

You have been provided with three data files. They are

- customers.csv
- products.csv
- sales.csv

You need to write a spark program in Scala using Dataframe/Dataset API's for below tasks. You need to use below versions of tools

- Spark 2.1
- Scala 2.11
- Sbt/Maven

Note:

- **Answers should be submitted within 48 hours from questions received. Otherwise answers will not be considered.**
- **We accept partial answers too, if you are not able to complete it within given time limit.**
- **Project you created should be sbt or maven project. You need to mail the whole project so that it can be directly imported and checked if the results are correct.**

Task 1 : Count of each files

Load all the above files as spark dataframes and print the count of each of files.

Task 2 : Handle Nulls

sales.csv contains null values in amount column. Replace those null values with mean of column in spark dataframe.

Task 3 : Join with Projection

Join all the above files on **customerId** and **itemId**. Select below columns from the joined data
date, customerId, itemId, name (from customers.csv with column name as customerName), category, amount

Task 4 : Add a Sequential Row ID

For a joined dataframe, add an id column which contains sequence numbers from 1 to number of rows

Task 5 : Sales By Week

Print Total sales for every week. Use **amount** for sales and **date** column for calculating week.

Task 6 : Sales with 5% Discount

Add a column **discount_amount** to joined dataframe, which holds 5% discounted amount.