

Deploying a recommender system on the cloud



Emilio Macias · [Follow](#)

9 min read · Jul 1, 2021

16

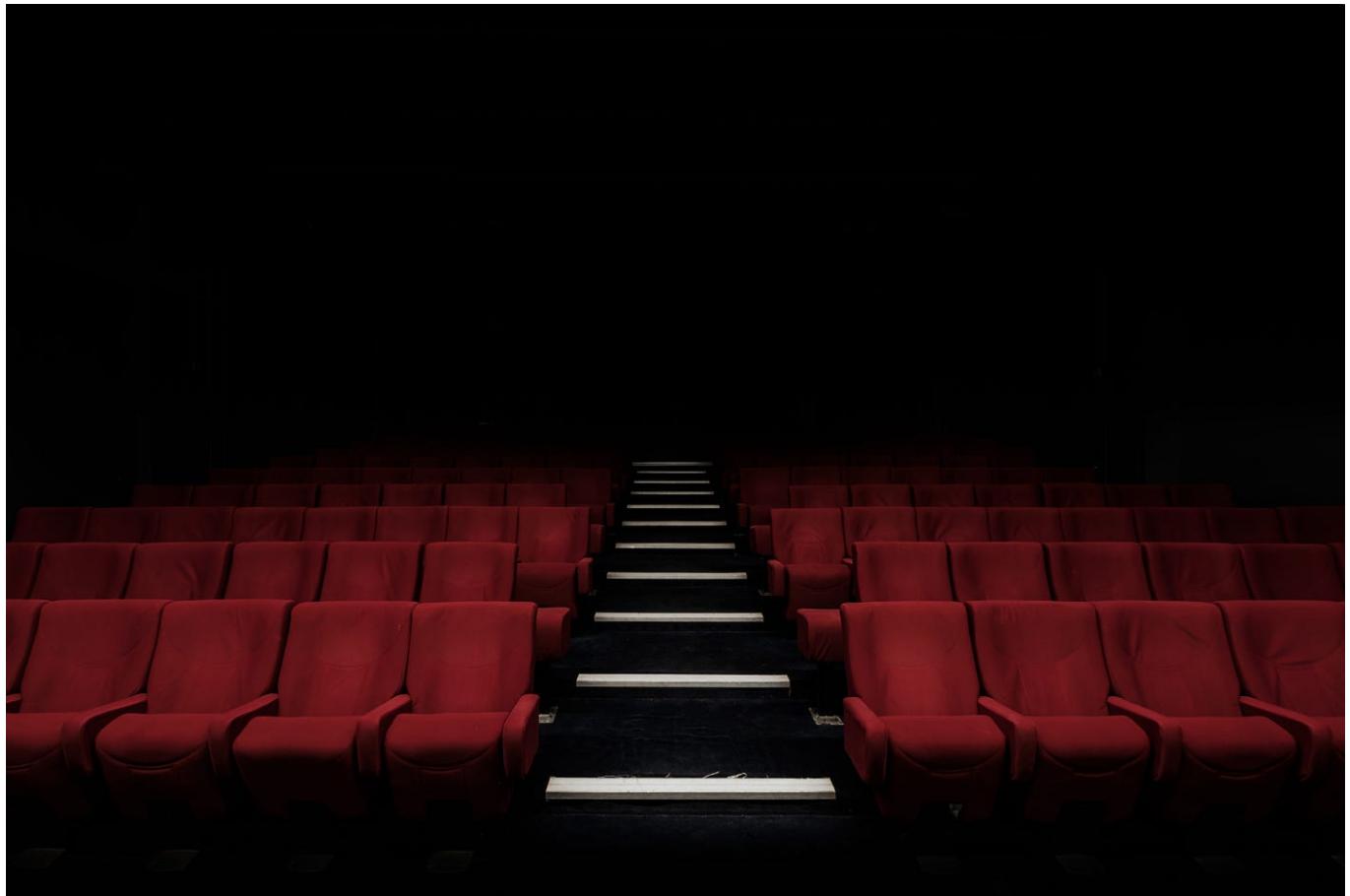


Photo by [Felix Mooneeram](#) on [Unsplash](#)

Abstract

This project covers the end-to-end process of training and deploying a recommender system on AWS, following a typical use case such as predicting movie ratings from users. Although a basic introduction to recommender systems is given, the mathematical reasoning is out of the scope of this article, focusing primarily on the practical side of model building and service deployment.

Background

Gone are the days when we would visit a video store to get some movies to watch. Users now have access to a great deal of movies with the click of a button. However there is something that they lack more and more everyday, and that's spare time. Most users are not willing to spend an hour searching for the ideal movie to watch at a certain time. Personal recommendations disappeared when the video stores closed.

