

# **TRABAJO PRACTICO**

## **N°2 – SPARK**

# **ORGANIZACIÓN DE**

## **DATOS 75.06**

**Alumno: Gonzalez, Agustin Nicolas**

**Padrón: 106086**

**Corrector: Matias.**

Link al colab:

<https://colab.research.google.com/drive/1JSeQLfMFR-CtBhXqrr0kcy7ZWFDvZL-N?usp=sharing>

```
[ ] #1
#Considerando los logs de acciones realizadas sobre ítems, mostrar el top 10 de ids
#de ítems que fueron afectados por mayor cantidad usuarios distintos
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/logs.csv', header=True, inferSchema=True)
rdd = df.rdd

rdd.filter(lambda x: x.title != None).map(lambda x: (x.title, 1)).reduceByKey(lambda x,y: x+y).takeOrdered(10, lambda x: -x[1])

[('Especial:Userlogin', 75423),
 ('delete', 13761),
 ('move', 5972),
 ('Usuario:Ontzak', 4993),
 ('create', 4940),
 ('Usuario:Marcelo', 4383),
 ('Usuario:Taichi', 2807),
 ('Usuario:Geom', 2805),
 ('Usuario:Fixertool', 2788),
 ('Usuario:Strakhov', 2705)]
```

```
#34
#¿Quién es el usuario que más ha bloqueado a otros?
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/logs.csv', header=True, inferSchema=True)
rdd = df.rdd

rdd.filter(lambda x: x.action=="block").map(lambda x: (x.contributor_username, 1)).reduceByKey(lambda x,y: x+y).reduce(lambda x,y: x if x[1] > y[1] else y)

('Magister Mathematicae', 29226)
```

```
[ ] #14
#El Top 5 de contenidos que tienen la mayor cantidad de redirecciones que apuntan a
#ellos
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/redirect_list.csv', header=True, inferSchema=True)
rdd = df.rdd

rdd.map(lambda x: (x.rd_title, 1)).reduceByKey(lambda x,y: x+y).takeOrdered(5, lambda x: -x[1])

/usr/local/lib/python3.7/dist-packages/pyspark/sql/context.py:79: FutureWarning: Deprecated in 3.0.0. Use SparkSession
FutureWarning
[('Artículo_futuro', 1161),
 ('IV_milenio', 1001),
 ('V_milenio', 998),
 ('Estaciones_de_Metrobús_de_la_Ciudad_de_México', 396),
 ('163.117.0.0', 310)]
```

```
[ ] #32
#¿Cuál es el segundo contenido con más referencias geográficas asignadas?
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/geo_tags.csv', header=True, inferSchema=True)
rdd = df.rdd

#sqlContext = SQLContext(sc)
#df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/contents', header=True, inferSchema=True)
#rdd_contenidos = df.rdd

rdd.map(lambda x: (x.gt_page_id, 1)).reduceByKey(lambda x,y: x+y).takeOrdered(2, lambda x: -x[1])[1][0]
#rdd_contenidos.filter(lambda x: x.id == id_contenido).collect()
```

FutureWarning  
7421943

*Aclaracion sobre el 32: El resultado que se muestra es el segundo ID con mas referencias geográficas asignadas, lo que esta comentado era lo que había realizado en primera instancia, que era asignar ese numero como "id\_contenido" y luego encontrar en el rdd\_contenidos cual era el contenido con ese id, pero lo comente porque no encontraba nada*

```

#20
Cantidad de Stubs por categoría en la Wikipedia.
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/categorylinks.csv', header=True, inferSchema=True)
rdd = df.rdd

texts = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/contents_text_sample.csv')
texts.to_parquet('contents_text_sample.parquet')
del texts
df = sqlContext.read.parquet('contents_text_sample.parquet')
rdd_contenidos = df.rdd

stubs = rdd_contenidos.filter(lambda x: x.text != None and "url" not in x.text and len(x.text) < 1500).map(lambda x: (str(x.id), x.text))
rdd_filtrado = rdd.map(lambda x: (x.cl_from, x.cl_to)).join(stubs)
rdd_filtrado.map(lambda x: (x[1][0], 1)).reduceByKey(lambda x,y: x+y).collect()

