-init[915]: total 12K
8 [ 25.949367] cloud-init[915]: -rw-r--r-- 1 ubuntu ubuntu 561 Oct 28 12:15 modelo.
 .pkl
9 [ 25.949582] cloud-init[915]: -rw-r--r-- 1 ubuntu ubuntu 673 Oct 28 12:15
  modelv1prod.py
10 [ 25.949784] cloud-init[915]: -rw-r--r-- 1 ubuntu ubuntu 896 Oct 28 12:15
  modelv2canary.py
```

Los endpoints también se están ejecutando adecuadamente

Bloque 7: consola endpoints

```
1 ubuntu@ip-10-0-20-16:~$ curl -X POST http://127.0.0.1:8000/api/v1/recommendation -H
  "Content-Type: application/json" -d '{"user_id": "test"}'
2 {"event_log":"User test purchase processed.","latency_ms":0,"model_id":"V1 - Stable"
  ,"recommendation":["itemX_oldmodel","itemY_oldmodel"]}
3 ubuntu@ip-10-0-20-16:~$
```

5 Configuración del Balanceador de Carga (ALB)

5.1 Target groups

Para esta parte debemos de tener corriendo 3 instancias con la versión de python_v_1 y 1 instancia con la versión de python_v_2. Esto lo haremos para crear los dos target groups, el de la V1 y el de la V2. Procederemos a crear 4 instancias privadas a partir de nuestra AMI y de nuestro user data, una vez las tengamos ejecutando nos conectamos una a una y vamos ejecutando el archivo python que corresponda a cada una, a las V1 `modelv1prod.py` y a la V2 `modelv2canary.py`, lo hacemos con 'nohup' para poder cerrar la terminal y que sigan funcionando en segundo plano. Debemos de tener estas instancias corriendo escuchando en el puerto correspondiente para poder configurar el ALB con estos target groups.

Instances (5) [Info](#)

🔍 Find Instance by attribute or tag (case-sensitive)

| <input type="checkbox"/> | Name 🔗 | ▲ | Instance ID | Instance state ▼ | Instance type |
|--------------------------|------------------------|---|-------------------------------------|---|---------------|
| <input type="checkbox"/> | bastion | | i-092e5cf4c8de61945 | ✔ Running 🔍 🔍 | t3.micro |
| <input type="checkbox"/> | V1_01 | | i-0500ef74baa09a6ea | ✔ Running 🔍 🔍 | t3.micro |
| <input type="checkbox"/> | V1_02 | | i-095393fa86993292c | ✔ Running 🔍 🔍 | t3.micro |
| <input type="checkbox"/> | V1_03 | | i-0b6e3868eda0fc3bb | ✔ Running 🔍 🔍 | t3.micro |
| <input type="checkbox"/> | V2_01 | | i-006549b59f57aa04e | ✔ Running 🔍 🔍 | t3.micro |

Figura 4: Instancias

Configuramos los dos targets al puerto 8000 con las instancias que hemos creado, de esta forma:

The image shows two screenshots of the AWS Management Console. The left screenshot shows the 'Registered targets' for a target group with 3 targets. The right screenshot shows the 'Registered targets' for a target group with 1 target.

| Instance ID | Name | Port |
|-------------------------------------|-------|------|
| i-0b6e3868eda0fc3bb | V1_03 | 8000 |
| i-095393fa86993292c | V1_02 | 8000 |
| i-0500ef74baa09a6ea | V1_01 | 8000 |

| Instance ID | Name | Port |
|-------------------------------------|-------|------|
| i-006549b59f57aa04e | V2_01 | 8000 |

Figura 5: TG-V1 y TG-V2

5.2 ALB y pruebas

Cuando tengamos los grupos creamos la ALB de forma pública usando nuestra vpc de la práctica, el SG ALB y las subredes públicas. El listener lo ponemos en el puerto 80 como indica el enunciado y ajustamos los pesos de 90 % para V1 y 10 % para V2

The image shows the configuration of an ALB listener in the AWS Management Console. The listener is named 'Listener HTTP:80' and is configured with the following settings:

- Protocol: HTTP
- Port: 80
- Default action: Forward to target groups
- Routing action: Forward to target groups
- Target group: V1 (Target type: Instance, IPv4 | Target stickiness: Off) with a weight of 9 and 90% percent.
