

## Módulo 12: Arquitecturas y procesos Big Data

### Capstone 12. Parte 1: Modelo de *sentiment* sobre Amazon Reviews

Enrique González, Jacinto Arias  
Máster en Ciencia de Datos e Ingeniería de Datos en la Nube  
Universidad de Castilla-La Mancha

Marta Bellón Castro  
Curso 2022-2023

## Índice

- [1. Introducción](#)
- [2. Análisis exploratorio](#)
- [3. Modelado](#)

```
In [1]: # Instalamos algunas librerías útiles para la práctica
```

```
import pyspark.sql.functions as sqlf
from pyspark.ml.pipeline import PipelineModel
from pyspark.ml.evaluation import BinaryClassificationEvaluator
```

```
VBox()
```

Starting Spark application

ID	YARN Application ID	Kind	State	Spark UI	Driver log	User	Current session?
2	application_1691082518949_0003	pyspark	idle	<a href="#">Link</a>	<a href="#">Link</a>	None	✓

```
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
SparkSession available as 'spark'.
```

```
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
```

```
In [2]: sc.install_pypi_package('pandas')
sc.install_pypi_package('seaborn')
sc.install_pypi_package('tabulate')
```

```
VBox()
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
Collecting pandas
  Downloading pandas-1.3.5-cp37-cp37m-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (11.3 MB)
Collecting python-dateutil>=2.7.3
  Downloading python_dateutil-2.8.2-py2.py3-none-any.whl (247 kB)
Requirement already satisfied: numpy>=1.17.3; platform_machine != "aarch64" and platform_machine != "arm64" and python_version < "3.10" in /usr/local/lib64/python3.7/site-packages
s (from pandas) (1.20.0)
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/site-packages (from pandas) (2023.3)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.7.3->pandas) (1.13.0)
Installing collected packages: python-dateutil, pandas
Successfully installed pandas-1.3.5 python-dateutil-2.8.2
```

```
Collecting seaborn
  Downloading seaborn-0.12.2-py3-none-any.whl (293 kB)
Requirement already satisfied: pandas>=0.25 in ./tmp/1691088765070-0/lib/python3.7/site-packages (from seaborn) (1.3.5)
Requirement already satisfied: numpy!=1.24.0,>=1.17 in /usr/local/lib64/python3.7/site-packages (from seaborn) (1.20.0)
Collecting typing_extensions; python_version < "3.8"
  Downloading typing_extensions-4.7.1-py3-none-any.whl (33 kB)
Collecting matplotlib!=3.6.1,>=3.1
  Downloading matplotlib-3.5.3-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.whl (11.2 MB)
Requirement already satisfied: python-dateutil>=2.7.3 in ./tmp/1691088765070-0/lib/python3.7/site-packages (from pandas>=0.25->seaborn) (2.8.2)
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/site-packages (from pandas>=0.25->seaborn) (2023.3)
Collecting packaging>=20.0
  Downloading packaging-23.1-py3-none-any.whl (48 kB)
Collecting pyparsing>=2.2.1
  Downloading pyparsing-3.1.1-py3-none-any.whl (103 kB)
Collecting cycler>=0.10
  Downloading cycler-0.11.0-py3-none-any.whl (6.4 kB)
Collecting kiwisolver>=1.0.1
  Downloading kiwisolver-1.4.4-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.whl (1.1 MB)
Collecting pillow>=6.2.0
  Downloading Pillow-9.5.0-cp37-cp37m-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (3.3 MB)
Collecting fonttools>=4.22.0
  Downloading fonttools-4.38.0-py3-none-any.whl (965 kB)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.7.3->pandas>=0.25->seaborn) (1.13.0)
Installing collected packages: typing-extensions, packaging, pyparsing, cycler, kiwisolver, pillow, fonttools, matplotlib, seaborn
Successfully installed cycler-0.11.0 fonttools-4.38.0 kiwisolver-1.4.4 matplotlib-3.5.3 packaging-23.1 pillow-9.5.0 pyparsing-3.1.1 seaborn-0.12.2 typing-extensions-4.7.1
```

```
Collecting tabulate
  Downloading tabulate-0.9.0-py3-none-any.whl (35 kB)
Installing collected packages: tabulate
Successfully installed tabulate-0.9.0
```

WARNING: The directory '/home/.cache/pip' or its parent directory is not owned or is not writable by the current user. The cache has been disabled. Check the permissions and owner of that directory. If executing pip with sudo, you may want sudo's -H flag.

WARNING: The directory '/home/.cache/pip' or its parent directory is not owned or is not writable by the current user. The cache has been disabled. Check the permissions and owner of that directory. If executing pip with sudo, you may want sudo's -H flag.

WARNING: The directory '/home/.cache/pip' or its parent directory is not owned or is not writable by the current user. The cache has been disabled. Check the permissions and owner of that directory. If executing pip with sudo, you may want sudo's -H flag.

```
In [3]: # Los siguientes paquetes están disponibles en el cluster
sc.list_packages()
```

```
VBox()
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
```

Package	Version
aws-cfn-bootstrap	2.0
beautifulsoup4	4.9.3
boto	2.49.0
click	8.1.3
cycler	0.11.0
docutils	0.14
fonttools	4.38.0
jmespath	1.0.1
joblib	1.2.0
kiwisolver	1.4.4
lockfile	0.11.0
lxml	4.9.2
matplotlib	3.5.3
mysqlclient	1.4.2
nltk	3.8.1
nose	1.3.4
numpy	1.20.0
packaging	23.1
pandas	1.3.5
Pillow	9.5.0
pip	20.2.2
py-dateutil	2.2
pyparsing	3.1.1
pystache	0.5.4
python-daemon	2.2.3
python-dateutil	2.8.2
python37-sagemaker-pyspark	1.4.2
pytz	2023.3
PyYAML	5.4.1
regex	2021.11.10
seaborn	0.12.2
setuptools	28.8.0
simplejson	3.2.0
six	1.13.0
tabulate	0.9.0
tqdm	4.65.0
typing-extensions	4.7.1
wheel	0.29.0
windmill	1.6

WARNING: The directory '/home/.cache/pip' or its parent directory is not owned or is not writable by the current user. The cache has been disabled. Check the permissions and owner of that directory. If executing pip with sudo, you may want sudo's -H flag.

## 1. Introducción

En este capstone vamos a aprender un modelo de detección del sentimiento utilizando MLlib y EMR. Una vez aprendido ampliaremos el proyecto serializando este modelo y comparándolo con el modelo de detección de sentimiento disponible en AWS.

Para ello utilizaremos el dataset de **amazon reviews** que está disponible de manera pública

<https://s3.amazonaws.com/amazon-reviews-pds/readme.html>

Este dataset tiene dos versiones una en tsv y otra en parquet. Nosotros usaremos la que está en parquet que está disponible a través de la ruta de s3: `s3://amazon-reviews-pds/parquet`.

