

UNIVERSITÀ TELEMATICA INTERNAZIONALE
UNINETTUNO

INGEGNERIA INFORMATICA
Corsi di Laurea Magistrale - Big Data

Introduzione ai Big Data



Progetto su R, MapReduce e Hadoop "Contabilità ditta"

Docente:
Prof. Pirrone Daniele

Studente:
Emanuele Coltro
mat: 671HHHINGINFOR

Anno Accademico 2022-2023

Abstract

Il presente paper ha l'obiettivo di esaminare le potenzialità e le performance degli strumenti utilizzati per la gestione e l'analisi dei Big Data. Viene approfondito il processo che porta alla produzione di report grafici, tabelle e schemi a partire dai dati di un problema o di una esigenza. La relazione fornisce informazioni dettagliate su come impostare un progetto, dalla scelta e impostazione dell'ambiente di sviluppo alla contestualizzazione delle decisioni prese. Inoltre, vengono allegati i codici per la generazione dei risultati e dei grafici.

Il progetto in esame mira a soddisfare l'esigenza dell'azienda (fittizia) STEC-CAPARAPETUTTI S.R.L. di risolvere un problema di contabilità mediante l'estrazione di alcune feature dai dati. A tal fine, sono stati impiegati strumenti avanzati per l'elaborazione dei dati, tra cui algoritmi di MapReduce e tecniche di visualizzazione avanzata.

L'analisi dei dati è stata eseguita utilizzando il linguaggio R per una prima analisi, seguito dall'utilizzo di un tool scritto in JavaScript per prototipare il sistema di MapReduce e risolvere il problema.

Indice

1	Introduzione	3
1.1	Richieste	3
1.2	Introduzione all'ambiente	3
1.2.1	MapReduce	3
1.2.2	Hadoop	4
1.2.3	R e RStudio	6
1.2.4	Python	6
1.2.5	Tool Javascript-html per MapReduce	6
1.3	Preparazione dell'ambiente	7
1.3.1	Preparazione dell'ambiente virtuale	7
1.3.2	Installazione di Hadoop	8
1.3.3	Integrazione di RStudio	8
1.3.4	Configurazione di Python	8
1.3.5	Conclusione configurazione	8
2	Svolgimento	10
2.1	Analisi del problema	10
2.2	Analisi del Dataset	10
2.3	Risoluzione tramite codice R	11
2.4	MapReduce usando Javascript	11
2.5	MapReduce usando R e Hadoop	13
2.5.1	Scripting	13
2.5.2	Implementazione in R	14
2.5.3	MapReduce con Python e Hadoop	15
2.6	Rapporto finale	15
2.6.1	Job 1 - Media di vendita per ogni mese	16
2.6.2	Job 2 - Varianza di vendita per ogni mese	16
2.6.3	Job 3 e 4 - Mesi migliori e peggiori di ogni anno	17
2.6.4	Analisi	18
3	Conclusioni	19
	Allegato	20

1 Introduzione

In questa fase, forniremo le basi dei concetti teorici e degli strumenti utilizzati nel corso del progetto. Vedremo in dettaglio come sia stato impostato l'ambiente di sviluppo e come sia stato affrontato il problema in analisi.

1.1 Richieste

La società fittizia STECCAPARAPETUTTI S.R.L. richiede l'analisi dei dati delle vendite del sistema di contabilità dell'anno precedente. I dati sono contenuti in un file CSV strutturato in modo omogeneo. L'analisi richiede il calcolo della media e della varianza delle vendite per ogni mese di ogni anno, l'identificazione del mese di ogni anno con la maggiore e minore vendita. Per ogni lavoro, è necessaria la documentazione dell'implementazione in R e dell'implementazione del paradigma MapReduce di Hadoop. Il risultato finale deve essere presentato sotto forma di un rapporto completo.

1.2 Introduzione all'ambiente

1.2.1 MapReduce

L'algoritmo MapReduce è un modello di programmazione e un'infrastruttura di elaborazione distribuita utilizzata per elaborare e generare informazioni da enormi quantità di dati in modo efficiente, sicuro e scalabile. È stato introdotto da Google[1] e ha giocato un ruolo fondamentale nel campo del data processing su larga scala.

L'idea alla base di MapReduce è suddividere un grande task di elaborazione dei dati in due fasi principali: la fase di "map" e la fase di "reduce" (più una fase intermedia totalmente automatizzata):

- **Map:** Durante la fase di MapReduce, i dati di input vengono suddivisi in piccoli frammenti (split) e passati a un insieme di processi paralleli chiamati "Worker". Ognuno di essi esegue una funzione di mapping definita dall'utente che prende in ingresso un dato e produce una serie di coppie chiave-valore intermedie. Queste coppie rappresentano il risultato dell'elaborazione dei dati nella fase di mapping.
- **Shuffle e Sort:** Dopo la fase di map, il sistema raggruppa le coppie chiave-valore in base alle chiavi e le ordina nella fase di Shuffle e Sort. Questo è un passaggio cruciale poiché consente ai dati correlati di essere inviati allo stesso processo di "reduce".

- **Reduce:** Nell'ultima fase il sistema assegna le coppie chiave-valore raggruppate nuovamente ai processi worker. Ogni processo esegue una funzione di "reduce" definita dall'utente che agisce sulle coppie chiave-valore correlate e produce i risultati finali dell'elaborazione.

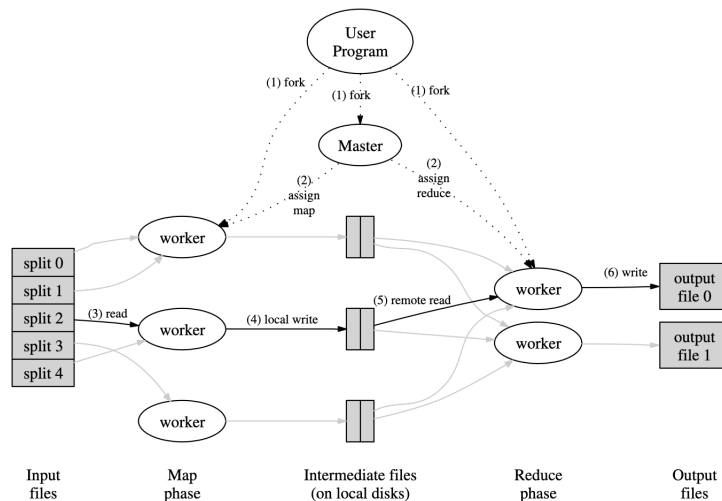


Figura 1: Schema del processo di MapReduce tratto dal paper di Google

L'obiettivo di MapReduce è la parallelizzazione e la distribuzione del carico di lavoro su un cluster di computer, consentendo di elaborare grandi quantità di dati in modo più veloce ed efficiente rispetto a un approccio sequenziale tradizionale. Questo modello di programmazione nasconde molti dettagli complessi dell'elaborazione distribuita, semplificando la creazione di applicazioni che possono sfruttare l'elaborazione su larga scala senza doversi preoccupare delle questioni di basso livello.

Questo ha portato negli anni alla nascita di tecnologie e framework, come Apache Hadoop e Apache Spark, che forniscono funzionalità di MapReduce ma più avanzate e versatili per l'elaborazione distribuita dei dati.

1.2.2 Hadoop

Apache Hadoop è un framework open-source sviluppato per consentire l'elaborazione distribuita di grandi quantità di dati su cluster di computer. È basato sulla paper originale di Google su MapReduce e il file system distribuito Google File System (GFS). Hadoop è progettato per gestire l'elaborazione di dati in parallelo su macchine commodity (hardware relativamente economico e accessibile).

Il cuore di Hadoop[3] è composto da due componenti principali: Hadoop Distributed File System (HDFS) e il framework di elaborazione MapReduce.

1. Hadoop Distributed File System (HDFS): HDFS è il sistema di archiviazione distribuito di Hadoop. Si basa su un modello di architettura master-slave in cui un nodo master, chiamato "NameNode", gestisce i metadati del file system e tiene traccia di dove sono archiviati i dati. I nodi slave, chiamati "DataNode", contengono i blocchi di dati reali. I dati vengono suddivisi in blocchi e replicati su vari DataNode per garantire la disponibilità e l'affidabilità.

