

PROGETTO: *Comparatore di Applicazioni*

ESAME: *Piattaforme per i Big Data*

Alessio Cimino

a.a. 2024

Il progetto sviluppato durante questo corso ha come obiettivo fornire agli utenti uno strumento che consenta loro di filtrare le applicazioni specificando parametri in base alla richiesta. Ogni giorno, infatti, vengono create e aggiunte nuove applicazioni con una facilità sempre crescente.

È stato utilizzato un archivio che contiene più di 2,3 milioni di applicazioni sviluppate per il sistema operativo Android e scaricabili attraverso la piattaforma Google Play.

Il file è aggiornato fino a giugno 2021 e rappresenta un insieme di dati di Big Data indispensabile per elaborare una grande quantità di informazioni, come in questo caso specifico.

Questo file contiene, per ogni riga, le seguenti 23 caratteristiche:

- App Name
- App Id
- Category
- Rating
- Rating Count
- Installs
- Minimum Installs
- Maximum Installs

- Free
- Price
- Currency
- Size
- Minimum Android
- Developer Id
- Developer Website
- Developer Email

- Released
- Privacy Policy
- Last Updated
- Content Rating
- Ad Supported
- In app purchases
- Editor Choice

L'interfaccia grafica vuole essere semplice ed intuitiva a causa dell'enorme quantità di dati; perciò ho diviso i metodi in due parti, come mostrato nella figura:

```
Menu:
[1] Info App
[2] Lista Top App
[0] Quit
Inserisci la tua scelta (0-2): 1
Sotto-menu:
[1] Info App
[2] ContentRating and Price
[3] Info Developer x App
[4] Info App x Developer
[0] Torna al Menu Principale
Inserisci la tua scelta (0-4): █
```

La prima parte dei metodi genera un output direttamente sul monitor, selezionando una delle opzioni nel sotto-menu visualizzato.

La seconda opzione, al contrario, genera, una volta eseguita, un file .csv con un numero di righe specificato dall'utente in un intervallo che va da 1 a 1000. L'utente può specificare direttamente le sotto-classi di interesse, applicando filtri prima di lanciare la query.

Tale decisione è volta a sviluppare in modo ottimale le linee di codice implementate, rendendo più efficiente l'esecuzione ed evitando di utilizzare inutilmente risorse computazionali e spazio di memoria fisica.

I metodi sono stati inizialmente testati su una macchina virtuale Kali Linux su Oracle installata su un laptop con le seguenti caratteristiche hardware:

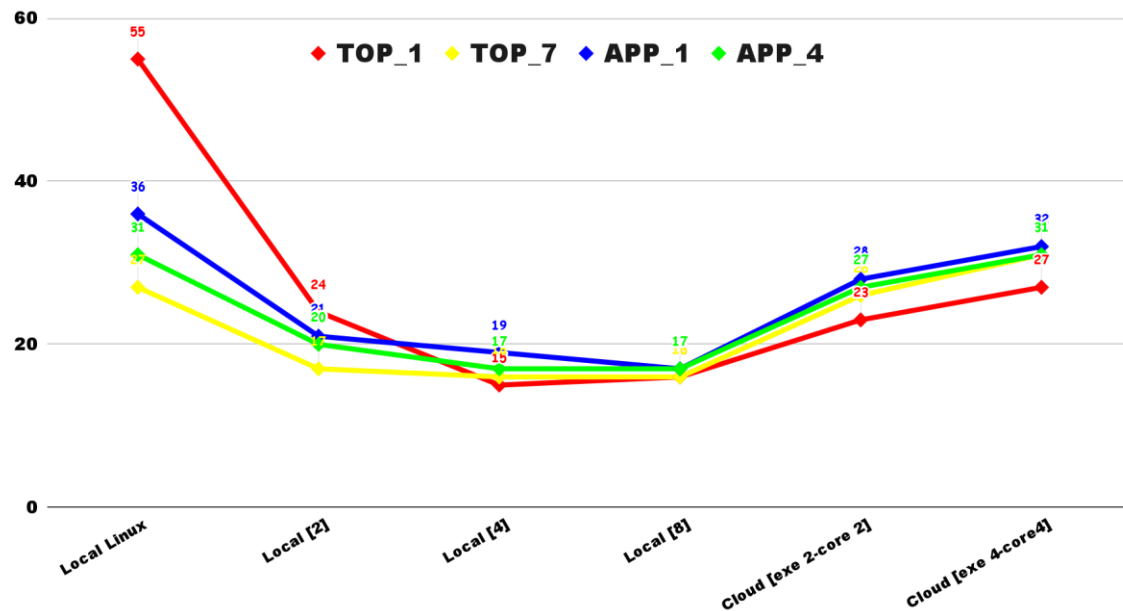
- Processore AMD 3020e (2C / 2T, 1.2 / 2.6GHz, 1MB L2 / 4MB L3)
- Memoria RAM da 8GB SO-DIMM DDR4-2400
- Archiviazione SSD da 256GB M.2 2242 PCIe NVMe 3.0x2
- Scheda grafica integrata AMD Radeon Graphics
- Sistema operativo: Windows 10

Successivamente, modificando una piccola parte del codice per eliminare la parte interattiva del menu, ho testato i vari metodi in modalità Cloud per confrontarne le prestazioni, ottenendo i grafici qui di seguito rappresentati.

