

Spark RDD(16/11/2018)

Carga de datos

In []:	<pre>lista=[3,6,2,7,9] listaRDD=sc.parallelize(lista)</pre>	Conjunto de datos
In []:	<pre>ventas=[("moto",2000), ("coche",10000), ("coche",12000), ("moto",3000), ("moto",2000), ("bici",200)] ventasRDD=sc.parallelize(ventas)</pre>	Similar a tabla. RDD de tuplas donde cada tupla representa una fila y cada elemento de la tupla se considera una columna en función de la posición.

Datos externos

In []:	<pre>lineas=sc.textFile("usuarios.csv")</pre>	RDD de líneas del archivo
In []:	<pre>archivos=sc.wholeTextFiles("/carpeta")</pre>	RDD de pares nombre y contenido

Recuperando información

In []:	<pre>listaRDD.collect()</pre>	Devuelve el conjunto como lista.
Out []:	<pre>[3,6,2,7,9]</pre>	
In []:	<pre>ventasRDD.first()</pre>	Primer elemento del RDD
Out []:	<pre>("moto",2000)</pre>	
In []:	<pre>listaRDD.take(3)</pre>	3 primeros elementos como lista
Out []:	<pre>[3,6,2]</pre>	
In []:	<pre>ventasRDD.count()</pre>	Número de elementos del RDD
Out []:	<pre>6</pre>	
In []:	<pre>listaRDD.sum()</pre>	Suma los elementos del RDD
Out []:	<pre>27</pre>	
In []:	<pre>listaRDD.max()</pre>	Máximo del RDD
Out []:	<pre>9</pre>	
In []:	<pre>listaRDD.min()</pre>	Mínimo del RDD
Out []:	<pre>2</pre>	
In []:	<pre>listaRDD.mean()</pre>	Media del RDD
Out []:	<pre>5.4</pre>	
In []:	<pre>listaRDD.stdev()</pre>	Máximo del RDD
Out []:	<pre>2.5768197453450252</pre>	
In []:	<pre>listaRDD.stats()</pre>	Conjunto de valores estadísticos
Out []:	<pre>(count: 5, mean: 5.4, stdev: 2.5768197453450252, max: 9.0, min: 2.0)</pre>	
In []:	<pre>listaRDD.variance()</pre>	Varianza del RDD
Out []:	<pre>6.640000000000001</pre>	

Transformaciones del RDD

In []:	<pre>doblerRDD=listaRDD.map(lambda a:a*2) doblerRDD.collect()</pre>	Devuelve un RDD transformado según la función lambda
Out []:	<pre>[6,12,4,14,18]</pre>	
In []:	<pre>filtroRDD=listaRDD.filter(lambda a:a>2) filtroRDD.collect()</pre>	Devuelve un RDD filtrado según la condición de la función lambda

Out []:	<pre>[3,6,7,9]</pre>	
In []:	<pre>txt=sc.parallelize(["Hola a todos"])</pre>	
In []:	<pre>palabras=txt.flatMap(lambda a:a.split(" ")) palabras.collect()</pre>	Devuelve un RDD en el que cada elemento de la lista forma parte del nuevo RDD
Out []:	<pre>["Hola","a","todos"]</pre>	

Claves y valores

In []:	<pre>claves = ventasRDD.keys() claves.collect()</pre>	Devuelve la primera columna (claves) como RDD
Out []:	<pre>["moto","coche","coche","moto","moto","bici"]</pre>	
In []:	<pre>valores = ventasRDD.values() valores.collect()</pre>	Devuelve la segunda columna (valores) como RDD
Out []:	<pre>[2000,10000,12000,3000,2000,200]</pre>	
In []:	<pre>nColumna=ventasRDD.map(lambda a:a[1])</pre>	Devuelve la columna indicada como RDD (empezando en 0)
In []:	<pre>listaRDD.reduce(lambda a,b:a+b)</pre>	Realiza la operación indicada en lambda para todos los elementos
Out []:	<pre>27</pre>	
In []:	<pre>res1=ventasRDD.reduceByKey(lambda a,b:a+b) res1.collect()</pre>	Agrupar por clave y realiza la operación indicada en lambda con los valores
Out []:	<pre>[('moto', 7000), ('coche', 22000), ('bici', 200)]</pre>	
In []:	<pre>porClave=res1.sortByKey() porClave.collect()</pre>	Ordena por clave
Out []:	<pre>[('bici', 200), ('coche', 22000), ('moto', 7000)]</pre>	
In []:	<pre>porValor=res1.sortBy(lambda a:-a[1]) porValor.collect()</pre>	Ordena por valor, de mayor a menor
Out []:	<pre>[('coche', 22000), ('moto', 7000), ('bici', 200)]</pre>	

Python. Librerías útiles con RDD

In []:	<pre>from datetime import datetime fecha=datetime.strptime('2/15/1999','%m/%d/%Y') fecha</pre>	Convierte una cadena de texto a objeto fecha
Out []:	<pre>datetime.datetime(1999, 2, 15, 0, 0)</pre>	
In []:	<pre>fecha.weekday()</pre>	Día de la semana(0=lunes,...)
Out []:	<pre>0</pre>	
In []:	<pre>import json texto='{ "Nombre": "Pedro", "Edad": 25, "Aficiones": ["Cine", "Lectura"] }' datos=json.loads(texto) datos.get("Nombre"),datos.get("Aficiones")[1]</pre>	Interpreta json y convierte a diccionario. Con datos.items() se obtendría una lista de tuplas con los valores del diccionario.
Out []:	<pre>('Pedro', 'Lectura')</pre>	
In []:	<pre>import re patronDNI=re.compile(r"(\d{8})([A-Z])") texto="El DNI de Felipe es 45678912W y \ el de María es 56432135A" patronDNI.findall(texto)</pre>	Expresiones regulares para encontrar información
Out []:	<pre>[('45678912', 'W'), ('56432135', 'A')]</pre>	