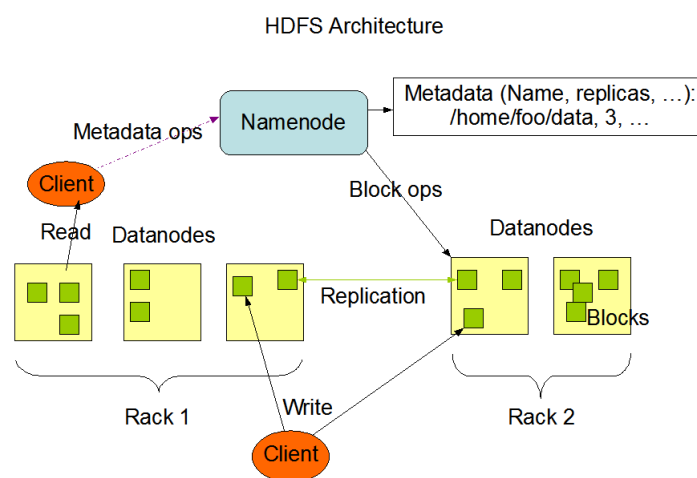


Figura 2: Schema dell'architettura HDFS tratto dal sito di Hadoop

2. MapReduce: Come descritto in precedenza, MapReduce è un modello di programmazione per l'elaborazione distribuita. Hadoop implementa questo modello consentendo agli sviluppatori di scrivere programmi MapReduce per l'elaborazione dei dati su cluster. Il framework si occupa della distribuzione delle attività di map e reduce su diversi nodi del cluster, dell'ordinamento e dell'aggregazione dei risultati intermedi e della gestione dei fallimenti dei nodi.

I punti di forza di Hadoop includono la scalabilità orizzontale, la tolleranza ai guasti, l'economia hardware, l'ecosistema di strumenti, l'elaborazione su larga scala, la flessibilità nei tipi di dati e l'elaborazione di dati grezzi. Tuttavia, Hadoop ha anche alcuni svantaggi come la complessità nella configurazione e gestione di un cluster, la latenza, la memoria limitata, l'approccio orientato ai batch, la complessità della programmazione e la concorrenza da parte di framework alternativi come Apache Spark.

1.2.3 R e RStudio

R è un linguaggio di programmazione e un ambiente di sviluppo utilizzati principalmente per l'analisi statistica e la visualizzazione dei dati, progettato per lavorare con dataset di varie dimensioni e complessità. Offre un'ampia gamma di funzioni statistiche, algoritmi e pacchetti per l'analisi dei dati, il machine learning, l'analisi delle serie temporali e altro ancora. Inoltre, R fornisce potenti strumenti di visualizzazione che consentono di creare grafici e grafici per rappresentare i dati in modo efficace.

RStudio è un ambiente di sviluppo integrato (IDE) progettato specificamente per lavorare con il linguaggio di programmazione R. Come R[5] anche RStudio[4] è un progetto open source che ha lo scopo di offrire un'interfaccia utente intuitiva e ben organizzata che semplifica la scrittura del codice R, l'analisi dei dati e la creazione di visualizzazioni.

1.2.4 Python

Python[2] è un linguaggio di programmazione che ha catturato l'attenzione di sviluppatori di tutto il mondo grazie alla sua natura versatile e alla sua sintassi chiara e leggibile. Creato da Guido van Rossum[6] e presentato nel lontano 1991, Python ha continuato a guadagnare popolarità nel corso degli anni, diventando uno dei linguaggi più utilizzati e apprezzati nella comunità dello sviluppo software.

Ciò che distingue Python è il suo approccio alla scrittura del codice. La sua sintassi è strutturata in modo simile al linguaggio naturale, rendendo il codice scritto in Python quasi come una conversazione tra lo sviluppatore e la macchina. Questa semplicità e leggibilità non solo agevolano la creazione del codice, ma anche la comprensione dello stesso da parte di altri sviluppatori, favorendo una collaborazione più agevole.

Python (come anche R) è un linguaggio interpretato, il che significa che non richiede una fase di compilazione separata. Questo aspetto favorisce un approccio di sviluppo più rapido e interattivo, consentendo agli sviluppatori di scrivere e testare il codice in modo immediato, senza dover attendere processi di compilazione lungo.

1.2.5 Tool Javascript-html per MapReduce

Per evitare possibili complicazioni dovute alla configurazione di Hadoop, il Professore Daniele Pirrone ha creato uno script web chiamato "Tool MapReduce". Questo strumento consente di simulare un programma MapReduce attraverso un browser generico su Internet.

Lo strumento è sviluppato in Javascript con interfaccia HTML, rendendolo compatibile con i principali sistemi operativi e browser. Il codice è rilasciato sotto licenza GPL.

Lo strumento è semplicemente uno strumento didattico di supporto per gli studenti e quindi:

- Consente l'utilizzo di file di input che sono limitati in dimensione alla capacità del computer su cui viene eseguito.
- Accetta solo file di testo come input (la maggior parte dei quali hanno un'estensione .txt).
- Simula il paradigma di programmazione parallela MapReduce, ma in realtà funziona in modalità stand-alone, non distribuita.

1.3 Preparazione dell'ambiente

Affrontiamo ora il processo di installazione e configurazione dell'ambiente di sviluppo su macchina virtuale, che si è rivelato necessario per una corretta comunicazione tra Hadoop, RStudio e Python.

Le principali motivazioni di questa scelta risiedono nel fatto che i packages di RHadoop per una corretta comunicazione tra RStudio e Hadoop non sono aggiornati alle ultime versioni ¹.

La scelta di utilizzare la versione di Hadoop all'interno di una macchina virtuale anziché installarla direttamente sul sistema locale è stata influenzata da problematiche di gestione del DFS estendendo tutto il disco a HDFS andando quindi a creare conflitti.

Sfruttando il contesto isolato che offre la virtualizzazione si ha potuto sperimentare con un approccio trial and error diverse configurazioni che permettessero il corretto funzionamento di tutte le componenti. Sono state così fatte scelte specifiche per quanto riguarda le versioni del sistema operativo e dei software da utilizzare, tenendo conto della compatibilità tra di essi e delle raccomandazioni della comunità.

Si è scelto Virtualbox come sistema di virtualizzazione in quanto software open source di ampia distribuzione.

1.3.1 Preparazione dell'ambiente virtuale

Il primo passo è stato creare una nuova macchina virtuale su VirtualBox dal sistema host assegnando risorse hardware adeguate.

¹l'ultima release di RHadoop risale al 2015

Ruolo	Sistema Operativo	CPU	RAM	Spazio su disco
Host	macOS Ventura 13.4	2.3 GHz Intel Core i5 dual-core	16 GB 2133 MHz LPD-DR3	500 GB
VM	Ubuntu 16.04	1 CPU	4 GB	80 GB

Tabella 1: Specifiche di sistema

È stata scelta una distribuzione Linux come sistema operativo ospite per sfruttare la flessibilità e le prestazioni offerte da questa piattaforma open source. Essendo infatti Apache Hadoop un sistema utilizzato principalmente lato server, la sua implementazione più diffusa e stabile è in ambiente UNIX.

1.3.2 Installazione di Hadoop

Per garanzie di compatibilità si è scelto la versione 2.6.5 di Apache Hadoop, andandolo a configurare in modalità Single Node Cluster. Successivamente si è anche cercato di aggiungere nodi al cluster ma con risultati poco stabili.²

1.3.3 Integrazione di RStudio

L'integrazione di RStudio è stata relativamente più semplice utilizzando la versione 1.0.153 del 2017 e versione di R 3.2.3 che ha permesso una installazione delle librerie di RHadoop senza conflitti.

1.3.4 Configurazione di Python

La configurazione di Python presente già nel sistema operativo offre le versioni 2.7 e 3.5.2 con i relativi comandi python e python3. L'unica configurazione necessaria è stata l'installazione del gestore di pacchetti Python "pip" per l'installazione delle librerie utilizzate poi negli script.

1.3.5 Conclusione configurazione

L'installazione e la configurazione dell'ambiente di sviluppo su VirtualBox con Hadoop, RStudio e Python hanno richiesto un approccio meticoloso e la riso-

²In una configurazione iniziale si era utilizzata sia una soluzione "fisica" con due Raspberry Pi 3 Model B come nodi dati, sia macchine virtuali Ubuntu Server 22.04. Entrambe le soluzioni hanno presentato problematiche, quindi si è preferito tornare alla configurazione Single Node.

luzione di diverse sfide tecniche. Le scelte fatte durante il processo sono state guidate dalla necessità di garantire l'integrazione corretta dei componenti e l'ottimizzazione delle risorse disponibili.

2 Svolgimento

2.1 Analisi del problema

Abbiamo iniziato analizzando le richieste del problema da un punto di vista matematico. Per risolverle, bisogna essere in grado di trovare i documenti che contribuiscono all'incremento del fatturato e discriminare gli altri. In un secondo momento, è necessario calcolare la media (Eq.1) e la varianza (Eq.2) utilizzando le formule sottostanti.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2)$$

Successivamente, sarà necessario calcolare l'ammontare del fatturato diviso per ogni mese di ogni anno. Una volta elaborate queste informazioni, si può procedere ad identificare il mese di ogni anno che ha il valore più alto e più basso, analizzandoli tramite le funzioni di minimo e massimo e raggruppando per anno.

2.2 Analisi del Dataset

Procediamo con l'analisi del dataset con i comandi

```
1 summary(df)
2 table(df$tipo)
```

che vengono utilizzati per valutare la tipologia dei dati e per estrarre i valori della prima colonna riferiti al tipo di documento. In prima analisi, ci rendiamo

Tipo di documento	Quantità
BUONO.PRELIEVO	2146
DDT	15
FATTURA	556
INVENTARIO	1
NOTA.DI.CREDITO	9
OFFERTA	622
PREVENTIVO	240
RICEVUTA	11

Tabella 2: Analisi delle tipologie di documento

conto che i documenti che apportano incremento al fatturato sono della tipologia

fattura e ricevuta. Andiamo quindi ad impostare un filtro che ci aiuta ad estrarre dal dataset i documenti interessati.