Thanks to software and statistics, this problem can be solved with a so-called *recommender system* which is nothing but a machine learning system that can predict users' interests and ratings they would give to an item with the aim of helping the users discover new products they might enjoy. Basically the more data we have collected about an user (*movie ratings*) the better we know their preferences and so the more accurate the *recommendations* will be.



Photo by [Sean Benesh](#) on [Unsplash](#)

Recommender systems

Coming back to the age of movie stores, the store assistant would come to you and ask you for your preferences, maybe they'd ask you for your favourite movies and recommend similar ones that you might not have watched. This is what we would call *content-based* recommendations since they are only based on the similarity of the items, and that can be easily solved with recommender systems.

Another way of recommending movies is called *collaborative filtering* by which items are recommended based on user profiles. Imagine the scenario where a store assistant would personally know everyone that visits their store, including their interests. When customer X arrives at the store, the assistant knows customer X has similar preferences to customer Y and so he/she will recommend movies that customer Y enjoyed and haven't been watched by customer X. As you can imagine, this approach was not feasible from a human perspective, however it can be solved with linear algebra and matrices by collecting rating and watching information in order to find similarities between users.

One type of collaborative filtering systems are model-based algorithms which make use of machine learning methods to find these user similarities. For this project we are using an algorithm called Singular Value Decomposition (SVD) which is available in [scikit-surprise](#), a Python library for recommender systems.

If you are interested in the mathematical details of these methods I suggest you reading the [Introduction to Recommender System](#).

ML Pipeline

One of the most common ways to deploy an AI service is through an API that can receive HTTP requests and uses a machine learning model to make predictions (in this case, a rating from a given user for a particular movie). Here we are going to build a ML pipeline using different services from AWS.

Below is the list of AWS services and external tools we'll be using for our pipeline:

- Flask, Pandas and Scikit-Surprise (Python libraries) for coding the ML algorithm.
- Docker for packing our algorithm into a container.
- Amazon S3 for storing the datasets and the trained model.
- Amazon ECR for hosting our custom algorithm in a Docker container
- Amazon SageMaker for training the model and for making predictions
- AWS Lambda function for invoking the SageMaker Endpoint and enriching the response.
- API Gateway for publishing our API service to the Internet.
- Amazon CloudWatch for event logging.
- Postman for sending HTTP requests to our public API.