- Target group: V2 (Target type: Instance, IPv4 | Target stickiness: Off) with a weight of 1 and 10% percent.

Figura 6: ALB

Al tener el ALB corriendo podemos probar nuestros endpoints desde el nombre de dns que nos proporciona AWS. Comprobamos como una de cada diez veces aproximadamente nos sale que se está usando el modelo V2.

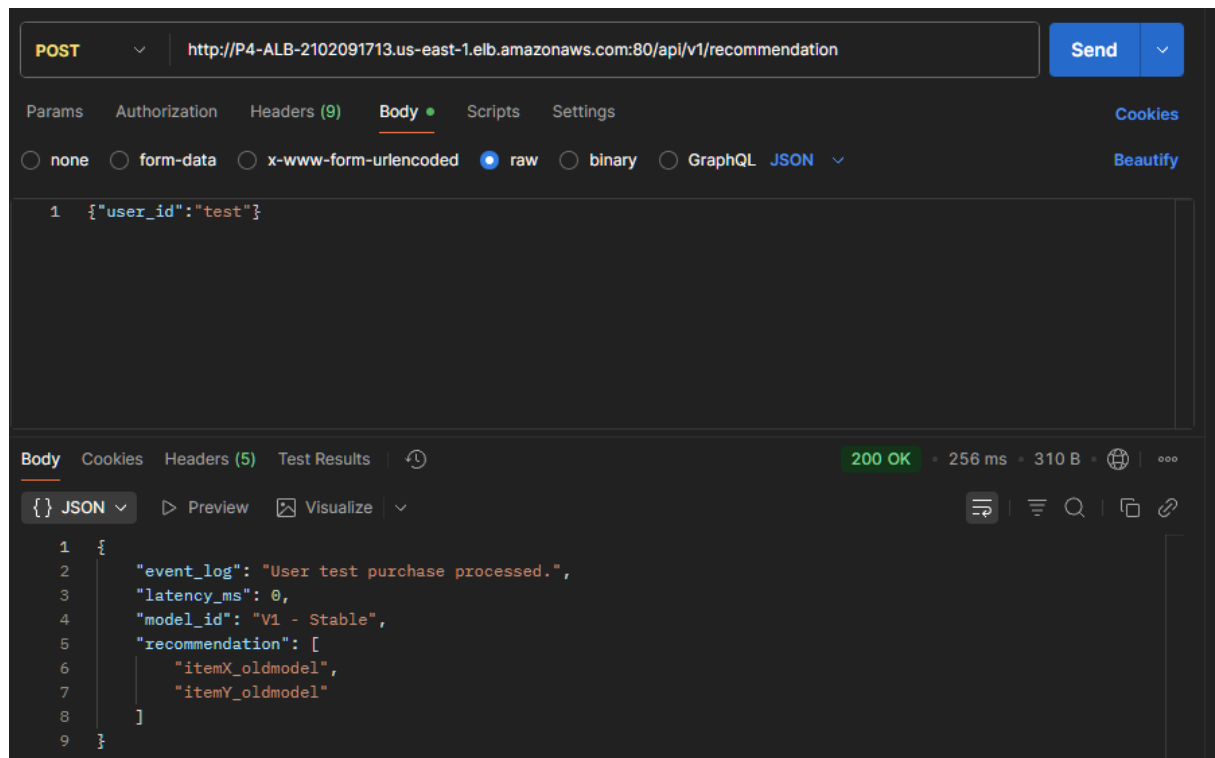


Figura 7: prueba V1

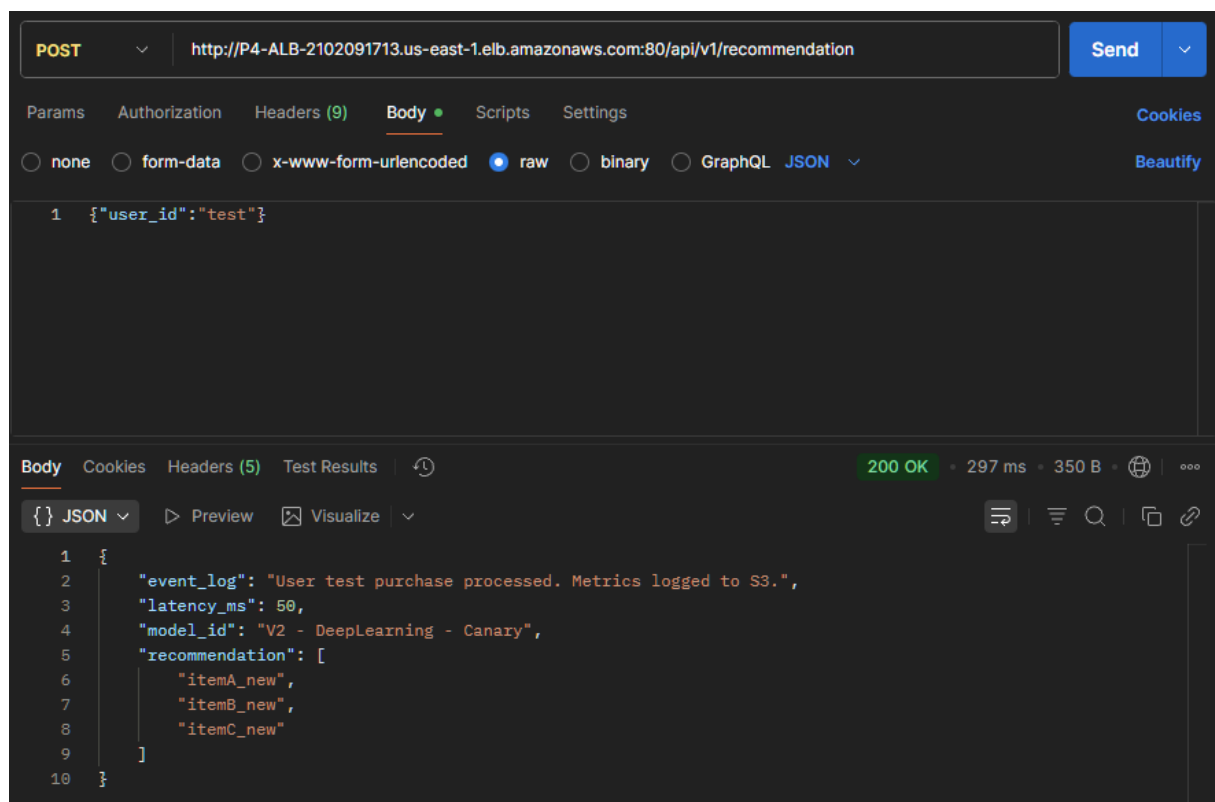


Figura 8: prueba V2

6 Transición a contenedores

Para este apartado ejecutaremos nuestra aplicación canary dentro de un contenedor docker. Lanzamos una instancia, esta vez puede ser sin la AMI ya que debemos de instalar por nuestra cuenta el docker únicamente. Instalaremos el docker, descargaremos nuestro modelv2canary.py desde el s3, o para mayor control, y así darnos cuenta de que proviene de la version con docker, podemos copiar y pegar una version del .py modificada que contenga el nombre de modelo como V2-docker. Cuando tengamos el archivo .py dentro de la instancia crearemos un dockerfile que copiara el archivo .py en el contenedor y lo ejecutará. Finalmente construiremos la imagen a partir del dockerfile y la lanzaremos mapeando el puerto 80 de la instancia con el puerto 8000 del contenedor.

Bloque 8: Consola instancia docker

```
1 ubuntu@ip-10-0-20-71:~$ nano modelv2canary.py
2 (copypaste)
3 ubuntu@ip-10-0-20-71:~$ nano Ddockerfile
4 FROM python:3.10-slim
5 WORKDIR /app
6 COPY modelv2canary.py .
7 RUN pip install flask
8 EXPOSE 8000
9 CMD ["python", "modelv2canary.py"]
10 ubuntu@ip-10-0-20-71:~$ sudo docker build -t modelv2canaryimage .
11 ...
12 ...