2.3 Risoluzione tramite codice R

A questo punto, procediamo con il calcolo della media (Allegato contabilita.r) utilizzando la funzione `aggregate`, aggregando in base ai mesi di ogni anno e applicando la funzione di calcolo della media (`FUN = mean`). Vengono fatti in seguito delle trasformazioni per estrarre dalla data le informazioni di anno e mese e renderle più facilmente leggibili. Successivamente, andiamo a calcolare la varianza utilizzando la stessa funzione `aggregate` e impostando `FUN = var`.

vendite_medie	data_formattata
179.93081	2016-Gen
131.92288	2016-Feb
222.55800	2016-Mar
182.44100	2016-Apr
200.41270	2016-Mag
143.58000	2016-Giu
185.62951	2016-Lug
66.66857	2016-Ago
119.33484	2016-Set
174.44197	2016-Ott

Figura 3: Estratto del dataset con la media

data_formattata	varianza
2016-Gen	110972.691
2016-Feb	38351.044
2016-Mar	295872.781
2016-Apr	191570.685
2016-Mag	196337.279
2016-Giu	69653.870
2016-Lug	160827.843
2016-Ago	5731.737
2016-Set	35471.358
2016-Ott	152986.941

Figura 4: Estratto del dataset con la varianza

Per la soluzione dei job 3 e 4, si è reso necessario andare a fare un calcolo preliminare, ovvero aggregare le informazioni delle vendite in base al mese e all'anno. Successivamente, si è estratto il valore maggiore (per il job 3) e il valore minore (per il job 4) e visto che si è persa l'informazione specifica del mese, è stato necessario andare a recuperarla facendo un merge con la tabella aggregata precedente e le due tabelle risultanti dall'estrazione dei valori di minimo e massimo.

```
1 mesi_max = merge(mesi_max_solo_anno, somma_mesi, by="
vendite_totali")
```

2.4 MapReduce usando Javascript

Una volta affrontato il problema con lo scritto in R, andiamo ad analizzare il problema sotto il paradigma di MapReduce. Per una prima analisi, utilizziamo

lo strumento del professor Pirrone che permette un rapido approccio a questo modo di scrivere codice. I primi step intrapresi sono andare a dividere il file di input per riga e darlo in pasto allo script di mapping. La fase di mapping prevede la divisione della linea letta in input e sui tre valori che porta (tipologia ordine, data e costo). Successivamente, andiamo (come fatto in precedenza) a filtrare per la tipologia di ordine fattura e ricevuta: nel caso in cui la tipologia d'ordine appartenga a queste due categorie, viene creata una coppia che ha come chiave la data e come valore l'importo, o costo, del documento.

```
1      return V_In_Map.map(function(item){
2          //scrivo il codice di map
3          var tipoOrdine = item.split(",")[0]
4          var data = item.split(",")[1]
5          var costo = item.split(",")[2]
6
7          var chiave;
8          var valore;
9          if(tipoOrdine === "FATTURA" ||
10             tipoOrdine === "RICEVUTA"){
11              chiave = data.slice(0,6);
12              valore = costo;
13          }
14          else
15          {
16              chiave = "NULL";
17              valore = 0;
18          }
19          return keyVal(chiave, valore);
20      });
```

Listing 1: Script di mapping

Serve fare una precisazione: in un paradigma di MapReduce normale, è possibile scartare dei dati in input. In questo tool, è sempre necessario fornire un output di mapping. Di conseguenza, invece di scartare i documenti che non sono utili al fine del programma, vengono mappati con una chiave "null", che nelle elaborazioni successive non verrà considerata.

Per il primo job (Allegato job1.js), lo script di reduce andrà a recuperare tutti i valori relativi ad una determinata chiave, ovvero ad una data, e a calcolare la media (andandoli a sommare e successivamente a dividere per il numero di valori associati a quella data).

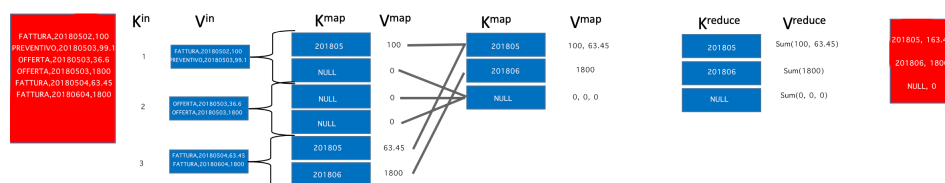


Figura 5: Schema di MapReduce per il primo job

Nel secondo job (Allegato job2.js), andiamo a riciclare lo stesso script di mapping. Per la fase di reduce, andiamo prima a calcolarci la media di valori associati a quella data (esattamente come il job precedente). Il secondo passaggio è quello di andare a calcolarci per ogni valore lo scarto quadratico rispetto alla media, utilizzando la formula (Eq.1, Eq.2). Come ultimo step, ogni valore di scarto quadratico viene sommato alla media.

Per il job 3 e il job 4, risultano necessari due passaggi del processo di MapReduce:

- Il primo passaggio (Allegato prejob34.js) prevede di raccogliere la somma degli importi di ogni mese dell'anno (viene utilizzato parte del codice del primo job senza andare a calcolare la media, ma solo facendo la somma)
- Il secondo passaggio prevede un mapping leggermente diverso, dove come chiave viene inserito l'anno e come valore viene inserito sia il mese che la media degli importi. Lo script di reduce andrà a cercare per ogni anno il valore massimo (per il job 3, Allegato job3.js) e il valore minimo (per il job 4, Allegato job4.js) dell'importo in maniera iterativa. Di conseguenza, questo ha permesso di salvare da parte le informazioni relative al mese e di riportarla a fine elaborazione.

2.5 MapReduce usando R e Hadoop

2.5.1 Scripting

Si è quindi passati ad implementare il MapReduce utilizzando script in R sulla piattaforma Hadoop. Per iniziare ad utilizzare la piattaforma, è fondamentale inizializzare le variabili d'ambiente e far partire i servizi per la gestione delle risorse (Yarn) e il file system condiviso (DFS).

Per l'esecuzione del singolo script, vengono fatte altre operazioni preliminari: per ogni script viene prima eliminata e poi ricreata sul file system condiviso la cartella di lavoro (andando così a pulire le tracce delle precedenti elaborazioni). In seguito, viene caricato all'interno del DFS il file che andremo ad utilizzare come input. Finalmente, possiamo andare a lanciare il nostro script utilizzando lo

streaming di Hadoop. Come ultimo passaggio, andiamo a visualizzare a schermo l'output prodotto, lo copiamo nell'hard disk locale e lo rinominiamo.

Queste operazioni sono state automatizzate da uno script in Bash (Allegato `all_jobs.sh`) e sono state inserite delle chiamate per monitorare le tempistiche e le performance (per salvarle poi in un file).

2.5.2 Implementazione in R

Gli script in R prendono spunto da quelli utilizzati prima nel linguaggio Javascript. Lo script di mapping (Allegato `mapper.r`), come prima, prende il file in input e per ogni riga, estrae il tipo di ordine, la data e il costo, mandando in output solo le tipologie di ordine presenti nel filtro. Come detto già in precedenza, lo script di mapping è univoco per tutti e quattro i job.

Nel primo script di reduce (Allegato `reducer_job1.r`), andiamo a sfruttare quelle che sono le peculiarità del paradigma di mappare Douce, andando a creare un nuovo Environment. Successivamente, andiamo a leggere l'output dello script precedente, estrarre per ogni riga la data e l'importo, e controllare se all'interno del nuovo ambiente quella data è già stata inserita:

- nel caso il valore della data (mese anno) non fosse presente, viene inserito come importo complessivo l'importo attuale e il conteggio degli elementi viene messo a uno. Il tutto poi viene inserito all'interno dell'ambiente.
- Nel caso il valore se è già presente, viene recuperato l'importo complessivo, al quale si somma l'importo attuale, aumentato il conteggio di uno e inserito il tutto nuovamente all'interno dell'ambiente.

Una volta ultimato questo processo, si passa ad estrarre dall'Environment tutti i valori inseriti e si procede con un processo di "abbellimento" prima di mandare i dati in output.

Il job 2 (Allegato `reducer_job2.r`), nel processo di reduce, ha una peculiarità: oltre ad implementare l'algoritmo per il calcolo della varianza (rispetto a quello di calcolo della media), per recuperare le informazioni relative alle medie, va a leggere all'interno del file system di Hadoop il file di output precedente.

I job 3 e 4 (Allegato `reducer_job3.r` e `reducer_job4.r`), invece, nella fase di reduce, salvano all'interno dell'ambiente la somma degli importi raggruppati per data (mese dell'anno). In un secondo momento, queste informazioni vengono recuperate, aggregate per anno ed estratti i valori di massimo e minimo.

