

Implementacja operacji relacyjnych w Apache Spark z wykorzystaniem MapReduce

Streszczenie

W niniejszym dokumencie przedstawiono implementację podstawowych operacji relacyjnych w środowisku Apache Spark z wykorzystaniem modelu MapReduce. Analizie poddano zbiory danych z plików `relation.txt` i `join.txt`. Zaprezentowano procesy przetwarzania i wyniki operacji takich jak suma zbiorów, część wspólna oraz złączenia relacji, wzbogacone o odpowiednie fragmenty kodu w języku Python.

Spis treści

1	Wstęp	1
2	Operacje na zbiorach z pliku <code>relation.txt</code>	1
2.1	Tworzenie zbiorów A i B	2
2.2	Implementacja w Apache Spark	2
2.3	Wyniki operacji	3
3	Złączenie relacji z pliku <code>join.txt</code>	3
3.1	Opis problemu	3
3.2	Dane wejściowe	3
3.3	Implementacja w Apache Spark	4
3.4	Wyniki złączenia	5
3.5	Analiza wyników	5
4	Wnioski	5

1 Wstęp

Przetwarzanie dużych zbiorów danych wymaga wydajnych i skalowalnych narzędzi. Apache Spark, oparty na modelu MapReduce, umożliwia efektywne wykonywanie operacji na danych rozproszonych. W niniejszym dokumencie skupimy się na implementacji podstawowych operacji relacyjnych z użyciem Spark, ilustrując procesy na przykładach.

2 Operacje na zbiorach z pliku `relation.txt`

Plik `relation.txt` zawiera pary składające się z nazwy relacji (A lub B) oraz wartości liczbowej. Celem jest wykonanie podstawowych operacji na zbiorach:

- Suma zbiorów $A \cup B$
- Część wspólna zbiorów $A \cap B$
- Różnica zbiorów $A - B$ oraz $B - A$

2.1 Tworzenie zbiorów A i B

Z pliku `relation.txt` wyodrębniamy wartości przypisane do relacji A i B :

$$A = \{0, 1, 4, 8, 9, 10, 17, \dots, 98\}$$

$$B = \{0, 2, 3, 5, 8, 9, 10, \dots, 99\}$$

2.2 Implementacja w Apache Spark

Poniżej przedstawiono implementację powyższych operacji z wykorzystaniem języka Python i biblioteki PySpark.

```
1 from pyspark import SparkContext
2
3 sc = SparkContext.getOrCreate()
4 data = sc.textFile('relation.txt')
5
6 # Mapowanie każdej linii na parę (RelationName, Value)
7 lines = data.map(lambda line: line.strip().split())
8
9 # Utworzenie RDD dla relacji A i B
10 rddA = lines.filter(lambda x: x[0] == 'A').map(lambda x: int(x[1]))
11 rddB = lines.filter(lambda x: x[0] == 'B').map(lambda x: int(x[1]))
12
13 # Mapowanie elementów w A i B do par (Value, RelationName)
14 mapA = rddA.map(lambda x: (x, 'A'))
15 mapB = rddB.map(lambda x: (x, 'B'))
16
17 # Połączenie zmapowanych RDD
18 unionMap = mapA.union(mapB)
19
20 # Grupowanie po kluczu (Value)
21 grouped = unionMap.groupByKey().mapValues(list)
22
23 # Operacje na zbiorach
24 unionValues = grouped.map(lambda x: x[0]) # Suma zbiorów A i B
25 intersectionValues = grouped.filter(lambda x: 'A' in x[1] and 'B' in
26     x[1]).map(lambda x: x[0]) # A - B
27 differenceAB = grouped.filter(lambda x: 'A' in x[1] and 'B' not in
28     x[1]).map(lambda x: x[0]) # A - B
29 differenceBA = grouped.filter(lambda x: 'B' in x[1] and 'A' not in
30     x[1]).map(lambda x: x[0]) # B - A
31
32 # Zbieranie i wyświetlanie wyników
33 print("Suma zbiorów A i B:")
34 print(unionValues.collect())
35
36 print("\nCzęść wspólna A i B:")
37 print(intersectionValues.collect())
38
39 print("\nRóżnica A - B:")
40 print(differenceAB.collect())
41
42 print("\nRóżnica B - A:")
43 print(differenceBA.collect())
```

Listing 1: Przetwarzanie pliku `relation.txt`

2.3 Wyniki operacji

Na podstawie powyższego kodu otrzymujemy następujące wyniki:

- Suma zbiorów $A \cup B$:

$$A \cup B = \{0, 1, 2, 3, 4, 5, \dots, 99\}$$

- Część wspólna zbiorów $A \cap B$:

$$A \cap B = \{0, 8, 9, 10, 17, \dots, 98\}$$

- Różnica zbiorów $A - B$:

$$A - B = \{1, 4, 18, \dots, 96\}$$

- Różnica zbiorów $B - A$:

$$B - A = \{2, 3, 5, \dots, 99\}$$

3 Złączenie relacji z pliku join.txt

3.1 Opis problemu

Celem jest wykonanie złączenia (*join*) między relacjami **Customers** i **Orders** na podstawie wspólnego atrybutu **CustomerID**, wykorzystując model MapReduce.

3.2 Dane wejściowe

Klienci (**Customers**):

CustomerID	Informacje o kliencie
1	CompanyName: Alfreds Futterkiste, ContactName: Maria Anders, Country: Germany
2	CompanyName: Ana Trujillo Emparedados y helados, ContactName: Ana Trujillo, Country: Mexico
3	CompanyName: Antonio Moreno Taquería, ContactName: Antonio Moreno, Country: Mexico
4	CompanyName: NO ORDERS, ContactName: NO ORDERS, Country: NO ORDERS

Zamówienia (**Orders**):

OrderID, CustomerID	OrderDate
10308, 1	'2016-09-18'
10309, 2	'2016-09-19'
10310, 3	'2016-09-20'
10311, 3	'2016-09-20'
10312, 3	'2016-09-20'
10313, 5	'2016-09-20'

3.3 Implementacja w Apache Spark

```
1 # Wczytanie danych z pliku join.txt
2 order_data = sc.textFile('join.txt')
3 # Podzielenie każdej linii na pola
4 order_lines = order_data.map(lambda line: line.strip().split('\t'))
5 # Funkcja przetwarzająca każdą linię i mapująca ją na parę (CustomerID, Record)
6 def process_order_line(line):
7     if line[0] == 'Orders':
8         # Mapowanie Orders na (CustomerID, ('Orders', OrderID, OrderDate))
9         return (line[2], ('Orders', line[1], line[3]))
10    elif line[0] == 'Customers':
11        # Mapowanie Customers na (CustomerID, ('Customers', CompanyName,
12            ContactName, Country))
13        return (line[1], ('Customers', line[2], line[3], line[4]))
14    else:
15        return None
16
17 # Mapowanie linii i filtrowanie wartości None
18 order_mapped = order_lines.map(process_order_line).filter(lambda x: x is not None)
19 # Grupowanie po CustomerID
20 order_grouped = order_mapped.groupByKey().mapValues(list)
21 # Funkcja łącząca informacje o kliencie z zamówieniami
22 def join_orders(customers_orders):
23     customer_id, records = customers_orders
24     customer_info = None
25     orders = []
26     for record in records:
27         if record[0] == 'Customers':
28             customer_info = record # ('Customers', CompanyName, ContactName,
29                 Country)
30         elif record[0] == 'Orders':
31             orders.append(record) # ('Orders', OrderID, OrderDate)
32     # Jeśli są dostępne informacje o kliencie, połącz z każdym zamówieniem
33     if customer_info:
34         return [
35             (customer_id, customer_info[1:], order[1:])
36             for order in orders
37         ]
38     else:
39         # Brak informacji o kliencie, nie można wykonać złączenia
40         return []
41
42 # Wykonanie operacji złączenia
43 joined_orders = order_grouped.flatMap(join_orders)
44 # Zbieranie i wyświetlanie wyników
45 print("\nPołączone zamówienia i klienci:")
46 for record in joined_orders.collect():
47     customer_id = record[0]
48     company_name, contact_name, country = record[1]
49     order_id, order_date = record[2]
50     print(f"CustomerID: {customer_id}, CompanyName: {company_name}, "
51         f"ContactName: {contact_name}, Country: {country}, "
52         f"OrderID: {order_id}, OrderDate: {order_date}")
```

Listing 2: Złączenie relacji z pliku join.txt

3.4 Wyniki złączenia

W wyniku wykonania powyższego kodu otrzymujemy następujące połączone rekordy:

- **CustomerID:** 1
 - **CompanyName:** Alfreds Futterkiste
 - **ContactName:** Maria Anders
 - **Country:** Germany
 - **OrderID:** 10308
 - **OrderDate:** 2016-09-18
- **CustomerID:** 2
 - **CompanyName:** Ana Trujillo Emparedados y helados
 - **ContactName:** Ana Trujillo
 - **Country:** Mexico
 - **OrderID:** 10309
 - **OrderDate:** 2016-09-19
- **CustomerID:** 3 (trzykrotnie dla różnych zamówień)
 - **CompanyName:** Antonio Moreno Taquería
 - **ContactName:** Antonio Moreno
 - **Country:** Mexico
 - **OrderID:** 10310, 10311, 10312
 - **OrderDate:** 2016-09-20

3.5 Analiza wyników

- Klient o **CustomerID** = 4 nie posiada zamówień, w związku z czym nie został uwzględniony w wynikach złączenia.
- Zamówienie o **OrderID** = 10313 z **CustomerID** = 5 nie zostało uwzględnione, ponieważ brak jest informacji o kliencie o **CustomerID** = 5 w relacji **Customers**.
- Operacja złączenia działa jako *inner join*, łącząc tylko te rekordy, dla których klucz **CustomerID** występuje w obu relacjach.

4 Wnioski

Zaprezentowane implementacje podstawowych operacji relacyjnych w środowisku Apache Spark pokazują efektywność modelu MapReduce w przetwarzaniu dużych zbiorów danych. Wykorzystanie RDD (Resilient Distributed Datasets) umożliwia łatwe i skalowalne operacje na danych, co jest kluczowe w analizie big data.

Literatura

- [1] Apache Spark – <https://spark.apache.org/>
- [2] J. Dean and S. Ghemawat, *MapReduce: Simplified Data Processing on Large Clusters*, Communications of the ACM, 51(1):107-113, 2008.