#### NOTA IMPORTANTE:

**El link para descargar los archivos de reviews no funciona, con lo cual realizaré esa parte del código como creo que debería ser (apartados 1-Carga de datos y 2-Filtrado) pero no lo descargará realmente.**

**Para poder continuar con la actividad, cargaré directamente en mi S3 los ficheros .parquet que componen el dataset de electronics facilitados por otro compañero que ya terminó el capstone y los ha descargado de su S3.**

Este dataset tiene las siguientes columnas (de su diccionario de datos):

marketplace	- 2 letter country code of the marketplace where the review was written.
customer_id	- Random identifier that can be used to aggregate reviews written by a single author.
review_id	- The unique ID of the review.
product_id	- The unique Product ID the review pertains to. In the multilingual dataset the reviews for the same product in different countries can be grouped by the same product_id.
product_parent	- Random identifier that can be used to aggregate reviews for the same product.
product_title	- Title of the product.
product_category	- Broad product category that can be used to group reviews (also used to group the dataset into coherent parts).
star_rating	- The 1-5 star rating of the review.
helpful_votes	- Number of helpful votes.
total_votes	- Number of total votes the review received.
vine	- Review was written as part of the Vine program.
verified_purchase	- The review is on a verified purchase.
review_headline	- The title of the review.
review_body	- The review text.
review_date	- The date the review was written.

De estas, la columna `product_category` se usa como clave de partición. Podéis encontrar toda la información en el enlace que os proporcionamos más arriba.

---

## 2. Análisis exploratorio

Antes de empezar con el modelado exploraremos los datos minimamente para poder estudiar sus propiedades.

---

### Tarea 1: Carga de datos

*Carga el dataset completo en formato parquet y cuenta sus registros. De momento, no lo persistas.*

In [4]:

```
from pyspark.sql import SparkSession
```

```
# Creamos la instancia SparkSession
```

```
spark = SparkSession.builder \
    .appName("SentimentAnalysis") \
    .getOrCreate()
```

```
# Define la ruta del archivo parquet en S3
```

```
parquet_path = "s3://amazon-reviews-pds/parquet-data"
```

```
# Lee el archivo parquet como un DataFrame
```

```
data_frame = spark.read.parquet(parquet_path)
```

```
# Cuenta el número de registros en el DataFrame
```

```
record_count = data_frame.count()
```

```
# Print the number of records
```

```
print("Total Records:", record_count)
```

```
VBox()
```

```
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
```

An error was encountered:

An error occurred while calling o140.parquet.

: java.io.IOException: com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.model.AmazonS3Exception: Access Denied (Service: Amazon S3; Status Code: 403; Error Code: AccessDenied; Request ID: B3MW5AVFZ6A6Z98K; S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=; Proxy: null), S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=

```
at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.list(Jets3tNativeFileSystemStore.java:423)
at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.isFolderUsingFolderObject(Jets3tNativeFileSystemStore.java:249)
at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.isFolder(Jets3tNativeFileSystemStore.java:212)
at com.amazon.ws.emr.hadoop.fs.s3n.S3NativeFileSystem.getFileStatus(S3NativeFileSystem.java:523)
at org.apache.hadoop.fs.FileSystem.exists(FileSystem.java:1767)
at com.amazon.ws.emr.hadoop.fs.EmrFileSystem.exists(EmrFileSystem.java:402)
at org.apache.spark.sql.execution.datasources.DataSource$.anonfun$checkAndGlobPathIfNecessary$4(DataSource.scala:784)
at org.apache.spark.sql.execution.datasources.DataSource$.anonfun$checkAndGlobPathIfNecessary$4$adapted(DataSource.scala:782)
at org.apache.spark.util.ThreadUtils$.anonfun$parmap$2(ThreadUtils.scala:372)
at scala.concurrent.Future$.anonfun$apply$1(Future.scala:659)
at scala.util.Success$.anonfun$map$1(Try.scala:255)
at scala.util.Success.map(Try.scala:213)
at scala.concurrent.Future$.anonfun$map$1(Future.scala:292)
at scala.concurrent.impl.Promise.LiftedTree1$1(Promise.scala:33)
at scala.concurrent.impl.Promise$.anonfun$transform$1(Promise.scala:33)
at scala.concurrent.impl.CallbackRunnable.run(Promise.scala:64)
at java.util.concurrent.ForkJoinTask$RunnableExecuteAction.exec(ForkJoinTask.java:1402)
at java.util.concurrent.ForkJoinTask.doExec(ForkJoinTask.java:289)
at java.util.concurrent.ForkJoinPool$WorkQueue.runTask(ForkJoinPool.java:1056)
at java.util.concurrent.ForkJoinPool.runWorker(ForkJoinPool.java:1692)
at java.util.concurrent.ForkJoinWorkerThread.run(ForkJoinWorkerThread.java:175)
```

Caused by: com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.model.AmazonS3Exception: Access Denied (Service: Amazon S3; Status Code: 403; Error Code: AccessDenied; Request ID: B3MW5AVFZ6A6Z98K; S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=; Proxy: null), S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=

```
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.handleErrorResponse(AmazonHttpClient.java:1879)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.handleServiceErrorResponse(AmazonHttpClient.java:1418)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeOneRequest(AmazonHttpClient.java:1387)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeHelper(AmazonHttpClient.java:1157)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.doExecute(AmazonHttpClient.java:814)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeWithTimer(AmazonHttpClient.java:781)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.execute(AmazonHttpClient.java:755)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.access$500(AmazonHttpClient.java:715)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutionBuilderImpl.execute(AmazonHttpClient.java:697)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:561)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:541)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5467)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5414)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5408)
at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.listObjectsV2(AmazonS3Client.java:971)
at com.amazon.ws.emr.hadoop.fs.s3.lite.call.ListObjectsV2Call.perform(ListObjectsV2Call.java:26)
at com.amazon.ws.emr.hadoop.fs.s3.lite.call.ListObjectsV2Call.perform(ListObjectsV2Call.java:12)
at com.amazon.ws.emr.hadoop.fs.s3.lite.executor.GlobalS3Executor$CallPerformer.call(GlobalS3Executor.java:111)
at com.amazon.ws.emr.hadoop.fs.s3.lite.executor.GlobalS3Executor.execute(GlobalS3Executor.java:138)
at com.amazon.ws.emr.hadoop.fs.s3.lite.AmazonS3LiteClient.invoke(AmazonS3LiteClient.java:191)
at com.amazon.ws.emr.hadoop.fs.s3.lite.AmazonS3LiteClient.invoke(AmazonS3LiteClient.java:186)
at com.amazon.ws.emr.hadoop.fs.s3.lite.AmazonS3LiteClient.listObjectsV2(AmazonS3LiteClient.java:75)
at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.list(Jets3tNativeFileSystemStore.java:414)
... 20 more
```