In questa fase, ho preferito non dilungarmi nell'analisi dell'algoritmo (visto che gran parte della complessità si è affrontata con i linguaggi precedenti), ma ho preferito dare risalto a quegli aspetti tipici di questo paradigma di programmazione che Hadoop offre.

2.5.3 MapReduce con Python e Hadoop

Visto che negli ultimi anni Python si sta affermando come linguaggio di programmazione per la data Science, molti programmi e piattaforme si sono aperti alla compatibilità con esso. Su Hadoop, è infatti possibile sfruttare il paradigma MapReduce facendo scripting, oltre che con R, anche con Java e Python. Ho voluto così tradurre il lavoro fatto fino adesso in R anche in Python (Allegato mapper.py, reducer_job1.py, reducer_job2.py, reducer_job3.py, reducer_job4.py) e fare una piccola comparazione delle performance di questi due linguaggi. Anche qui, non entro nel merito dell'implementazione, in quanto l'algoritmo è esattamente lo stesso. Cambia solo il linguaggio in cui è scritto. È stato scritto uno script Bash che, a fini statistici, esegue i job di MapReduce circa cento volte e li salva su file.

```
1 #!/bin/bash
2 for ((i=1; i<=100; i++))
3 do
4     ./all_jobs.sh
5 done
```

Listing 2: Script per il salvataggio dei tempi

Successivamente, lo script è stato eseguito sia con uno script R che con uno script Python, salvando i tempi su due file distinti. Successivamente, è stata calcolata la media di esecuzione utilizzando lo strumento JavaScript di MapReduce, come segue:

Tabella 3: Confronto delle Tempistiche tra R e Python per Job

Job	Tempistica R (sec)	Tempistica Python (sec)
1	6.45	5.36
2	8.32	7.21
3	7.54	6.39
4	6.12	5.42
Totale	28.43	24.38

2.6 Rapporto finale

L'analisi delle vendite fornisce una panoramica dettagliata dell'andamento del fatturato dell'azienda con una granularità mensile. Tale analisi può aiutare l'azienda

a comprendere i modelli di vendita, prendere decisioni informate e pianificare per il futuro.³

2.6.1 Job 1 - Media di vendita per ogni mese

Dalla Figura 6 si può avere una visione globale delle vendite medie mensili. Tale dato consente di identificare i periodi in cui le vendite sono state particolarmente robuste o al di sotto delle aspettative.

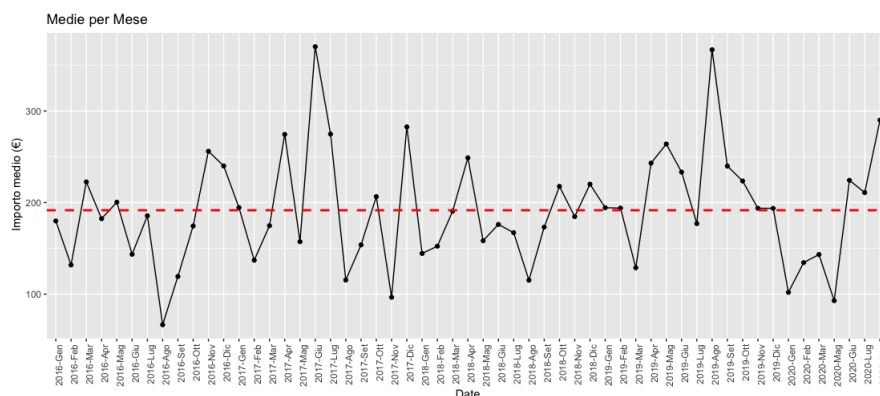


Figura 6: Grafico di elaborazione delle medie mensili

L'andamento di crescita o di stagnazione delle vendite può confermare e influenzare le strategie aziendali e gli sforzi di marketing in determinati mesi e su determinate campagne.

2.6.2 Job 2 - Varianza di vendita per ogni mese

La Figura 7 evidenzia l'ampiezza delle fluttuazioni mensili rispetto alla media. Un valore più alto indica una maggiore volatilità nelle vendite, mentre un valore più basso suggerisce una maggiore stabilità.

³È stato notato come in alcuni mesi non sia stato prodotto del fatturato (ad esempio, aprile 2020). Si sarebbe potuto decidere di andare a forzare il valore a 0, ma è stata scelta l'opzione di escludere tale dato poiché è possibile che coincida con un periodo di chiusura aziendale e un valore a 0 non porterebbe informazioni aggiuntive al fine di fare una analisi strategica con le strategie implementate.

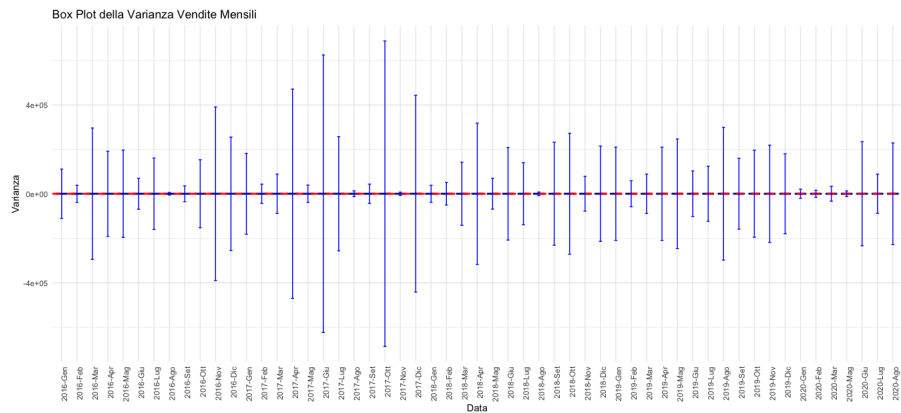


Figura 7: Grafico di elaborazione delle varianze mensili

Queste informazioni sono preziose per comprendere l'incertezza e l'irregolarità delle vendite durante alcuni periodi dell'anno.

2.6.3 Job 3 e 4 - Mesi migliori e peggiori di ogni anno

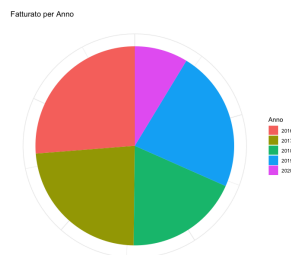


Figura 8: Suddivisione fatturato per anno

porta alla luce fattori di debolezza.

La Figura 8 fornisce una visione globale del fatturato dell'azienda in questi anni, contestualizzando i mesi migliori e peggiori di ogni anno (sempre in riferimento alla media globale delle vendite).

L'individuazione del mese con le maggiori vendite (Figura 9) è fondamentale per focalizzare l'attenzione sui fattori che hanno contribuito a questo successo, mentre l'analisi dei mesi con vendite più basse (Figura 10)

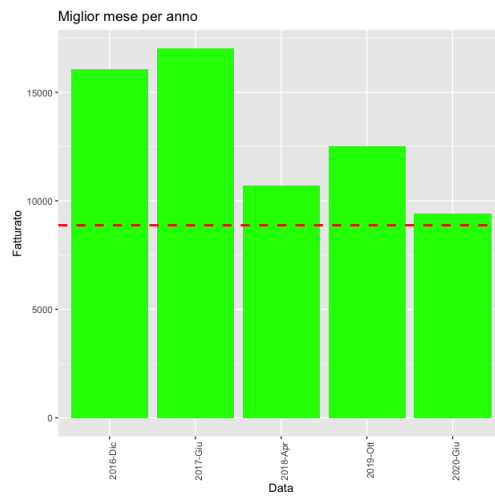


Figura 9: Migliori mesi

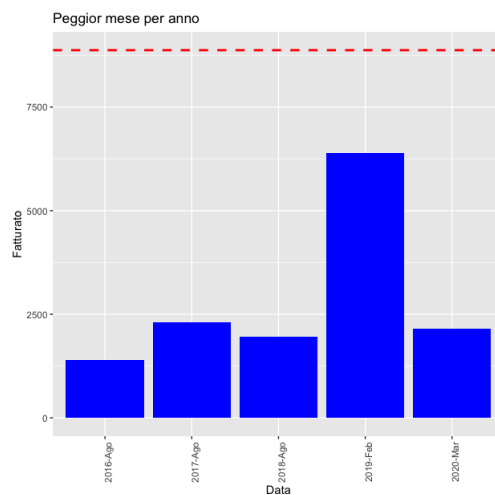


Figura 10: Peggiori mesi

Fattori come la stagionalità, le tendenze di mercato o le concomitanze avverse possono aver influenzato le vendite, così come potrebbero essere state implementate strategie specifiche o eventi particolari che hanno spinto le vendite al di sopra della media.

2.6.4 Analisi

Analizzare le condizioni e le tattiche associate ai mesi in questione può fornire lezioni preziose per guidare le future iniziative aziendali e aiutare a mitigare gli effetti negativi in situazioni simili in futuro.

