

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

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduzione</b>                           | <b>3</b>  |
| 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         |
| <b>2</b> | <b>Svolgimento</b>                            | <b>10</b> |
| 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 . . . . .                | 13        |
| 2.5.3    | MapReduce con Python e Hadoop . . . . .       | 14        |
| 2.6      | Analisi dei risultati . . . . .               | 14        |
| <b>3</b> | <b>Conclusioni</b>                            | <b>15</b> |
|          | <b>Allegato 1</b>                             | <b>16</b> |

# 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):

- Durante la fase di MapReduce, i dati di input vengono suddivisi in piccoli frammenti e passati a un insieme di processi paralleli chiamati "mapper". Ogni mapper esegue una funzione di mappatura 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 da parte dei mapper.
- 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".

- In questa fase, il sistema assegna le coppie chiave-valore raggruppate ai processi "reducer". Ogni processo reducer esegue una funzione di riduzione definita dall'utente che agisce sulle coppie chiave-valore correlate e produce i risultati finali dell'elaborazione.

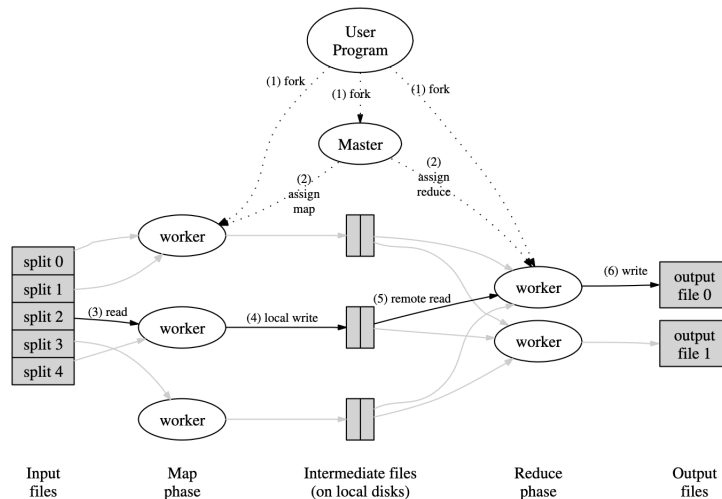


Figura 1: Schema del processo di MapReduce tratto dal paper di Google

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.

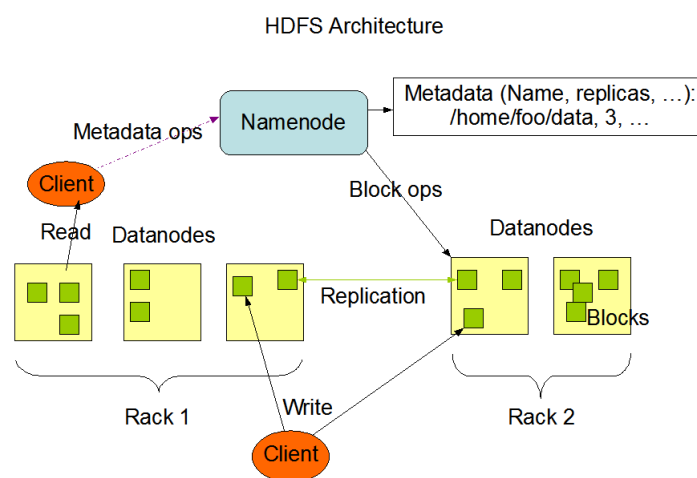


Figura 2: Schema dell'architettura HDFS tratto dal sito di Hadoop

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 <sup>1</sup>.

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.

---

<sup>1</sup>l'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.<sup>2</sup>

### 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-

---

<sup>2</sup>In 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 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, 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).

Nel secondo job, 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 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) e il valore minimo (per il job 4) 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 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, 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, 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, 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, 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 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.

## 2.6 Analisi dei risultati

### 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 1

## Fattori di standardizzazione in base alle tematiche

### Energia Primaria

```
1 # Esempio di codice Python
2 import numpy as np
3
4 def incmatrix(genl1,genl2):
5     m = len(genl1)
6     n = len(genl2)
7     M = None #to become the incidence matrix
8     VT = np.zeros((n*m,1), int) #dummy variable
9
10    #compute the bitwise xor matrix
11    M1 = bitxormatrix(genl1)
12    M2 = np.triu(bitxormatrix(genl2),1)
13
14    for i in range(m-1):
15        for j in range(i+1, m):
16            [r,c] = np.where(M2 == M1[i,j])
17            for k in range(len(r)):
18                VT[(i)*n + r[k]] = 1;
19                VT[(i)*n + c[k]] = 1;
20                VT[(j)*n + r[k]] = 1;
21                VT[(j)*n + c[k]] = 1;
22
23            if M is None:
24                M = np.copy(VT)
25            else:
26                M = np.concatenate((M, VT), 1)
27
28            VT = np.zeros((n*m,1), int)
29
30    return M
```

Listing 2: Python exampl

```
1
2 #JOB 1
3 media_mesi <- aggregate(documenti_vendita$costo, by = list(
4     mesi), FUN = mean)
5 colnames(media_mesi) <- c("data","vendite_medie")
6
7
```

```

8 nomi_mesi <- c("Gen", "Feb", "Mar", "Apr", "Mag", "Giu",
9             "Lug", "Ago", "Set", "Ott", "Nov", "Dic")
10
11 media_mesi <- media_mesi %>%
12   mutate(mese = substring(media_mesi$data, 5, 6))
13 media_mesi$anno <- substr(media_mesi$data, 1, 4)
14
15 media_mesi <- media_mesi %>%
16   mutate(mese_testuale = nomi_mesi[as.numeric(mese)])
17
18 media_mesi <- media_mesi %>%
19   mutate(data_formattata = paste(anno, mese_testuale, sep = "
20     -"))
21
22 media_mesi <- media_mesi %>%
23   select(data, vendite_medie, data_formattata)
24
25
26 #JOB 2
27 varianza_mesi = aggregate(documenti_vendita$costo, by = list(
28   mesi), FUN = var)
29 colnames(varianza_mesi) <- c("data", "varianza")
30 varianza_mesi$anno <- substr(varianza_mesi$data, 1, 4)
31 varianza_mesi$mese <- substr(varianza_mesi$data, 5, 6)

```

Questo indicatore considera la richiesta di energia primaria per l'intero ciclo di vita del prodotto considerato, tenendo conto, ad esempio, della trasformazione dei materiali combustibili in energia elettrica.

A questo indicatore contribuiscono quindi i materiali combustibili con il loro contenuto di energia primaria.

Il fattore di caratterizzazione è in questo caso il potere calorifico del materiale considerato.

## Effetto Serra

L'indicatore effetto serra viene calcolato considerando, tra le sostanze emesse in aria, quelle che contribuiscono al potenziale riscaldamento globale del pianeta Terra.

La quantità in massa di ciascuna sostanza, calcolata sull'intero ciclo di vita del prodotto, viene moltiplicata per un coefficiente di peso chiamato potenziale di riscaldamento globale (GWP, Global Warming Potential). Sommando poi i contributi delle varie sostanze, si ottiene il valore aggregato dell'indicatore.

Le sostanze che contribuiscono all'effetto serra sono principalmente: CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O, CFC, gli HCFC e gli HFC.

La CO<sub>2</sub> è la sostanza di riferimento per questo indicatore, vale a dire che il suo coefficiente di peso è uguale a 1 e i valori dell'indicatore sono espressi in kg di CO<sub>2</sub> equivalente (kg CO<sub>2</sub> eq).