Traceback (most recent call last):

```
File "/mnt/yarn/usercache/livy/appcache/application_1691082518949_0003/container_1691082518949_0003_01_000001/pyspark.zip/pyspark/sql/readwriter.py", line 364, in parquet
return self._df(self._jreader.parquet(_to_seq(self._spark._sc, paths)))
```

```
File "/mnt/yarn/usercache/livy/appcache/application_1691082518949_0003/container_1691082518949_0003_01_000001/py4j-0.10.9.5-src.zip/py4j/java_gateway.py", line 1322, in __call__
-
    answer, self.gateway_client, self.target_id, self.name)
File "/mnt/yarn/usercache/livy/appcache/application_1691082518949_0003/container_1691082518949_0003_01_000001/pyspark.zip/pyspark/sql/utls.py", line 190, in deco
    return f(*a, **kw)
File "/mnt/yarn/usercache/livy/appcache/application_1691082518949_0003/container_1691082518949_0003_01_000001/py4j-0.10.9.5-src.zip/py4j/protocol.py", line 328, in get_return_v
alue
    format(target_id, ".", name), value)
py4j.protocol.Py4JJavaError: An error occurred while calling o140.parquet.
: java.io.IOException: com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.model.AmazonS3Exception: Access Denied (Service: Amazon S3; Status Code: 403; Error Code: Acces
sDenied; Request ID: B3MW5AVFZ6A6Z98K; S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=; Proxy: null), S3 Extended Request I
D: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=
    at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.list(Jets3tNativeFileSystemStore.java:423)
    at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.isFolderUsingFolderObject(Jets3tNativeFileSystemStore.java:249)
    at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.isFolder(Jets3tNativeFileSystemStore.java:212)
    at com.amazon.ws.emr.hadoop.fs.s3n.S3NativeFileSystem.getFileStatus(S3NativeFileSystem.java:523)
    at org.apache.hadoop.fs.FileSystem.exists(FileSystem.java:1767)
    at com.amazon.ws.emr.hadoop.fs.EmrFileSystem.exists(EmrFileSystem.java:402)
    at org.apache.spark.sql.execution.datasources.DataSource$.anonfun$checkAndGlobPathIfNecessary$4(DataSource.scala:784)
    at org.apache.spark.sql.execution.datasources.DataSource$.anonfun$checkAndGlobPathIfNecessary$4$adapted(DataSource.scala:782)
    at org.apache.spark.util.ThreadUtils$.anonfun$parmap$2(ThreadUtils.scala:372)
    at scala.concurrent.Future$.anonfun$apply$1(Future.scala:659)
    at scala.util.Success$.anonfun$map$1(Try.scala:255)
    at scala.util.Success.map(Try.scala:213)
    at scala.concurrent.Future$.anonfun$map$1(Future.scala:292)
    at scala.concurrent.impl.Promise.LiftedTree1$1(Promise.scala:33)
    at scala.concurrent.impl.Promise$.anonfun$transform$1(Promise.scala:33)
    at scala.concurrent.impl.CallbackRunnable.run(Promise.scala:64)
    at java.util.concurrent.ForkJoinTask$RunnableExecuteAction.exec(ForkJoinTask.java:1402)
    at java.util.concurrent.ForkJoinTask.doExec(ForkJoinTask.java:289)
    at java.util.concurrent.ForkJoinPool$WorkQueue.runTask(ForkJoinPool.java:1056)
    at java.util.concurrent.ForkJoinPool.runWorker(ForkJoinPool.java:1692)
    at java.util.concurrent.ForkJoinWorkerThread.run(ForkJoinWorkerThread.java:175)
Caused by: com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.model.AmazonS3Exception: Access Denied (Service: Amazon S3; Status Code: 403; Error Code: AccessDenied; Re
quest ID: B3MW5AVFZ6A6Z98K; S3 Extended Request ID: V6YRZFzLFt5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=; Proxy: null), S3 Extended Request ID: V6YRZFzLFt
5Kb5UPpS9taMXau92MVpmVMOh8DduyZYj6DNNQ8kjGEa+NP+dA7syuW9c6y/U8WdE=
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.handleErrorResponse(AmazonHttpClient.java:1879)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.handleServiceErrorResponse(AmazonHttpClient.java:1418)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeOneRequest(AmazonHttpClient.java:1387)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeHelper(AmazonHttpClient.java:1157)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.doExecute(AmazonHttpClient.java:814)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.executeWithTimer(AmazonHttpClient.java:781)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.execute(AmazonHttpClient.java:755)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutor.access$500(AmazonHttpClient.java:715)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient$RequestExecutionBuilderImpl.execute(AmazonHttpClient.java:697)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:561)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:541)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5467)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5414)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:5408)
    at com.amazon.ws.emr.hadoop.fs.shaded.com.amazonaws.services.s3.AmazonS3Client.listObjectsV2(AmazonS3Client.java:971)
    at com.amazon.ws.emr.hadoop.fs.s3lite.call.ListObjectsV2Call.perform(ListObjectsV2Call.java:26)
    at com.amazon.ws.emr.hadoop.fs.s3lite.call.ListObjectsV2Call.perform(ListObjectsV2Call.java:12)
    at com.amazon.ws.emr.hadoop.fs.s3lite.executor.GlobalS3Executor$CallPerformer.call(GlobalS3Executor.java:111)
    at com.amazon.ws.emr.hadoop.fs.s3lite.executor.GlobalS3Executor.execute(GlobalS3Executor.java:138)
    at com.amazon.ws.emr.hadoop.fs.s3lite.AmazonS3LiteClient.invoke(AmazonS3LiteClient.java:191)
    at com.amazon.ws.emr.hadoop.fs.s3lite.AmazonS3LiteClient.invoke(AmazonS3LiteClient.java:186)
    at com.amazon.ws.emr.hadoop.fs.s3lite.AmazonS3LiteClient.listObjectsV2(AmazonS3LiteClient.java:75)
```

```
at com.amazon.ws.emr.hadoop.fs.s3n.Jets3tNativeFileSystemStore.list(Jets3tNativeFileSystemStore.java:414)
... 20 more
```

**Resultado esperado:** 160796570 registros

**Como expliqué anteriormente, haré solo la parte de código pero no podré ejecutarlo correctamente.**

---

## Tarea 2: Filtrado

Como el dataset es masivo para entrenar el modelo de sentiment vamos a trabajar únicamente con una partición. Concretamente utilizaremos la partición de `Electronics`. Filtra los datos para quedarte con esta partición y cuenta ahora el total de elementos de este nuevo dataset. No cachees este dataset.

