

Analisi Spark Core Covid-19 Opendata Vaccini

Roberto Pavia

Corso di laurea magistrale in Ingegneria Informatica

Università degli studi di Roma "Tor Vergata"

roberto.pavia@alumni.uniroma2.eu

Sommario—Questo documento descrive l'implementazione Spark Core per l'analisi di datasets opensource riguardanti le vaccinazioni Covid-19 in Italia. Vengono riportati in dettaglio i diversi frameworks utilizzati ed una visione dettagliata dell'architettura.

I. INTRODUZIONE

Lo studio effettuato, risponde a due queries usando tre datasets differenti. Il primo dataset `punti-somministrazione-tipologia.csv` contiene i dati sui punti di somministrazione per ciascuna regione italiana. Questo dataset è composto da diversi campi e l'unico di interesse a quest'analisi è `nome_area`. Tipo di dato string che indica la denominazione standard dell'area dov'è presente tale hub vaccinale. Il secondo dataset `somministrazioni-vaccini-latest.csv` contiene dati sulle somministrazioni giornaliere dei vaccini suddivisi per regioni, fasce d'età e categorie di appartenenza dei soggetti vaccinati. Sono stati considerati i campi d'interesse `data_somministrazione`, `fascia_anagrafica`, `Sesso_femminile` e `nome_area`. Il terzo dataset `somministrazioni-vaccini-summary-latest.csv` contiene dati sul totale delle somministrazioni giornaliere per regioni e categorie di appartenenza dei soggetti vaccinati. In questo caso i campi d'interesse sono `data_somministrazione`, `totale`, `nome_area`. Di seguito le motivazioni della scelta di questi campi e l'implementazione delle diverse queries.

II. QUERIES

A. Q1

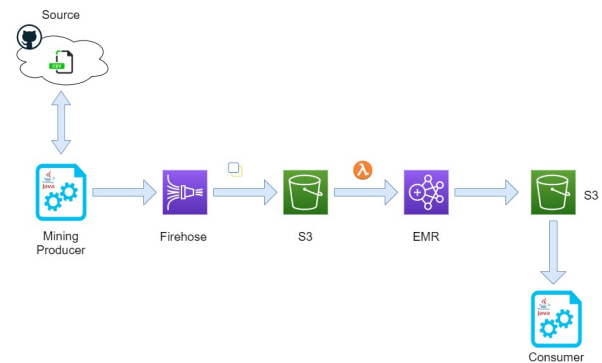
Nella prima query, per ogni mese solare e per ciascuna area si deve calcolare il numero medio di somministrazioni che è stato effettuato giornalmente in un centro vaccinale generico in quell'area e durante quel mese. Come da specifica, vanno considerati i dati a partire dal 01/01/2021. Tenendo quindi in considerazione i datasets `punti-somministrazione-tipologia.csv` e `somministrazioni-vaccini-summary-latest.csv` sono stati rispettivamente calcolati:

- Numero di hubs vaccinali totali per ogni regione sfruttando il numero di occorrenze del dato `nome_area`
- Numero medio di vaccinazioni totali in un mese raggruppando il totale filtrato sulla `data_somministrazione` e mappato su `nome_area`

B. Q2

Nella seconda query, per le donne, per ogni fascia anagrafica e per ogni mese solare, bisogna determinare le prime 5 aree per le quali è previsto il maggior numero di vaccinazioni il primo giorno del mese successivo. Per determinare la classifica mensile e prevedere il numero di vaccinazioni si è dovuto considerare la retta di regressione che approssima l'andamento delle vaccinazioni giornaliere. Per ogni mese e categoria si è calcolata la classifica a partire dai dati raccolti dal 01/02/2021. Si è tenuto in considerazione il dataset `somministrazioni-vaccini-latest.csv`

III. ARCHITETTURA



A. Pulizia dei dati e Ingestion

I datasets proposti sono in generale molto puliti, tuttavia è stato necessario prendere delle accortezze. Inoltre, come già notato, per rispondere alle diverse queries serve considerare un numero ridotto di colonne. Per sfruttare queste caratteristiche naturali dei dati si sarebbe potuto procedere in diversi modi:

- Parsing dei dati tramite un crawler Github con filtering and fix prima della fase di ingestion
- Filtering and fix durante la fase di ingestion
- Subito prima del processamento all'interno del cluster Spark

