

Текстовая расшифровка видео:

RESILIENT DISTRIBUTED DATASETS

План:

- Базовое API;
- Что делать с RDD;
- Простейшие трансформации;
- Распределенка;
- Как все сломать;
- Стадии вычислений;
- Действия (action);
- Частичная сортировка и «разравнивание»;
- Промежуточные итоги.

Базовое API

У нас есть данные, которые мы закидываем на кластер. Эти данные нарезаются на «кусочки» – Partition'ы. На каждый Partition попадает порция из этих изначальных данных. Когда мы запускаем распределенные вычисления, каждый из этих кусочков обрабатывается независимо (если получается; если не получается, то можно использовать горизонтальные пересылки данных).

RDD – это датасет, «размазанный» по кластеру в виде Partition'ов.

Рассмотрим код:



▶ Запустить стенд



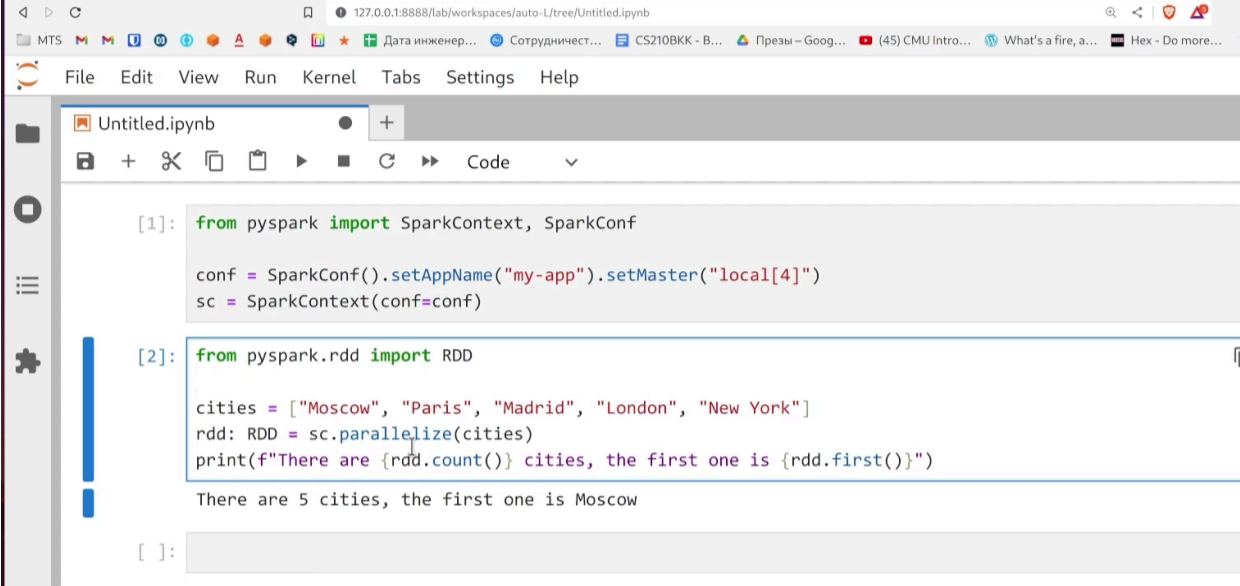
Дедлайн 07 июля, 23:59 Мск



```
1. from pyspark.rdd import RDD
2.
3. cities = ["Moscow", "Paris", "Madrid", "London", "New York"]
4. rdd: RDD = sc.parallelize(cities)
5. print(f"There are {rdd.count()} cities, the first one is {rdd.first()}")
```

Здесь специально добавлен импорт RDD для того, чтобы показать из каких вызовов возвращается объект типа RDD. Это важно, так как иногда может возвращаться RDD при запуске кода, а может возвращаться нечто иное (это довольно опасная ситуация).

Ранее мы создали notebook и запустили контекст. Теперь мы можем скопировать код, вставить в notebook и запустить:



```
[1]: from pyspark import SparkContext, SparkConf
conf = SparkConf().setAppName("my-app").setMaster("local[4]")
sc = SparkContext(conf=conf)

[2]: from pyspark.rdd import RDD

cities = ["Moscow", "Paris", "Madrid", "London", "New York"]
rdd: RDD = sc.parallelize(cities)
print(f"There are {rdd.count()} cities, the first one is {rdd.first()}")

There are 5 cities, the first one is Moscow

[ ]:
```

Основной метод, который эту коллекцию превращает в RDD, – **sc.parallelize**.

SC – это SparkContext.

Когда вызываем `sc.parallelize`, мы получаем RDD. Поверх RDD можем вызывать разные методы для того, чтобы что-либо посчитать. Например, мы можем посмотреть, сколько всего объектов в этом RDD, достать первый объект и тд.

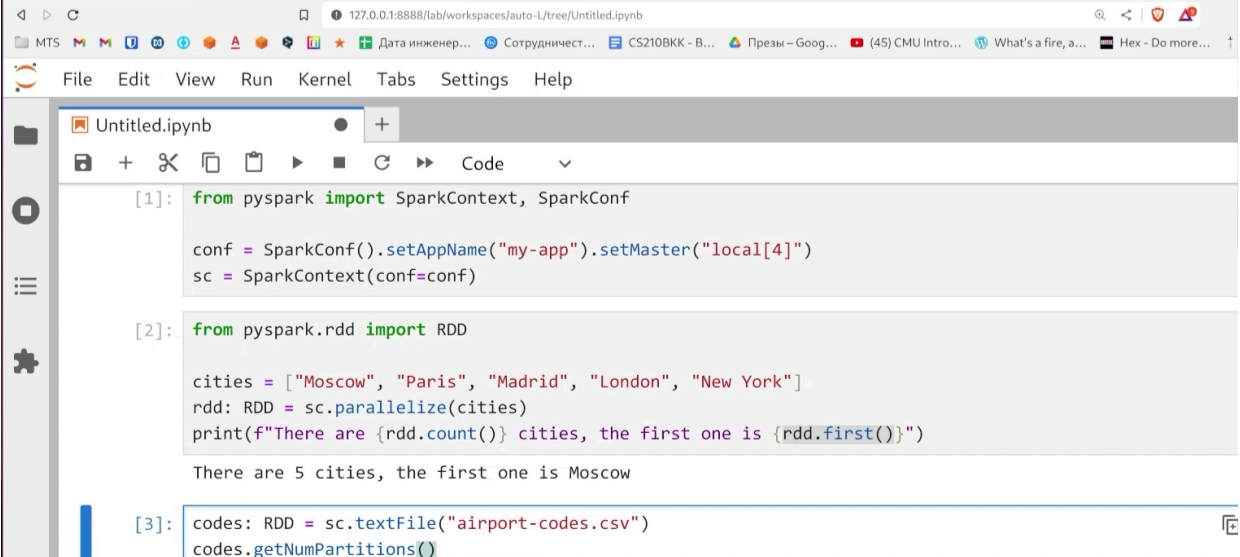
Если мы хотим загрузить данные не из стандартного контейнера Python, а из файла, например, HDFS, то можно использовать путь к файлу (можно его сначала прочитать, а потом «скормить»):

```
1. codes: RDD = sc.textFile("/tmp/data/airport-codes.csv")
2. codes.getNumPartitions()
3. # 2, может быть другим у вас
```

Можно использовать **sc.textFile**. Он прочитает нужный файл, а потом превратит его в RDD.

В качестве примера будем использовать тот же файл, который мы использовали на прошлых занятиях, а именно – коды аэропортов.

Метод **getNumPartitions** покажет на сколько «кусочков» разделится файл:



```
[1]: from pyspark import SparkContext, SparkConf
conf = SparkConf().setAppName("my-app").setMaster("local[4]")
sc = SparkContext(conf=conf)

[2]: from pyspark.rdd import RDD

cities = ["Moscow", "Paris", "Madrid", "London", "New York"]
rdd: RDD = sc.parallelize(cities)
print(f"There are {rdd.count()} cities, the first one is {rdd.first()}")

There are 5 cities, the first one is Moscow

[3]: codes: RDD = sc.textFile("airport-codes.csv")
codes.getNumPartitions()

2
```

Видим, что весь файл разбит на два Partition'a.

У нас есть переменная RDD с типом RDD из нашей распределенной коллекции, и есть codes, которые тоже RDD, но другого типа. Мы добавили данные в кластер. Что мы можем с ними сделать?

Что делать с RDD?

Есть два термина:

- Трансформация (Transformations);
- Действия (Actions).

Хитрость RDD:

Слово «Resilient» в аббревиатуре переводится как «устойчивый», «неизменный». Вспомним, что Spark написан на Scala с применением функционального подхода. Как мы знаем, все объекты неизменны. Здесь во многом используется тот же подход.

Трансформация – это набор методов, который мы можем применить к объекту RDD. Результатом выполнения мы не поменяем существующий RDD, а получим новый, который будет представлять из себя «ленивую» обертку.

У нас есть **способ вычисления**: для того, чтобы получить новый RDD на основе предыдущего, нужно к предыдущему применить какое-то преобразование.

Трансформации не запускают никакие дополнительные вычисления на кластере. Яркий пример: мы хотим привести все слова к верхнему регистру. Мы добавляем трансформацию, однако пока не будет запроса на развертывание графа вычислений, никаких реальных вычислений не будет.

Кратко о Трансформации: обычно ничего не считают, просто добавляют в граф вычислений. Результат – новый RDD.

Действия (Actions) используются для развертывания графа. Это другой набор вызовов, приводящий к тому, что цепочка трансформаций, написанная нами, раскручивается и схлопывается к какому-то результату.

Нюанс: это может быть опасно, так как при выполнении действия данные, которые мы преобразовали, посчитали и т.д. пересылаются на нашу машину, с которой мы запускаем обработку. Если этих данных много, то они могут упасть. **Принимайте решения взвешенно и осознавайте возможные последствия.**

Кратко о Действиях: раскручивают граф вычислений и запускают их распределённо. Результат зависит от того, что попросили.

Простейшие трансформации

Рассмотрим на примере:

```
def plus_one(i):  
    return i + 1  
  
plus_one = lambda i: i + 1
```

У нас есть RDD, к которому мы применяем примитивную функцию с помощью **map()/filter()**.

Map()/Filter() – это уже готовые и уже имеющиеся в стандартной библиотеке Python методы. С «**Map()**» мы «лениво» применяем одну и ту же функцию ко всем элементам. С помощью «**Filter()**» мы выкидываем те элементы, для которых предикат возвращает false, оставляя те, для которых он

Поработаем с данным кусочком кода:

```
filtered: RDD = rdd.filter(lambda city: city.startswith("M"))
```

У нас есть города. Мы можем применить фильтр и отфильтровать те города, названия которых начинаются на букву «М». Обратите внимание, что вычисления не запустились:

```
[6]: filtered: RDD = rdd.filter(lambda city: city.startswith("M"))
[7]: filtered
[7]: PythonRDD[5] at RDD at PythonRDD.scala:53
[ ]:
```

Переменная «Filtered» все еще представляет из себя просто RDD и для того, чтобы получить результат, нужно вызвать action. Забегая наперед скажем, что можно вызвать collect (самый базовый, но и самый опасный action, вызывающий раскручивание всей цепочки и возврат всего результата). Если мы вызовем collect, у нас вернутся два города, начинающихся на букву «М»:

```
[6]: filtered: RDD = rdd.filter(lambda city: city.startswith("M"))
[8]: filtered.collect()
[8]: ['Moscow', 'Madrid']
[ ]:
```

Помимо базовых операций map()/filter() есть большое количество других, которые тоже представляют собой трансформации, то есть никакие вычисления на кластере по умолчанию не запускают. Например:

- **Distinct()**, когда мы выкидываем все дубликаты;
- **Sample()**, когда мы берем по определенной логике из всех элементов какую-то подгруппу;
- **Операции по работе с сетями** (с множествами), когда мы считаем объединение множеств, пересечение множеств и т.д.:

```
union()          subtract()
intersection()   cartesian()
```

Распределенка

Если мы хотим в примере привести города к верхнему регистру, то в коде мы можем сделать следующее:

```
[5]: airport-codes.csv MapPartitionsRDD[4] at textFile at NativeMethodAccessorImpl.java:0
[6]: filtered: RDD = rdd.filter(lambda city: city.startswith("M"))
[8]: filtered.collect()
[8]: ['Moscow', 'Madrid']
[9]: rdd.map(lambda city: city.upper()).collect()
[9]: ['MOSCOW', 'PARIS', 'MADRID', 'LONDON', 'NEW YORK']
[ ]:
```

Когда мы вызвали parallels, данные по кластеру уже были «размазаны». Мы предположили, что три города из пяти попали в один partition, два других – в другой:


```
1. # Трансформация map: не запускает вычислений,
2. # не изменяет изначальный RDD
3. upper_rdd: RDD[str] = rdd.map(lambda city: city.upper())
4. # Метод take возвращает N первых элементов RDD
5. upper_lst: list = upper_rdd.take(3)
6. old_lst: list = rdd.take(3)
7. print(f"New RDD: {' '.join(upper_lst)}")
8. # New RDD: MOSCOW, PARIS, MADRID
9. print(f"Old RDD: {' '.join(old_lst)}")
10. # Old RDD: Moscow, Paris, Madrid
```

Это иллюстрация в коде, где мы приходим к верхнему регистру и берем три элемента в новом RDD, который уже «лениво» посчитан. Нам ничего не мешает достать старый RDD, спросить у него. Там все также будут неизменные элементы.

Нюанс: стадии по умолчанию будут пересчитываться.

Действия (actions)

Когда мы хотим получить какой-либо результат, мы хотим схлопнуть RDD к какому-то конкретному значению и переместить данные на основную ноду, на которой мы все запускаем.

Одним из типичных actions является **reduce()**, который вынесли в functools:

```
from functools import reduce
reduce(lambda a, b: a + b, [21, 7, 14])
```

Он берет элемент и аккумулятор по умолчанию и применяет функцию до тех пор, пока не схлопнет весь массив. Каждый раз **в первый аргумент попадает следующий элемент из изначального контейнера, а во второй – результат применения предыдущей операции.**

В Spark мы можем использовать action reduce(). Обратите внимание, что **он возвращает уже не RDD, а int:**

```
1. # Действие reduce применяет функцию f к промежуточному результату
2. # от предыдущей итерации со следующим элементом коллекции
3. count: int = rdd.map(lambda x: len(x)).reduce(lambda a, b: a + b)
4. print(f"The RDD contains {count} letters")
5. # The RDD contains 31 letters
```

Таким образом мы пройдемся по всем городам, у каждого посчитаем длину и посмотрим общее количество символов в датасете:

```
[12]: count: int = rdd.map(lambda x: len(x)).reduce(lambda a, b: a + b)
      print(f"The RDD contains {count} letters")
```

```
The RDD contains 31 letters
```

Action'ов не так много, но все они достаточно простые, так, например, take() берет один или несколько элементов, стараясь не раскручивать по максимуму граф вычисления для остальных.

Count() – это внезапный подсчет количества элементов. Это тоже action, так как в случае с RDD мы можем не знать точное количество записей.

Top() – это action, позволяющий доставать элемент по функции сравнения.

Частичная сортировка и «разравнивание»

TakeOrdered – комбинация из sort и take:

```
1. sc.parallelize(  
2.   [10, 1, 2, 9, 3, 4, 5, 6, 7]  
3. ).takeOrdered(6, key=lambda x: -x)  
4. # [10, 9, 7, 6, 5, 4]
```

Рассмотрим пример:

```
1. mapped_rdd = rdd.map(lambda x: list(x))  
2. print(mapped_rdd.take(2))  
3. # [['M', 'o', 's', 'c', 'o', 'w'], ['P', 'a', 'r', 'i', 's']]  
4. flatmapped_rdd = rdd.flatMap(lambda x: x.lower())  
5. print(flatmapped_rdd.take(4))  
6. # ['m', 'o', 's', 'c']
```

Представим, что каждый элемент RDD – список. Нас интересует возможность пройтись по всем элементам, распаковав все списки. Это позволяет сделать специальная версия `map`'а – **flatMap**, он разравнивает промежуточные результаты.

Мы берем датасет с городами, каждый город мы разбиваем по буквам. Мы можем захотеть составить RDD, которая будет представлять из себя коллекцию букв в нижнем регистре. Мы применяем `flatMap`, вызываем метод приведения к нижнему регистру и получаем набор RDD, где каждый элемент – буква.

Промежуточные итоги

[RDD – неизменяемый \(иммутабельный\) распределенный набор данных.](#)

Трансформации «лениво» создают новый RDD и не изменяют существующий.

Только **действия** приводят к запуску реальных вычислений.

Как вам урок?



Изучил, далее >

Слёрм ©

[+7 \(495\) 248-05-80](#)

[Лицензия №ДЛ-1368 от 22.08.2019](#)

[Политика конфиденциальности](#)

[Публичная оферта](#)