

Big Data e Big Problem

Normalmente si può trovare come definizione di big data la seguente:

I dati sono troppo grandi per entrare in memoria/disco rigido

Esiste però un altro problema Big:

Anche su dati piccoli, il mio modello potrebbe richiedere più risorse di quelle a mia disposizione

The Cancer Genome Atlas

<https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>
(<https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>)

portal.gdc.cancer.gov/repository?files_sort=%5B%7B"field"%3A"file_size"%2C"order"%3A"desc"%7D%5D

NIH NATIONAL CANCER INSTITUTE GDC Data Portal Home Projects Exploration Analysis **Repository** Quick Search Manage Sets Login Cart 0 GDC Apps

Search: e.g. 142682.bam, 4f6e2e7a-b...

Data Category

- ☐ simple nucleotide variation 237,055
- ☐ copy number variation 101,617
- ☐ transcriptome profiling 90,057
- ☐ sequencing reads 79,897
- ☐ biospecimen 55,598
- 5 More...

Data Type

- ☐ Annotated Somatic Mutation 100,858
- ☐ Raw Simple Somatic Mutation 84,132
- ☐ Aligned Reads 79,897
- ☐ Gene Expression Quantification 54,701
- ☐ Masked Annotated Somatic Mutation 44,755
- 17 More...

Experimental Strategy

- ☐ WXS 185,836
- ☐ Targeted Sequencing 118,616
- ☐ RNA-Seq 89,962
- ☐ Genotyping Array 68,455
- ☐ miRNA-Seq 45,210
- 5 More...

Files (596,758) Cases (84,392) Add All Files to Cart Manifest View 84,392 Cases in Exploration View Images

Primary Site Project Data Category Data Type Data Format

Showing 1 - 20 of 596,758 files 1.57 PB

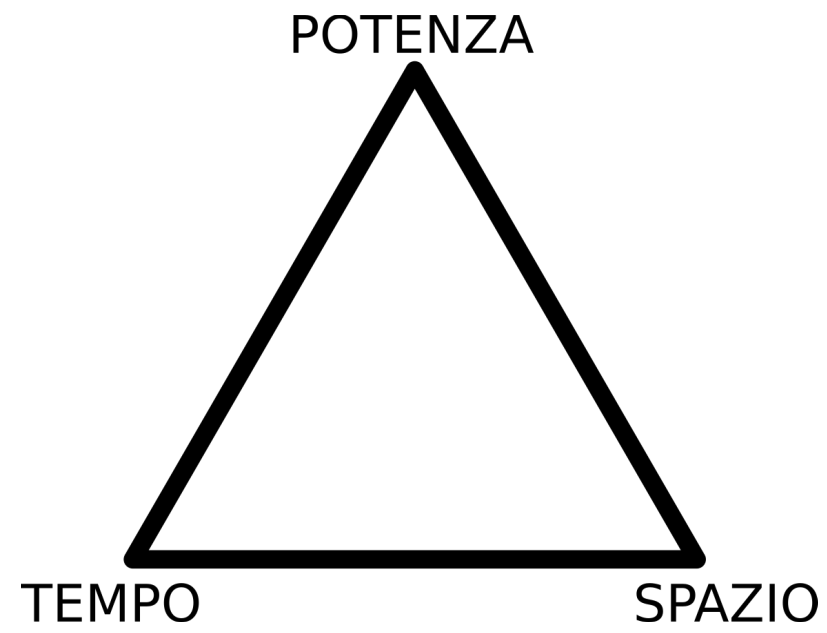
Access	File Name	Cases	Project	Data Category	Data Format	File Size	Annotations
controlled	259956ef-c652-42d4-b020-8d96a9af9dc7_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	750.32 GB	0
controlled	553c2bfe-7385-424f-a06a-0432172c24d4_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	736.26 GB	0
controlled	0b3b6347-7999-40c5-a74b-e5b3590894e1_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	683.73 GB	0
controlled	400315c6-e11e-40c5-b7d4-d5a90ef6541d_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	675.48 GB	0
controlled	780dec29-17f2-42e6-8193-95beedd3c786_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	674.62 GB	0
controlled	a5a18558-d8a1-4a3d-9f25-d44f4d19b801_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	670.74 GB	0
controlled	db1937ae-8a25-46dd-8e62-edda641c7ca4_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	666.01 GB	0
controlled	72cd6cc1-a15a-4aab-a7b8-cac497422a03_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	665.46 GB	0
controlled	1a1ea6a7-e4ff-4949-9db9-041c6f69df86_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	664.98 GB	0
controlled	58f9142d-6c91-41d1-910b-9db7365772a2_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	661.1 GB	0
controlled	216acf24-234a-44a4-8f7b-a95406550f1a_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	660.24 GB	0
controlled	23aa94b8-a430-48c8-b582-75529473336b_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	659.36 GB	0
controlled	c4835dd6-74c8-44b3-8b20-bf7de7d72043_wgs_gdc_realn.bam	1	WCDT-MCRPC	Sequencing Reads	BAM	655.77 GB	0

