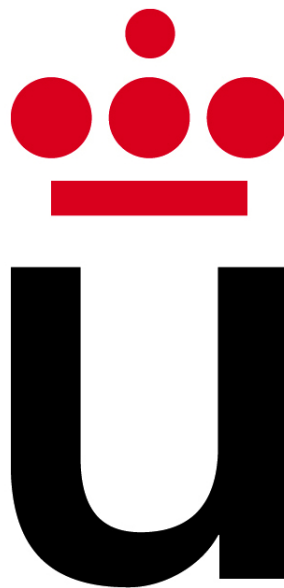


PRÁCTICA FINAL

Análisis de Big Data

ANÁLISIS DE UNA COMUNIDAD VIRTUAL
DE STACK EXCHANGE

José Ignacio Escribano



MÓSTOLES, 10 DE ABRIL DE 2016

Índice de tablas

Índice

1. Introducción	1
2. Carga de los datos	1
3. Ejercicio 1	1
4. Ejercicio 2	4
5. Conclusiones	9
6. Código Python	10

1. Introducción

En esta práctica final resolveremos una serie de cuestiones relativas a los usuarios que han participado en la comunidad <http://scifi.stackexchange.com>, dedicada a resolver preguntas sobre temas de ciencia ficción y fantasía.

El archivo `users.xml` contiene información sobre el perfil de un grupo del total de usuarios de la comunidad que han participado en ella, hasta marzo de 2015.

Utilizaremos Apache Spark para realizar el análisis de este archivo, utilizando un `DataFrame`.

2. Carga de los datos

Para cargar los datos, usamos el siguiente comando:

```
df = sqlContext.read.format('com.databricks.spark.xml').options(rowTag='row').load()
```

Esto creará un `DataFrame` a partir del fichero `users.xml`, que se leerá como archivo `xml`, y se guardará los hijos de un elemento `row` como una columna, quedando como una tabla.

3. Ejercicio 1

Inspeccionamos las primeras filas del archivo para ver cuáles son los elementos de los que consta cada uno de ellos.

```
<users>
  <row>
    <Id>-1</Id>
    <Reputation>1</Reputation>
    <CreationDate>2011-01-11T19:19:36.483</CreationDate>
    <DisplayName>Community</DisplayName>
    <LastAccessDate>2011-01-11T19:19:36.483</LastAccessDate>
    <Location>on the server farm</Location>
    <Views>0</Views>
    <UpVotes>2587</UpVotes>
    <DownVotes>3953</DownVotes>
    <AccountId>-1</AccountId>
  </row>
  <row>
    <Id>2</Id>
    <Reputation>101</Reputation>
    <CreationDate>2011-01-11T19:50:40.620</CreationDate>
```

```

    <DisplayName>Geoff Dalgas</DisplayName>
    <LastAccessDate>2015-02-05T00:03:28.030</LastAccessDate>
    <Location>Corvallis, OR</Location>
    <Views>21</Views>
    <UpVotes>1</UpVotes>
    <DownVotes>0</DownVotes>
    <Age>38</Age>
    <AccountId>2</AccountId>
</row>
...
</users>

```

Observamos que tenemos 11 elementos en cada fila:

- AccountId: Id de la cuenta
- Age: edad del usuario
- CreationDate: fecha de creación de la cuenta
- DisplayName: nombre de usuario
- DownVotes: votos a favor
- Id: Id de usuario
- LastAccessDate: fecha del último acceso
- Location: localización
- Reputation: puntos de reputación
- UpVotes: votos a favor
- Views: visitas a su página de perfil

Otra forma de comprobar el número de columnas de nuestro fichero es imprimir el esquema usando el comando

```
df.printSchema()
```

La salida de este comando es

```

root
 |-- AccountId: long (nullable = true)
 |-- Age: long (nullable = true)
 |-- CreationDate: string (nullable = true)
 |-- DisplayName: string (nullable = true)

```

```
|-- DownVotes: long (nullable = true)
|-- Id: long (nullable = true)
|-- LastAccessDate: string (nullable = true)
|-- Location: string (nullable = true)
|-- Reputation: long (nullable = true)
|-- UpVotes: long (nullable = true)
|-- Views: long (nullable = true)
```

Además de las columnas, este comando nos devuelve el tipo de cada variable.

Si queremos ver el contenido del fichero, usamos el comando, que muestra las 10 primeras filas del fichero.

```
df.show(10)
```

El comando devuelve la Tabla ??.

Ac-	Age	Dis-	Down- Id		Repu-	Up-	Views
coun-		play-	Vo-		tation	Vo-	
Id	CreationDate	Name	tes	LastAccessDate		tes	

	null	2011-		-1	2011-	on the	2587	0
	38	01-11T19:19:...		2	01-11T19:19:	server	1	3 21
-1	30	2011-	Community	3	2015-	farm	1	0 10
2	29	01-11T19:50:	Geoff	4	02-05T00:03:	Corvallis,	101	32 7
7598	null	2011-	Dalgas	5	2015-	OR	101	63 11
1998	30	01-11T19:55:	Nick	6	03-01T13:49:	Winston-	101	5 13
29738	30	2011-	Craver	7	2013-	Salem,	101	19 15
32917	38	01-11T20:17:	Em-	8	05-06T20:52:	NC San	99	211 98
33603	48	2011-	mett	9	2015-	Francisco,	146	78 85
51549	42	01-11T20:18:	Kevin	10	02-15T05:27:	CA New	103	48
3874		2011-	Mon-		2015-	York City,	1816	
15651		01-11T20:29:	trose		03-06T22:57:	Ne... New	1787	
		2011-	David		2014-	York, NY		
		01-11T20:37:	Fuller-		07-23T12:00:	Sweden		
		2011-	ton		2015-	United		
		01-11T20:37:	Soran-		02-18T23:17:	Kingdom		
		2011-	tis		2015-	Houston,		
		01-11T20:38:	GATh-		03-04T19:47:	TX		
		2011-	rawn		2015-	United		
		01-11T20:38:	Rodger		03-06T01:30:	States		
			Cooley					
			Matt-					
			hew-					
			Martin					

4. Ejercicio 2

A continuación, resolveremos una serie de consultas sobre el fichero XML.

1. Número total de registros de usuarios en el fichero.

La consulta que realizamos es la siguiente:

```
totalUsers = df.count()
print(totalUsers)
```

La consulta devuelve 20333. Es decir, tenemos que el fichero tiene 20333 usuarios.

2. Todos los datos disponibles sobre el usuario con `DisplayName` igual a `ambient_memory`.

La consulta que realizamos es la siguiente:

```
df.filter(df["DisplayName"] == "ambient_memory").show()
```

La consulta busca las filas que tienen `ambient_memory` como `DisplayName` y mostramos el resultado en forma tabular. La Tabla ?? muestra los resultados.

AccountId	Age	CreationDate	DisplayName	Down-Votes	Id	LastAccessDate	Location	Reputation	Up-Votes	Views
4529331	11	2014-06-20T12:01:00	ambient_memory	0	28494	2014-06-20T12:01:00	New York, United ...	1	0	0

3. Los 10 usuarios con mayor reputación (`Reputation`).

La consulta es la siguiente:

```
df.sort(df.Reputation.desc()).limit(10).show()

''' 0 de forma equivalente '''
df.sort("Reputation", ascending=False).limit(10).show()
```

La consulta ordena la columna `Reputation` de forma descendente, selecciona las 10 primeras filas y muestra el resultado en forma tabular. El resultado de la consulta se puede ver en la Figura 3.

AccountId	Age	CreationDate	DisplayName	Down-Votes	Id	LastAccessDate	Location	Reputation	Up-Votes	Views
-----------	-----	--------------	-------------	------------	----	----------------	----------	------------	----------	-------

41067	43	2011-		976	2015-		148259	5503	6224
893673	null	02-28T04:00:...		2765	03-08T02:21:...		118763	1934	3101
377643	null	2011-	DVK 1806	20774	2015-	New York, NY	92783	3878	7577
1057622	null	09-06T19:47:...	Thaddeu41	3500	03-07T21:21:...	Hayward, CA	67739	841	3562
12203	32	2013-	Richard2288	656	2015-	UK	62833	2230	1362
48648244		12-26T15:54:...		8719	03-08T02:36:...	Azkaban	57681	910	1515
375283	null	2011-	Slyt- 44	1027	2015-	Cincinnati, OH	49791	4875	2560
35978853		11-22T14:47:...	herin- 192	1693	03-08T00:08:...	null	48075	1919	1517
10264336		2011-	cess 1618	45	2015-	null	47505	3799	1671
1170648	null	02-01T22:00:...	Jeff 7	5184	03-08T02:28:...	Rich-	47257	2003	1771
		2012-	Darth 242		2015-	mond,			
		09-10T13:52:...	Satan 406		03-07T21:49:...	VA, USA			
		2011-	Keen		2015-	Salaberry-			
		03-07T02:09:...	Tango		03-08T01:44:...	de-Vall...			
		2011-	Da-		2015-	Orlando,			
		04-28T17:25:...	vRob60		03-06T02:46:...	FL			
		2011-	phan-		2015-				
		01-11T20:56:...	tom42		03-06T11:41:...				
		2012-			2015-				
		03-06T21:45:...			03-08T03:15:...				

4. Los usuarios con la fecha de creación más antigua y la más reciente, respectivamente.

La consulta es la siguiente:

```
''' Usuario más antiguo '''
df.sort("CreationDate", ascending=False).limit(1).show()

''' Usuario más reciente '''
df.sort("CreationDate", ascending=True).limit(1).show()
```

La consulta ordena descendientemente (ascentemente) por la columna `CreationDate`, se selecciona la primera fila y se muestra de forma tabular. La primera consulta devuelve el usuario más antiguo y la segunda, el usuario más reciente.

