

March 3, 2025

0.1 Películas más votadas

0.1.1 Sin nombre, solo tabla Ratings

```
[0]: from pyspark.sql import functions as func
from pyspark.sql.types import StructType, StructField, IntegerType, LongType

# Definir el esquema del DataFrame de ratings
schema = StructType([ \
    StructField("userID", IntegerType(), True), \
    StructField("movieID", IntegerType(), True), \
    StructField("rating", IntegerType(), True), \
    StructField("timestamp", LongType(), True)])

# Cargar el archivo de ratings u-data en un DataFrame, con separador el ↵
↵ tabulador \t
moviesDF = spark.read.option("sep", "\t").schema(schema).csv("dbfs:/FileStore/u.
↵ data")

# Funciones tipo SQL para agrupar y ordenar
topMovieIDs = moviesDF.groupBy("movieID").count().orderBy(func.desc("count"))

# Top 10
topMovieIDs.show(10)
```

```
+-----+-----+
|movieID|count|
+-----+-----+
|      50|  583|
|     258|  509|
|     100|  508|
|     181|  507|
|     294|  485|
|     286|  481|
|     288|  478|
|        1|  452|
|     300|  431|
|     121|  429|
+-----+-----+
```

only showing top 10 rows

0.1.2 Con nombre, haciendo inner join con películas (u.item)

```
[0]: from pyspark.sql.types import StructType, StructField, IntegerType, StringType
from pyspark.sql import functions as F

# Definir el esquema del DataFrame de ratings
esquemaRatings = StructType([
    StructField("UserID", IntegerType(), True),
    StructField("MovieID", IntegerType(), True),
    StructField("Rating", IntegerType(), True),
    StructField("Timestamp", IntegerType(), True)
])

# Cargar el archivo de ratings u-data en un DataFrame, con separador el
↳ tabulador \t
dfRatings = spark.read.csv("dbfs:/FileStore/u.data", sep="\t",
↳ schema=esquemaRatings, header=False)

# Definir esquema para el DataFrame de películas
esquemaPelículas = StructType([
    StructField("MovieID", IntegerType(), True),
    StructField("Title", StringType(), True),
    StructField("ReleaseDate", StringType(), True),
    StructField("EmptyColumn", StringType(), True),
    StructField("IMDB_URL", StringType(), True),
    StructField("Unknown", IntegerType(), True),
    StructField("Action", IntegerType(), True),
    StructField("Adventure", IntegerType(), True),
    StructField("Animation", IntegerType(), True),
    StructField("Children", IntegerType(), True),
    StructField("Comedy", IntegerType(), True),
    StructField("Crime", IntegerType(), True),
    StructField("Documentary", IntegerType(), True),
    StructField("Drama", IntegerType(), True),
    StructField("Fantasy", IntegerType(), True),
    StructField("FilmNoir", IntegerType(), True),
    StructField("Horror", IntegerType(), True),
    StructField("Musical", IntegerType(), True),
    StructField("Mystery", IntegerType(), True),
    StructField("Romance", IntegerType(), True),
    StructField("SciFi", IntegerType(), True),
    StructField("Thriller", IntegerType(), True),
    StructField("War", IntegerType(), True),
    StructField("Western", IntegerType(), True)
```

```

])

# Cargar el archivo de películas en un DataFrame, con separador /
dfPelículas = spark.read.csv("dbfs:/FileStore/u.item", sep="|",
    ↪ schema=esquemaPelículas, header=False)

# Mostrar las 10 películas con más votos
dfRatingsNombres = dfRatings.join(dfPelículas,on="MovieID",how="inner")

dfRatingsAgrupados = dfRatingsNombres.groupBy("Title").agg(F.count("Title").
    ↪ alias("Ratings")).orderBy("Ratings",ascending=False)
dfRatingsAgrupados.show(10)

```

```

+-----+-----+
|          Title|Ratings|
+-----+-----+
|   Star Wars (1977)|   583|
|   Contact (1997)|   509|
|   Fargo (1996)|   508|
|Return of the Jed...|   507|
|   Liar Liar (1997)|   485|
|English Patient, ...|   481|
|   Scream (1996)|   478|
|   Toy Story (1995)|   452|
|Air Force One (1997)|   431|
|Independence Day ...|   429|
+-----+-----+
only showing top 10 rows

```

0.1.3 Con un diccionario, UDF y broadcast

```

[0]: from pyspark.sql import functions as func
from pyspark.sql.types import StructType, StructField, IntegerType, LongType
import codecs

# Genera un diccionario con código y nombre de película
def loadMovieNames():
    movieNames = {}
    lines = sc.textFile("dbfs:/FileStore/u.item")
    collected_lines = lines.collect() # Convierte el RDD en una lista
    for line in collected_lines:
        fields = line.split('|')
        movieNames[int(fields[0])] = fields[1]
    return movieNames