```

```

[('Wikipedia:Artículo con identificadores Microsoft Academic', 423),
 ('Años 1210', 1),
 ('Armas de las artes marciales de Japón', 2),
 ('Física de la materia condensada', 1),
 ('Poetas', 1),
 ('Terminología de arqueología', 2),
 ('Emiratos Árabes Unidos', 2),
 ('Lenguaje y comunicación sonora', 1),
 ('Localidades de Río Grande del Sur', 17),
 ('Portal:Colombia', 1),
 ('Nacidos en 1951', 9),
 ('Wikipedia:Artículo con identificadores Open Library', 87),
 ('Provincia de Fermo', 1),
 ('Wikipedia:Categorías de artículos por plantillas', 48),
 ('Deportistas de Cantabria', 2),
 ('Distritos de Helsinki', 2),
 ('Científicos de Escocia', 4),
 ('Viola', 1),
 ('Wikipedia:Consultas de borrado con resultado por determinar', 25),
 ('Álbumes de rap', 5),
 ('Jardines botánicos por país', 2),
 ('Finales', 1),
 ('Ópera de Canadá', 1),
 ('Terroristas de Colombia', 1),
 ('Siglas', 46),
 ('Transporte de la Unión Europea por país', 1),
 ('Estados y territorios fundados en 1802', 1),
 ('Reyes de Arda', 1),
 ('Fallecidos en los años 1480', 2),
 ('Infraestructuras de los Países Bajos', 3),
 ('Intérpretes de música clásica de España', 2),
 ('Wikipedia:Artículo con identificadores Persa', 7),
 ('Pandaceae', 1),
 ('Localidades de Belice', 2),
 ('Novelas en inglés', 8),
 ('Premio Locus a la mejor novela de fantasía', 1),
 ('520', 1),
 ('Música clásica por país', 5),

```

*Aclaración sobre el 20: En primera instancia que el resultado continua, pero solo muestro esta parte, y luego es sobre el criterio que use para determinar cuando un artículo es un stub, tome como referencia lo que decía en la página de Wikipedia, menos de 1500 caracteres, y que no sea links a otras páginas, por eso también bloquee la palabra "url" que observe que se presentaba en muchos casos que eran simplemente para contener el link*

```

#8
#Considerando el pagelink_sample.csv, usando una representación de grafos obtener
#aquellos contenidos que tienen "relaciones no correspondidas". Entendemos como
#funciona una relación correspondida con un ejemplo: Si el contenido A tiene un link
#al B, pero B no tiene un link a A, podemos decir que B tiene una relación no
#correspondida con A
sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/pagelinks_sample.csv', header=True, inferSchema=True)
rdd = df.rdd

sqlContext = SQLContext(sc)
df = sqlContext.read.csv('/content/drive/MyDrive/Colab Notebooks/contents.csv', header=True, inferSchema=True)
rdd_contenidos = df.rdd

union = rdd.map(lambda x: (x.pl_title, x.pl_from)).join(rdd_contenidos.map(lambda x: (x.title, x.id))).map(lambda x: (str(x[1][0]), str(x[1][1]))).cache()
relaciones = union.map(lambda x: tuple(sorted(x))).map(lambda x: (x,1)).reduceByKey(lambda x,y: x+y).filter(lambda x:x[1]==1).map(lambda x: (x[0][0], x[0][1]))
relaciones_filtradas = relaciones.map(lambda x: (x[0], 1)).reduceByKey(lambda x,y: x+y) + relaciones.map(lambda x: (x[1], 1)).reduceByKey(lambda x,y: x+y)
relaciones_filtradas.map(lambda x: (x[0])).collect()

```

▶ [ '2280364',  
'5155470',  
'7819264',  
'8163946',  
'1307772',  
'1085272',  
'233143',  
'20654',  
'26877',  
'1534295',  
'206641',  
'194658',  
'316264',  
'3742767',  
'3138775',  
'3060109',  
'4264338',  
'4570784',  
'4844762',  
'2607991',  
'2365473',  
'2404060',  
'1749293',  
'4206602',  
'3170114',  
'2195925',  
'1919920',  
'1419338',  
'1934334',  
'4843955',  
'3105521',  
'1661496',  
'1515208',  
'2793612',  
'3189741',  
'3082926',  
'1623525',  
'2108232',  
'.....']

*Aquí también el resultado continua.*