



Relazione Progetto: amazonAnalyzerKubernetes

Federico Canzonieri 1000024369

Luglio 2022

Contents

1	Introduzione	1
2	Componenti	2
2.1	Python	2
2.2	S3	3
2.3	Apache Spark	3
2.4	Lambda	4
2.5	OpenSearch	5
2.6	Grafana	5

1 Introduzione

Questo progetto fa uso di diverse tecnologie per costruire una data pipeline che ha come sorgente dati **Amazon** (le recensioni di un prodotto) , al fine di eseguire della **sentiment analysis** per capire come variano le recensioni degli utenti.

In particolare le tecnologie usate sono **Kubernetes** e qualche servizio di **AWS** come **S3**, **Lambda**, **Amazon Opensearch Service** e **Amazon Managed Grafana** . Il repository è disponibile al seguente indirizzo: <https://github.com/federicocanzonieri/amazonAnalyzerKubernetes>.

Vediamo l'architettura dell'applicazione per capire il flusso dati.

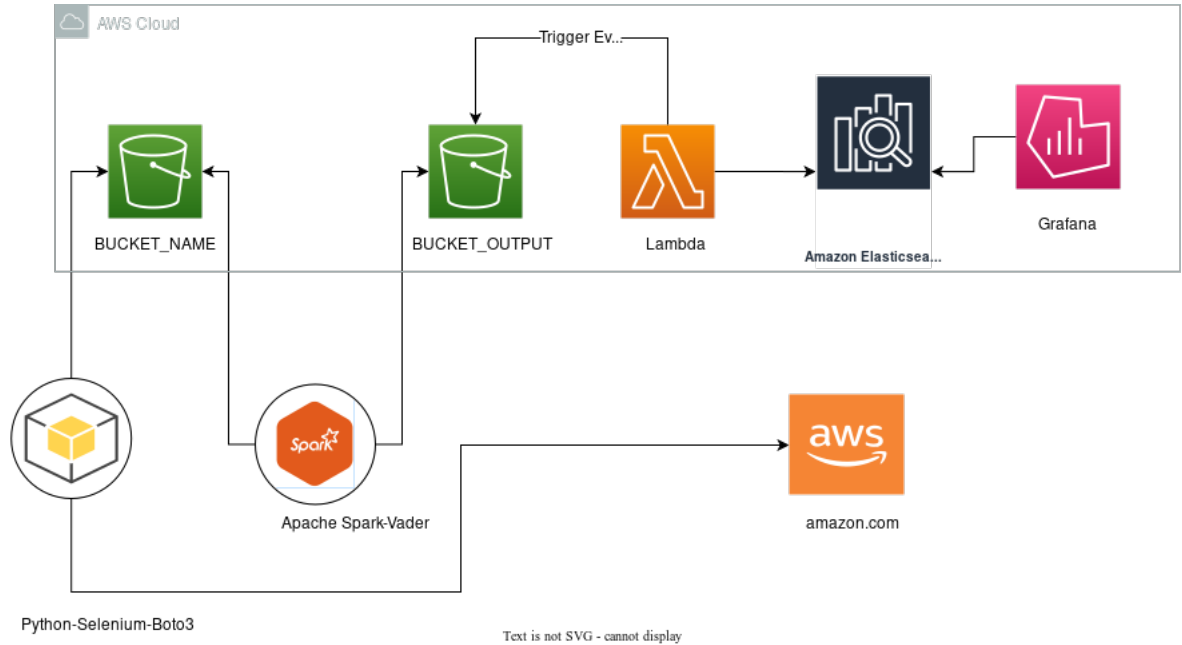


Figure 1: I Cerchi sono Pod di K8S

Come si può vedere dal grafico il flusso dei dati è il seguente: Python → S3 → Spark → S3 → Lambda → Amazon Opensearch Service/Grafana

2 Componenti

2.1 Python

Il primo componente in questa pipeline è il servizio di **scraping** dei dati da amazon. Tale servizio è costituito da una pod con al suo interno **Python** e **Selenium**, quest'ultimo è richiesto per eseguire il Javascript delle pagine web. Ci sono differenti parametri che si possono settare sul **python-deploment.yaml**

- **CODE_PRODUCT**: codice del prodotto da analizzare
- **START_PAGE**: numero pagina da cui iniziare
- **END_PAGE**: numero pagina finale
- **MODE_REVIEWS**: "recent" or "useful"
- **DOMAIN_URL**: "it", "co.uk"

Talvolta può capitare che le richieste vengano bloccate tramite un **ReCaptcha**, per minimizzare la possibilità che ciò avvenga viene usata una libreria python **fake_useragent**, usata per arricchire le richieste e renderle più credibili. Per il parsing delle pagine viene usata la libreria **BeautifulSoup**, per estrarre i dati interessanti:

- title, titolo della recensione
- rating, le stelle date al prodotto (1 a 5)
- body, la recensione
- date, data della recensione
- name, nome utente del utente
- verified_buy, se l'utente è verificato o meno
- helpful_vote, se la recensione ha avuto voti positivi
- country, nazione

Estratti i dati e trasformati in JSON, viene usato **boto3**, per inviare i dati sul bucket S3 [2].

```
s3 = boto3.client('s3',aws_access_key_id=AWS_ACCESS_KEY_ID,
                  aws_secret_access_key=AWS_SECRET_ACCESS_KEY )
# print(dir(s3))
# print(review)
global counter
s3.put_object(
    Body=json.dumps(review),
    Bucket=BUCKET_NAME,
    Key=NAME_FILES_S3+str(counter)
)
```

Figure 2: Boto3

2.2 S3

I dati dal python vengono memorizzati su un bucket S3 [3].

2.3 Apache Spark

Il servizio di Apache Spark che gira su un altro pod, va a prendersi i dati dal bucket ed effettua delle operazioni su di esso tra cui il calcolo dello score/polarità della recensione, per capire se l'accezione è negativa, positiva o neutra, infine li va a mettere i dati processati su un altro bucket [4].

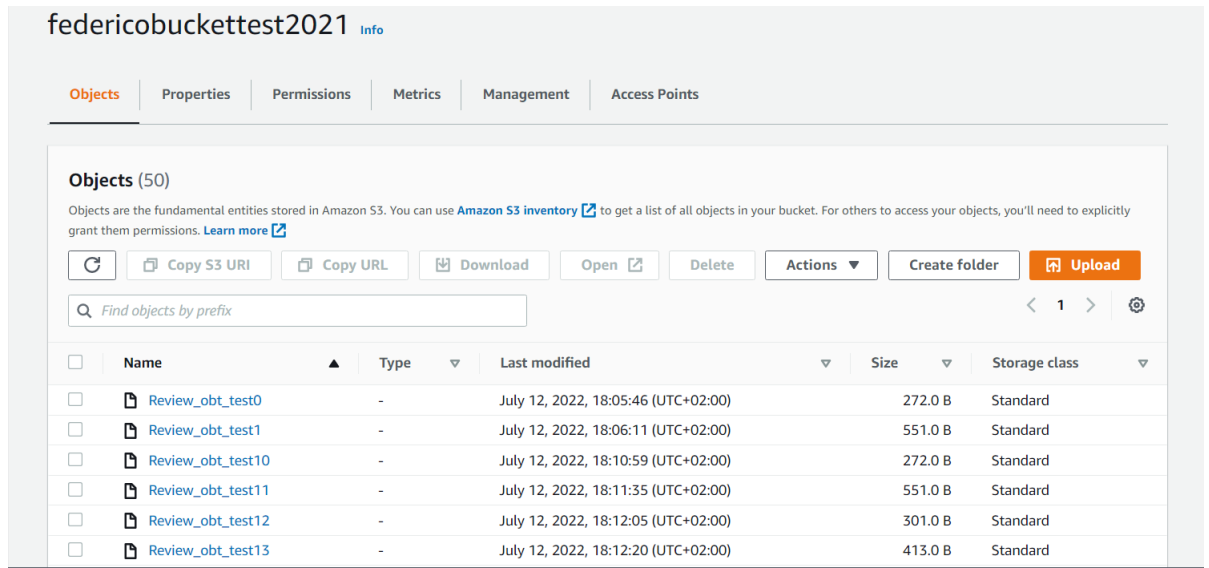


Figure 3: S3 contenente i dati originali

```
sentimen=udf(get_sentiment,DoubleType())
splitt=udf(splitting,ArrayType(StringType()))
df=df.withColumn("sentiment",sentimen("title"))
df=df.withColumn("words",splitt("title"))

df.writeStream\
  .format("json")\
  .option("checkpointLocation", "./checkpoints") \
  .option("path", "s3a:///"+BUCKET_OUTPUT +"/")\
  .start()\
  .awaitTermination()
```

Figure 4: Spark code

2.4 Lambda

Quando Spark mette i dati processati sull'altro bucket, viene **triggerata** una lambda function, che prende i dati e li mette su **opensearch** ⁵. Bisogna configurare alcune variabili d'ambiente ⁷.

```

1 import boto3
2 import re
3 import requests
4 from requests_aws4auth import AWS4Auth
5 import os
6 import json
7
8
9 region = 'eu-west-3' # e.g. us-west-1
10 service = 'es'
11 # credentials = boto3.Session().get_credentials()
12 awsauth = AWS4Auth(os.getenv("aws_access_key_id"), os.getenv("aws_secret_access_key"), region, service)
13 host = 'https://search-test-opensearch-progetto-sis-gwmm1rrvnf6xb33q5tmkdj7yq1.eu-west-3.es.amazonaws.com' # the OpenSearch Servi
14 index = 'reviews_amazon'
15 type = '_doc'
16 url = host + '/' + index + '/' + type
17 headers = { "Content-Type": "application/json" }
18 s3 = boto3.client('s3',aws_access_key_id=os.getenv("aws_access_key_id"),aws_secret_access_key=os.getenv("aws_secret_access_key"))
19 # Regular expressions used to parse some simple log lines
20 # Lambda execution starts here
21 def handler(event, context):
22     print("PROVA\n")
23     for record in event['Records']:
24
25         bucket = record['s3']['bucket']['name']
26         key = record['s3']['object']['key']
27
28         # Get, read, and split the file into lines
29         obj = s3.get_object(Bucket=bucket, Key=key)
30
31         body=obj['Body'].read()
32         lines=body.splitlines()
33         for line in lines:
34             r = requests.post(url, auth=awsauth, json=json.loads(line) , headers=headers)

```

Figure 5: Editor variabili lambda

General configuration	Environment variables (6) Edit														
Triggers	The environment variables below are encrypted at rest with the default Lambda service key.														
Permissions															
Destinations															
Function URL															
Environment variables	<table> <tr> <th>Key</th><th>Value</th></tr> <tr> <td>aws_access_key_id</td><td>AKIARLBECYJEWZCKTXQS</td></tr> <tr> <td>aws_secret_access_key</td><td>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</td></tr> <tr> <td>domain_url_opensearch</td><td>https://search-cluster-progetto-amk-cloud-7towj5eggx2razpypiw3uzyloe.us-east-1.es.amazonaws.com</td></tr> <tr> <td>index_opensearch</td><td>reviews_amazon</td></tr> <tr> <td>password_opensearch</td><td>p*xXar8UP828</td></tr> <tr> <td>username_opensearch</td><td>federico-master</td></tr> </table>	Key	Value	aws_access_key_id	AKIARLBECYJEWZCKTXQS	aws_secret_access_key	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	domain_url_opensearch	https://search-cluster-progetto-amk-cloud-7towj5eggx2razpypiw3uzyloe.us-east-1.es.amazonaws.com	index_opensearch	reviews_amazon	password_opensearch	p*xXar8UP828	username_opensearch	federico-master
Key	Value														
aws_access_key_id	AKIARLBECYJEWZCKTXQS														
aws_secret_access_key	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														
domain_url_opensearch	https://search-cluster-progetto-amk-cloud-7towj5eggx2razpypiw3uzyloe.us-east-1.es.amazonaws.com														
index_opensearch	reviews_amazon														
password_opensearch	p*xXar8UP828														
username_opensearch	federico-master														
Tags															
VPC															
Monitoring and operations tools															

Figure 6: Editor variabili lambda

2.5 OpenSearch

L'ultimo componente è Opensearch di AWS, che permette di creare e visualizzare dashboard per fare delle analisi sui dati ad esempio delle **World Cloud**.

2.6 Grafana

Si possono visualizzare i dati anche tramite Grafana.

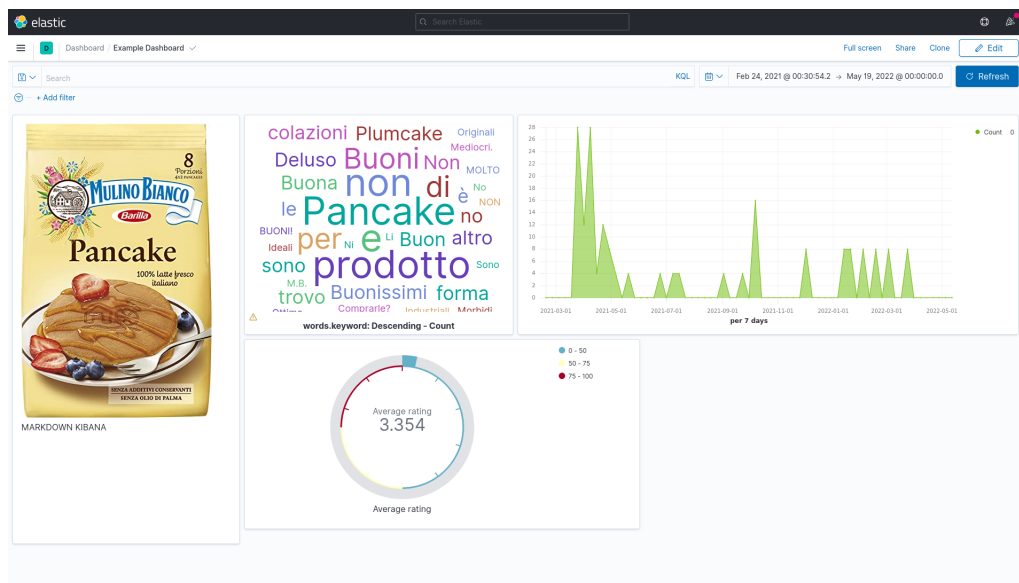


Figure 7: Dashboard