13 Successfully built 5865039f3e08
14 Successfully tagged modelv2canaryimage:latest
15 ubuntu@ip-10-0-20-71:~$ sudo docker run -d -p 80:8000 modelv2canaryimage
16 02af51a5fb4c5d18145061323f93872ed9311caa210b0e1d2db38dc9e0158140
17 ubuntu@ip-10-0-20-71:~$ sudo docker ps
18 CONTAINER ID   IMAGE                                COMMAND                                CREATED        STATUS
19 02af51a5fb4c   modelv2canaryimage                 "python ...modelv2canar"            55 seconds ago Up 54
20 seconds       0.0.0.0:80->8000/tcp, [::]:80->8000/tcp    gifted_chatterjee
```

6.1 Nueva regla SG Backend

Como estamos conectando el puerto 80 en vez de el 8000, debemos de cambiar la configuración del SG Backend para que permita la entrada del puerto 80 desde el SG ALB. También al crear el target group seleccionamos el puerto 80 en vez del 8000 y conectamos la instancia docker.

| Targets | Monitoring | Health checks | Attributes | Tags |
|--|---------------------|---------------|--------------|------|
| Registered targets (1) Info | | | | |
| Target groups route requests to individual registered targets using the protocol and automatically applied to HTTP/HTTPS target groups with at least 3 healthy targets | | | | |
| <input type="text" value="Filter targets"/> | | | | |
| <input type="checkbox"/> | Instance ID | ▼ | Name | ▼ |
| <input type="checkbox"/> | i-03ef5e3a06a7a6cb0 | | V2_02_docker | 80 |

Figura 9: TG V2 docker

6.2 Editar listener

Cuando tengamos el TG editamos el listener añadiendo otro grupo. No se si debía incluir TG V2 docker nuevo o sustituirlo por el anterior TG V2, pero optaré por incluir el nuevo TG de forma que tengamos el 80 % de trafico para V1, el 10 % para V2 normal y el ultimo 10 % para el V2 con docker. Al incluirlo lo probamos con el postman.

| Target group | Weight | Percent | |
|---|--------|---------|-------------------------|
| V2 Target type: Instance, IPv4 Target stickiness: Off | 1 | 10% | <button>Remove</button> |
| V1 Target type: Instance, IPv4 Target stickiness: Off | 8 | 80% | <button>Remove</button> |
| V2-docker Target type: Instance, IPv4 Target stickiness: Off | 1 | 10% | <button>Remove</button> |

0-999

Figura 10: ALB final

6.3 Pruebas endpoint del ALB

Y como podemos observar en la última captura, nos redirige de forma correcta al endpoint del contenedor.

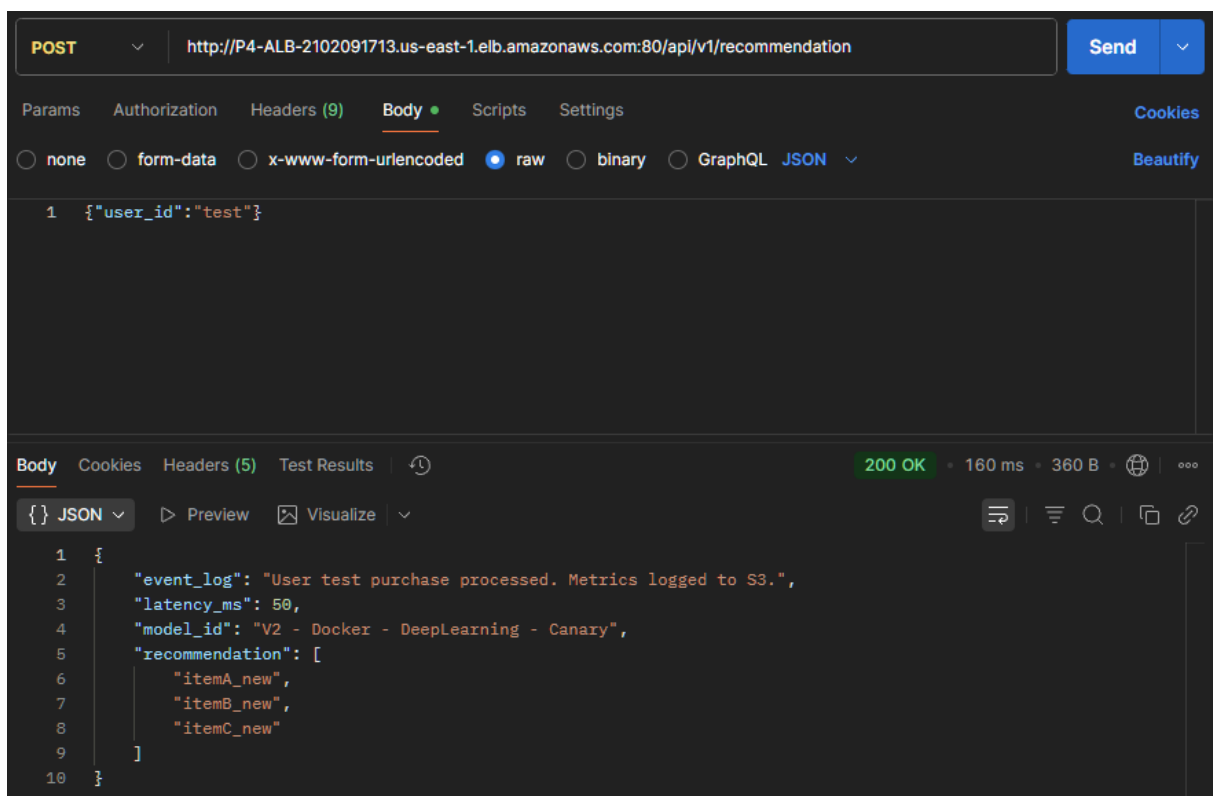


Figura 11: post docker

7 Preguntas de reflexión

7.1 Seguridad y NACLs