In [5]:

```
# Solución
# Proceso
from pyspark.sql import SparkSession

# Inicia una instancia de SparkSession
spark = SparkSession.builder \
    .appName("SentimentAnalysis") \
    .getOrCreate()

#Esto es lo que habría que hacer si funcionase el link
# Ruta del archivo parquet en S3
ruta_parquet = "s3://amazon-reviews-pds/parquet-data"

# Lee el archivo parquet y crea un DataFrame
data_frame = spark.read.parquet(ruta_parquet)

# Filtra los datos para obtener la categoría de Electrónicos
data_frame_electronics = data_frame.filter(data_frame["product_category"] == "Electronics")

#Utilizaré directamente los parquet de electronics que he subido a S3

# Ruta del archivo "Electronics" parquet en S3
ruta_parquet = "s3://capstone12bucket2/electronics"

# Lee el archivo parquet y crea un DataFrame
data_frame_electronics = spark.read.parquet(ruta_parquet)

# Cuenta el número de registros en el nuevo DataFrame
conteo_electronics = data_frame_electronics.count()

# Imprime el conteo de registros
print("Registros en electronics:", conteo_electronics)
```

VBox()



```
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
Registros en electroncis: 3120938
```

**Resultado esperado:** 3120938 registros

---

## Tarea 3: Almacenamiento

Para no seguir trabajando con los datos públicos, vamos a escribir los datos en un bucket de S3 en nuestra cuenta. Para ello, crea un bucket the S3 para este capstone y escribe los datos dentro del bucket en el directorio `eElectronics`. Utiliza `repartition` para tener 32 particiones. Tras esto, vuelve a cargar el dataset y cachéalo.

**No los cargo porque ya los tengo subidos, comento el código**

In [6]:

```
# Solución

#####

# Ruta S3 para almacenar los datos
s3_bucket = "s3://capstone12bucket2"

# Escribe los datos en el depósito S3 en la subcarpeta "eElectronica" con 32 particiones
data_frame_electronics.repartition(32).write.parquet(s3_bucket + "/eElectronica", mode="overwrite")

# Carga nuevamente el conjunto de datos desde S3
data_frame_electronics = spark.read.parquet(s3_bucket + "/eElectronica").cache()

#####
```

```
VBox()
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
'\n# Ruta S3 para almacenar los datos\ns3_bucket = "s3://capstone12bucket2"\n\n# Escribe los datos en el dep?sito S3 en la subcarpeta "eElectronica" con 32 particiones\ndata_frame_electronics.repartition(32).write.parquet(s3_bucket + "/eElectronica", mode="overwrite")\n\n# Carga nuevamente el conjunto de datos desde S3\ndata_frame_electronics = spark.read.parquet(s3_bucket + "/eElectronica").cache()\n\n'
```

---

## Tarea 4: Almacenamiento

Obten los siguiente resultados del dataset que acabáis de cargar:

1. Muestra el total de reviews para cada posible número de estrellas recibidas (`star_rating`)
2. Obtén los 10 productos con mayor número de votos (`total_votes`) mostrando su nombre, numero de votos y valoración media (`star_rating`)
3. Obtén la cantidad de reviews (1 registro de dataset -> 1 review) y la valoración media (`star_rating`) por mes y año. Obten los últimos 15 registros ordenador por año y mes.

In [7]:

```
from tabulate import tabulate
import pandas as pd
```

```

from pyspark.sql.functions import desc
from pyspark.sql.functions import month, year, avg
from pyspark.sql.window import Window
from pyspark.sql.functions import row_number
from pyspark.sql import functions as F

# Análisis de reseñas por número de estrellas

from pyspark.sql.functions import count

# Calcular el total de reseñas por número de estrellas
reviews_per_star = data_frame_electronics.groupBy("star_rating").agg(count("*").alias("total_reviews"))

review_data = reviews_per_star.toPandas()
print(tabulate(review_data, headers="keys", tablefmt="fancy_grid", floatfmt=".0f"))

# 2-10 productos con mayor número de votos

# 10 productos con mas votos
top_ten_products = data_frame_electronics.groupBy("product_id", "product_title") \
    .agg(F.sum("total_votes").alias("total_votes"), F.avg("star_rating").alias("star_rating")) \
    .orderBy(F.desc("total_votes")) \
    .limit(10)

top_ten_products.show()

# 3-Cantidad de reviews

# Cantidad de reviews y valoración media por mes y año
reviews_avg_rating = data_frame_electronics.select("review_date", "star_rating") \
    .withColumn("year", year("review_date")) \
    .withColumn("month", month("review_date")) \
    .groupBy("year", "month") \
    .agg(count("*").alias("reviews_count"), avg("star_rating").alias("mean_star_rating")) \
    .orderBy(desc("year"), desc("month"))

# Últimos 15 registros por mes y año
window_spec = Window.orderBy(desc("year"), desc("month"))
last_15_records = reviews_avg_rating.withColumn("row_number", row_number().over(window_spec)) \
    .filter("row_number <= 15") \
    .orderBy(desc("year"), desc("month"))

# Obtener Los últimos 15 registros ordenados por año y mes
last_15_records_selected = last_15_records.select("year", "month", "reviews_count", "mean_star_rating")

# Mostrar Los resultados seleccionados
last_15_records_selected.show()

```

VBox()

FloatProgress(value=0.0, bar\_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...

	star_rating	total_reviews
0	1	360558
1	3	240859
2	4	542181
3	5	1796672
4	2	180668

+-----+			
product_id	product_title total_votes	star_rating	
+-----+			
B000I1X6PM	Denon AKDL1 Dedic...	46348	3.4917491749174916
B000J36XR2	AudioQuest K2 Ter...	20515	3.9357429718875503
B004QK7HI8	Mohu Leaf 30 TV A...	18198	4.091145435081817
B0001FTVEK	Sennheiser On-Ear...	18031	4.034208432776452
B000EPLP3C	Zune 30 GB Digita...	17598	3.7341513292433537
B001FA1018	Apple iPod touch ...	17103	4.384232365145229
B00D5Q75RC	Bose SoundLink Mi...	16028	4.721810699588477
B0054JJ0QW	Bose QuietComfort...	13383	4.439661515820457
B0002L5R78	High Speed HDMI C...	12859	4.462750716332378
B000WYVBR0	VideoSecu ML531BE...	11102	4.579711991844017
+-----+			

+-----+			
year month reviews_count	mean_star_rating		
+-----+			
2015	8	103336	4.0944104668266625
2015	7	100128	4.086149728347715
2015	6	91815	4.093296302347111
2015	5	89676	4.100483964494402
2015	4	93469	4.102825535739122
2015	3	109175	4.1156858255095035
2015	2	107623	4.117920890515967
2015	1	120852	4.152268890874789
2014	12	108294	4.120302140469462
2014	11	77844	4.107843892914033
2014	10	78519	4.113870528152422
2014	9	78126	4.115928116120114
2014	8	82550	4.116026650514839
2014	7	79816	4.117582940763756
2014	6	48707	4.015726692261892
+-----+			