Name	Status	Creation time	Source	Data transformation	Destination
stream-query1-file1	Active	2021-06-04T12:40+0200	Direct PUT and other sources	Disabled	Amazon S3 file1-query1
stream-query1-file2	Active	2021-06-04T12:41+0200	Direct PUT and other sources	Disabled	Amazon S3 file1-query2

La scelta è stata quella di sfruttare un *Mining Producer* per parsare i dati dalla repository Github via HTTP request per poi sfruttare *Kinesis Firehose*. Trattasi di un servizio Amazon completamente gestito e di facile utilizzo per caricare flussi di dati in data lake, datastore e servizi di analisi in modo

```

graph LR
    SR[Source records] --> T[Transform source records  
Invoke Lambda function]
    T --> C[Convert record format  
Refer to Glue table for schema]
    C --> PR[Processed records]
  
```

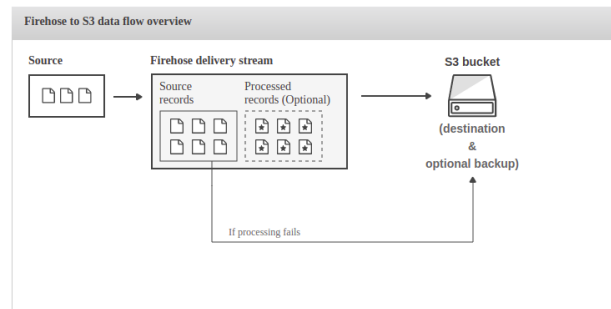
The diagram shows a workflow starting with 'Source records' (four document icons). An arrow points to a dashed box containing two main steps: 'Transform source records' (with subtext 'Invoke Lambda function') and 'Convert record format' (with subtext 'Refer to Glue table for schema'). An arrow from the second step points to 'Processed records' (three star icons). Below the dashed box, the text '----- Optional' is shown.

RegExp: `(?=[^"]*"*)*[""]*`

Sample: `col1, col2, "col3, col", col4` ✓

B. App Spark e storage

- spark-submit via lambda function
- spark-submit manuale via SSH



- Registrazione di log su bucket S3
- La modalità d'avvio è impostata su Cluster e non esecuzione fase (cosa che sarebbe stata prevista tramite lambda)
- La versione di EMR è la 6.3.0 per poter implementare il codice in Spark 3.1.1
- Il tipo di istanza utilizzata, con pattern general purpose, è una `m3.xlarge`. Un nodo master e due principali.

[illegible]

In parallelo all'avvio del CovidCluster si procede a generare il jar dell'applicazione Spark eseguendo una build artifacts del codice. Sarà necessario specificare la main class prima di fare la build del file jar da distribuire secondo lo schema SPMD. Nel caso in cui questo non venga fatto è comunque possibile specificare la main class tramite il comando di `spark-submit` usando il parametro `--class`. Dopo aver caricato il jar in S3, dalla console EMR acceduta via SSH copiamo in memoria il jar.

C. MLib RDD-based

Per l'implementazione dell'algoritmo di regressione lineare si sono riscontrati problemi con l'utilizzo della classe `StreamingLinearRegressionWithSGD()` nella versione 3.1.1 di Spark. In particolare l'integrazione in Scala, seppur fatta seguendo la [documentazione ufficiale](#) è risultata problematica a livello di librerie. Viene infatti restituito un problema con l'inizializzazione del regressore. Sono state provate più versioni della stessa classe ma il risultato è stato sempre lo stesso anche in diversi ambienti di sviluppo. Si è deciso quindi di implementare l'algoritmo di regressione lineare classica direttamente in Scala nell'applicazione Spark. Nella classe `utils.Helper.scala` si può trovare il calcolo dei coefficienti della regressione mediante una soluzione veloce del calcolo della media. Questa soluzione prevede il calcolo della media `as we go`, ovvero aggiungendo il nuovo valore calcolato ogni volta alla media calcolata fino a quel punto, sfruttando la logica delle tuple.