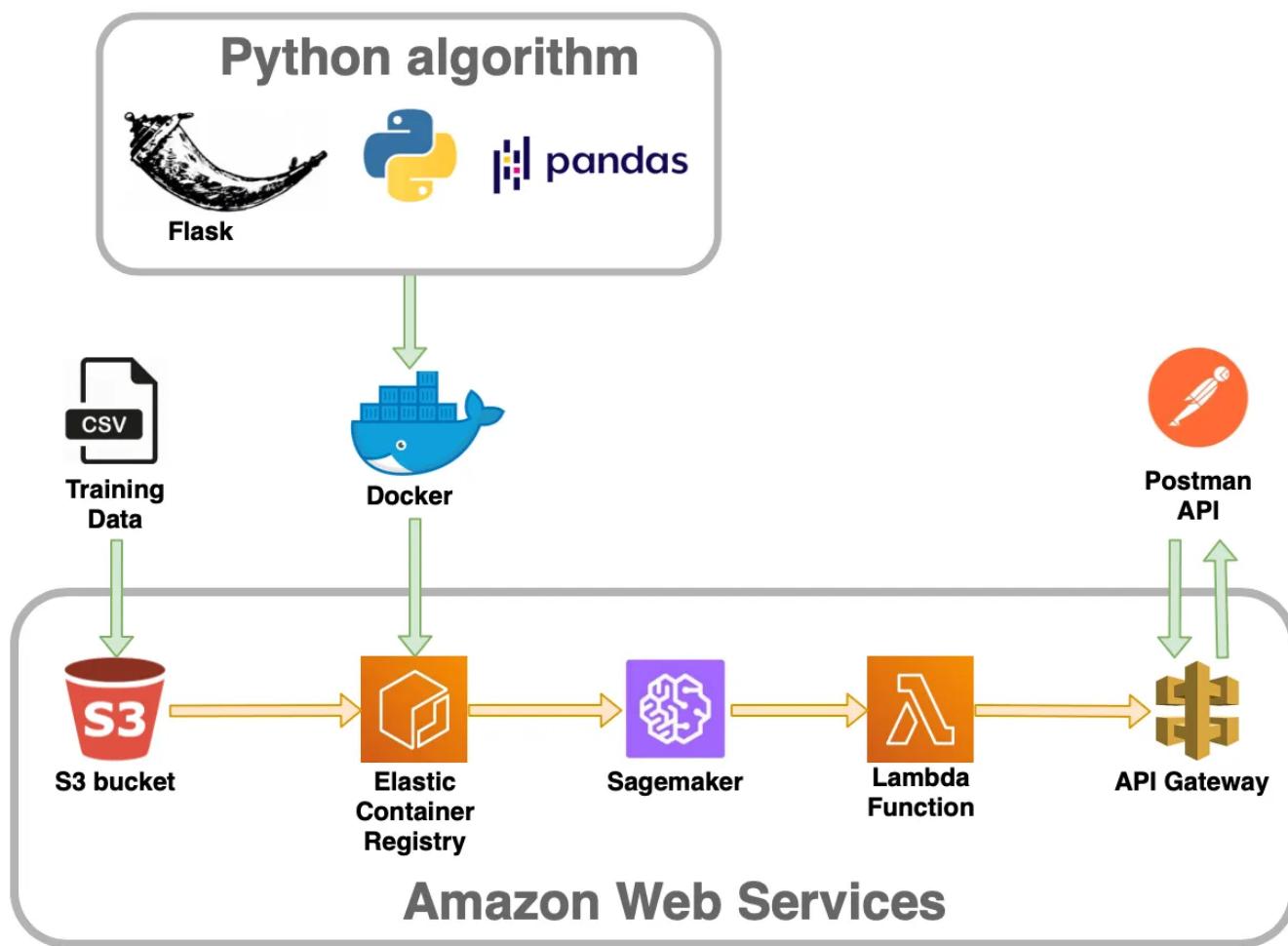


Image by Author

Preparing Data

The first step is to collect the data that we need for our recommender system. In particular, we are interested in movie ratings information that will help our model find similarities between movies and users. One of the largest movie databases is MovieLens, a dataset collected by the GroupLens Research and available [here](#).

AWS provides a service named S3 where we can upload our datasets into object buckets that can be accessed via other AWS services, custom models or even containers.

emacias-movielens

Objects (3)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

[Actions ▾](#) [Create folder](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	movies.csv	csv	June 24, 2021, 09:18:14 (UTC+02:00)	482.8 KB	Standard
<input type="checkbox"/>	output/	Folder	-	-	-
<input type="checkbox"/>	ratings.csv	csv	June 20, 2021, 15:25:24 (UTC+02:00)	2.4 MB	Standard

Image by Author

Containerising our algorithm

Even though SageMaker provides a large list of built-in ML algorithms, building your own custom algorithm will provide you better flexibility for particular scenarios.

AWS expects your ML algorithm to have a *train* and a *serve* component, as explained in this post: <https://aws.amazon.com/blogs/machine-learning/train-and-host-scikit-learn-models-in-amazon-sagemaker-by-building-a-scikit-docker-container/>

An SVD recommender system can be trained with Scikit-Surprise library using a grid search to find the hyper-parameters that minimise the RMSE metric, as follows:

```
1 # Process ratings with Surprise Scikit
2 reader = Reader(line_format='user item rating timestamp', sep=',', rating_scale = (0.5,5)
3 ratings = Dataset.load_from_file(ratings_path, reader=reader)
4 trainset = ratings.build_full_trainset()
5
6 # Find the best hyperparameters for an SVD model via grid search cross-validation
7 param_grid = {
8     'n_factors': [10, 100, 500],
9     'n_epochs': [5, 20, 50],
10    'lr_all': [0.001, 0.005, 0.02],
11    'reg_all': [0.005, 0.02, 0.1]}
12
13 gs_model = GridSearchCV(
14     algo_class = SVD,
15     param_grid = param_grid,
16     n_jobs = -1,
17     joblib_verbose = 5)
18
19 gs_model.fit(ratings)
20
21 # Train the SVD model with the parameters that minimise the root mean squared error
22 best_SVD = gs_model.best_estimator['rmse']
23 best_SVD.fit(trainset)
```

train.py hosted with ❤ by GitHub

[view raw](#)

The way to pack your algorithm is by using Docker containers that you can then push into AWS. Below you can see the Dockerfile I used for my model:

```
1 # Build an image that can do training and inference in SageMaker
2 # This is a Python 3 image that uses the nginx, gunicorn, flask stack
3 # for serving inferences in a stable way.
4
5 FROM ubuntu:18.04
6
7 MAINTAINER Amazon AI <sage-learner@amazon.com>
8
9 RUN apt-get -y update && apt-get install -y --no-install-recommends \
10      wget \
11      python3-pip \
12      python3-setuptools \
13      nginx \
14      ca-certificates \
15      gcc \
16      libpq-dev \
17      python3-dev \
18      python3-venv \
19      python3-wheel \
20      && rm -rf /var/lib/apt/lists/*
21
22 RUN ln -s /usr/bin/python3 /usr/bin/python
23 RUN ln -s /usr/bin/pip3 /usr/bin/pip
24
25 # surprise won't install without numpy preinstalled
26 RUN pip3 install wheel
27
28 # Here we get all python packages.
29 # There's substantial overlap between scipy and numpy that we eliminate by
30 # linking them together. Likewise, pip leaves the install caches populated which uses
31 # a significant amount of space. These optimizations save a fair amount of space in the
32 # image, which reduces start up time.
33 RUN pip --no-cache-dir install numpy==1.16.2 flask gunicorn
34
35 # surprise won't install without numpy preinstalled
36 RUN pip install scikit-surprise
37
38 # Set some environment variables. PYTHONUNBUFFERED keeps Python from buffering our stand-
39 # output stream, which means that logs can be delivered to the user quickly. PYTHONDONTW-
40 # keeps Python from writing the .pyc files which are unnecessary in this case. We also u-
41 # PATH so that the train and serve programs are found when the container is invoked.
42
43 ENV PYTHONUNBUFFERED=TRUE
44 ENV PYTHONDONTWRITEBYTECODE=TRUE
45 ENV PATH="/opt/program:${PATH}"
```

```
46
47 # Set up the program in the image
48 COPY src /opt/program
49 WORKDIR /opt/program
```