Resultados esperados:

1. Muestra el total de reviews para cada posible número de estrellas recibidas (star\_rating)

star_rating	count
1	360558
3	240859

<i>star_rating</i>	<i>count</i>
4	542181
5	1796672
2	180668

1. Obtén los 10 productos con mayor número de votos (*total\_votes*) mostrando su nombre, numero de votos y valoración media (*star\_rating*)

<i>product_title</i>	<i>total_votes</i>	<i>star_rating</i>
Denon AKDL1 Dedicated Link Cable (Discontinued by Manufacturer)	12944	3
AudioQuest K2 Terminated Speaker Cable - UST 2.44 m Plugs 8' Pair (Discontinued by Manufacturer)	9072	1
Panasonic ErgoFit In-Ear Earbud Headphone	8680	5
Apple iPod touch 8GB (4th Generation)	6353	5
Denon AKDL1 Dedicated Link Cable (Discontinued by Manufacturer)	5546	1
Apple iPod touch 8 GB 2nd Generation	4595	5
Bose QuietComfort 15 Acoustic Noise Cancelling Headphones (Discontinued by Manufacturer)	4556	4
Panasonic ErgoFit In-Ear Earbud Headphone	4341	5
X-Mini II XAM4-B Portable Capsule Speaker, Mono	4260	1
Denon AKDL1 Dedicated Link Cable (Discontinued by Manufacturer)	4242	2

1. Obtén la cantidad de reviews (1 registro de dataset -> 1 review) y la valoración media (*star\_rating*) por mes y año. Obten los últimos 15 registros ordenador por año y mes.

<i>year</i>	<i>month</i>	<i>review_count</i>	<i>mean_star_rating</i>
2015	8	103336	4.09441
2015	7	100128	4.08615
2015	6	91815	4.0933
2015	5	89676	4.10048
2015	4	93469	4.10283
2015	3	109175	4.11569
2015	2	107623	4.11792
2015	1	120852	4.15227
2014	12	108294	4.1203
2014	11	77844	4.10784
2014	10	78519	4.11387
2014	9	78126	4.11593
2014	8	82550	4.11603
2014	7	79816	4.11758

year	month	review_count	mean_star_rating
2014	6	48707	4.01573

### 3. Modelado

Como paso previo al modelado realizaremos dos procesos de limpieza sobre los datos:

#### Tarea 6: Preparación del texto

Limpiad el texto de las reviews ( `review_body` ) utilizando expresiones sobre strings o expresiones regulares

- Pasar todo el texto a minúsculas.
- Eliminar números y signos de puntuación.
- Si existen, elimina los registros con valores nulos en el body con las transformaciones anteriores.

Muestra los resultados para las primeras 10 filas del dataframe ordenadas por `review_id`

```
In [8]: # Solución
from pyspark.sql import functions as F
from pyspark.sql.functions import lower, regexp_replace, col

# Pasar todo el texto a minúsculas.
df_lowercase = data_frame_electronics.withColumn("cleaned_review_body", lower(col("review_body")))

# Eliminar números y signos de puntuación
df_cleaned = df_lowercase.withColumn("cleaned_review_body", regexp_replace(col("cleaned_review_body"), "[^a-zA-Z\s]", ""))

# Eliminar los registros con valores nulos
df_filtered = df_cleaned.filter(df_cleaned["cleaned_review_body"] != "")

# Mostrar los resultados de las primeras 10 filas ordenadas por review_id
df_filtered.select("review_id", "review_body", "cleaned_review_body").orderBy("review_id").show(10, truncate=False)
```

```
VBox()
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
```

review_id	review_body	cleaned_review_body
R10000WMGX51T	Great little emergency radio. Very good reception. The weather band is a feature. Can't beat the quality for this price/	great little emergency radio very good reception thebr weather band is a feature cant beat the quality for this price
R10001L4QTCA84	Lives up to its claim, and really does fit bulky phone cases. Braided cable is sturdy but flexible. I think it stays a little more flexible in the cold weather, which is nice. Definitely getting a few more in the future!	lives up to its claim and really does fit bulky phone cases braided cable is sturdy but flexible i think it stays a little more flexible in the cold weather which is nice definitely getting a few more in the future
R100030LR2P5UE	I've gone through three pairs of these in the last two years. I am in love with the sound quality, and even though I know it's not the best I particularly love how the bass sounds. They're comfortable to wear and very isolating. With these headphones, you don't even need noise canceling. There is very little sound leak, unless you like to listen to music ridiculously loud. All in all, I was very impressed with these. They're without a doubt the best sounding headphones I've ever owned.	ive gone through three pairs of these in the last two years i am in love with the sound quality and even though i know its not the best i particularly love how the bass sounds theyre comfortable to wear and very isolating with these headphones you dont even need noise canceling there is very little sound leak unless you like to listen to music ridiculously loud all in all i was very impressed with these theyre without a doubt the best sounding headphones ive ever owned
R100050193PJ6W	stopped working after a while, changed batteries, it worked for a few days, then it quit	stopped working after a while changed batteries it worked for a few days then it quit

[illegible]

**Resultado esperado:**

review_id	review_body	clean_review_body
R10000WMGX5S1T	Great little emergency radio. Very good reception. The weather band is a feature. Can't beat the quality for this price/	great little emergency radio very good reception thebr weather band is a feature cant beat the quality for this price
R10001L4QTCA84	Lives up to its claim, and really does fit bulky phone cases. Braided cable is sturdy but flexible. I think it stays a little more flexible in the cold weather, which is nice. Definitely getting a few more in the future!	lives up to its claim and really does fit bulky phone cases braided cable is sturdy but flexible i think it stays a little more flexible in the cold weather which is nice definitely getting a few more in the future
R10003OLR2P5UE	I've gone through three pairs of these in the last two years. I am in love with the sound quality, and even though I know it's not the best I particularly love how the bass sounds. They're comfortable to wear and very isolating. With these headphones, you don't even	ive gone through three pairs of these in the last two years i am in love with the sound quality and even though i know its not the best i particularly love how the bass sounds theyre comfortable to wear and very isolating with these headphones you dont even need noise canceling there is very little sound leak unless you like to listen to music ridiculously loud all in all i was very impressed with these theyre without a doubt the best sounding headphones ive ever ownedbr now the problem the wires are thin and stringy and do not last on my first pair the part of the wire that connected

