



Esempi di Apache Spark

Programmazione

Apache Spark è un framework open-source per l'elaborazione distribuita dei dati, progettato per gestire in modo estremamente rapido grandi volumi di informazioni su cluster di macchine.

Nasce come evoluzione del modello MapReduce di Hadoop, migliorandone drasticamente le prestazioni grazie all'uso massiccio della memoria RAM: invece di scrivere continuamente su disco tra un'operazione e l'altra, Spark conserva i dati in memoria, riducendo la latenza e rendendo possibile l'analisi interattiva, il machine learning e lo streaming in tempo reale.

La sua architettura si basa su un driver che coordina il lavoro e su più executor che processano i dati in parallelo, distribuendo automaticamente il carico e gestendo la tolleranza ai guasti.



0. SparkSession + DataFrame

La SparkSession è l'interfaccia principale per lavorare con Apache Spark: rappresenta il punto d'ingresso per creare DataFrame, leggere file, configurare l'applicazione e utilizzare tutti i moduli (SQL, Streaming, MLlib).

Un DataFrame è una collezione distribuita di dati organizzati in colonne, con schema e tipi, ottimizzato tramite il motore Catalyst e progettato per gestire grandi volumi in modo parallelo e fault-tolerant.



0. Setup: creare la SparkSession e un DataFrame di esempio

```
1.from pyspark.sql import SparkSession
2.
3.# Creazione della SparkSession (punto di ingresso di Spark)
4.spark = SparkSession.builder \
5.  .appName("EsempiSpark") \
6.  .getOrCreate()
7.
8.# Dati di esempio in memoria (lista di tuple)
9.data = [("Mario", 30, "IT"),
10.   ("Luca", 25, "FR"),
11.   ("Anna", 40, "IT")]
12.
13.# Creazione di un DataFrame con schema esplicito
14.df = spark.createDataFrame(data, ["nome", "eta", "paese"])
15.
16.df.show()
17.df.printSchema()
```

Cosa fa:

- **Crea la sessione Spark.**
- **Costruisce un DataFrame a partire da una lista Python.**
- **Mostra i dati e lo schema (tipi delle colonne).**



1. Select e Filter

Le operazioni di selezione (select) e filtraggio (filter/where) permettono di estrarre specifiche colonne o righe da un DataFrame.

Sono trasformazioni lazy: Spark non le esegue subito, ma costruisce un piano logico ottimizzato che verrà calcolato solo quando servirà un'azione.

Queste operazioni sono fondamentali per pulizia dati, pre-processing e creazione di viste ridotte.



1. Selezionare colonne e filtrare righe (select, filter / where)

```
1.from pyspark.sql.functions import col  
2.  
3.# Selezionare solo alcune colonne  
4.df_nome_eta = df.select("nome", "eta")  
5.  
6.# Filtrare le righe con eta > 30  
7.df_over_30 = df.filter(col("eta") > 30)  
8.  
9.df_nome_eta.show()  
10.df_over_30.show()  
11.
```

Cosa fa:

- **select** sceglie un **sottoinsieme** di colonne.
- **filter (o where)** applica una **condizione** sulle righe tramite `col("eta") > 30`.



2. **withColumn (nuove colonne e trasformazioni)**

withColumn permette di aggiungere una nuova colonna o sovrascriverne una esistente.

È essenziale per creare feature derivate, normalizzazioni, conversioni di tipi e calcoli vettoriali.

Spark applica queste trasformazioni in modo distribuito su tutte le partizioni, senza modificare il DataFrame originale (immutabilità funzionale).



2. Aggiungere una colonna calcolata (withColumn)

Cosa fa:

- **withColumn crea (o sovrascrive) una colonna.**
- **lit crea una costante (qui +10).**
- **L'operazione è vettoriale su tutte le righe.**

```
1. from pyspark.sql.functions import lit  
2.  
3. # Aggiungere una colonna "eta_tra_10_anni"  
4. df_eta_future = df.withColumn("eta_tra_10_anni", col("eta") + lit(10))  
5.  
6. df_eta_future.show()  
7.
```



3. GroupBy e Aggregazioni

Il groupBy permette di raggruppare i dati in base al valore di una o più colonne, mentre agg consente di applicare funzioni di aggregazione (conteggio, media, somma, min/max).

Si tratta di un'operazione chiave nell'analisi dei dati e nei sistemi distribuiti, perché richiede uno shuffle: Spark redistribuisce le righe in base ai gruppi e combina i risultati parallelamente.



3. Raggruppare e aggregare (groupBy, agg, funzioni di aggregazione)

Cosa fa:

- **groupBy("paese")** crea gruppi per valore di paese.
- **agg** applica funzioni di aggregazione (**avg**, **count**).
- **alias** rinomina le colonne risultanti.

```
1. from pyspark.sql.functions import avg, count
2.
3. # Calcolare l'eta media per paese e quanti utenti per paese
4. df_stats_paese = df.groupBy("paese") \
5.   .agg(
6.     avg("eta").alias("eta_media"),
7.     count("*").alias("num_persone")
8.   )
9.
10. df_stats_paese.show()
11.
```



4. Join tra DataFrame

Le join consentono di combinare due DataFrame su una colonna comune, analogamente ai database relazionali.

