

Oozie

Hadoop Streaming

1. wrzucić mapper.py i reducer.py na hdfs
2. rozwinąć na HUE dropdown menu przy przycisku 'Query', wejść w menu 'Scheduler' i kliknąć 'Workflow'
3. kliknąć w dropdown 'DOCUMENTS' i wybrać 'Actions'
4. Przeciągnąć akcję streaming do workflowu
5. W pole Mapper wpisać 'mapper.py', a w pole Reducer wpisać 'reducer.py'
6. Kliknąć 'Add'
7. Kliknąć w przycisk 'FILES+' i znaleźć plik 'mapper.py', następnie kliknąć jeszcze raz i znaleźć plik 'reducer.py'
8. Kliknąć w prawym górnym rogu akcji w przycisk ustawień
9. Kliknąć w przycisk 'PROPERTIES+' i wpisać w pierwsze pole 'mapred.input.dir', a w drugim podać ścieżkę do katalogu / pliku wejściowego
10. Kliknąć jeszcze raz i dodać property 'mapred.output.dir' z zamiarem na katalog wyjściowy
11. Zapisz workflow klikając save w prawym górnym rogu

Mapreduce

1. wrzucić jar z jobem na hdfs
2. dodać akcję 'Java program' do workflowu
3. w polu 'Jar name' znaleźć jar z jobem na hdfsie
4. w polu 'Main class' wpisać nazwę klasy razem z pakietem
5. Kliknąć 2 razy w 'ARGUMENTS+'. W pierwszym polu wpisać plik/katalog wejściowy, a w drugim wyjściowy
6. Zapisz workflow klikając save w prawym górnym rogu

Sqoop

1. wrzuć 'mysql-connector-java-5.1.46-bin.jar' na hdfs
2. dodaj akcję 'sqoop1' do workflowu
3. w polu 'Sqoop command' wpisz komendę do importu sqoop (bez polecenia sqoop) i kliknij add
4. Kliknij w prawym górnym rogu akcji w przycisk ustawień
5. kliknij na 'ARCHIVES+'
6. znajdź na 'hdfsie mysql-connector-java-5.1.46-bin.jar'
7. Zapisz workflow klikając save w prawym górnym rogu

Hive

1. stwórz skrypt 'query.sql' zawierający komenty sqlowe do wykonania
2. umieść skrypt 'query.sql' na hdfsie
3. dodaj akcję 'HiveServer2' do workflowu
4. podaj namiary na skrypt 'query.sql' i kliknij 'add'
5. Zapisz workflow klikając save w prawym górnym rogu

Koordinator

1. rozwiń dropdown menu przy przycisku 'Query', wejdź w menu 'Scheduler' i kliknij 'Schedule'
2. Kliknij na dropdown 'Choose a workflow...' i wybierz workflow, który chciałbyś zaschedulować.
3. Wybierz jak często job ma się wykonywać. Możesz kliknąć w Options i zaznaczyć 'Advanced syntax', żeby wpisać wyrażenie crona.
4. Wybierz zakres dat w jakim scheduler ma się wykonywać
5. kliknij save, żeby zapisać koordynatora

Dodanie parametru do joba

`${nazwa_parametru}`

Zadania

1. Utwórz workflow z joba streaming sliczającego słowa i joba mapreduce sortującego po liczbie wystąpień
2. Utwórz workflow ściągający dane ze sqoopu, ładujący je do tabeli 'owners' na hivie i wyliczający tabelę wynikową będącą połączeniem tabeli 'transfers' z tabelą 'owners'