3 Conclusioni

In questo paper abbiamo visto come l'analisi dei dati sia un processo fondamentale per l'ottimizzazione di qualsiasi attività. Attraverso l'analisi di un dataset relativo al fatturato di un'azienda abbiamo mostrato come sia possibile utilizzare gli strumenti più efficaci per la comprensione dei dati.

In particolare, abbiamo utilizzato R come strumento di analisi dei dati e JavaScript, Hadoop e R per l'implementazione di algoritmi di mapreduce. Grazie all'utilizzo di R, abbiamo creato grafici chiari ed efficaci per la comprensione dei dati, come grafici lineari, box plot e grafici a barre con grafico a torta.

Abbiamo anche illustrato il processo di data science, partendo dall'analisi del problema e del dataset, passando per l'elaborazione di algoritmi di mapreduce fino alla produzione di grafici per rispondere alle necessità del problema iniziale.

Inoltre, abbiamo visto come lo strumento di JavaScript per l'implementazione di algoritmi di mapreduce sia un ottimo strumento didattico per approcciarsi al paradigma di programmazione di mapreduce, soprattutto per chi è alle prime armi e non conosce R o Python. Tuttavia, l'utilizzo di Hadoop come software open source è risultato alquanto complicato da configurare correttamente, soprattutto se si vuole andare ad aggiungere nodi al cluster.

In conclusione, abbiamo apprezzato la possibilità di scalare la propria potenza di calcolo e parallelizzare le informazioni, abbiamo visto come la scelta del linguaggio e delle strategie di implementazione di un algoritmo possano incidere sulle performance e di conseguenza sui costi computazionali di un processo.

La gestione di grandi quantità di dati è una disciplina cruciale in un mondo in cui i dati sono sempre più presenti e utilizzati. Gli strumenti e le metodologie utilizzati in questo documento rappresentano i principi base per l'approccio a questa materia, evidenziando come sia necessario adattarsi alle diverse problematiche proposte.

Allegato

Risoluzione problema tramite R

contabilita.r

```
1 path_attuale <- dirname(rstudioapi::getActiveDocumentContext
2   ($path)
3   setwd(path_attuale)
4
5 library(lubridate)
6 library(dplyr)
7
8 #impostazione dati
9 nomeFile="Ordini.csv"
10 df=read.csv(nomeFile, sep= ",", header= FALSE)
11 colnames(df) <- c("tipo","data","costo")
12 summary(df)
13 table(df$tipo)
14 documenti_vendita <- subset(df, tipo=="FATTURA" | tipo=="
15   RICEVUTA")
16 mesi <- substring(documenti_vendita$data, 1, 6)
17
18 #JOB 1
19 media_mesi <- aggregate(documenti_vendita$costo, by = list(
20   mesi), FUN = mean)
21 colnames(media_mesi) <- c("data","vendite_medie")
22 nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
23   "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
24 media_mesi <- media_mesi %>%
25   mutate(mese = substring(media_mesi$data, 5, 6))
26 media_mesi$anno <- substr(media_mesi$data, 1, 4)
27
28 media_mesi <- media_mesi %>%
29   mutate(mese_testuale = nomi_mesi[as.numeric(mese)])
30
31 media_mesi <- media_mesi %>%
32   mutate(data_formattata = paste(anno, mese_testuale, sep = "
33   -"))
34
35 media_mesi <- media_mesi %>%
36   select(data ,vendite_medie, data_formattata)
37
38 #JOB 2
39 varianza_mesi = aggregate(documenti_vendita$costo, by = list(
40   mesi), FUN = var)
41 colnames(varianza_mesi) <- c("data","varianza")
42 varianza_mesi$anno <- substr(varianza_mesi$data, 1, 4)
```

```

38 varianza_mesi$mese <- substr(varianza_mesi$data, 5, 6)
39 varianza_mesi <- varianza_mesi %>%
40   mutate(mese_testuale = nomi_mesi[as.numeric(mese)])
41 varianza_mesi <- varianza_mesi %>%
42   mutate(data_formattata = paste(anno, mese_testuale, sep = "
   -"))
43 varianza_mesi <- varianza_mesi %>%
44   select(data, varianza, data_formattata)
45
46 #JOB3
47 somma_mesi = aggregate(documenti_vendita$costo, by = list(
   mesi), FUN = sum)
48 colnames(somma_mesi) <- c("data", "vendite_totali")
49 somma_mesi$anno <- substr(somma_mesi$data, 1, 4)
50 anni <- somma_mesi$anno
51 mesi_max_solo_anno = aggregate(somma_mesi$vendite_totali, by
   =list(anni), FUN=max)
52 colnames(mesi_max_solo_anno) <- c("anno_di_riferimento", "
   vendite_totali")
53 mesi_max = merge(mesi_max_solo_anno, somma_mesi, by="vendite_
   totali")
54 mesi_max$anno <- substr(mesi_max$data, 1, 4)
55 mesi_max$mese <- substr(mesi_max$data, 5, 6)
56 mesi_max <- mesi_max %>%
57   mutate(mese_testuale = nomi_mesi[as.numeric(mese)])
58 mesi_max <- mesi_max %>%
59   mutate(data_formattata = paste(anno, mese_testuale, sep = "
   -"))
60
61 #JOB 4
62 mesi_min_solo_anno = aggregate(somma_mesi$vendite_totali, by
   =list(anni), FUN=min)
63 colnames(mesi_min_solo_anno) <- c("anno_di_riferimento", "
   vendite_totali")
64 mesi_min = merge(mesi_min_solo_anno, somma_mesi, by="vendite_
   totali")
65 mesi_min$anno <- substr(mesi_min$data, 1, 4)
66 mesi_min$mese <- substr(mesi_min$data, 5, 6)
67 mesi_min <- mesi_min %>%
68   mutate(mese_testuale = nomi_mesi[as.numeric(mese)])
69 mesi_min <- mesi_min %>%
70   mutate(data_formattata = paste(anno, mese_testuale, sep = "
   -"))
71
72 #creazione grafici
73 library(ggplot2)
74
75 media_generale <- mean(media_mesi$vendite_medie)
76 media_mesi <- media_mesi[order(media_mesi$data),]

```

```

77 media_mesi$data_formattata <- factor(media_mesi$data_
    formattata, levels = unique(media_mesi$data_formattata))
78
79 #medie per mese
80 ggplot(media_mesi, aes(x = data_formattata, y = vendite_medie
    , group = 1,)) +
81   geom_line() + geom_point() + ggtitle("Medie per Mese") +
82   geom_hline(yintercept = media_generale, color = "red",
    linetype = "dashed", size = 1) +
83   theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
84   labs(x = "Date", y = " Importo medio")
85
86 #varianza per mese
87 dati_vendite <- merge(media_mesi, varianza_mesi, by = "data",
    all = TRUE)
88 dati_vendite <- dati_vendite[order(dati_vendite$data),]
89 dati_vendite$data_formattata.x <- factor(dati_vendite$data_
    formattata.x, levels = unique(dati_vendite$data_formattata
    .x))
90 #boxplot varianza
91 ggplot(dati_vendite, aes(x = data_formattata.x, y = vendite_
    medie)) +
92   geom_boxplot(fill = "lightblue", color = "darkblue") +
93   geom_hline(aes(yintercept = mean(vendite_medie)), color = "
    red", linetype = "dashed", linewidth = 1) +
94   geom_errorbar(aes(ymin = vendite_medie - varianza, ymax =
    vendite_medie + varianza), width = 0.2, color = "blue") +
95   labs(title = "Box Plot della Varianza Vendite Mensili",
    x = "Data",
96     y = "Varianza") +
97   theme_minimal() +
98   theme(axis.text.x = element_text(angle = 90, hjust = 1))
99
100
101 fatturato_per_anno <- aggregate(vendite_totali ~ anno, data =
    somma_mesi, sum)
102 # Creazione del grafico a torta
103 ggplot(fatturato_per_anno, aes(x = "", y = vendite_totali,
    fill = as.factor(anno))) +
104   geom_bar(stat = "identity", width = 1) +
105   coord_polar(theta = "y") +
106   labs(title = "Fatturato per Anno",
    fill = "Anno") +
107   theme_minimal() +
108   theme(axis.title = element_blank(),
    axis.text = element_blank(),
109     axis.ticks = element_blank(),
110     legend.position = "right")
111
112 media_generale <- mean(somma_mesi$vendite_totali)
113 #migliori mesi

```

```

115 ggplot(mesi_max, aes(x = data_formattata, y = vendite_totali)
116 ) +
117 geom_bar(stat = "identity", fill = "green", position =
118 position_dodge(width = 0.8)) +
119 labs(title = "Miglior mese per anno",
120 x = "Data",
121 y = "Fatturato") + theme(axis.text.x = element_text(
122 angle = 90, hjust = 1)) +
123 geom_hline(yintercept = media_generale, color = "red",
124 linetype = "dashed", size = 1)
125 #peggiori mesi
126 ggplot(mesi_min, aes(x = data_formattata, y = vendite_totali)
127 ) +
128 geom_bar(stat = "identity", fill = "blue", position =
129 position_dodge(width = 0.8)) +
130 labs(title = "Peggior mese per anno",
131 x = "Data",
132 y = "Fatturato") + theme(axis.text.x = element_text(
133 angle = 90, hjust = 1)) +
134 geom_hline(yintercept = media_generale, color = "red",
135 linetype = "dashed", size = 1)