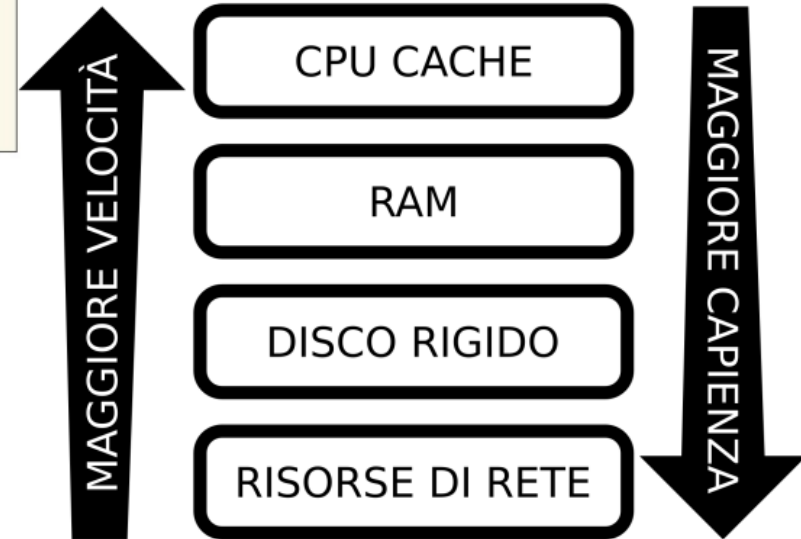
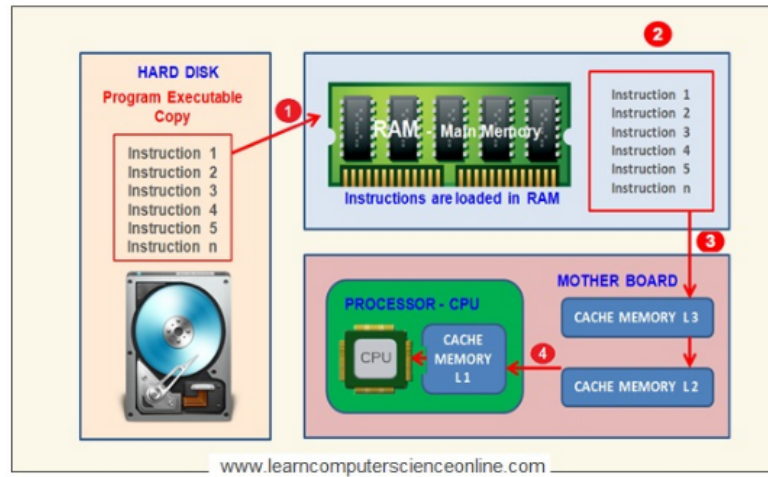
Nei big data dobbiamo fare delle scelte

Semplificando, ci sono 3 variabili nel nostro sistema:

- il **tempo** che siamo disponibili ad aspettare
- lo **spazio** su RAM che siamo disponibili ad allocare
- l'ammontare di informazioni che vogliamo ottenere (statistical **power**)

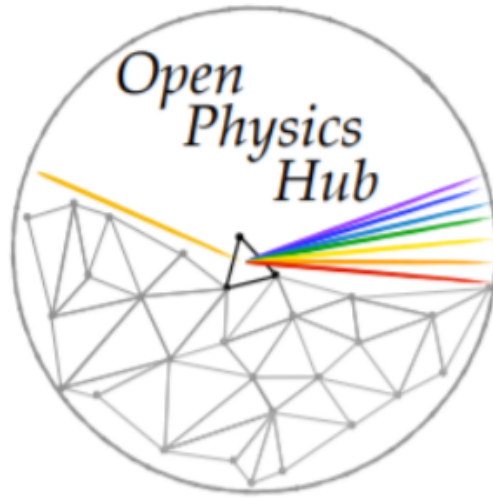
Tanto più ottimizziamo uno di questi parametri, tanto più ci rimettiamo negli altri due.