<i>review_id</i>	<i>review_body</i>	<i>clean_review_body</i>
	<p>need noise canceling. There is very little sound leak, unless you like to listen to music ridiculously loud. All in all, I was very impressed with these. They're without a doubt the best sounding headphones I've ever owned.</p> <p>Now, the problem: The wires are thin and stringy, and do NOT last. On my first pair, the part of the wire that connected to the left cup came apart. I'm not an abuser of headphones, either. On the other two pairs, they wire at the base next to the adapter came apart. I went at them with a soldering iron, desperately trying to make them last as long as I could, but they'd always crap out on me again. The sound quality distorts over time, and the foam around the cups is cheap and wears out quickly.</p> <p>They aren't worth the price for such bad quality. I'd suggest looking around for other pairs, Sony, Denon, and Sennheiser all have superior headphones for a similar price. I myself just ordered a pair of Denon AHD1001's, and here's hoping they last longer!</p>	<p>to the left cup came apart im not an abuser of headphones either on the other two pairs they wire at the base next to the adapter came apart i went at them with a soldering iron desperately trying to make them last as long as i could but theyd always crap out on me again the sound quality distorts over time and the foam around the cups is cheap and wears out quicklybr br they arent worth the price for such bad quality id suggest looking around for other pairs sony denon and sennheiser all have superior headphones for a similar price i myself just ordered a pair of denon ahds and heres hoping they last longer</p>
R10005O193PJ6W	stopped working after a while, changed batteries, it worked for a few days, then it quit	stopped working after a while changed batteries it worked for a few days then it quit
R10008LR7CU84N	I ordered this cable and it doesn't work when I contacted them they told me I was doing something wrong. I then had my dad who it a certified computer tech look at it and there is something wrong with the cable. When I told them they never responded to me again.	i ordered this cable and it doesnt work when i contacted them they told me i was doing something wrong i then had my dad who it a certified computer tech look at it and there is something wrong with the cable when i told them they never responded to me again
R10009JN2UWOJC	Have not owned it that long however it has the features , feel and works like a quality unit that would be at a much higher price point	have not owned it that long however it has the features feel and works like a quality unit that would be at a much higher price point
R1000AMVKPW32O	Bought for a gift and it is just what was needed to mount the new 32" TV outdoors. The fact that it has full motion swing makes it even better because we can move it around to see it from different angles and still have a sturdy mount.	bought for a gift and it is just what was needed to mount the new tv outdoors the fact that it has full motion swing makes it even better because we can move it around to see it from different angles and still have a sturdy mount
R1000CJMO2L8X4	Perfect for the gym	perfect for the gym
R1000EDGJU3CU	Love these !!! The sound quality is amazing ! The price was amazing especially for the quality.	love these the sound quality is amazing the price was amazing especially for the quality
R1000EG9XXBLXT	I have had good success with these disks, and have used hundreds of them successfully on both computers and a dedicated Panosonic DVD recorder. They seem very reliable, and the lines on the disk label help to keep labeling neat and straight.	i have had good success with these disks and have used hundreds of them successfully on both computers and a dedicated panosonic dvd recorder they seem very reliable and the lines on the disk label help to keep labeling neat and straight

## Tarea 7: Obtención del sentiment

Cread la variable `sentiment` en función del número de estrellas asumiendo que una review de menos (<) de 3 estrellas es negativa, usando 1 para el sentiment positivo y 0 para el negativo. Para poder generar la variable que determine el sentiment a partir del número de estrellas podéis utilizar la función de spark `when` . Muestra el resultado para las primeras 10 reviews ordenadas por `review_id` .



```
In [11]: # Solución
from pyspark.sql.functions import when

# Agregar la columna "sentiment" basada en la valoración
df_sentiment = data_frame_electronics.withColumn("sentiment", when(data_frame_electronics["star_rating"] < 3, 0).otherwise(1))

# Mostrar los primeros 10 resultados con las columnas relevantes
df_result = df_sentiment.select("review_id", "review_body", "star_rating", "sentiment").orderBy("review_id").limit(10)
df_result.show(truncate=False)
```

VBox()

FloatProgress(value=0.0, bar\_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...

review_id	review_body	
star_rating	sentiment	
R1000WVGXSS51T	Great little emergency radio. Very good reception. The weather band is a feature. Can't beat the quality for this price!	
5	1	
R10001L4QTCA84	Lives up to its claim, and really does fit bulky phone cases. Braided cable is sturdy but flexible. I think it stays a little more flexible in the cold weather, which is nice. Definitely getting a few more in the future!	
5	1	
R10003OLR2P5UE	I've gone through three pairs of these in the last two years. I am in love with the sound quality, and even though I know it's not the best I particularly love how the bass sounds. They're comfortable to wear and very isolating. With these headphones, you don't even need noise canceling. There is very little sound leak, unless you like to listen to music ridiculously loud. All in all, I was very impressed with these. They're without a doubt the best sounding headphones I've ever owned.	
	The wires are thin and stringy, and do NOT last. On my first pair, the part of the wire that connected to the left cup came apart. I'm not an abuser of headphones, either. On the other two pairs, they wire at the base next to the adapter came apart. I went at them with a soldering iron, desperately trying to make them last as long as I could, but they'd always crap out on me again. The sound quality distorts over time, and the foam around the cups is cheap and wears out quickly.	
	They aren't worth the price for such bad quality. I'd suggest looking around for other pairs, Sony, Denon, and Sennheiser all have superior headphones for a similar price. I myself just ordered a pair of Denon AHD1001's, and here's hoping they last longer!	3
		1
R10005O193PJ6W	Stopped working after a while, changed batteries, it worked for a few days, then it quit	
3	1	
R10008LR7CU84N	I ordered this cable and it doesn't work when I contacted them they told me I was doing something wrong. I then had my dad who is a certified computer tech look at it and there is something wrong with the cable. When I told them they never responded to me again.	
1	0	
R10009JN2UWOJC	Have not owned it that long however it has the features , feel and works like a quality unit that would be at a much higher price point	
5	1	
R1000AMVKPW32O	Bought for a gift and it is just what was needed to mount the new 32" TV outdoors. The fact that it has full motion swing makes it even better because we can move it around to see it from different angles and still have a sturdy mount.	
5	1	
R1000CJM02L8X4	Perfect for the gym	
5	1	
R1000EDGJU03CU	Love these !!! The sound quality is amazing ! The price was amazing especially for the quality.	
5	1	
R1000EG9XXBLXT	I have had good success with these disks, and have used hundreds of them successfully on both computers and a dedicated Panasonic DVD recorder. They seem very reliable, and the lines on the disk label help to keep labeling neat and straight.	
5	1	