Dockerfile hosted with ❤️ by GitHub

[view raw](#)

Once we have the Dockerfile ready we can build the image and push it into Amazon Elastic Container Registry (ECR). You can automate this process with the following script that is provided by AWS to automatically connect with your username and region, build the Docker image and make it available for SageMaker:

```
1 #!/usr/bin/env bash
2
3 # This script shows how to build the Docker image and push it to ECR to be ready for use
4 # by SageMaker.
5
6 # The argument to this script is the image name. This will be used as the image on the l
7 # machine and combined with the account and region to form the repository name for ECR.
8 image=$1
9
10 if [ "$image" == "" ]
11 then
12     echo "Usage: $0 <image-name>"
13     exit 1
14 fi
15
16 chmod +x src/train
17 chmod +x src/serve
18
19 # Get the account number associated with the current IAM credentials
20 account=$(aws sts get-caller-identity --query Account --output text)
21
22 if [ $? -ne 0 ]
23 then
24     exit 255
25 fi
26
27
28 # Get the region defined in the current configuration (default to us-west-2 if none defi
29 region=$(aws configure get region)
30 region=${region:-us-west-2}
31
32
33 fullname="${account}.dkr.ecr.${region}.amazonaws.com/${image}:latest"
34
35 # If the repository doesn't exist in ECR, create it.
36
37 aws ecr describe-repositories --repository-names "${image}" > /dev/null 2>&1
38
39 if [ $? -ne 0 ]
40 then
41     aws ecr create-repository --repository-name "${image}" > /dev/null
42 fi
43
44 # Get the login command from ECR and execute it directly
45 aws ecr get-login-password --region "${region}" | docker login --username AWS --password
```

```

46
47 # Build the docker image locally with the image name and then push it to ECR
48 # with the full name.
49
50 docker build -t ${image} .
51 docker tag ${image} ${fullname}
52
53 docker push ${fullname}

```

[build_and_push.sh](#) hosted with ❤ by GitHub

[view raw](#)

Once you run this script you can go to ECR and verify your image has been loaded in the registry.

Images (10)							View push commands	Edit
<input type="checkbox"/>	Image tag	Pushed at	Size (MB)	Image URI	Digest	Scan status	Vulnerabilities	
<input type="checkbox"/>	latest	Jun 24, 2021 09:24:08 AM	266.36	Copy URI	sha256:fca7a132e9efd77...	-	-	

Image by Author

Training the model

AWS provides a great tool for building machine learning models, and this is SageMaker. This tool will run our Docker container, first for training the model with the movie ratings dataset and then, when our API is invoked, it will make rating inferences with the model, given a user and a movie IDs.

Creating a training job with SageMaker is very simple. We start by pointing to the path of our custom container in ECR. Having our own algorithm helps us to find the best model hyperparameters via a grid search in the code.

However, SageMaker also gives you the choice to use specific values in the training job configuration. The next step is to provide the channels which basically means configuring the input data sources, making sure you are using the same channel name in the AWS training job as in the model algorithm. Finally we just need to specify the output directory where the model will be saved.

The screenshot shows the 'Training jobs' section of the Amazon SageMaker console. A specific training job is selected, with its name 'emacias-svd-movielens-copy-06-24' displayed prominently at the top. Below the job name, there are several tabs: 'Job settings' (selected), 'Algorithm', 'Outputs', 'Logs', and 'Metrics'. The 'Job settings' tab displays detailed information about the job, including its ARN, creation time, last modified time, status (Completed), and IAM role ARN. The 'Algorithm' tab shows the algorithm ARN, instance type (ml.m4.xlarge), and other parameters like additional volume size and maximum runtime. At the bottom of the page, the text 'Image by Author' is visible.

