

PROYECTO : *Recomendador de App*

ASIGNATURA: *Cloud y Big Data*



Alessio Cimino, Luca Neve

El proyecto que desarrollamos durante este curso tiene como objetivo proporcionar a los usuarios una herramienta que les permita filtrar las aplicaciones especificando parámetros según la solicitud. Cada día, de hecho, se crean y añaden nuevas aplicaciones con una siempre más creciente facilidad.

Hemos utilizado un archivo que contiene más de 2,3 millones de aplicaciones, desarrolladas para el sistema operativo Android y descargables a través de la plataforma Google Play.

El archivo está actualizado hasta junio de 2021 y representa un conjunto de datos de Big Data indispensable para procesar una gran cantidad de información, como en este caso específico.

Este archivo contiene, para cada línea, las siguientes 23 características:

- | | | |
|----------------|----------------|------------------|
| • App Name | • Free | • Developer |
| • App Id | • Price | • Email |
| • Category | • Currency | • Released |
| • Rating | • Size | • Privacy Policy |
| • Rating Count | • Minimum | • Last Updated |
| • Installs | • Android | • Content Rating |
| • Minimum | • Developer Id | • Ad Supported |
| • Installs | • Developer | • In app |
| • Maximum | • Website | • purchases |
| • Installs | | • Editor Choice |

Decidimos desarrollar una interfaz gráfica simple e intuitiva, a causa de la enorme cantidad de datos. Por esta razón, dividimos los métodos en dos partes, como se muestra en la 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 primera parte de los métodos genera una salida directamente en el monitor, seleccionando una de las opciones del submenú que se muestran.

La segunda opción, en cambio, genera, una vez ejecutada, un archivo .csv con un número de líneas que debe ser especificado por el usuario

en un rango que va de 1 a 1000. El usuario puede especificar directamente las subclases de interés, aplicando filtros antes de lanzar la consulta. Hemos tomado esta decisión tanto para desarrollar de manera óptima las líneas de código implementadas, como para hacer más eficiente la ejecución, evitando ocupar innecesariamente recursos computacionales y espacio de memoria física.

Los métodos fueron probados inicialmente en una máquina virtual Kali Linux en Oracle instalada en un computador portátil con las siguientes características hardware:

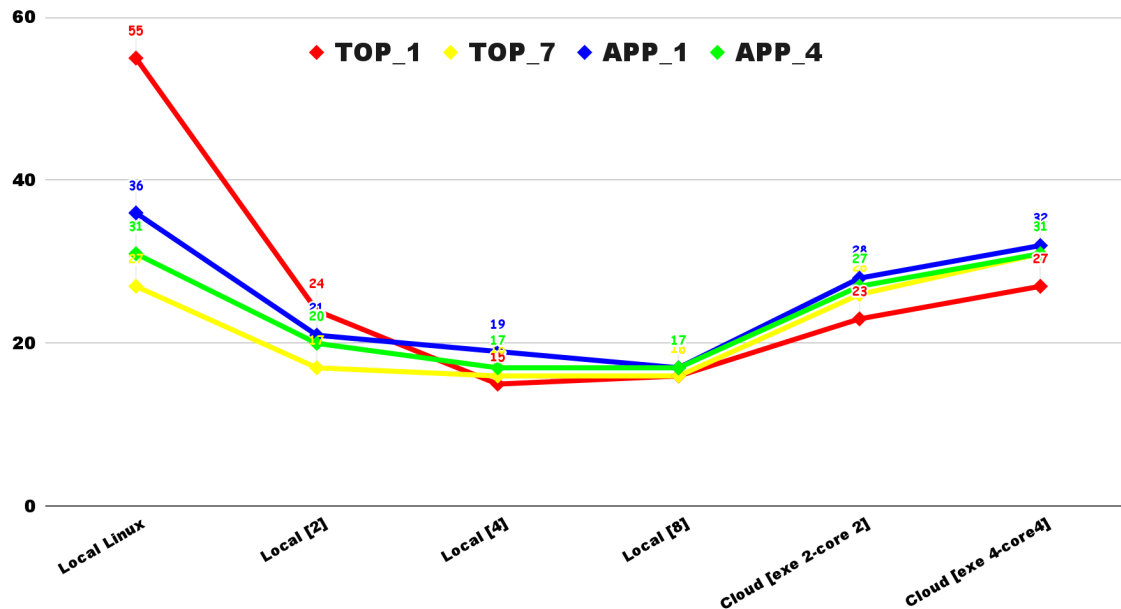
- Procesador AMD 3020e (2C / 2T, 1.2 / 2.6GHz, 1MB L2 / 4MB L3)
- Memoria RAM de 8GB SO-DIMM DDR4-2400
- Almacenamiento SSD de 256GB M.2 2242 PCIe NVMe 3.0x2
- Tarjeta gráfica integrada AMD Radeon Graphics
- Sistema operativo: Windows 10

Posteriormente, utilizamos las herramientas presentadas durante el curso para ejecutar los diversos métodos directamente en modo Cloud; para lograrlo, tuvimos que modificar una pequeña parte del código, eliminando la parte de menú interactivo.

Rendimiento:

tiempos de ejecución - Local Virtual Machine y Google Cloud

SpeedTest



Como se puede observar en el gráfico presentado, los tiempos de ejecución en la forma local disminuyen significativamente a medida que aumenta el número de núcleos utilizados. Naturalmente, los tiempos representados por la ejecución en la nube son más altos en comparación con el local debido al tiempo necesario para cargar los datos.

```
*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
```

La infraestructura que aprovechamos para desarrollar el proyecto incluye varias herramientas informáticas, como PySpark con muchas bibliotecas esenciales para analizar diferentes tipos de datos (cadenas de texto, enteros, fechas, boolean, etc.), Google Cloud y, por supuesto, Excel para verificar la precisión de los datos extraídos.

```

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}' generado con sucesso.")
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)

```

El menú de usuario que decidimos implementar se basa en una interfaz simple e intuitiva, ligera y esencial, como se muestra en la figura, donde es posible especificar la/las categoría/s de interés y el número de líneas de salida deseadas (max 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

```

```
(kali@kali) - [~/Desktop/CLO/PROGETTO]
$ python recomendador_de_app.py

RECOMENDADOR DE APP

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|
+-----+-----+-----+-----+-----+-----+-----+-----+

RECOMENDADOR DE APP

Eseguito: Info App
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):
```

Los resultados obtenidos, más allá de los tiempos de ejecución que pueden variar según la configuración de hardware, muestran un resultado acorde con nuestras expectativas iniciales de proyecto: las salidas procesadas, ya sea en formato .csv o directamente en monitor, constituyen una base sobre la cual será posible implementar, en un segundo momento, métodos adicionales que se vuelvan necesarios para optimizar un sistema de consulta de este tipo, actualizando, por supuesto, el DataSet utilizado.