**Resultado esperado:**

review_id	review_body	star_rating	sentiment
R10000WMGX51T	Great little emergency radio. Very good reception. The weather band is a feature. Can't beat the quality for this price/	5	1
R10001L4QTCA84	Lives up to its claim, and really does fit bulky phone cases. Braided cable is sturdy but flexible. I think it stays a little more flexible in the cold weather, which is nice. Definitely getting a few more in the future!	5	1
R10003OLR2P5UE	I've gone through three pairs of these in the last two years. I am in love with the sound quality, and even though I know it's not the best I particularly love how the bass sounds. They're comfortable to wear and very isolating. With these headphones, you don't even need noise canceling. There is very little sound leak, unless you like to listen to music ridiculously loud. All in all, I was very impressed with these. They're without a doubt the best sounding headphones I've ever owned.  Now, the problem: The wires are thin and stringy, and do NOT last. On my first pair, the part of the wire that connected to the left cup came apart. I'm not an abuser of headphones, either. On the other two pairs, they wire at the base next to the adapter came apart. I went at them with a soldering iron, desperately trying to make them last as long as I could, but they'd always crap out on me again. The sound quality distorts over time, and the foam around the cups is cheap and wears out quickly.  They aren't worth the price for such bad quality. I'd suggest looking around for other pairs, Sony, Denon, and Sennheiser all have superior headphones for a similar price. I myself just ordered a pair of Denon AHD1001's, and here's hoping they last longer!	3	1
R10005O193PJ6W	stopped working after a while, changed batteries, it worked for a few days, then it quit	3	1
R10008LR7CU84N	I ordered this cable and it doesn't work when I contacted them they told me I was doing something wrong. I then had my dad who it a certified computer tech look at it and there is something wrong with the cable. When I told them they never responded to me again.	1	0
R10009JN2UWOJC	Have not owned it that long however it has the features , feel and works like a quality unit that would be at a much higher price point	5	1
R1000AMVKPW32O	Bought for a gift and it is just what was needed to mount the new 32" TV outdoors. The fact that it has full motion swing makes it even better because we can move it around to see it from different angles and still have a sturdy mount.	5	1
R1000CJMO2L8X4	Perfect for the gym	5	1
R1000EDGJU03CU	Love these !!! The sound quality is amazing ! The price was amazing especially for the quality.	5	1
R1000EG9XXBLXT	I have had good success with these disks, and have used hundreds of them successfully on both computers and a dedicated Panosonic DVD recorder. They seem very reliable, and the lines on the disk label help to keep labeling neat and straight.	5	1

## Tarea 8: División del conjunto de datos

Divide el conjunto de datos en entrenamiento (70% de los datos) y test (30% de los datos). Una vez hecho esto, guarda los datos de test en el bucket the s3 previamente creado (usa el nombre `electronics_test` )

```
In [12]: # Solución

#División del conjunto de datos en entrenamiento y prueba
train_data, test_data = df_sentiment.randomSplit([0.7, 0.3], seed=5)
train_data = train_data.na.drop(subset=["review_body"])
test_data = test_data.na.drop(subset=["review_body"])

# Guardar Los datos de prueba en el bucket S3
test_data.write.parquet("s3://capstone12bucket2/electronics_test", mode="overwrite")
train_data.show()
```

VBox()

Amazon.de - Reviews für Philips Fidelio L...									
/marketplace/customer_id/	review_id/product_id/product_parent/	product_title/star_rating/helpful_votes/total_votes/vine/verified_purchase/	review_headline/						
review_body/review_date/year/product_category/sentiment/									
DE  10873130  RV8R4ZISL74R0 B0066XYIJ4	437018546 Philips Fidelio L...	4  0  1  N	Y Very good headpho... Very sati						
sified bu...  2013-01-11 2013  Electronics  1									
DE  16052556 R310SD8WTTVAEG B000065BPB	114634020 Sennheiser HD280P...	4  0  0  N	Y  nice sound quality even thou						
gh they ...  2014-12-05 2014  Electronics  1									
DE  16725375 R356QYZ9KZVW2G B002SJKWZ4	428789313 Magic Wand Button...	3  0  2  N	Y  gedult Es ist ec						
ht cool ...  2015-01-02 2015  Electronics  1									
DE  16729634 R13NFYSFEZKTBR B009346RSS	375251831 FiiO Taishan D03K...	5  5  5  N	Y FiiO D03K mit Phi... verbindet						
bei mir...  2014-01-05 2014  Electronics  1									
DE  16780689 R327JAJ6DQGSUI B005HJWW8	724543729 FiiO Fujiyama E06...	2  0  0  N	Y FiiO Fujiyama E06... Leider na						
ch einem...  2014-09-15 2014  Electronics  0									
DE  16855834 R302T4TNX6BDZT B00PCY000Q	175237145 Apple MD820ZM/A A...	5  1  1  N	Y  sehr gut geht alles keinerlei						
Problem...  2015-08-12 2015  Electronics  1									
DE  17411119  R75U0KQJSJA75 B002C1BHIO	872923735 AmazonBasics HDMI...	5  0  0  N	N  robistu55 Hallo !						
Absolute...  2010-11-07 2010  Electronics  1									
DE  17758998 R3767SKRX4YH0Y B0077JREJC	129618844 Panasonic TY-CC20...	4  0  1  N	Y fits nice, design... fits nic						
e, design...  2015-07-04 2015  Electronics  1									
DE  17906058  REQPYSEK5CYRU B00579Z0ZY	407582114 Diablo III Mouse ...	5  3  4  N	Y  Nur geil^^ Das ist a						
uf jeden...  2011-11-10 2011  Electronics  1									
DE  18076227 R1MZ2M2HN65S4D B000065BPB	114634020 Sennheiser HD280P...	5  0  1  N	N Super Abschirmung... Hab die K						
opfh&oum...  2011-08-19 2011  Electronics  1									
DE  18277633 R1KEGV179ZRQG B002C1BHIO	872923735 AmazonBasics HDMI...	5  0  0  N	Y Tolles Kabel wo P... Um es kur						
z zu mac...  2013-01-04 2013  Electronics  1									
DE  18664601 R35F571EBS4ZZI B004XBDRKC	685653138 Onkyo UWF-1 WIFI ...	4  0  0  N	Y Bei mir läuft all... Also mein						
Receive...  2012-12-05 2012  Electronics  1									
DE  18698259  RRDMYHXCONAW B0067TQQI8	441125258 Anker 2. Gen Astr...	5  0  0  N	Y  top Akku Lädt mit						
1000mAh...  2013-09-11 2013  Electronics  1									
DE  18843374  RG3S2B4WZ1JGN B00E906C96	986493541 FiiO X3 Digital M...	1  3  5  N	N In one year, the ... After a y						
ear of m...  2015-03-06 2015  Electronics  0									
DE  19024734 R22R10BTQZ6PGG B005HJWW8	724543729 FiiO Fujiyama E06...	5  1  1  N	Y  Kleines Monster! Hab ihn e						
rst seit...  2014-02-21 2014  Electronics  1									
DE  19072669 R1GY056FO316MA B004RE90NI	452956477 Philips EXP2546 t...	1  1  1  N	Y spielt nur mit im... siehe übe						
rschrift...  2014-09-17 2014  Electronics  0									
DE  19346997 R2HDK2U007W4GU B002C1BHIO	872923735 AmazonBasics HDMI...	5  0  0  N	Y  Hammerteil Das HDM						
Kabel ka...  2010-08-14 2010  Electronics  1									
DE  19350553 R36NRBAVANLIBY B0067TQQI8	441125258 Anker 2. Gen Astr...	5  0  0  N	N  super Lieferung						
super s...  2013-10-23 2013  Electronics  1									
DE  19382923 R1GBKM6DATP5BJ B0077JREJC	129618844 Panasonic TY-CC20...	3  19  20  N	N  Helligkeitsprobleme Die Kamer						
a Liefer...  2012-10-30 2012  Electronics  1									
DE  19437520 R17BAAOIHXHPF3S B002C1BHIO	872923735 AmazonBasics HDMI...	5  1  1  N	Y  Absolut zufrieden Ich habe						
das Amaz...  2010-09-29 2010  Electronics  1									
only showing top 20 rows									

