

Age Group Sales Dominance

Background: Segment-dominated spending refers to situations where a particular demographic group—such as an age bracket, gender, or region—contributes disproportionately to total sales. Rather than consumer spending being evenly distributed, one segment often stands out as the primary revenue driver. Identifying these dominant segments enables businesses to optimize marketing strategies, personalize offerings, and allocate resources more effectively by focusing on the groups that generate the most value.

Problem Definition: You are given an E-commerce website logs data (*sales.csv*) created to practice exploratory data analysis. Each record in the dataset consists of a country code, user's language, user's age, and the sales value. A sample input file (*sample.txt*) has been provided. Your task is to utilize MRJob to find out the top-k age-sales for each country based on the following steps:

- For each country, calculate the country average sales (**country_avg_sales**).
- For each country, calculate the average sales in each age group (**country_age_avg_sales**).
- For each country, calculate the age dominance score (**ads**) for each age group where $\text{ads} = \text{country_age_avg_sales} / \text{country_avg_sales}$
- Report the age groups in each country whose **ads** value is not smaller than a given threshold τ .
- Records with invalid age column values should be ignored when computing average sales.
- You don't need to round anything during the computation.

An **age group** refers to a 10-year interval defined by rounding down the age to the nearest multiple of 10. For example, ages 0–9 are in group **0**, 10–19 in group **10**, 20–29 in group **20**, and so on.

Output Format: The output should contain three fields: the country code, the age group, and the ads value, in the format of `<the country code>**\t** <the age group>, < the ads value>`. The results should be sorted by country code alphabetically first and then by age group value in ascending order. Given the sample input file and the threshold $\tau=1.1$, the result should be like:

```
"CA"      "20,1.264528441270623"
"CN"      "40,1.1028398125172318"
"CN"      "60,1.5839536807278742"
```

Code Format: The code template has been provided. Your code should take three

parameters: the input file, the output folder on HDFS, and the threshold value k . We will also use more than 1 reducer to test your code. Assuming $\tau=1.2$ and using 2 reducers, you need to use the command below to run your code:

```
$ python3 project1.py -r hadoop sample.txt -o hdfs_output --jobconf  
myjob.settings.tau=1.2 --jobconf mapreduce.job.reduces=2
```

Note: You can access the value of τ in your program like `τ =`

```
jobconf_from_env('myjob.settings.tau') , and you need to import jobconf_from_env  
by from mrjob.compat import jobconf_from_env (see the code template)
```