Aunque el SG-BACKEND ya filtra el tráfico no deseado al puerto 8000, la NACL añade una capa adicional de defensa a nivel de red (network layer) para todos los recursos de la subred, independiente de los SG de cada instancia. Así, mitiga riesgos ante errores en configuración de SGs y protege ante tráfico malicioso de otras fuentes dentro de la VPC. Un atacante no podría interceptar la respuesta de la API usando solo la NACL, ya que una NACL únicamente permite o deniega paquetes; no tienen capacidad de leer, modificar ni reenviar tráfico como lo haría un proxy. La confidencialidad y protección de datos requieren cifrado (HTTPS) además del filtrado de red.

7.2 WAF vs. NACL

Un ataque de Inyección de Código SQL sería detenido por el WAF (Web Application Firewall) y no por la NACL. El motivo es que el WAF inspecciona y filtra el contenido de las cargas útiles HTTP, detectando patrones peligrosos en la capa de aplicación. En cambio, la NACL solo filtra por IP, puerto y protocolo, sin analizar el contenido, por lo que nunca detectaría ataques específicos de aplicaciones.

7.3 VPC Endpoints

El VPC Endpoint permite que el tráfico de logging del Modelo V2 (por ejemplo, guardar logs en S3) viaje directamente entre la VPC y el servicio de AWS sin atravesar el Internet Gateway (IGW). Esto significa que los datos nunca salen a Internet, manteniendo toda la transferencia dentro de la red privada de AWS y, por tanto, aumentando la seguridad y reduciendo la latencia.

7.4 Redes de Contenedores

El Balanceador de Carga debe dirigir el tráfico a la IP privada de la VM Host Docker y al puerto 80 (el puerto mapeado externamente). Docker se encarga de reenviar ese tráfico interno al puerto 8000 del proceso Flask en el contenedor. El componente concreto de red de Docker responsable de esta traducción de puertos se llama Docker bridge network + Port Mapping (NAT); es decir, Docker redirige automáticamente las peticiones del puerto publicado en el host al puerto interno del contenedor.

7.5 Rollback Rápido

En un canary deployment, el empleo de contenedores Docker permite una reversión a la versión V1 mucho más rápida y segura que si el V2 estuviera en una VM tradicional. Con Docker, basta detener y eliminar el contenedor V2 para cortar completamente el tráfico nuevo, sin eliminar la VM ni reconfigurar el sistema operativo. Esto reduce a segundos el tiempo de rollback y minimiza los riesgos de “código residual”, facilitando además automatización y consistencia en los cambios.