Job name	Status	SageMaker metrics time series	IAM role ARN
emacias-svd-movielens-copy-06-24	Completed View history	Disabled	arn:aws:iam::567169250895:role/service-role/AmazonSageMaker-ExecutionRole-20210608T194844
ARN	Creation time	Training time (seconds)	
arn:aws:sagemaker:us-east-1:567169250895:training-job/emacias-svd-movielens-copy-06-24	Jun 24, 2021 07:24 UTC	169	
	Last modified time	Billable time (seconds)	
	Jun 24, 2021 07:30 UTC	169	
		Managed spot training savings	
		0%	
		Tuning job source/parent	-

Algorithm			
Algorithm ARN	Instance type	Additional volume size (GB)	Volume encryption key
-	ml.m4.xlarge	1	-
Training image	Instance count	Maximum runtime (s)	
567169250895.dkr.ecr.us-east-	1	86400	

Image by Author

Once the training job is complete you can proceed to create an inference model from the job. Remember it is always a good practice to include some logs in your algorithm for debugging the different stages. The training logs can be found under the AWS CloudWatch service.

Great! Our movie recommender model is ready now. So how can we deploy our service?

Deploying the model

We'll keep on using the SageMaker service since we'll be creating a SageMaker Endpoint directly from our model. An endpoint configuration is required so you'll need to either choose an existing configuration or create a new one. What you need to be aware of in order to avoid a surprise in your bill is that you'll be charged for the time your SageMaker Endpoint is up, so make sure you remember to delete the endpoint when you no longer need it.

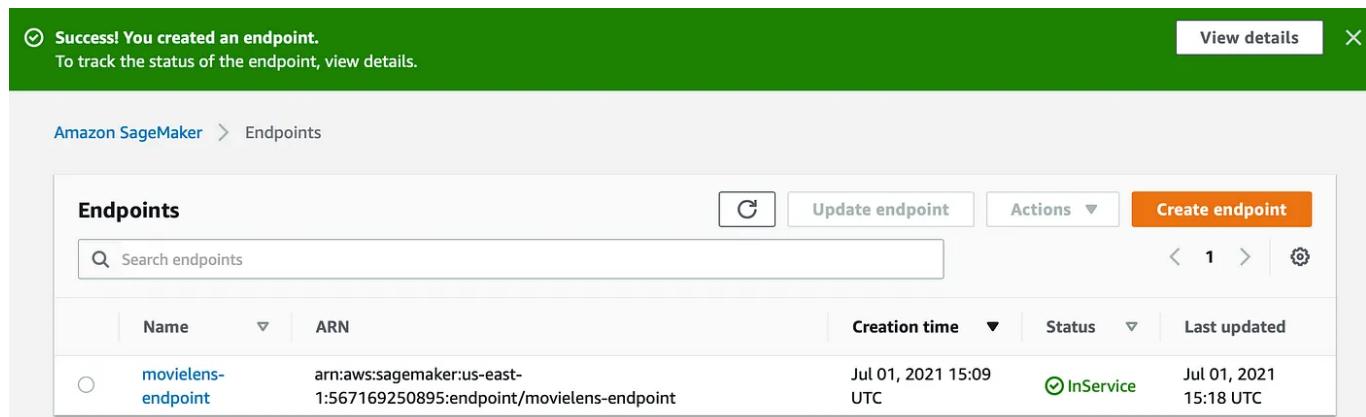


Image by Author

Once the endpoint is in service (it might take some time to launch), the next step in the deployment is to create a Lambda function that will invoke your endpoint, passing in a user and a movie, and expecting to receive a predicted rating. The way we link this function to our endpoint is by creating an environment variable whose value is the name of our SageMaker Endpoint and then loading this variable in the Lambda code.

One of the key benefits of using a Lambda function is that we can enrich the HTTP response returned by the SageMaker Endpoint. In our scenario, apart from returning the predicted movie rating, we'll also transform the movie id into the corresponding title. We can achieve this by making the Lambda function load a movie dataset (CSV file) that relates the movie IDs to titles. One thing to remember though is that any required Python library apart from the built-in ones (os, boto3, json) will have to be added as a Lambda layer. You can use the following lambda handler as a starting point:

```
1 ...
2 This code should be copied into a Lambda function of AWS.
3 Make sure to create an 'ENDPOINT_NAME' variable to points to your SageMaker Endpoint.
4 ...
5
6 import os
7 import boto3
8 import json
9 import pandas as pd
10
11 # grab environment variables
12 ENDPOINT_NAME = os.environ['ENDPOINT_NAME']
13
14 runtime = boto3.client('runtime.sagemaker')
15 s3_client = boto3.client('s3')
16 bucket = <YOUR_BUCKET_NAME>
17 movies_filename = 'movies.csv'
18
19 def lambda_handler(event, context):
20
21     payload = json.dumps(event)
22
23     response = runtime.invoke_endpoint(EndpointName=ENDPOINT_NAME,
24                                         ContentType='application/json',
25                                         Body=payload)
26
27     file_contents = s3_client.get_object(Bucket=bucket, Key=movies_filename)
28     movies = pd.read_csv(file_contents['Body'], usecols=[0,1], index_col=0, squeeze=True)
29
30     movie_id = event['movieId']
31     predicted_value = movies[int(movie_id)] + ' : ' + response['Body'].read().decode()
32
33     return predicted_value
```

lambda_function.py hosted with ❤ by GitHub

[view raw](#)

