

# Reducción de costos de Machine Learning con Endpoint Serverless

We Love the Cloud

#### Resumen

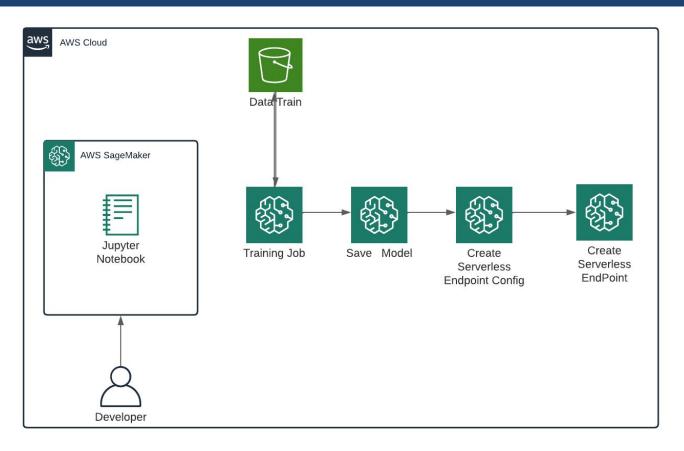


Al momento de incorporar un modelo de aprendizaje automático dentro de tu organización, un concepto a definir es el tipo de consumo de esta solución. De acuerdo al contexto del negocio, una alternativa es desplegar una API (EndPoint) que de respuesta en tiempo real, lo que requiere habilitar una instancia 24/7 en el ambiente de producción, que por defecto, se traduce en un cobro por todo el tiempo que esté habilitado.

A continuación exploramos la solución desarrollada por AWS a ésta problemática, gracias a el despliegue de un End-Point Serverless. Con este tipo de puntos de enlace (EndPoint), podemos reducir los costos relacionados a las inferencias, donde el aprovisionamiento se hará en el momento que esté sea llamado, es decir, bajo demanda, y totalmente administrado por AWS.

# **DIAGRAMA DE ARQUITECTURA**





#### **DIAGRAMA DE ARQUITECTURA**



Servicios de AWS se utilizan para qué fines en esta arquitectura:

- Sagemaker Studio: Desde aquí se escribirá el código relacionado a todo el despliegue de la solución mostrada, todo será escrito desde una Jupyter Notebook
- SageMaker(TrainingJob): El encargado de levantar la instancia de entrenamiento, el resultado del entrenamiento nos entregará nuestro artefacto.

#### **DIAGRAMA DE ARQUITECTURA**



- SageMaker(ModelRegistry): El encargado de registrar nuestro modelo dentro de la nube, para posteriormente utilizarlo en levantamiento de inferencias.
- SageMaker(Serverless Endpoint): Aquí levantaremos nuestro Endpoint que no necesitará aprovisionamiento de servidor
- Simple Storage Service: Desde aquí obtendremos nuestra data de entrenamiento.

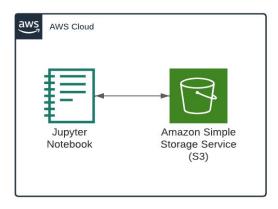


1- **Preparación:** Nuestro primer paso será preparar el ambiente a utilizar, en éste caso usaremos las Jupyter Notebook proporcionadas por SageMaker como nuestro entorno de desarrollo, por lo cual lo primero será importar librerías a usar, además del rol de ejecución para realizar correctamente el levantamiento de tareas.

```
import json
import pathlib
import pickle
import tarfile
import joblib
import numpy as np
import pandas as pd
import xgboost
import argparse
import os
```



2- **Data:** Para empezar a prepararnos para nuestro entrenamiento, necesitaremos un set de datos que esté alojado en S3 para la ejecución de este. Puedes fácilmente leer un set de datos en tu Jupyter Notebook, prepararlo y luego dejarlo en un bucket, incluso leer desde un bucket para luego editar y subir a un destino diferente. Siempre teniendo en cuenta que el resultado final debe ser un bucket.





3- Entrenamiento: Una vez que nuestros pasos anteriores fueron realizados, pasaremos a realizar nuestro entrenamiento con los datos anteriormente almacenados en S3. El resultado de éste entrenamiento será nuestro artefacto el cual deberemos trabajar en el paso siguiente.

```
train input = TrainingInput("
validation input = TrainingInpu' "
estimator.fit({'train': train input, 'validation': validation input})
2022-07-11 01:19:56 Starting - Starting the training job...ProfilerReport-1657502396: InProgress
                            Amazon SageMaker > Training jobs >
                            sagemaker-xgboost-
                               Job settings
                                 Job name
```

sagemaker-xgboost

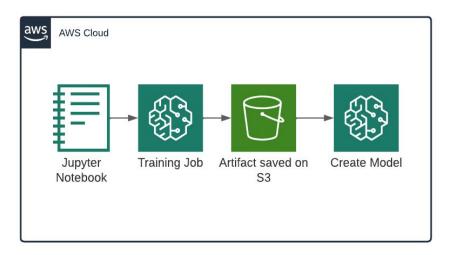
Status

 ○ Completed View history



4- **Crear modelo:** El resultado de nuestro entrenamiento nos arroja un artefacto que es almacenado en S3, para poder utilizarlo en los pasos posteriores; deberemos registrar éste artefacto en el directorio de modelos de AWS.