Spark ottimizza automaticamente il tipo di join (broadcast, shuffle) in base alla dimensione delle tabelle.

Questo permette di collegare dataset diversi, come tabelle anagrafiche, log, lookup mantenendo coerenza con i paradigmi SQL ma con esecuzione distribuita.



4. Join tra DataFrame

Cosa fa:

- Crea un secondo DataFrame con la “dimensione” paesi.
- Esegue una join (inner) sulla colonna comune paese.
- Aggiunge al DataFrame originale il nome esteso del paese.

```
1. # Secondo DataFrame con info sui paesi  
2. data_paesi = [("IT", "Italia"), ("FR", "Francia")]  
3. df_paesi = spark.createDataFrame(data_paesi, ["paese", "nome_paese"])  
4.  
5. # Join sui paesi (inner join)  
6. df_join = df.join(df_paesi, on="paese", how="inner")  
7.  
8. df_join.show()  
9.
```



5. Lettura CSV e scrittura Parquet

Spark supporta la lettura e scrittura di dati in molti formati (CSV, JSON, Parquet, Avro).

La lettura da CSV è utile per ingest iniziali e dati provenienti da sistemi legacy; il Parquet, invece, è un formato colonnare compresso progettato per query analitiche ad alte prestazioni.

Spark gestisce automaticamente inferenza dello schema, partizionamento e parallelizzazione dell'I/O.



5. Lettura da file CSV e scrittura su disco

Cosa fa:

- **spark.read.csv legge un file CSV, con intestazione e schema inferito.**
- **write.parquet salva i dati in formato Parquet, più efficiente di CSV per analisi.**

```
1. # Lettura da CSV (schema inferito)
2. df_csv = spark.read \
3.   .option("header", "true") \
4.   .option("inferSchema", "true") \
5.   .csv("path/dati_clienti.csv")
6.
7. df_csv.show(5)
8.
9. # Scrittura in formato Parquet (colonnare)
10. df_csv.write \
11.   .mode("overwrite") \
12.   .parquet("path/output_clienti_parquet")
13.
```



6. Cache e Persistenza

La cache permette di memorizzare in RAM il risultato di una trasformazione per evitare ricalcoli.

È fondamentale per algoritmi iterativi (machine learning, analisi ricorsive) e per workflow che usano più volte lo stesso DataFrame.

Spark conserva i dati nelle partizioni dei nodi finché servono, migliorando le performance e riducendo i costi di computazione.



6. Cache / persistenza in memoria (cache)

Cosa fa:

- Costruisce un DataFrame filtrato e con colonna calcolata.
- `cache()` chiede a Spark di mantenerlo in memoria dopo la prima azione.
- Successive action (come `count()` e `show()`) saranno più veloci.

```
1.# Preparare un DataFrame derivato "pesante" e metterlo in cache  
2.df_derivato = df.filter(col("paese") == "IT") \  
3.      .withColumn("eta_x2", col("eta") * 2) \  
4.      .cache() # o .persist()  
5.  
6.# Prima azione: esegue la pipeline e mette in cache  
7.df_derivato.count()  
8.  
9.# Le azioni successive riusano i dati in RAM  
10.df_derivato.show()
```



7. Structured Streaming

Structured Streaming abilita l'elaborazione di flussi di dati in tempo reale usando la stessa API dei DataFrame batch.

Il flusso è trattato come una tabella che si aggiorna continuamente man mano che arrivano nuovi eventi.

Spark gestisce checkpointing, fault-tolerance e aggiornamento incrementale, consentendo casi d'uso come conteggio di parole, ingest real-time, dashboard e pipeline di log/telemetria.



7. Esempio super base di Structured Streaming

Solo per dare l'idea del pattern, non da lanciare così com'è in produzione.

Cosa fa:

- Legge uno stream di testo da una socket (host:port).
- Divide le righe in parole, le raggruppa e le conta.
- Stampa le word count aggiornate continuamente sulla console.

```
1. from pyspark.sql.functions import explode, split
2.
3. # Lettura di uno stream di testo (es. socket)
4. lines = spark.readStream \
5.   .format("socket") \
6.   .option("host", "localhost") \
7.   .option("port", 9999) \
8.   .load()
9.
10. # Split delle righe in parole
11. words = lines.select(
12.   explode(split(col("value"), " ")).alias("word"))
13.)
14.
15. # Conteggio parole in tempo quasi-reale
16. word_counts = words.groupBy("word").count()
17.
18. # Scrittura su console dello stream
19. query = word_counts.writeStream \
20.   .outputMode("complete") \
21.   .format("console") \
22.   .start()
23.
24. # Attesa della terminazione (bloccante)
25. query.awaitTermination()
```





Buon Davante a tutti