Before moving on to the next step, you should make sure your Endpoint and Lambda function are working correctly. Lambda service provides an easy tool to test your function, all you need to do is configuring a test event where you pass in a sample input to verify the correct result of the query. In our

[Open in app ↗](#)

[Sign up](#)

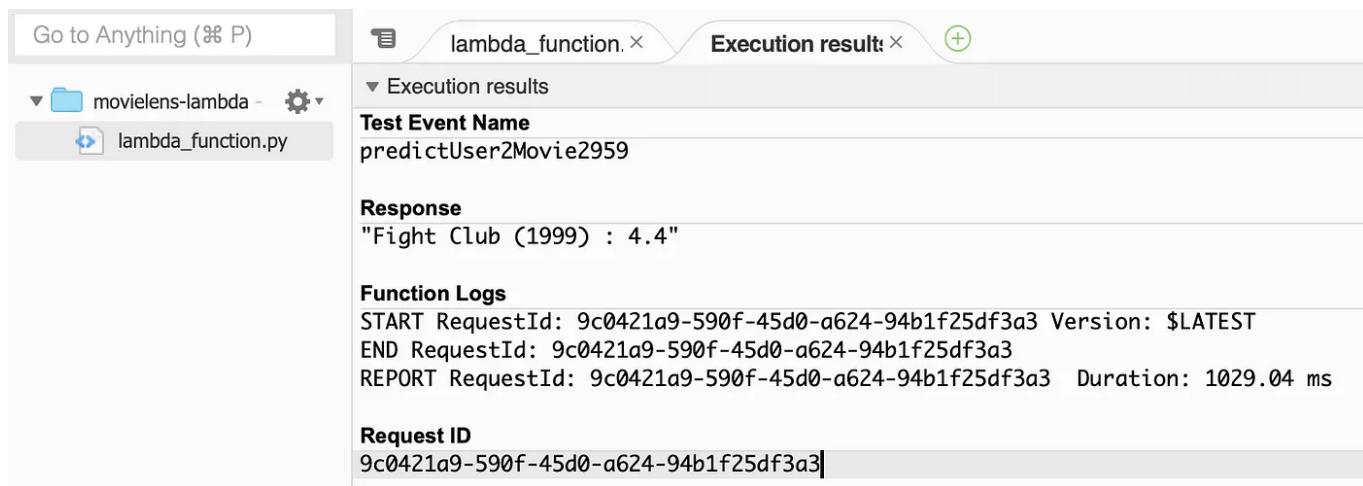
[Sign In](#)

```

  "userId": "2",
  "movieId": "2959"
}

```

After running the test, and once the Endpoint response is enriched by the Lambda function, we'll obtain the following response:



The screenshot shows the AWS Lambda Test Execution interface. On the left, there's a sidebar with a 'Go to Anything' search bar and a folder named 'movielens-lambda' containing 'lambda_function.py'. The main area has tabs for 'lambda_function.' and 'Execution result'. Under 'Execution result', it shows a 'Test Event Name' of 'predictUser2Movie2959'. The 'Response' field contains the output: "'Fight Club (1999) : 4.4'". Below that, 'Function Logs' show the execution details: START RequestId: 9c0421a9-590f-45d0-a624-94b1f25df3a3 Version: \$LATEST, END RequestId: 9c0421a9-590f-45d0-a624-94b1f25df3a3, REPORT RequestId: 9c0421a9-590f-45d0-a624-94b1f25df3a3 Duration: 1029.04 ms. At the bottom, the 'Request ID' is listed as '9c0421a9-590f-45d0-a624-94b1f25df3a3'.

Image by Author

The last step of our ML deployment is to create a public API that can be invoked from outside our AWS private virtual cloud. API Gateway is the service that will let us build a public REST API from scratch. All we need to do inside this service is creating a resource and a method (in our case, POST, since we'll need to provide the user ID and movie ID in the body of the request), making sure to choose the option to integrate this REST method with our Lambda function.