## Tarea 8: Modelo de sentiment

Construye un **pipeline** de entrenamiento de un modelo de sentiment a partir de los datos preparados anteriormente. Deberas utilizar una secuencia de diferentes **transformadores** y **estimadores**:

- `Tokenizer` nos permitirá construir un vector de palabras a partir de nuestras sentencias
- `StopWordsRemover` nos permitirá limpiar de nuestros vectores de palabras las de menor significado
- Construcción de características dos alternativas:
  - Modelo TF-IDF usando `HashingTF` e `IDF`
  - `Word2Vec` nos permitirá crear un vector a partir de la lista de palabras
- Clasificación binaria, basada en la variable sentiment que hemos utilizado, aplica un clasificador (`LogisticRegression`, `DecisionTree`) evita ensembles por su alto tiempo de aprendizaje.

Buscando en la documentación, encuentra los distintos elementos y conéctalos en un pipeline junto a un algoritmo de clasificación

- Recomendamos utilizar una muestra (método `sample`) pues el tiempo puede ser excesivo
- Es posible ajustar hiperparámetros, pero igualmente puede ser bastante lento

**Valida el modelo con el conjunto de test anterior usando el area bajo la curva ROC**

In [14]:

```
# Solución

from pyspark.ml import Pipeline
from pyspark.ml.feature import Tokenizer, StopWordsRemover, HashingTF, IDF
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.ml.classification import LogisticRegression

# Creamos objeto Tokenizer
tokenizer = Tokenizer(inputCol="review_body", outputCol="words")

# Creamos objeto StopWordsRemover
remover = StopWordsRemover(inputCol="words", outputCol="filtered_words")
# Datos
df_tokenized = tokenizer.transform(train_data)
df_clean = remover.transform(df_tokenized)
df_clean.select("review_body", "words", "filtered_words").show()

# Vectores de frecuencia
hashingTF = HashingTF(inputCol="filtered_words", outputCol="rawFeatures", numFeatures=8000)
idf = IDF(inputCol="rawFeatures", outputCol="features")

df_features = hashingTF.transform(df_clean)
df_rescaled = idf.fit(df_features).transform(df_features)

# Clasificación binaria y creación del pipeline
classifier = LogisticRegression(featuresCol="features", labelCol="sentiment")
pipeline = Pipeline(stages=[tokenizer, remover, hashingTF, idf, classifier])
```

```
# Entrenamiento del modelo
model = pipeline.fit(train_data)

# Realizamos Las precciones
predictions = model.transform(test_data)
predictions.select("sentiment", "prediction").show()

# Crear un evaluador de clasificación
evaluator = MulticlassClassificationEvaluator(labelCol="sentiment", predictionCol="prediction", metricName="accuracy")

# Calcular la precisión del modelo
accuracy = evaluator.evaluate(predictions)
print("Accuracy:", accuracy)

# Calcular la recuperación (recall)
evaluator.setMetricName("weightedRecall")
recall = evaluator.evaluate(predictions)
print("Recall:", recall)

# Calcular la puntuación F1
evaluator.setMetricName("f1")
f1 = evaluator.evaluate(predictions)
print("F1 Score:", f1)
```

VBox()

FloatProgress(value=0.0, bar\_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...

review_body	words	filtered_words
Very satisfied bu...	[very, satisfied,...]	[satisfied, beat,...]
even though they ...	[even, though, th...]	[even, though, st...]
Es ist echt cool ...	[es, ist, echt, c...]	[es, ist, echt, c...]
verbindet bei mir...	[verbindet, bei, ...]	[verbindet, bei, ...]
Leider nach einem...	[leider, nach, ei...]	[leider, nach, ei...]
keinerlei Problem...	[keinerlei, probl...]	[keinerlei, probl...]
Hallo ! Absolute...	[hallo, !, , abso...]	[hallo, !, , abso...]
fits nice, design...	[fits, nice,, des...]	[fits, nice,, des...]
Das ist auf jeden...	[das, ist, auf, j...]	[das, ist, auf, j...]
Hab die Kopfh&oum...	[hab, die, kopfh&...]	[hab, die, kopfh&...]
Um es kurz zu mac...	[um, es, kurz, zu...]	[um, es, kurz, zu...]
Also mein Receive...	[also, mein, rece...]	[also, mein, rece...]
Lädt mit 10000mAh...	[lädt, mit, 10000...]	[lädt, mit, 10000...]
After a year of m...	[after, a, year, ...]	[year, moderate, ...]
Hab ihn erst seit...	[hab, ihn, erst, ...]	[hab, ihn, erst, ...]
siehe überschrift...	[siehe, überschri...]	[siehe, überschri...]
Das HDMI Kabel ka...	[das, hdmi, kabel...]	[das, hdmi, kabel...]
Lieferung super s...	[lieferung, super...]	[lieferung, super...]
Die Kamera Liefer...	[die, kamera, lie...]	[die, kamera, lie...]
Ich habe das Amaz...	[ich, habe, das, ...]	[ich, habe, das, ...]

only showing top 20 rows

sentiment	prediction
1	1.0
1	1.0
1	1.0
1	1.0
1	1.0
1	1.0
1	1.0
1	1.0
0	0.0
1	1.0
1	1.0
1	1.0
1	1.0
0	1.0
1	1.0
1	0.0
1	1.0
1	1.0
1	1.0
1	1.0

only showing top 20 rows

Accuracy: 0.8959088643584634  
Recall: 0.8959088643584635  
F1 Score: 0.8881667544596061

## Tarea 9: Serialización

Guarda el modelo entrenado en S3 (al bucket que creaste anteriormente) utilizando la opción nativa de Spark.

```
In [15]: # Solución
# Guardo el modelo en mi S3
model.save("s3://capstone12bucket2/trained_model")

VBox()
FloatProgress(value=0.0, bar_style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
```

```
In [ ]:
```