| Composto                      | Formula                          | GWP100 [kg CO <sub>2</sub> /kg gas] |
|-------------------------------|----------------------------------|-------------------------------------|
| Diossido di carbonio          | CO <sub>2</sub>                  | 1                                   |
| Ossido di carbonio            | CO                               | 2                                   |
| Metano                        | CH <sub>4</sub>                  | 11                                  |
| Ossido di azoto               | N <sub>2</sub> O                 | 320                                 |
| CFC-11                        | CFCI <sub>3</sub>                | 4.000                               |
| CFC-12                        | CF <sub>2</sub> CI <sub>2</sub>  | 8.500                               |
| Clorotrifluorometano (CFC-13) | CF <sub>3</sub> CI               | 11.700                              |
| Tetrafluorometano (CFC-14)    | CF <sub>4</sub>                  | 9.300                               |
| HCFC-22                       | CHF <sub>2</sub> CI              | 1.700                               |
| HCFC-125                      | CHF <sub>2</sub> CF <sub>3</sub> | 3400                                |
| Halon-1301                    | CF <sub>3</sub> Br               | 5.600                               |
| Diclorometano                 | CH <sub>2</sub> CI <sub>2</sub>  | 25                                  |
| Cloroformio                   | CHCI <sub>3</sub>                | 15                                  |

Tabella 3: Fattori di standardizzazione per i principali responsabili dell'effetto serra, basati sul loro diretto contributo al riscaldamento globale con un tempo-orizzonte di 100 anni.

### Assottigliamento della fascia di ozono stratosferico

La riduzione della fascia di ozono stratosferico viene calcolata come l'indicatore precedente, ma utilizzando diverse sostanze (CFC, HCFC) e un diverso coefficiente di peso, chiamato potenziale di riduzione dell'ozono (ODP, Ozone Depletion Potential).

La sostanza di riferimento in questo caso è un clorofluorocarburo, precisamente il CFC-11.

### Eutrofizzazione

Questo indicatore valuta l'effetto dell'eutrofizzazione, ovvero l'aumento della concentrazione di sostanze nutritive negli ambienti acquatici. Le sostanze che contribuiscono al fenomeno dell'eutrofizzazione sono i composti a base di fosforo e azoto.

La sostanza di riferimento è il fosfato (PO<sub>4</sub>) e il coefficiente di peso prende il nome di potenziale di nutrizione (NP, Nutrifcation Potential).

| Formula           | NEP [kg NO <sub>3</sub> -/kg compost] |
|-------------------|---------------------------------------|
| NO <sub>3</sub> - | 1                                     |
| NO <sub>2</sub>   | 1.35                                  |
| NO <sub>x</sub>   | 1.35                                  |
| NO                | 2.07                                  |
| N <sub>2</sub> O  | 2.82                                  |
| NH <sub>3</sub>   | 3.64                                  |
| HCN               | 2.29                                  |
| N                 | 4.43                                  |
| PO <sub>4</sub> — | 10.45                                 |
| P                 | 32.03                                 |

Tabella 4: Fattori di standardizzazione per i principali responsabili dell'effetto serra, basati sul loro diretto contributo al riscaldamento globale con un tempo-orizzonte di 100 anni.

### Formazione di smog fotochimico (photo-smog)

Il termine "smog estivo" si riferisce a tutte le sostanze organiche volatili che portano alla formazione di ozono troposferico attraverso reazioni fotochimiche (in presenza di radiazione solare).

Il fattore di caratterizzazione utilizzato è chiamato "potenziale di formazione di ozono fotochimico" (POCP, Photochemical Ozone Creation Potential) e la sostanza di riferimento è l'etilene (C<sub>2</sub>H<sub>4</sub>).

| Composto | POCP [g C <sub>2</sub> H <sub>4</sub> /g di composto] |
|----------|---|
| metano   | 0,007   |
| etano    | 0,100   |
| propano  | 0,500   |
| aldeidi  | 0,3±0,2   |
| CO       | 0,040   |
| metanolo | 0,123   |
| etanolo  | 0,268   |

Tabella 5: Fattori di standardizzazione per i principali responsabili dello smog fotochimico.

## **Rifiuti Solidi**

L'indicatore in questione raggruppa tutti i rifiuti di tipo solido generati in qualsiasi attività nel ciclo di vita di un prodotto, ad esempio durante la generazione di energia elettrica necessaria per una lavorazione o durante la produzione delle lamiere di acciaio.

Non esistono fattori di caratterizzazione per questo indicatore, e ogni sostanza viene sommata alle altre tenendo semplicemente conto della quantità emessa in massa.

## 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 |

## 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>.