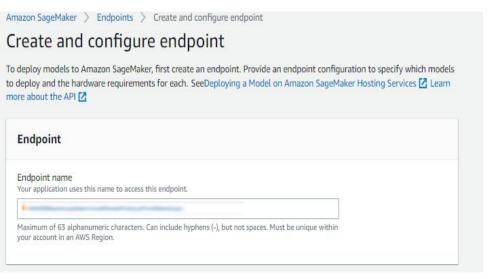


5- Crear configuración endpoint: Una vez que nuestro modelo fue registrado, podemos pasar a crear nuestra configuración de EndPoint. Aquí le podemos definir si consumirá alguna instancia en específico o será Serverless. En éste caso estamos lanzando una configuración Serverless.

mazon SageMaker > Endpoint configuration erverless-configuration	
Endpoint configuration settings	
Name serverless-configuration	ARN
Production variants	
Model name	Training job
Serverless-EndPoint-Model	



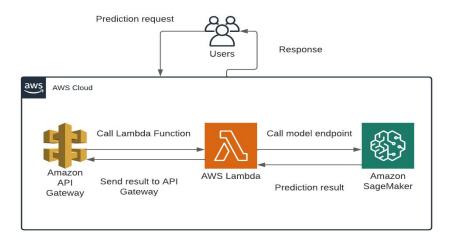
6- **Levantamiento de endpoint:** Una vez que la configuración ya está arriba, pasaremos a levantar el enlace. Esto se puede realizar fácilmente tanto desde la consola, como desde nuestra Jupyter Notebook.



A the se suistine and sist seefin	tion	O Create a new and solution	Samuel and	
<ul> <li>Use an existing endpoint configuration</li> <li>Use an existing endpoint configuration or clone an endpoint configuration.</li> </ul>		<ul> <li>Create a new endpoint configuration</li> <li>Add models and configure the instance and initial weight for each model.</li> </ul>		
indpoint configuration			C	
indpoint configuration  Q Search resources				
			< 1	
	▼ ARN			
Q Search resources	∇ ARN		< 1	



7- Invocar endpoint: Luego de haber levantado el endpoint llega el momento de consumir éste mismo. Puedes consumir directamente desde SageMaker o si ya confías totalmente en tu solución, puedes pasarla a producción. A continuación mostraremos un caso de uso a la hora de consumir EndPoints.





6- **Borrar endpoint:** Recuerda que un endpoint es un recurso que está siendo levantado y queda prendido a menos que no le indiques su finalización. Evita siempre incurrir en gastos no esperados y recuerda eliminar tu endpoint. A continuación se mostrará dos maneras.

Endpoints		C	Update endpoint	t Actions ▲
Q Search endpoints				Add/Edit tags
				Delete
Name	Creation time	•	Status	▼ Last updated
• serverless-endpoint	Jul 11, 2022 03:54 UTC		<b>⊘</b> InService	Jul 11, 2022 03:

#### Conclusion



En este ejercicio, se logró levantar un Endpoint serverless, demostrando el paso a paso, desde preparar el ambiente hasta consumir dicho endpoint. Se ocupó una máquina para la notebook y el entrenamiento, pero a diferencia de un Endpoint en real-time, ésta vez nuestra inferencia no le especificamos máquina, haciendo de éste método una alternativa muy amigable a la hora de reducir costos en el ámbito de Machine Learning.

Gracias a la flexibilidad entregada por AWS, no se debe hacer mayores cambios a la arquitectura si actualmente estás usando EndPoint real-time, únicamente deberás apegarte a tu regla de negocio. Puesto que si necesitas una capacidad totalmente inmediata y de gran potencia, real-time deberá ser la opción. Si tus predicciones pueden esperar un tiempo de "en frío" y quieres reducir gastos de operación, ServerLess Endpoint es la alternativa.

#### LINKS RELACIONADOS



- https://docs.aws.amazon.com/sagemaker/latest/dg/serverless-endpoints.
   html
- https://docs.aws.amazon.com/sagemaker/latest/dg/serverless-endpoints
   -create-invoke-update-delete.html
- https://docs.aws.amazon.com/sagemaker/latest/dg/serverless-endpoints -monitoring.html
- https://docs.aws.amazon.com/sagemaker/latest/dg/serverless-endpoints -troubleshooting.html
- https://docs.aws.amazon.com/sagemaker/latest/dg/deploy-model.html