```

Codice MapReduce

Javascript per Mapreduce Tool

Script inseriti nel Tool Js+html

job1.js

```

1
2 function jobInputSplit ( input_str ) {
3     return input_str.split('\n');
4 }
5 function jobMap ( V_In_Map ) {
6
7     return V_In_Map.map(function(item){
8         //scrivo il codice di map
9         var tipoOrdine = item.split(",")[0]
10        var data = item.split(",")[1]
11        var costo = item.split(",")[2]
12
13        var chiave;
14        var valore;
15        if(tipoOrdine === "FATTURA" || tipoOrdine === "
16        RICEVUTA"){
17            chiave = data.slice(0,6);
18            valore = costo;

```



```

18     }
19     else
20     {
21         chiave = "NULL";
22         valore = 0;
23     }
24     return keyVal(chiave, valore);
25 });
26 }
27
28 function jobReduce ( K_In_Reduce_V_In_Reduce ) {
29
30     return K_In_Reduce_V_In_Reduce.map(function (items){
31         var K_In_Reduce = items.split(S[0])[0];
32         var V_In_Reduce = items.split(S[0])[1].split(S[1]);
33         //scrivo il codice di split
34         var Reduce = V_In_Reduce.reduce(function (accumulator
35 , item,) {
36             if(K_In_Reduce !== "NULL"){
37                 return parseInt(accumulator) + parseInt(item)
38             };
39             return 0;
40         },0);
41         var k_out = K_In_Reduce;
42         var v_out = (Reduce/V_In_Reduce.length).toFixed(2) ;
43         return keyVal (k_out,v_out);
44     });
45 }

```

job2.js⁴

```

1 function jobReduce ( K_In_Reduce_V_In_Reduce ) {
2
3     return K_In_Reduce_V_In_Reduce.map(function (items){
4         var K_In_Reduce = items.split(S[0])[0];
5         var V_In_Reduce = items.split(S[0])[1].split(S[1]);
6         //scrivo il codice di reduce
7         var totale = V_In_Reduce.reduce(function (accumulator,
8 item,) {
9             return parseInt(accumulator) + parseInt(item);
10         },0);
11         var media = totale/V_In_Reduce.length;
12         console.log(media);
13         var scarti_quadratici = V_In_Reduce.map( item =>
14             Math.pow(parseInt(media) - parseInt(item),2)
15         );
16     });
17 }

```

⁴NB: per il job2 e il prejob34.js non è riportato lo script di mapping in quanto è uguale a quello del job1

```

14     );
15     var scarto_quadratico_medio = scarti_quadratici.reduce(
16         function (accumulator, item,) {
17             return parseInt(accumulator) + parseInt(item)
18         },0)/scarti_quadratici.length;
19     console.log(scarto_quadratico_medio);
20     var k_out = K_In_Reduce;
21     var v_out = scarto_quadratico_medio.toFixed(2) ;
22     return keyVal (k_out,v_out);
23 }

```

prejob34.js⁵

```

1 function jobReduce ( K_In_Reduce_V_In_Reduce ) {
2
3     return K_In_Reduce_V_In_Reduce.map(function (items){
4         var K_In_Reduce = items.split(S[0])[0];
5         var V_In_Reduce = items.split(S[0])[1].split(S[1]);
6         //scrivo il codice di reduce
7         var Reduce = V_In_Reduce.reduce(function (accumulator,
8             item,) {
9             return parseInt(accumulator) + parseInt(item);
10        },0);
11        var k_out = K_In_Reduce;
12        var v_out = Reduce.toFixed(2);
13        return keyVal (k_out,v_out);
14    });
15 }

```

job3.js

```

1 function jobMap ( V_In_Map ) {
2
3     return V_In_Map.map(function(item){
4         //scrivo il codice di map
5         var data = item.split("|")[0]
6         var anno = data.slice(0,4);
7         var mese = data.slice(4,6);
8         var media = item.split("|")[1]
9         var chiave= anno;
10        var valore = mese+"-"+media;
11        return keyVal(chiave, valore);
12    });
13 }
14

```

⁵Questo script è stato utilizzato come pre processazione per i job3 e job4

```

15 function jobReduce ( K_In_Reduce_V_In_Reduce ) {
16
17     return K_In_Reduce_V_In_Reduce.map(function (items){
18         var K_In_Reduce = items.split(S[0])[0];
19         var V_In_Reduce = items.split(S[0])[1].split(S[1]);
20         //scrivo il codice di reduce
21         var mese_massimo;
22         var Reduce = V_In_Reduce.reduce(function (massimo, item,)
23         {
24             var mese_item = item.split("-")[0];
25             var media_item = item.split("-")[1];
26             if((parseFloat(media_item)>parseFloat(massimo))||
27             massimo==0){
28                 mese_massimo = mese_item;
29                 return media_item;
30             } else {
31                 return massimo;
32             }
33         },0);
34         var k_out = K_In_Reduce;
35         var v_out = mese_massimo+"-"+ Reduce;
36         return keyVal (k_out,v_out);
37     });
38 }

```

job4.js

```

1 function jobMap ( V_In_Map ) {
2
3     return V_In_Map.map(function(item){
4         //scrivo il codice di map
5         var data = item.split("|")[0]
6         var anno = data.slice(0,4);
7         var mese = data.slice(4,6);
8         var media = item.split("|")[1]
9         var chiave= anno;
10        var valore = mese+"-"+media;
11        return keyVal(chiave, valore);
12    });
13 }
14
15 function jobReduce ( K_In_Reduce_V_In_Reduce ) {
16
17     return K_In_Reduce_V_In_Reduce.map(function (items){
18         var K_In_Reduce = items.split(S[0])[0];
19         var V_In_Reduce = items.split(S[0])[1].split(S[1]);
20         //scrivo il codice di reduce
21         var mese_minimo;
22         var Reduce = V_In_Reduce.reduce(function (minimo, item,)

```

```

23     {
24         var mese_item = item.split("-")[0];
25         var media_item = item.split("-")[1];
26         if((parseFloat(media_item)<parseFloat(minimo))||
27         minimo==0){
28             mese_minimo = mese_item;
29             return media_item;
30         } else {
31             return minimo;
32         }
33     },0);
34     var k_out = K_In_Reduce;
35     var v_out = mese_minimo+"-"+ Reduce;
36     return keyVal (k_out,v_out);
37 }

```

Script Bash

all_jobs.sh⁶

```

1  #!/bin/bash
2  export HADOOP_OPTS=-Djava.library.path=/usr/local/hadoop/lib/
   native
3  export HADOOP_HOME=/usr/local/hadoop
4  export HADOOP_CMD=/usr/local/hadoop/bin/hadoop
5  export HADOOP_STREAMING=/usr/local/hadoop/share/hadoop/tools/
   lib/hadoop-streaming-2.6.5.jar
6  export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
7  work_folder=$PWD
8  cd $HADOOP_HOME/sbin
9  start-dfs.sh
10 start-yarn.sh
11
12 cd work_folder
13 rm -f performance.txt
14 hdfs dfs -rm -r /homework
15 hdfs dfs -mkdir /homework
16
17 # Job1
18 hdfs dfs -rm -r /homework/job1
19 hdfs dfs -mkdir /homework/job1
20 hdfs dfs -copyFromLocal $work_folder/Ordini.csv /homework/
   job1
21 rm -f medie_calcolate