```

```

# Carga el diccionario en una variable distribuida a todos los nodos con
↳broadcast
nameDict = spark.sparkContext.broadcast(loadMovieNames())

# Crea el esquema de rating
schema = StructType([ \
    StructField("userID", IntegerType(), True), \
    StructField("movieID", IntegerType(), True), \
    StructField("rating", IntegerType(), True), \
    StructField("timestamp", LongType(), True)])

# Carga df de ratings
moviesDF = spark.read.option("sep", "\t").schema(schema).csv("dbfs:/FileStore/u.
↳data")

movieCounts = moviesDF.groupBy("movieID").count()

# Crea una función definida por usuario (UDF) para buscar nombres de películas
↳en el diccionario distribuido a partir del código
def lookupName(movieID):
    return nameDict.value[movieID]

lookupNameUDF = func.udf(lookupName)

# Añade la columna nombre de película usando la función UDF
moviesWithNames = movieCounts.withColumn("movieTitle", lookupNameUDF(func.
↳col("movieID")))

# Ordena los resultados
sortedMoviesWithNames = moviesWithNames.orderBy(func.desc("count"))

# Muestra los 10 primeros
sortedMoviesWithNames.show(10, False)

```

```

+-----+-----+-----+
|movieID|count|movieTitle|
+-----+-----+-----+
|50      |583  |Star Wars (1977)|
|258     |509  |Contact (1997)|
|100     |508  |Fargo (1996)|
|181     |507  |Return of the Jedi (1983)|
|294     |485  |Liar Liar (1997)|
|286     |481  |English Patient, The (1996)|
|288     |478  |Scream (1996)|
|1       |452  |Toy Story (1995)|
|300     |431  |Air Force One (1997)|
|121     |429  |Independence Day (ID4) (1996)|

```

+-----+-----+-----+-----+
only showing top 10 rows

0.1.4 Obtener una lista de nombres de películas y un diccionario con el número de votos en cada puntuación:

```
[0]: # Nos quedamos con las columnas MovieID y Rating
dfRatings = dfRatings.select("MovieID", "Rating")

# Convertir el DataFrame de ratings a un RDD de filas
rddFilas = dfRatings.rdd

# Convertir el RDD de filas a un RDD de tuplas
rddTuplas = rddFilas.map(lambda fila: (fila[0], (fila[1],1)))

# Función para crear un diccionario con el número de votos para cada puntuación
def crearRatingDict(tuplas):
    RatingDict = {}
    for rating, cont in tuplas:
        if rating in RatingDict:
            RatingDict[rating] += cont
        else:
            RatingDict[rating] = cont
    return RatingDict

# Agrupar por MovieID y agregar las puntuaciones
rddRatingsAgrupados = rddTuplas.groupByKey().mapValues(crearRatingDict)

# Volver a convertir a dataframe y hacer join con películas para obtener el
↪ nombre

# Mostrar 10 películas con su nombre y puntuaciones
```

0.1.5 Superhéroe más relacionado

En cada línea aparece un código de superhéroe y los códigos de otros superhéroes que aparecen con él en algún comic. Puede aparecer repetido en más de una línea.

Ficheros: Marvel-names.txt, Marvel-graph.txt

Obtener cuál es el superhéroe que más relaciones tiene con otros superhéroes.

0.1.6 Distancia entre superhéroes

Grado de separación entre dos superhéroes, calculándolo a partir de las apariciones conjuntas en un comic.

Utilizamos algoritmo Breadth-first search: recorre un árbol o grafo nivel por nivel, comenzando

desde la raíz (o nodo inicial) y explorando todos los nodos vecinos en el nivel actual antes de moverse al siguiente nivel.

BFS es útil para encontrar la ruta más corta en grafos no ponderados y para explorar todos los nodos a una cierta “profundidad” del nodo inicial.

Pasos del algoritmo:

- Inicialización:
 - Coloca el nodo inicial en una cola (queue).
 - Marca el nodo inicial como visitado.
- Proceso de recorrido:
 - Mientras la cola no esté vacía:
 - * Saca (dequeue) el nodo al frente de la cola.
 - * Procesa el nodo (por ejemplo, imprime su valor).
 - * Para cada nodo vecino no visitado:
 - Marca el vecino como visitado.
 - Añade (enqueue) el vecino a la cola.

```
[0]: # Ejemplo en Python:
from collections import deque

def bfs(tree, start_node):
    visited = set()
    queue = deque([start_node])
    visited.add(start_node)

    while queue:
        node = queue.popleft()
        print(node) # Procesa el nodo

        for neighbor in tree[node]:
            if neighbor not in visited:
                visited.add(neighbor)
                queue.append(neighbor)

# Ejemplo de uso
tree = {
    'A': ['B', 'C'],
    'B': ['D', 'E'],
    'C': ['F', 'G'],
    'D': [],
    'E': [],
    'F': [],
    'G': []
}
```

```
bfs(tree, 'A')
```

A
B
C
D
E
F
G