Rendimento:

Tempi di esecuzione - Local Virtual Machine e Google Cloud

SpeedTest



Come si può osservare nel grafico presentato, i tempi di esecuzione nella modalità locale diminuiscono significativamente all'aumentare del numero di core utilizzati. Naturalmente, i tempi rappresentati dall'esecuzione in cloud sono più elevati rispetto a quelli locali a causa del tempo necessario per caricare i dati.

```
*tempistiche.txt: Bloc de notas
Archivo  Edición  Formato  Ver  Ayuda
TOP_1
top_1 local linux -> 55.56
top_1 local[2] -> 24
top_1 local[4] -> 15
top_1 local[8] -> 16
top_1 cloud --num-executors 2 --executive-cores 4 --> 23
top_1 cloud --num-executors 4 --executive-cores 4 --> 27

TOP_7
top_7 local linux -> 27
top_7 local[2] -> 17
top_7 local[4] -> 16
top_7 local[8] -> 16
top_7 cloud --num-executors 2 --executive-cores 4 --> 26
top_7 cloud --num-executors 4 --executive-cores 4 --> 31

APP_1
app_1 local linux -> 36.4
app_1 local[2] -> 21
app_1 local[4] -> 19
app_1 local[8] -> 17
app_1 cloud --num-executors 2 --executive-cores 4 --> 28
app_1 cloud --num-executors 4 --executive-cores 4 --> 32

APP_4
app_4 local linux -> 31.4
app_4 local[2] -> 20
app_4 local[4] -> 17
app_4 local[8] -> 17
app_4 cloud --num-executors 2 --executive-cores 4 --> 27
app_4 cloud --num-executors 4 --executive-cores 4 --> 31
```

L'infrastruttura utilizzata per sviluppare il progetto include diverse risorse informatiche, come PySpark con molte librerie essenziali per analizzare diversi tipi di dati (stringhe, interi, date, booleani, ecc.), Google Cloud e, naturalmente, Excel per verificare l'accuratezza dei dati estratti.

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, unix_timestamp, datediff, current_date
3 from pyspark.sql.types import TimestampType
4
5 import sys
6
7 def generate_top_lastUpdated_csv(listaGeneros, output_path, number_choice):
8     spark = SparkSession.builder.appName("TopLastUpdated").getOrCreate()
9
10    file_path = "Google-Playstore.csv"
11    df = spark.read.option("header", "true").csv(file_path)
12
13    filtered_df = df.filter(col("Category").isin(listaGeneros))
14    filtered_df = filtered_df.filter(col("Last Updated").isNotNull() & (col("Last Updated") != ""))
15    date_format = "MMM d, yyyy"
16    filtered_df = filtered_df.withColumn("Last Updated", unix_timestamp(col("Last Updated"), date_format).cast(TimestampType()))
17    filtered_df = filtered_df.withColumn("DateDifference", datediff(current_date(), col("Last Updated")))
18    sorted_df = filtered_df.orderBy("DateDifference")
19
20    top_df = sorted_df.limit(number_choice)
21    result_df = top_df.select(
22        "App Name",
23        "Released",
24        col("Last Updated").cast("date").alias("Last Updated"),
25        "DateDifference"
26    )
27
28    result_df.write.option("header", "true").csv(output_path, mode="overwrite")
29
30    spark.stop()
31    print(f"File CSV '{output_path}' generato con successo.")
32
33    print("Argomenti di input:", sys.argv)
34    listaGeneros = sys.argv[1].split(',')
35    number_choice = int(sys.argv[2])
36
37    if not listaGeneros:
38        print("Specifica almeno un genere.")
39        sys.exit(1)
40
41    output_path = "output_top/Top_lastUpdated.csv"
42    generate_top_lastUpdated_csv(listaGeneros, output_path, number_choice)
```

Il menù utente implementato si basa su una interfaccia semplice ed intuitiva, leggera ed essenziale, come mostrato nella figura qui accanto, nella quale è possibile specificare la categoria di interesse (o più di una) ed il numero di righe output desiderate per un massimo di 1000.

```
Sotto-menu:
[1] Top Download
[2] Top Rating
[3] Top Paid x Download
[4] Top Paid x Price
[5] Top Release Update
[6] Top Pegi 18
[7] Top best developers
[0] Torna al Menu Principale
Inserisci la tua scelta (0-7): 1
[?] Choose the Categories:
[ ] SelectAll
[X] Action
[ ] Adventure
> [ ] Arcade
[ ] Art
[X] Audio
[X] Auto
[ ] Beauty
[ ] Board
[ ] Books
[ ] Business
[ ] Card
[ ] Casino

How many occurrences do you want printed? (max 1000)
→100
```

```

Menu:
[1] Info App
[2] Lista Top App
[0] Quit
Inserisci la tua scelta (0-2): 1
Sotto-menu:
[1] Info App
[2] ContentRating and Price
[3] Info Developer x App
[4] Info App x Developer
[0] Torna al Menu Principale
Inserisci la tua scelta (0-4): 1
Inserisci il nome dell'app: Snapchat
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
23/12/12 08:04:02 WARN Utils: Your hostname, kali resolves to a loopback address: 127.0.1.1; using 10.0.2.15 instead (on interface eth0)
23/12/12 08:04:02 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
23/12/12 08:04:04 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
+-----+-----+-----+-----+-----+-----+-----+-----+
|App Name|App Id          |Category|Rating|Rating Count|Installs      |Size|Developer Id|Released   |Last Updated|
+-----+-----+-----+-----+-----+-----+-----+-----+
|Snapchat|com.snapchat.android|Social  |4.3   |26340056    |1,000,000,000+|72M |Snap Inc    |Oct 29, 2012|Jun 15, 2021|
+-----+-----+-----+-----+-----+-----+-----+-----+

```

I risultati ottenuti, al di là dei tempi di esecuzione che possono facilmente variare in base alla configurazione hardware utilizzata, mostrano un risultato in linea con le aspettative del progetto presentato per il quale sarà possibile, in un secondo momento, implementare metodi aggizionali che si rendano necessari per ottimizzare un sistema di consultazione di questo tipo, aggiornando ovviamente il DataSet utilizzato.