```

⁶il codice per eseguire i job in python è stato omissso in quanto è molto simile a quello in R, cambiano solo i nomi degli script in input

```

22 echo "Job1 - Somma iniziato"
23 start_time=$(date +%s)
24 hadoop jar $HADOOP_STREAMING -files $work_folder/mapper_job1.
    r,$work_folder/reducer_job1.r -mapper 'Rscript mapper_job1
    .r' -reducer 'Rscript reducer_job1.r' -input /homework/
    job1/Ordini.csv -output /homework/job1/output
25 end_time=$(date +%s)
26 hdfs dfs -cat /homework/job1/output/part-00000
27 hdfs dfs -copyToLocal /homework/job1/output/part-00000 ./
28 elapsed_time=$((end_time - start_time))
29 echo "Job1 : ${elapsed_time} seconds" >> performance.txt
30 mv ./part-00000 ./medie_calcolate
31 hdfs dfs -mv /homework/job1/output/part-00000 /homework/job1/
    output/medie_calcolate
32
33 # Job2
34 hdfs dfs -rm -r /homework/job2
35 hdfs dfs -mkdir /homework/job2
36 hdfs dfs -copyFromLocal $work_folder/Ordini.csv /homework/
    job2
37 rm -f varianze
38 echo "Job2 - Varianze iniziato"
39 start_time=$(date +%s)
40 hadoop jar $HADOOP_STREAMING -files $work_folder/mapper_job1.
    r,$work_folder/reducer_job2.r -mapper 'Rscript mapper_job1
    .r' -reducer 'Rscript reducer_job2.r' -input /homework/
    job2/Ordini.csv -output /homework/job2/output
41 end_time=$(date +%s)
42 hdfs dfs -cat /homework/job2/output/part-00000
43 hdfs dfs -copyToLocal /homework/job2/output/part-00000 ./
44 elapsed_time=$((end_time - start_time))
45 echo "Job2 : ${elapsed_time} seconds" >> performance.txt
46 mv ./part-00000 ./varianze
47 hdfs dfs -mv /homework/job2/output/part-00000 /homework/job2/
    output/varianze
48
49 # Job3
50 hdfs dfs -rm -r /homework/job3
51 hdfs dfs -mkdir /homework/job3
52 hdfs dfs -copyFromLocal $work_folder/Ordini.csv /homework/
    job3
53 rm -f max.txt
54 echo "Job 3 - Miglior mese per anno"
55 start_time=$(date +%s)
56 hadoop jar $HADOOP_STREAMING -files $work_folder/mapper_job1.
    r,$work_folder/reducer_job3.r -mapper 'Rscript mapper_job1
    .r' -reducer 'Rscript reducer_job3.r' -input /homework/
    job3/Ordini.csv -output /homework/job3/output
57 end_time=$(date +%s)

```

```

58 hdfs dfs -cat /homework/job3/output/part-00000
59 hdfs dfs -copyToLocal /homework/job3/output/part-00000 ./
60 elapsed_time=$((end_time - start_time))
61 echo "Job3: ${elapsed_time} seconds" >> performance.txt
62 mv ./part-00000 ./max.txt
63 hdfs dfs -mv /homework/job3/output/part-00000 /homework/job3/
    output/max
64
65 # Job4
66 hdfs dfs -rm -r /homework/job4
67 hdfs dfs -mkdir /homework/job4
68 hdfs dfs -copyFromLocal $work_folder/Ordini.csv /homework/
    job4
69 rm -f min.txt
70 echo "Job 4 - Peggior mese per anno"
71 start_time=$(date +%s)
72 hadoop jar $HADOOP_STREAMING -files $work_folder/mapper_job1.
    r,$work_folder/reducer_job4.r -mapper 'Rscript mapper_job1
    .r' -reducer 'Rscript reducer_job4.r' -input /homework/
    job4/Ordini.csv -output /homework/job4/output
73 end_time=$(date +%s)
74 hdfs dfs -cat /homework/job4/output/part-00000
75 hdfs dfs -copyToLocal /homework/job4/output/part-00000 ./
76 elapsed_time=$((end_time - start_time))
77 echo "Job4 : ${elapsed_time} seconds" >> performance.txt
78 mv ./part-00000 ./min.txt
79 hdfs dfs -mv /homework/job4/output/part-00000 /homework/job4/
    output/min
80
81 cd $HADOOP_HOME/sbin
82 stop-all.sh

```

Script R

mapper.r

```

1 filtro <- c("FATTURA", "RICEVUTA")
2
3 connessione <- file("stdin", open = "r")
4
5 while (length(line <- readLines(connessione, n = 1, warn =
    FALSE)) > 0) {
6     fields <- unlist(strsplit(line, ","))
7
8     tipoOrdine <- fields[1]
9     data <- substr(fields[2], 1, 6)
10    costo <- fields[3]
11

```

```

12   if (tipoOrdine %in% filtro){
13     cat(data, "\t", costo, "\n", sep = "")
14   }
15 }
16 close(connessione)

```

reducer_job1.r

```

1  nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
2                "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
3
4  env <- new.env(hash = TRUE)
5
6  connessione <- file("stdin", open = "r")
7
8  while (length(line <- readLines(connessione, n = 1, warn =
9    FALSE)) > 0) {
10    valore_linea <- unlist(strsplit(line, "\t"))
11    vendita <- list(mese = valore_linea[1], importo = as.
12      numeric(valore_linea[2]))
13    risultato <- list(importo = 0, count = 0)
14
15    if (exists(vendita$mese, envir = env, inherits = FALSE)) {
16      valore_corrente <- get(vendita$mese, envir = env)
17      risultato$count <- valore_corrente$count + 1
18      risultato$importo <- valore_corrente$importo + vendita$
19        importo
20    }
21    else {
22      risultato$count <- 1
23      risultato$importo <- vendita$importo
24    }
25
26    assign(vendita$mese, risultato, envir = env)
27  }
28  close(connessione)
29
30  for (data in ls(env, all = TRUE)) {
31    valore <- get(data, envir = env)
32    mese = substring(data, 5, 6)
33    anno = substring(data, 1, 4)
34    mese_testuale = nomi_mesi[as.numeric(mese)]
35    data_formattata = paste(anno, mese_testuale, sep = "-")
36    importo = round(valore$importo / valore$count, digits =
37      2)
38    cat(data, "\t", importo, "\t", data_formattata, "\n", sep =
39      "")
40  }

```

reducer_job2.r

```
1 library(data.table)
2 nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
3               "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
4
5 nome_file_precedente <- "hadoop fs -text /homework/job1/
6   output/medie_calcolate"
7
8 env <- new.env(hash = TRUE)
9 connessione <- file("stdin", open = "r")
10 while (length(line <- readLines(connessione, n = 1, warn =
11   FALSE)) > 0) {
12
13   valore_linea <- unlist(strsplit(line, "\t"))
14   vendita <- list(data = valore_linea[1], importo = as.
15     numeric(valore_linea[2]))
16   risultato <- list(totale = 0, count = 0)
17
18   if (exists(vendita$data, envir = env, inherits = FALSE)) {
19     valore_corrente <- get(vendita$data, envir = env)
20
21     risultato$count <- valore_corrente$count + 1
22     media_del_mese <- dati_precedenti[dati_precedenti$V1 ==
23       vendita$data]$V2
24     risultato$totale <- (valore_corrente$totale + (vendita$
25       importo - media_del_mese)^2)
26   } else {
27     risultato$count <- 1
28     media_del_mese <- dati_precedenti[dati_precedenti$V1 ==
29       vendita$data]$V2
30     risultato$totale <- ((vendita$importo - media_del_mese)
31       ^2)
32   }
33
34   assign(vendita$data, risultato, envir = env)
35 }
36 close(connessione)
37
38 for (data in ls(env, all = TRUE)) {
39   valore <- get(data, envir = env)
40   mese = substring(data, 5, 6)
41   anno = substring(data, 1, 4)
42   mese_testuale = nomi_mesi[as.numeric(mese)]
43   data_formattata = paste(anno, mese_testuale, sep = "-")
44   varianza = round(valore$totale / valore$count, digits = 2)
45   cat(data, "\t", varianza, "\t", data_formattata, "\n", sep =
46     "")
47 }
```


40

}

reducer_job3.r

```

1  library(stringr)
2  nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
3                "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
4
5  env <- new.env(hash = TRUE)
6
7  connessione <- file("stdin", open = "r")
8
9  while (length(line <- readLines(connessione, n = 1, warn =
10    FALSE)) > 0) {
11    valore_linea <- unlist(strsplit(line, "\t"))
12    vendita <- list(mese = valore_linea[1], importo = as.
13      numeric(valore_linea[2]))
14    risultato = 0
15
16    if (exists(vendita$mese, envir = env, inherits = FALSE)) {
17      valore_salvato <- get(vendita$mese, envir = env)
18      risultato <- valore_salvato + vendita$importo
19    }
20    else {
21      risultato = vendita$importo
22    }
23
24    assign(vendita$mese, risultato, envir = env)
25  }
26
27  close(connessione)
28
29  chiavi <- ls(env, all = TRUE)
30  valori <- sapply(chiavi, function(chiave) get(chiave, envir =
31    env))
32  df <- data.frame(Chiave = chiavi, Valore = valori)
33  df$anno <- substr(df$Chiave, 1, 4)
34  df$mese <- substr(df$Chiave, 5, 6)
35  df$mese_testuale <- nomi_mesi[as.numeric(df$mese)]
36  df$data_formattata <- paste(df$anno, df$mese_testuale, sep =
37    "-")
38
39  lista_anni = unique(df$anno)
40  for (year in lista_anni) {
41    values_in_year <- df[which(df$anno == year), ]
42    max <- values_in_year[which.max(values_in_year$Valore), ]
43    cat(max$anno, "\t", max$Valore, "\t", max$data_formattata, "
44      \n", sep = "")
45  }

