

# Spark / Scala Exercise 2

Create a spark application written in scala that processes the [data](#), and stores the entire output as specified below per the requirements and submission guidelines. Pay close attention to performance. Assume that the same code will be run on a much larger input dataset.

## Overview

When processing large amounts of data over time it's important to process only the new data in incremental batches for performance and cost efficiency. This exercise is designed to simulate processing an incremental and full batches.

The input data contains company attributes from many data sources. Each source has a numeric source id. In every input batch, data is loaded from a subset of sources and may only cover part of the companies in the dataset. In other words, each input batch contains parts of the final company records you will produce in the output.

## Requirements

### Project

- Create a Scala SBT project
- Use the Spark **Scala API** and Dataframes/Datasets
  - Do not use Spark SQL with a sql string!
  - **If you write `spark.sql("select ....")` you are doing it wrong!!**
- **This application should be close to production quality, with consideration of performance, ease of use, well organized and clean. Comment as necessary. Remove extra debug code or comment it out to avoid extra processing.**

### Processing

#### Application Parameters

inputPath: String

Full path to the input data for each run.

outputPath: String

Full path to the output location where data will be written.

*batchId: Long*

The batch id to process

*lastFullBatchId: Option[Long]*

The batch id indicating the most recent full batch available, None for first run.

*createFullBatch: Boolean*

If true, then create a full batch by combining previous full batch plus incremental batches since.

If false, only create an incremental batch output with the changes in the current batch.

### Attribute Rules

- For each company, take all attributes with highest probability.
  - We will identify maximum attribute\_prob for each company (across all attribute\_id) and then select only those attribute\_id & attribute\_value's which have atleast one attribute\_prob which matches with maximum attrib\_prob (maximum within company across attribute\_id)
- If probabilities are the same, then take the attribute value from the newer batch.
  - This rule states how we handle duplicates on attribute probability for the same company\_id, attribute\_id combination when found in 2 different batch files. It does not state any logic for handling duplicates within same batch but I believe that is not needed as I have verified in data and found no such batch.
  - In case there is an attribute\_id which has more than one row having the maximum attrib\_prob (maximum within company across attribute\_id) then we need to pick the values from the row which belongs to the latest batch\_id.

### Run 1

- Pass parameters inputPath="..input/run\_1", batchId=495, lastFullBatchId=None, createFullBatch= true.
- Since createFullBatch is true and lastFullBatchId is None, read in all input batches up to and including 495
- Produce a single batch output, applying the attribute rules
- Output format is Parquet and stored in the output directory provided. The output batch should be stored in its own partition with batch id 495 (same as input structure)

### Run 2

- Pass parameters inputPath="..input/run\_2", batchId=597, lastFullBatchId=495, createFullBatch= false.
- Read input batch 597 and the previous full batch output 495 you produced in run 1
  - My understanding is that we will need to pull the source data from all batch\_id > 495 upto the batch\_id = 597 (I believe that without this we may lose data for all intermediary batches so I am proceeding with this assumption). This data needs to be combined with the target data produced in run 1 for batch\_id=495
- Process batch 597 using the same rules you've already implemented and produce only the incremental batch changes as follows:
  - Include **all attributes** for only the companies that have **at least one** change.
    - My understanding is that by change we mean either a new attribute is added with higher probability or an existing attribute appears again with higher value
  - Store the incremental output in a partition, using batch id 597.

### Run 3

- Pass parameters `inputPath= "../input/run_3", batchId=648, lastFullBatchId=495, createFullBatch= true`.
- Read outputs for 495 (run 1) and incremental 597 (run 2) and combine with all data in this run to produce a full output with batch\_id 648 in the output folder
- My understanding here is again that we will be pulling all the source data for batch\_id > 597 and upto batch\_id=648. This data needs to be combined with the target data of partition 495 & 597 loaded in run 1 and run 2 respectively to produce the final output.

## Data

All data for the runs is available [here](#). The input batches are split into 3 groups for clarity and to simulate multiple runs. The format is parquet files under each partition.

### Input

#### input/run\_1:

batch\_id=73

batch\_id=75

batch\_id=76

...

batch\_id=495

#### input/run\_2:

batch\_id=508

batch\_id=514

...

batch\_id=597

#### input/run\_3:

batch\_id=601

batch\_id=605

...

batch\_id=648

### Output

output/batch\_id=495

output/batch\_id=597

output/batch\_id=648

## Schema

The input and output data schema should be the same:

```
batch_id: long  
company_id: Integer  
source_id: integer  
attribute_id: integer  
attribute_prob: double  
attribute_value: string
```

## Submission Guidelines

Please submit the entire scala sbt project with all code necessary to build and run the app. You can bundle it as a zip, tarball or github link. Also include a copy of the output folders that is generated from running the app.