Questo schema non tiene conto anche di altri fattori:

- la leggibilità e trasferibilità del codice (potremmo scrivere in codice binario)
- il tempo che serve all'analista per produrre il codice



HIGH PERFORMANCE COMPUTING CLUSTER

“MATRIX” features

- 32 multi-core Intel XeonGold 5120 processors with dual thread
- 8 GB of memory per physical core (4 GB per virtual dual threaded core)
- 100 Gb/s InfiniBand connection
- 240 TB of disk space on SAS disks

“BLADE RUNNER” specifications:

Infiniband connected nodes:

- 2 nodes with 24 virtual cores and 64 GB RAM
- 2 nodes with 32 virtual cores and 64 GB RAM
- 1 node with 32 virtual cores and 128 GB RAM
- 2 nodes with 16 virtual cores and 24 GB RAM

Ethernet connected nodes:

- 2 nodes with 12 virtual cores and 16 GB RAM
- 2 nodes with 32 virtual cores and 64 GB RAM

<https://site.unibo.it/openphysicshub/en>

Leggere files (grandi)

Quando abbiamo a che fare con files molto grandi, per leggerli è utile usare la funzione `open` .

Questa funzione può essere usata per delimitare un blocco di codice all'interno del quale il file è aperto.

Al termine delle operazioni, questo verrà chiuso in modo automatico e sicuro, evitando che rimanga aperto e venga corrotto da altri processi.


```
In [1]: %%file prova.txt  
        tonno  
        mandibola  
        rum  
        pinocchio  
        sigmoide
```

Overwriting prova.txt

```
In [2]: with open('prova.txt') as file:
```

```
    for line in file:
```

```
        print(repr(line))
```

```
        print(len(line))
```

```
'tonno\n'
```

```
6
```

```
'mandibola\n'
```

```
10
```

```
'rum\n'
```

```
4
```

```
'pinocchio\n'
```

```
10
```

```
'sigmoide\n'
```

```
9
```

Sulla RAM non viene mai caricato tutto il file ma solo una riga alla volta.

Ci sono tre operazioni fondamentali per l'analisi dati:

- iterazione lazy
- map (ripeti un'operazione su tutti gli elementi)
- filter (seleziona solo una parte degli elementi)

Altri tipi di operazioni che discuteremo saranno:

- reduce (comporre insieme gli elementi)
- functional programming

Iterazione Lazy

Che cosa intendiamo con iterazione **lazy**?

Le operazioni non vengono compiute finchè il risultato non è richiesto!

In Python questa cosa è gestita da degli oggetti chiamati **iteratori**.

Sono gli oggetti su cui faccio i cicli **for**.

→ Un iteratore può essere percorso una volta sola!

Questo è controintuitivo: se provo a fare un ciclo for su di una lista, lo posso fare quante volte voglio

```
In [12]: lista = [1, 2, 3]

print("--- prima iterazione ---")
for elemento in lista:
    print(elemento)

print("--- seconda iterazione ---")
for elemento in lista:
    print(elemento)
```

```
--- prima iterazione ---
1
2
3
--- seconda iterazione ---
1
2
3
```

Ma se provo a farlo su di un file, lo posso leggere una volta sola!

Se volessi rileggerlo, dovrei aprirlo di nuovo!

```
In [13]: # creiamo il file NONTOCCARE  
! echo "un testo inutile" > ./NONTTOCCARE.TXT
```

```
In [14]: import os # pip install os
```

```
In [15]: os.getcwd() # in quale cartella mi trovo?
```

```
Out[15]: '/home/claudia/Documenti/Didattica/AA2020_2021/2021_PLS/PLSBigDataNetworks'
```

```
In [16]: os.listdir() # quali files e sotto-cartelle ci sono?
```

```
Out[16]:
```

```
['Lezione 1a - Introduzione a Python.ipynb',  
 'Lezione 1b - Iteratori e Big Data.ipynb',  
 'README.md']
```



```
In [17]: directory = "./"
filename = "NONTOCCARE.TXT"

position = directory + filename

with open(position) as file:

    print("--- prima iterazione ---")
    for line in file:
        print(repr(line))

    print("--- seconda iterazione ---")
    for line in file:
        print(repr(line))
```

```
--- prima iterazione ---
'un testo inutile\n'
--- seconda iterazione ---
```

Python ce lo nasconde, ma in realtà ogni volta che iteriamo sulla lista lui crea un nuovo iteratore che scorre la lista e poi scompare.

Possiamo farlo esplicitamente con il comando **iter**

```
In [18]: lista = [1, 2, 3]

iteratore_lista = iter(lista)

print("prima iterazione")
for elemento in iteratore_lista:
    print(elemento)

print("seconda iterazione")
for elemento in iteratore_lista:
    print(elemento)
```

```
prima iterazione
1
2
3
seconda iterazione
```

Non sempre è possibile evitare di caricare l'intero data set.