D. Consumer, Monitoring e risultati

L'applicazione Consumer è speculare a quella del Producer. Allo stesso modo, sfruttando l'SDK di AWS disponibile sia per Kinesis Firehose che per S3, è stato simulato il download delle diverse partizioni create su un dato bucket dall'applicazione Spark. Le diverse partizioni vengono salvate all'interno di una cartella locale del Consumer, il quale potrà fare su queste monitoraggio, visualizzazione e statistiche varie. Non avendo la possibilità diretta in Spark Core di salvare un RDD come un file CSV, se non passando per il framework di SparkSQL, il consumer si occuperà di produrre i risultati da presentare lavorando sulle partizioni.

<input type="checkbox"/>	Nome	Tipo	Ultima modifica	Dimensioni	Classe di storage
<input type="checkbox"/>	..SUCCESS	-	04 Jun 2021 01:56:02 PM CEST	0 B	Standard
<input type="checkbox"/>	part-00000	-	04 Jun 2021 01:56:02 PM CEST	4.5 KB	Standard
<input type="checkbox"/>	part-00001	-	04 Jun 2021 01:56:02 PM CEST	3.0 KB	Standard

Come già spiegato in precedenza il cluster è composto da un nodo master e da due nodi core instance tutti su macchine `m3.xlarge`. Di seguito vengono riportati i relativi gruppi con le caratteristiche hardware specifiche.

Gruppi di istanze					
Filter: Filtra i gruppi di istanze (tutti caricati)					
ID	Stato	Tipo e nome nodo	Tipo di istanza	Numero istanze	
ig-2V3L0B6UD3LGT	Terminato (1 richiesto)	MASTER Master Instance Group	m3.xlarge 4 vCore, 15 GB di memoria, 80 SSD GB di storage Storage EBS: nessuno	0 Istanze	
ig-8GKP1SUED1W	Terminato (2 richiesto)	CORE Core Instance Group	m3.xlarge 4 vCore, 15 GB di memoria, 80 SSD GB di storage Storage EBS: nessuno	0 Istanze	

La configurazione prevede una chiave di accesso al nodo master mostrato sopra per l'accesso SSH necessario ad eseguire `spark-submit`

Per concludere riportiamo le tempistiche prese direttamente dalla console di AWS. I diversi processi hanno terminato con stato di `Riuscito`. L'applicazione esegue le due query, con output di log minimale per evitare rallentamenti. Nei tempi riportati, oltre all'esecuzione delle diverse fasi, è incluso anche il tempo di salvataggio dei dati sui bucket.

Processi Fasi Esecutori						
Utente: hadoop						
Tempo di attività totale: 46 s						
Processi completati: 5						
Tempistiche eventi						
Processi (5)						
Filter: Filtra i processi (tutti caricati)						
ID processo	Descrizione	Inviato (UTC+2)	Durata	Fasi	Attività	
4	Riuscito runJob at SparkHadoopWriter.scala:83	2021-06-04 13:56 (UTC+2)	0.9 s	2 / 5	4 / 10	
3	Riuscito sortBy at VaccinationCovid19.scala:95	2021-06-04 13:55 (UTC+2)	2 s	3 / 4	6 / 8	
2	Riuscito sortBy at VaccinationCovid19.scala:79	2021-06-04 13:55 (UTC+2)	10 s	2 / 2	4 / 4	
1	Riuscito runJob at SparkHadoopWriter.scala:83	2021-06-04 13:55 (UTC+2)	0.6 s	2 / 5	4 / 10	
0	Riuscito sortBy at VaccinationCovid19.scala:95	2021-06-04 13:55 (UTC+2)	6 s	4 / 4	8 / 8	