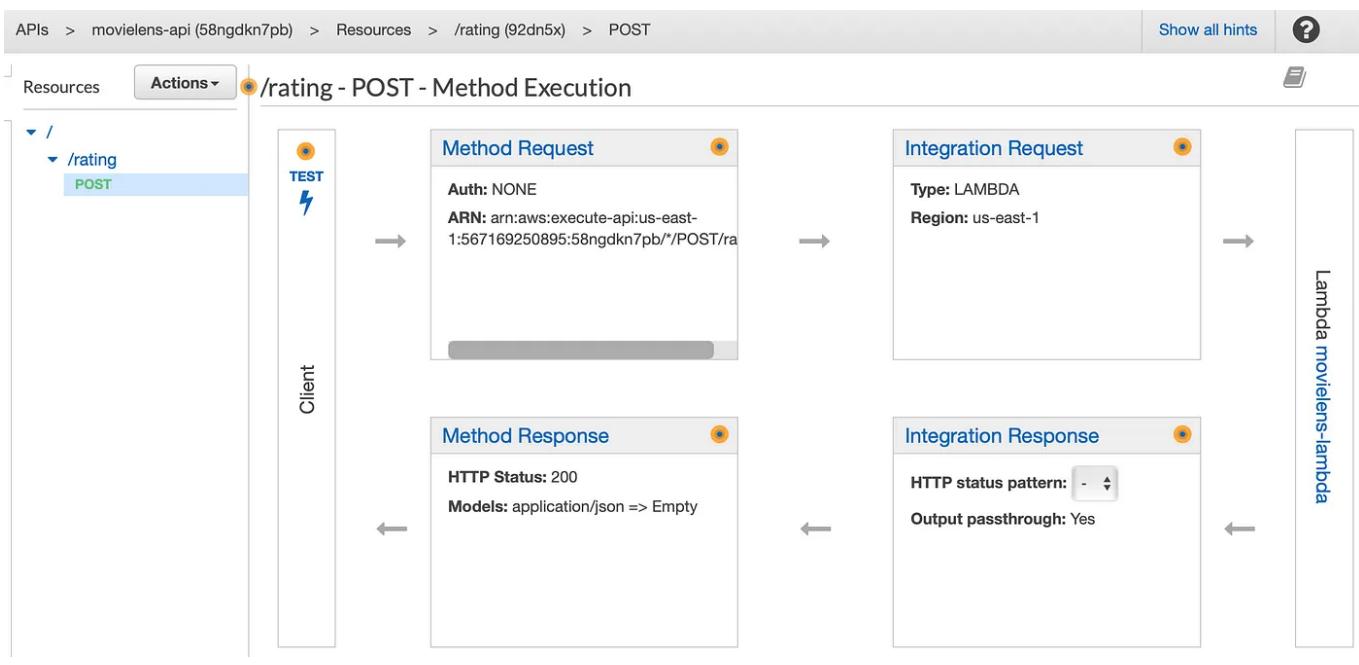


Image by Author

Our recommender system is now live! API Gateway provides the public URL we need to invoke our service. It's time to ask for some movie predictions...

Testing our AI service

There are different ways you can invoke our movie recommender. If you prefer to do it via the terminal, you can use a command-line utility such as [curl](#). Alternatively you can use a visual tool such as Postman, as we're doing here. Just remember we have configured our service with a POST method, so in Postman, we'll have to choose this type of method, enter the URL we got from API Gateway and include a JSON body in the HTTP request. The expected format of the JSON body is as follows:

```
{
  "userId": String,
```

```

“movieId”: String
}

```

And voilà! Below we can see the response from our API including the name of the movie (*Fight Club*) and the predicted rating from the user (4.4).

The screenshot shows a Postman interface with the following details:

- Request URL:** https://58ngdkn7pb.execute-api.us-east-1.amazonaws.com/movielens-stage/rating
- Method:** POST
- Body:** JSON (Pretty)


```

1 {
2   ...
3     "userId": "2",
4     "movieId": "2959"
5 }
```
- Response:**
 - Status: 200 OK
 - Time: 20.61 s
 - Size: 316 B
 - Content: { "Fight Club (1999) : 4.4" }

Image by Author

Conclusion

We have covered how we can take a custom machine learning algorithm into production via a public cloud provider such as AWS, applying it to the use case of a recommender system that can predict movie ratings from the users.

A similar outcome can be achieved with Google Cloud Platform or Microsoft Azure, so it is up to you to pick the public cloud provider that meets your requirements in the most efficient and economic way.

Feel free to use my code for your own projects or for learning purposes:
https://github.com/ejmacias/aws_movie_recommender

Happy deployment!

Recommender Systems

AWS

Sagemaker

Flask

Docker



Written by Emilio Macias

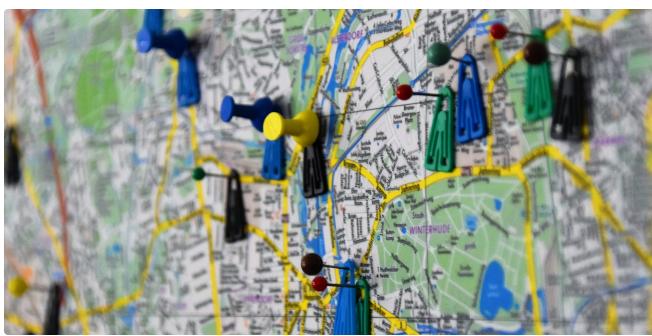
34 Followers

Follow



Software engineer with a passion for data science and its social impact.