Supponiamo di voler calcolare tutte le combinazioni di elementi di una sequenza: non possiamo risolvere questo problema senza tenere in memoria l'intera sequenza!

Map

Un tipo di operazione molto frequente sulle sequenze è il cosiddetto **mapping**, ovvero applicare una funzione a tutti gli elementi di una lista, uno alla volta ed indipendentemente dagli altri.

Ad esempio, avendo una serie di numeri, potrei voler prendere il quadrato di ciascuno.

```
In [19]: numeri = [0, 1, 2, 3, 4, 5, 6]

quadrati = []
for numero in numeri:
    quadrato = numero **2
    quadrati.append(quadrato)

print(quadrati)
```

```
[0, 1, 4, 9, 16, 25, 36]
```

Questo può essere espresso in modo più conciso con una *comprehension*, che è funzionalmente identica al ciclo visto prima, ma più sintetica

```
In [20]: numeri = [0, 1, 2, 3, 4, 5, 6]
         quadrati = [x**2 for x in numeri]
         print(quadrati)
[0, 1, 4, 9, 16, 25, 36]
```

Il concetto di **map** è un'astrazione di questo procedimento.

Python fornisce una funzione, chiamata appunto **map**, che prende in input una funzione ed un oggetto iterabile e ritorna un iteratore i cui elementi sono il risultato dell'operazione.

```
In [21]: def quadrato(n):  
         return n**2  
  
         numeri = [0, 1, 2, 3, 4, 5, 6]  
         quadrati = map(quadrato, numeri)  
         print(quadrati)
```

```
<map object at 0x7fa9266c3780>
```

Ricordiamoci che il risultato delle operazioni sugli iteratori, quando possibile, è a sua volta un iteratore!

Siamo noi che dobbiamo esplicitamente **concretizzare** l'iterazione

```
In [22]: list(quadrati)
```

```
Out[22]: [0, 1, 4, 9, 16, 25, 36]
```

ricordiamoci che l'iterazione è compiuta una volta sola, quindi se vogliamo il risultato dobbiamo salvarcelo alla prima concretizzazione!

```
In [23]: list(quadrati)
```

```
Out[23]: []
```

Filter

Le operazioni **filter** selezionano un sottoinsieme dei dati e risultano nella generazione di un secondo iteratore.

```
In [24]: numeri = [-2, -1, 0, 1, 2]
positivi = []
for numero in numeri:
    if numero > 0:
        positivi.append(numero)

print(positivi)
```

```
[1, 2]
```


in modo simile all'operazione di **map**, anche l'operazione di **filter** ha un costrutto nel linguaggio tramite le *comprehension*

```
In [25]: numeri = [-2, -1, 0, 1, 2]
          positivi = [x for x in numeri if x > 0]
          print(positivi)
```

```
[1, 2]
```

ed esattamente come prima, abbiamo una funzione che prende una funzione di filtro (che ci dice se l'elemento è accettabile o no) e la applica ad un operatore

```
In [26]: def is_positive(n):  
          return n > 0  
  
          positivi = filter(is_positive, numeri)  
          print(list(positivi))
```

```
[1, 2]
```

Reduce

questa operazione combina gli elementi di un iteratore in un elemento unico

```
In [27]: numeri = [1, 2, 3, 4]
         totale = 0
         for numero in numeri:
             totale += numero
         print(totale)
```

10

Come per casi precedenti, esiste una funzione preesistente per effettuare le operazioni di riduzione: la funzione **reduce**

```
In [28]: from functools import reduce

def somma(a, b):
    return a + b

numeri = [1, 2, 3, 4]

totale = reduce(somma, numeri)

print(totale)
```

10

Questo tipo di operazioni è così comune che ci sono una serie di operazioni predefinite:

- **sum** per la somma
- **min** e **max** per il minimo e massimo

e così via

Una tipica riduzione, che useremo molto, è la stima delle frequenze.

```
In [29]: from collections import Counter

numeri = [ 1, 1, 1, 2, 2, 3, 3, 4, 4, 4, 4, 4 ]

Counter( numeri )
```

```
Out[29]: Counter({1: 3, 2: 2, 3: 2, 4: 5})
```

Una proprietà importante delle riduzioni è che i risultati si possono combinare: dati i conteggi su due serie, posso sommare insieme i due conteggi ed ottenere i conteggi totali fra le due serie

Map - Reduce

il famoso metodo MAP-REDUCE è una combinazione di queste idee:

- prendo una sequenza, la divido in sottosequenze
- invio le sequenze a diversi computer
- compio una riduzione su ciascuna sottosequenza
- raccolgo le sottosequenze e le combino insieme
- tutto questo fatto in modo ricorsivo

In []: