E Correction
■ Moje kursy Moje

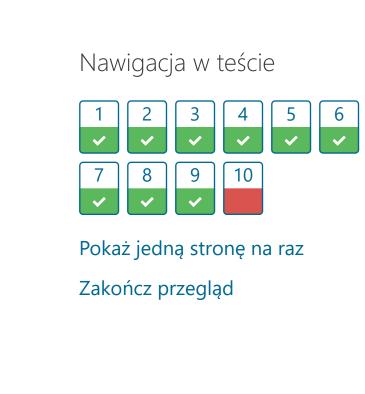
PO12: Przetwarzanie Big Data 22/23-Z

Strona główna / Moje kursy / WliT / Informatyka / Stacjonarne / I stopień / Semestr 7 [WliT-Inf-st-I] / PO12: Przetwarzanie Big Data 22/23-Z / 10. Spark - Delta Lake / Spark - Delta Lake - zadania - zarejestruj swoje rozwiązania tutaj

```
Rozpoczęto środa, 21 grudnia 2022, 20:13
            Stan Ukończone
     Ukończono środa, 21 grudnia 2022, 20:37
   Wykorzystany 23 min. 28 sek.
            czas
          Punkty 50,00/60,00
          Ocena 8,33 pkt. na 10,00 pkt. możliwych do uzyskania (83,33%)
                                                                                                                                                                                                       Pytanie 1
Zadanie 1
                                                                                                                                                                                                       Poprawnie
                                                                                                                                                                                                       Punkty: 4,00 z 4,00
Utwórz pustą tabelę Delta Lake o nazwie customers, której lokalizacją będzie /tmp/delta-customers.

♥ Oflaguj pytanie

Kolumny tabeli muszą odpowiadać kolumnom danych źródłowych.
Wszystkie kolumny powinny być ciągami znaków o długości do 200 znaków
DROP TABLE IF EXISTS ✓ customers;
CREATE TABLE customers
 id VARCHAR(200),
 name VARCHAR(200),
  address VARCHAR(200),
  zipcode VARCHAR(200),
 city VARCHAR(200),
  country VARCHAR(200),
  effectiveDate VARCHAR(200)
                    '/tmp/delta-customers';
```



USING DELTA LOCATION Pytanie **2** Zadanie 2 Poprawnie Punkty: 4,00 z 4,00 Wprowadź do utworzonej przez Ciebie tabeli dane o naszych klientach obowiązujące na dzień 2021-01-01. ♥ Oflaguj pytanie ✓ (id, name, address, zipcode, city, country, effectiveDate) INSERT INTO customers ✓ , city, country, effectiveDate SELECT ✓ id, name, address, zipcode from (select id, name, address, zipcode, city, country, effectiveDate, rank() over (partition by id order by to_date(effectiveDate,"dd-MM-yyyy") desc) as version source_data to_date(effectiveDate,"dd-MM-yyyy") < date "2021-01-01"</pre> where version = 1; Pytanie **3** Zadanie 3 Poprawnie Zmień kraj zamieszkania na wartość Poland u wszystkich klientów posiadających wartość id mniejszą niż 50. Punkty: 4,00 z 4,00 Zwróć uwagę, że obecnie id jest ciągiem znaków. Użyj wyrażenia cast(id as int) aby wykonać zadanie prawidłowo. ♥ Oflaguj pytanie update customers = 'Poland' set country cast(id as int) <</pre> Pytanie **4** Zadanie 4 Poprawnie Usuń z tabeli customers tych klientów, którzy zmienili swoje dane w styczniu 2021. Identyfikacja klientów odbywać się powinna za każdym razem w oparciu o atrybut id. Punkty: 4,00 z 4,00 ♥ Oflaguj pytanie Jeśli okaże się, że polecenie DELETE nie wspiera podzapytań, skorzystaj z interfejsu w Scali, lub za pomocą poniższego kodu, w dodatkowym paragrafie uzyskaj identyfikatory klientów do usunięcia, a następnie wkomponuj uzyskaną wartość w polecenie usuwające klientów. val usunIDS = spark.sql("""select id as usun from (select id, name, address, zipcode, city, country, effectiveDate, rank() over (partition by id order by to_date(effectiveDate,"dd-MM-yyyy") desc) as version from source_data where to_date(effectiveDate, "dd-MM-yyyy") >= date "2021-01-01" to_date(effectiveDate,"dd-MM-yyyy") < date "2021-02-01"</pre> where version = 1""").collect.map(_.toSeq).flatten.mkString(",") ✓ .sql(s""" delete from customers id in (\$usunIDS Pytanie **5** Zadanie 5 Poprawnie Korzystając z jednego polecenia merge jednocześnie Punkty: 10,00 z 10,00 wstaw nowych klientów, którzy pojawili się po raz pierwszy w lutym 2021 • zaktualizuj dane adresowe oraz effectiveDate tych klientów, którzy w tym samym czasie te dane zmienili.

♥ Oflaguj pytanie Poniższe polecenie uzyskuje dane z obu grup klientów: select id, name, address, zipcode, city, country, effectiveDate from (select id, name, address, zipcode, city, country, effectiveDate, rank() over (partition by id order by to_date(effectiveDate, "dd-MM-yyyy") desc) as version source_data to_date(effectiveDate,"dd-MM-yyyy") >= date "2021-02-01" to_date(effectiveDate,"dd-MM-yyyy") < date "2021-03-01"</pre> where version = 1MERGE INTO customers USING (select id, name, address, zipcode, city, country, effectiveDate from (select id, name, address, zipcode, city, country, effectiveDate, rank() over (partition by id order by to_date(effectiveDate, "dd-MM-yyyy") desc) as version from source_data where to_date(effectiveDate, "dd-MM-yyyy") >= date "2021-02-01" to_date(effectiveDate,"dd-MM-yyyy") < date "2021-03-01"</pre> where version = 1) new_customers = new_customers.id customers ✓ THEN MATCHED UPDATE ✓ SET customers.address = new_customers.address, customers.zipcode = new_customers.zipcode, customers.city = new_customers.city, customers.country = new_customers.country, customers.effectiveDate = new_customers.effectiveDate NOT MATCHED INSERT ✓ , zipcode, city, country, effectiveDate) (id, name, address VALUES (id, name, address, zipcode, city, country, effectiveDate)

Kilka pytań

WRITE✓

MERGE

DELETE

UPDATE

WRITE

Tak **♦ ✓**

Zadanie 7

select id,

from

kolumnę do schematu

val data032021 = spark.sql("""

effectiveDate,

source_data

where version = 1""")

data032021

name, address, zipcode, city, country,

select id, name, address, zipcode, city, country, effectiveDate,

✓ .write.option(" mergeSchema

to_date(effectiveDate,"dd-MM-yyyy") >= date "2021-03-01"

to_date(effectiveDate,"dd-MM-yyyy") < date "2021-04-01"</pre>

rank() over (partition by id order by to_date(effectiveDate,"dd-MM-yyyy") desc) as version

✓ ", "true").

cast(null as date) as endDate

✓ ("overwrite")

.option(" overwriteSchema

ADD COLUMNS

CREATE OR REPLACE TABLE AS SELECT

☑CREATE [OR REPLACE] TABLE ✔

Poprawna odpowiedź to:

• CREATE [OR REPLACE] TABLE

Punkty: 1,00 z 1,00

Zwróć uwagę, że niektóre z tych operacji mają flagę isBlindAppend zapaloną. Które to były operacje?

Czy liczba plików utworzonych przy pierwszej operacji (zaobserwowanych w katalogu tabeli) zgadza się z wpisem w uzyskanych powyżej metadanych?

Pytanie **7** Zadanie 6 Poprawnie Okazało się, że ta operacja zmiany kraju na wartość Poland była błędem. Punkty: 10,00 z 10,00 Twoim zadaniem jest przywrócić wartości, które zostały nadpisane przez tą modifikację. Nie naprawiaj danych jeśli pojawiły się późniejsze (po Twoim poleceniu update) aktualizacje adresu (skorzystaj z effectiveDate). ♥ Oflaguj pytanie Świetnie do takiej naprawy może się przydać operacja merge. Tym razem będzie ona miała tylko sekcję when MATCHED THEN. Oczywiście nie korzystaj z danych źródłowych - ich już nie ma. Jest tylko Twoja tabela Delta Lake. Skorzystaj z DataFrame API. import io.delta.tables._ val deltaTable = DeltaTable. forPath ✓ (spark, "/tmp/delta-customers") val poprawneDF = spark.read ✓ ("delta") format .option(" versionAsOf ✓ ("/tmp/delta-customers") . load deltaTable.as("old") ✓ .as("new"), "old.id = new.id and to_date(old.effectiveDate,'dd-MM-yyyy') < date '2021-01-01'")</pre> poprawneDF whenMatched (Map("country" -> " new updateExpr execute

Kolumna ta jest nam potrzebna i będzie wykorzystywana w następnym zadaniu, dlatego skorzystaj z poniższego kodu oraz wskazówki powyżej, aby podczas wstawiania tych danych wprowadzić brakującą

format(" delta mode(" append saveAsTable ✓ ("customers") Pytanie **9** Zadanie 8 Poprawnie Korzystając z tego mechanizmu zmień typy kolumn, o których wspominaliśmy. Punkty: 7,00 z 7,00 ♥ Oflaguj pytanie Przy okazji wycofaj ostatnią aktualizację. Nie była przemyślana. Pojawiły się duplikaty w naszych danych. Sprawdź, której wersji danych potrzebujesz. import io.delta.tables._ val fullHistDF = custDeltaTable.history().select(\$"version",\$"timestamp",\$"operation",\$"isBlindAppend",\$"operationMetrics") fullHistDF.show(false) val custDeltaTable = DeltaTable.forPath(spark, "/tmp/delta-customers") spark.read.format("delta") .option("versionAsOf", "10") .table("customers") ", col(" effectiveDate .withColumn(" effectiveDate ").cast("date")) withColumn ").cast("int")) .write .format("delta")

Zadanie 9

Zadanie 9

Zalmpiementuj funkcje, która na podstawie danych zrożkowych z kolejnego micriozo (parametr funkcji) będzie aktualizowała zawartość tabeli custowers zgodnie a regulami SCD Type 2.

Po zakończonej implementuji sprawdzi jej działanie.

Jeśli uwszasz. że obecne dane w tabeli customers powinny zostać poprawione, dokonaj wcześniej stosownych korekt.

Jeśli uwszasz. że obecne dane w tabeli customers powinny zostać poprawione, dokonaj wcześniej stosownych korekt.

Jeśli uwszasz. że obecne dane w tabeli customers powinny zostać poprawione, dokonaj wcześniej stosownych korekt.

Jeśli uwszasz. że obecne dane w tabeli customers powinny zostać poprawione, dokonaj wcześniej stosownych korekt.

Jeśli uwszasz. że obecne dane w tabeli customers powinny zostać poprawione, dokonaj wcześniej stosownych korekt.

Jeśli triczesz skorzystaj se stromy:

https://docs.delta.io/0.4.0/delta-update.htmli/siowly-changing-data-scd-type-2-operation-into-delta-tables

Pytanie **6**

Poprawnie

Pytanie **8**

Poprawnie

Punkty: 5,00 z 5,00

♥ Oflaguj pytanie

Punkty: 2,00 z 2,00

♥ Oflaguj pytanie