```

reducer_job4.r

```
1 library(stringr)
2 nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
3               "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
4
5 env <- new.env(hash = TRUE)
6
7 connessione <- file("stdin", open = "r")
8
9 while (length(line <- readLines(connessione, n = 1, warn =
10   FALSE)) > 0) {
11   valore_linea <- unlist(strsplit(line, "\t"))
12   vendita <- list(mese = valore_linea[1], importo = as.
13     numeric(valore_linea[2]))
14   risultato = 0
15
16   if (exists(vendita$mese, envir = env, inherits = FALSE)) {
17     valore_salvato <- get(vendita$mese, envir = env)
18     risultato <- valore_salvato + vendita$importo
19   }
20   else {
21     risultato = vendita$importo
22   }
23
24   assign(vendita$mese, risultato, envir = env)
25 }
26
27 close(connessione)
28
29 chiavi <- ls(env, all = TRUE)
30 valori <- sapply(chiavi, function(chiave) get(chiave, envir =
31   env))
32 df <- data.frame(Chiave = chiavi, Valore = valori)
33 df$anno <- substr(df$Chiave, 1, 4)
34 df$mese <- substr(df$Chiave, 5, 6)
35 df$mese_testuale <- nomi_mesi[as.numeric(df$mese)]
36 df$data_formattata <- paste(df$anno, df$mese_testuale, sep =
37   "-")
38
39 lista_anni = unique(df$anno)
40 for (year in lista_anni) {
41   values_in_year <- df[which(df$anno == year), ]
42   min <- values_in_year[which.min(values_in_year$Valore), ]
43   cat(min$anno, "\t", min$Valore, "\t", min$data_formattata, "
44     \n", sep = "")
45 }
```

Script Python

mapper.py

```
1 filtro = ["FATTURA", "RICEVUTA"]
2 import sys
3
4 for line in sys.stdin:
5     fields = line.strip().split(",")
6
7     tipoOrdine = fields[0]
8     data = fields[1][:6]
9     costo = fields[2]
10
11     if tipoOrdine in filtro:
12         print(data + "\t" + costo)
```

reducer_job1.py

```
1 nomi_mesi = ["Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
2             "Lug", "Ago", "Set", "Ott", "Nov", "Dic"]
3 env = {}
4 import sys
5
6 for line in sys.stdin:
7     valore_linea = line.strip().split("\t")
8     vendita = {"mese": valore_linea[0], "importo": float(
9     valore_linea[1])}
10    risultato = {"importo": 0, "count": 0}
11
12    if vendita["mese"] in env:
13        valore_corrente = env[vendita["mese"]]
14        risultato["count"] = valore_corrente["count"] + 1
15        risultato["importo"] = valore_corrente["importo"] +
16        vendita["importo"]
17    else:
18        risultato["count"] = 1
19        risultato["importo"] = vendita["importo"]
20
21    env[vendita["mese"]] = risultato
22
23 for data, valore in env.items():
24     mese = data[4:6]
25     anno = data[0:4]
26     mese_testuale = nomi_mesi[int(mese) - 1]
27     data_formattata = f"{anno}-{mese_testuale}"
28     importo = round(valore["importo"] / valore["count"], 2)
29     print(f"{data}\t{importo}\t{data_formattata}")
```

reducer_job2.py

```
1 import subprocess
2 import pandas as pd
3
4 nomi_mesi = ["Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
5             "Lug", "Ago", "Set", "Ott", "Nov", "Dic"]
6
7 nome_file_precedente = subprocess.check_output(['hadoop', 'fs',
8         '-text', '/homework/job1/output/medie_calcolate']).
9         decode('utf-8')
10 dati_precedenti = pd.read_csv(nome_file_precedente, sep='\t')
11
12 env = {}
13
14 while True:
15     line = input()
16     if not line:
17         break
18
19     valore_linea = line.strip().split("\t")
20     vendita = {'data': valore_linea[0], 'importo': float(
21         valore_linea[1])}
22     risultato = {'totale': 0, 'count': 0}
23
24     if vendita['data'] in env:
25         valore_corrente = env[vendita['data']]
26         media_del_mese = dati_precedenti[dati_precedenti['
27 data'] == vendita['data']]['media'].values[0]
28         risultato['count'] = valore_corrente['count'] + 1
29         risultato['totale'] = valore_corrente['totale'] + (
30 vendita['importo'] - dati_precedenti[dati_precedenti['data']
31 == vendita['data']]['media'].values[0]) ** 2
32     else:
33         media_del_mese = dati_precedenti[dati_precedenti['
34 data'] == vendita['data']]['media'].values[0]
35         risultato['count'] = 1
36         risultato['totale'] = (vendita['importo'] -
37 dati_precedenti[dati_precedenti['data'] == vendita['data']
38 ][ 'media'].values[0]) ** 2
39
40     env[vendita['data']] = risultato
41
42 for data, valore in env.items():
43     mese = data[5:7]
44     anno = data[0:4]
45     mese_testuale = nomi_mesi[int(mese) - 1]
46     data_formattata = f"{anno}-{mese_testuale}"
47     varianza = round(valore['totale'] / valore['count'], 2)
```

```
39 print(f"{data}\t{varianza}\t{data_formattata}")
```

reducer_job3.py

```
1 import pandas as pd
2 import sys
3
4 def get_month_name(month):
5     nomi_mesi = ["Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
6                 "Lug", "Ago", "Set", "Ott", "Nov", "Dic"]
7     return nomi_mesi[month - 1]
8
9 env = {}
10
11 for line in sys.stdin:
12     valore_linea = line.strip().split("\t")
13     vendita = {'mese': valore_linea[0], 'importo': float(
14               valore_linea[1])}
15     risultato = 0
16
17     if vendita['mese'] in env:
18         valore_salvato = env[vendita['mese']]
19         risultato = valore_salvato + vendita['importo']
20     else:
21         risultato = vendita['importo']
22
23     env[vendita['mese']] = risultato
24
25 chiavi = list(env.keys())
26 valori = list(env.values())
27
28 df = pd.DataFrame({'Chiave': chiavi, 'Valore': valori})
29 df['anno'] = df['Chiave'].str[:4]
30 df['mese'] = df['Chiave'].str[4:6]
31 df['mese_testuale'] = df['mese'].astype(int).apply(
32     get_month_name)
33 df['data_formattata'] = df['anno'] + '-' + df['mese_testuale']
34
35 lista_anni = df['anno'].unique()
36
37 for year in lista_anni:
38     values_in_year = df[df['anno'] == year]
39     max_value = values_in_year.loc[values_in_year['Valore'].
40                                   idxmax()]
41     print(f"{max_value['anno']}\t{max_value['Valore']}\t{
42           max_value['data_formattata']}")
```

reducer_job4.py

```

1 import pandas as pd
2 import sys
3
4 nomi_mesi = ["Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
5             "Lug", "Ago", "Set", "Ott", "Nov", "Dic"]
6 env = {}
7
8 for line in sys.stdin:
9     valore_linea = line.strip().split("\t")
10    vendita = {'mese': valore_linea[0], 'importo': float(
11              valore_linea[1])}
12    risultato = 0
13
14    if vendita['mese'] in env:
15        valore_salvato = env[vendita['mese']]
16        risultato = valore_salvato + vendita['importo']
17    else:
18        risultato = vendita['importo']
19
20    env[vendita['mese']] = risultato
21
22 chiavi = list(env.keys())
23 valori = list(env.values())
24
25 df = pd.DataFrame({'Chiave': chiavi, 'Valore': valori})
26 df['anno'] = df['Chiave'][:4]
27 df['mese'] = df['Chiave'][4:6]
28 df['mese_testuale'] = df['mese'].apply(lambda mese: nomi_mesi
29                                       [int(mese) - 1])
30 df['data_formattata'] = df['anno'] + '-' + df['mese_testuale']
31
32 lista_anni = df['anno'].unique()
33
34 for year in lista_anni:
35     values_in_year = df[df['anno'] == year]
36     min_value = values_in_year.loc[values_in_year['Valore'].
37                                   idxmin()]
38     print(f"{min_value['anno']}\t{min_value['Valore']}\t{
39           min_value['data_formattata']}")

```

Elenco delle figure

1	Schema del processo di MapReduce tratto dal paper di Google . . .	4
2	Schema dell'architettura HDFS tratto dal sito di Hadoop	5
3	Estratto del dataset con la media	11
4	Estratto del dataset con la varianza	11
5	Schema di MapReduce per il primo job	13
6	Grafico di elaborazione delle medie mensili	16
7	Grafico di elaborazione delle varianze mensili	17
8	Suddivisione fatturato per anno	17
9	Migliori mesi	18
10	Peggiori mesi	18

Riferimenti bibliografici

- [1] Jeffrey Dean and Sanjay Ghemawa. Mapreduce: Simplified data processing on large clusters. url: <https://static.googleusercontent.com/media/research.google.com/it/archive/mapreduce-osdi04.pdf>.
- [2] Python Software Foundation. Python. url: <https://www.python.org/>.
- [3] The Apache Software Foundation. Apache hadoop. url: <https://hadoop.apache.org/>.
- [4] PBC Posit. Rstudio ide. url: <https://www.rstudio.com/categories/rstudio-ide/>.
- [5] The CRAN team. R archive. url: <https://cran.r-project.org/>. E-mail: CRAN@R-project.org.
- [6] Wikipedia. Python. url: <https://it.wikipedia.org/wiki/Python>.