More from Emilio Macias



Emilio Macias in Towards Data Science

Store locations

Emilio Macias in Towards Data Science

Clustering news articles based on named entities

Finding the optimal location for a new low-cost supermarket in Madrid

8 min read · Dec 22, 2020

👏 32 💬 2

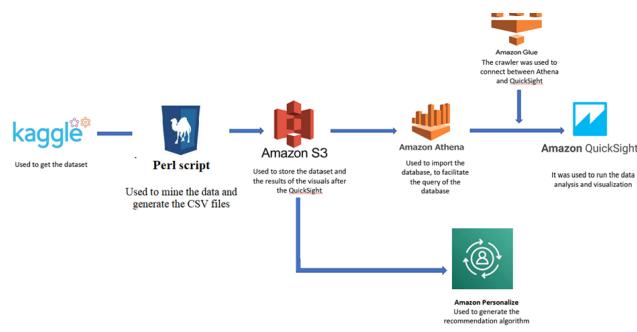
How to group real-time news in Spanish by people, locations and organisations

12 min read · May 2, 2021

Bookmark + 7 💬 1

[See all from Emilio Macias](#)

Recommended from Medium



Sahar Farhat

Recommendation system for e-commerce website

I. Scenario Definition

5 min read · Oct 13

👏 1 💬



ai geek (wishesh)

Best Practices for Deploying Large Language Models (LLMs) in...

Large Language Models (LLMs) have revolutionized the field of natural language...

10 min read · Jun 26

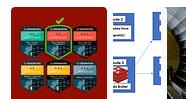
Bookmark + 209 💬 1

Lists



Coding & Development

11 stories · 258 saves



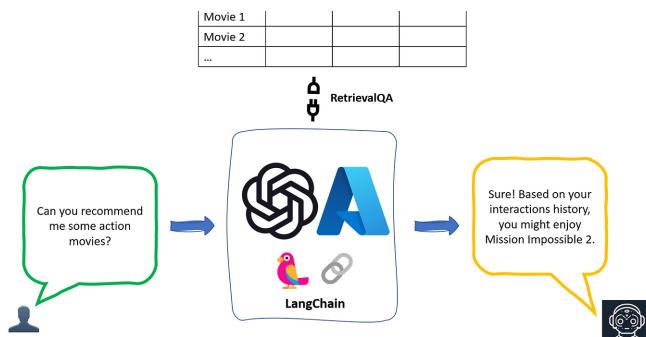
New_Reading_List

174 stories · 180 saves



Natural Language Processing

803 stories · 370 saves



Valentina Alto in Microsoft Azure

Enhancing Recommendation Systems with Large Language...

An implementation with LangChain and Azure OpenAI

★ · 13 min read · Aug 23

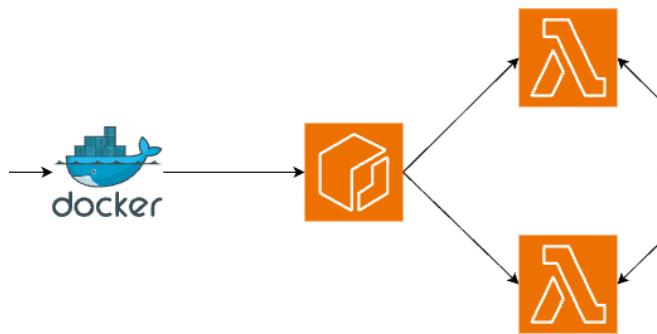


 Amy @GrabNGoInfo in GrabNGoInfo

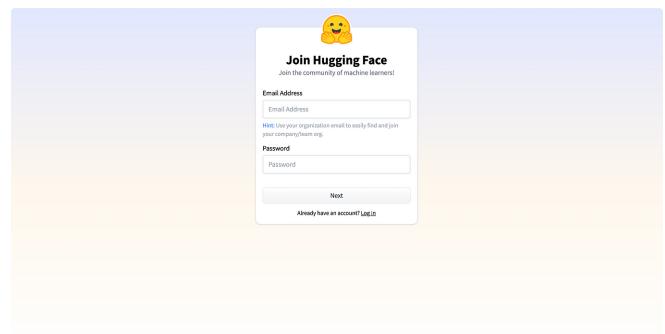
Hybrid Recommendation System Using User-based and Item-base...

4 hybrid methods to combine user-user and item-item collaborative filtering

★ · 8 min read · Jul 11



FanchenBao



sriram c in Technology at Nineleaps

API Service with FastAPI + AWS Lambda + API Gateway and Make ...

The goal of this article is to show one way to create a FastAPI app, deploy it to AWS...

14 min read · Sep 12



36



11



See more recommendations

How to deploy your own LLM(Large Language Models)

If you are new to LLM world, Then I would suggest you to go through previous articles t...

5 min read · Sep 17