Las Tablas 4 y 5 muestra los usuarios con fecha de creación de su cuenta más antigua y más reciente, respectivamente.

Las Tablas 6 y 7 muestran los usuarios más jóvenes y viejos, respectivamente.

Ac- coun- tId	Age	CreationDate	Dis- play- Name	Down- Vo- tes	Id	LastAccessDate	Repu- tation	Up- Vo- tes	Views
4529331	null	2014-06-20T12:01	am- bient_me- mory	0	28494	2014-06-20T12:01	New York, United ...	1 0	0

Ac- coun- tId	Age	CreationDate	Dis- play- Name	Down- Vo- tes	Id	LastAccessDate	Repu- tation	Up- Vo- tes	Views
-1	null	2011-01-11T19:19	Com- munity	3953	-1	2011-01-11T19:19	on the server farm	1	2587 0

5. El usuario más joven y el más viejo (de los que han indicado una edad válida en el campo `Age`).

La consulta es la siguiente:

```
''' Importamos las funciones `max` y `min` '''
from pyspark.sql.functions import min, max

''' Calculamos la edad máxima y mínima de la columna `Age` '''
ages = df.select(min("Age").alias("min"), max("Age").alias("max")).collect()
```

La consulta devuelve

```
[Row(min=14, max=95)]
```

Es decir, la edad mínima es 14 y la máxima es 95.

Para conocer la información de los usuarios que tienen estas edades, hacemos la siguiente consulta:

```
''' Usuarios con la edad mínima guardada en la variable `ages` '''
df.filter(df.Age == ages[0].min).show()
```

```
''' Usuarios con la edad máxima guardada en la variable `ages` '''
df.filter(df.Age == ages[0].max).show()
```

Ac- coun- tId	Age	CreationDate	Dis- play- Name	Down- Vo- tes	Id	LastAccessDate	Repu- tation	Up- Vo- tes	Views
-1	null	2011-01-11T19:19:00	Com- munity	3953	-1	2011-01-11T19:19:00	on the server farm	2587	0

Ac- coun- tId	Age	CreationDate	Dis- play- Name	Down- Vo- tes	Id	LastAccessDate	Repu- tation	Up- Vo- tes	Views
---------------------	-----	--------------	-----------------------	---------------------	----	----------------	-----------------	-------------------	-------

	95	2011-		367	2011-		0
	95	01-19T21:58:...		1266	02-13T13:03:...		31
2913	95	2011-	bill weaver	1761	2015-	United States	2 1 0
9606	95	03-20T00:47:...	Mark Bessley	1925	03-04T22:20:...	California	78 1 0
	95	2011-	Kevin 0	2115	2015-	Cincinnati OH	0 3
237100	95	05-03T23:27:...	Sam Meldrum	2181	02-20T01:48:...	Chesham,	101 0 22
8802	95	2011-	jufrepeji 0	2205	2011-	United	101 1 1
299848	95	05-20T08:27:...	Keng 0	2341	05-20T17:58:...	K... Spain	101 0 10
564	95	2011-	HackToHell	2472	2012-	Parts	101 9 0
413744	95	06-04T09:59:...	Crazy Eddie	3120	09-12T09:15:...	Unknown	161 0 29
114869	95	2011-	Quinn 0	3192	2011-	Mordor	103 0 2
388430	95	06-10T12:49:...	Culver 0	3213	06-10T12:49:...	Gliese	431 8 38
47040	95	2011-	xiaohou- 18	5323	2015-	581G	101 3 0 8
21774	95	06-13T08:12:...	zi79 0	6365	02-22T17:09:...	South	944 22 0 0
128120	95	2011-	1.01pm 11	6392	2014-	Bend, IN	101 2
150461	95	07-02T23:26:...	JK. 0	8091	08-03T18:08:...	Australia	521 352
1399708	95	2011-	iamse- 0	8531	2015-	Atlantis,	101 9 4
1196073	95	07-27T03:05:...	rious 0	8554	01-20T16:08:...	FL	188 9 0
427266	95	2011-	JNat 0	1161	2014-	Teresina,	101
286976	95	10-21T03:45:...	MrEngi- 0		01-16T23:46:...	Brazil	101 89
989850	2011-	10-29T00:25:...	Merlin	2013-	07-10T04:29:...	London	
2187757	2011-	11-02T00:13:...	Pubby	2015-	03-04T22:47:...	Europe	
	2012-	03-14T16:07:...	Majti	2015-	03-05T16:45:...	United States	
	2012-	05-16T19:20:...	Watch	2015-	03-02T16:56:...	United States	
	2012-	05-17T23:24:...		2015-	02-12T14:56:...	Celje,Slo-	
	2012-	08-02T19:32:...		2013-	11-06T02:00:...	venia On	
	2012-	08-29T05:03:...		2014-	01-12T13:06:...	someones	
	2012-	08-30T10:02:...			02-04T12:19:...	wrist	
	2012-	12-29T01:52:...			02-13T08:37:...		

5. Conclusiones

6. Código Python