L'applicazione Spark in output ritorna un messaggio di log per ogni nuovo stage creato nel DAG. Si è tenuto conto quindi di quando all'interno dello Spark Core avviene una fase di re-shuffle. Si può notare come la classe `MultipartUploadOutputStream` ritorna in output il momento di chiusura del canale verso i due bucket, sui quali scrivere le diverse parti dei risultati ottenuti. Come ultima riga il comando di `spark-submit` utilizzato con una jar prodotto senza la main class e quindi con la necessità di specificare come già spiegato in precedenza.

```
21/06/04 13:55:33 INFO utils: Using initial executors = 50, max of spark.dynamicAllocation.initialExecutors, spark.dynamicAllocation.minExecutors and spark.executor.instances
21/06/04 13:55:33 INFO SparkClientSchedulerBackend: SchedulerBackend is ready for scheduling beginning after reached minRegisteredResourceRatio: 0.0
21/06/04 13:55:33 INFO VACCINATIONCOVID19_QUERY1: started!
21/06/04 13:55:34 INFO ClientConfiguration: loaded native-gli library
21/06/04 13:55:34 INFO LocalCodecs: Successfully loaded & initialized native-ito library (hadoop-ito rev 049362b7cf33ff5f73d6d6153245772cd445e)
21/06/04 13:55:35 INFO ClientConfiguration: set initial getobject socket timeout to 2000 ms
21/06/04 13:55:37 INFO VACCINATIONCOVID19_QUERY1: shuffle -> new stage created!
21/06/04 13:55:37 INFO ClientConfiguration: set initial getobject socket timeout to 2000 ms
21/06/04 13:55:37 INFO VACCINATIONCOVID19_QUERY1: shuffle -> new stage created!
21/06/04 13:55:37 INFO VACCINATIONCOVID19_QUERY1: joining RDDs
21/06/04 13:55:45 INFO MultipartUploadOutputStream: close closed: false s3://file-query/2021-06-04_Q1/SUCCESS
21/06/04 13:55:47 INFO VACCINATIONCOVID19_QUERY1: ended
21/06/04 13:55:47 INFO VACCINATIONCOVID19_QUERY2: started
21/06/04 13:55:47 INFO ClientConfiguration: set initial getobject socket timeout to 2000 ms
21/06/04 13:55:48 INFO LocalCodecs: Successfully loaded & initialized native-ito library
21/06/04 13:55:48 INFO ClientConfiguration: set initial getobject socket timeout to 2000 ms
21/06/04 13:55:50 INFO VACCINATIONCOVID19_QUERY2: shuffle -> new stage created!
21/06/04 13:55:50 INFO ClientConfiguration: set initial getobject socket timeout to 2000 ms
21/06/04 13:55:50 INFO VACCINATIONCOVID19_QUERY2: regression -> R^2=0.95
21/06/04 13:56:01 INFO MultipartUploadOutputStream: close closed: false s3://file-query/2021-06-04_Q2/SUCCESS
21/06/04 13:56:02 INFO VACCINATIONCOVID19_QUERY2: ended
[hadoop@ip-172-31-87-2 ~]$ spark-submit --class VaccinationCovid19 --jars vaccinationCovid19.jar
```

E. Conclusioni e limitazioni

Rispetto ai test fatti in locale, notiamo un uptime del cluster di circa il doppio sulla piattaforma di Cloud Computing AWS. Per un'applicazione molto semplice è difficile notare grandi differenze, anzi sembrerebbe funzionare meglio in locale. Mantenendo le configurazioni tra locale e cloud, il più simili possibile, in locale si sta lavorando con un `intel-core i5 4200u` equipaggiato con 8 GB di memoria RAM. Con applicazioni più complesse e query computazionalmente più pesanti il gap creatosi tra le prestazioni su cloud e quelle in locale tenderà a svanire invertendosi. Inoltre è facile prevedere tempi migliori in cloud senza l'aggiunta del framework di SparkSQL sopra Spark Core, perdendo però la possibilità di trattare i dati in modo agile e strutturato. In questo caso le dimensioni ridotte dei datasets, hanno influito sui tempi presentati. Un limite riscontrato è quello del `vendor lock-in`. Bisogna infatti notare come i servizi di AWS tendono ad integrarsi alla perfezione tra di loro ma si è limitati nell'integrazione con altri servizi di cloud computing come ad esempio quello offerto da Google.

Completed Jobs (5)						
Filter: Filtra i gruppi di istanze (tutti caricati)						
ID	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total	
4	runJob at SparkHadoopWriter.scala:83	2021/06/07 17:30:18	0.2 s	2/2 (3 skipped)	2/4 (3 skipped)	
3	sortBy at VaccinationCovid19.scala:95	2021/06/07 17:30:16	2 s	3/3 (1 skipped)	6/8 (2 skipped)	
2	sortBy at VaccinationCovid19.scala:80	2021/06/07 17:30:12	3 s	2/2	4/4	
1	runJob at SparkHadoopWriter.scala:83	2021/06/07 17:30:12	0.5 s	2/2 (3 skipped)	4/8 (3 skipped)	
0	sortBy at VaccinationCovid19.scala:95	2021/06/07 17:30:09	3 s	4/4	8/8	