



# PARALLEL DISTRIBUTED PROCESSING (HADOOP)

## Assignment 2

Abdullah Ghayumi  
634072@student.inholland.nl

# Assignment 2

## Details

Name: Abdullah Ghayumi

Student number: 634072

GitHub URL: <https://github.com/Aras53/Hadoop>

## Table of Contents

Details .....	1
Problem .....	2
Solution.....	3
Uploading orders.csv .....	3
Script.....	3
Result .....	4

## Problem

- Make an alphabetic list from all locations from the orders.csv.
- Group by “location” with target “Holland”
- Count how many times Holland was the target from that location
- Code is executed from Pig View
- Upload a document on Moodle with the following information:
  - [Name and student number](#)
  - [URL from GitHub with your source code](#)
- Explain in the document what steps have to be taken to execute the code
- Explain always in your own words every step included source code.
- Make a screenshot from the result and include it in your document

## Solution

### Uploading orders.csv

The first step I took was to upload the file to the Files View, but I quickly found out that the file was too big and wouldn't fully upload. To mitigate this, I tried uploading the file via MobaXTerm.

I uploaded the file first to the /home/maria\_dev location that I could access. Then I used the sudo root command to be able to access and execute command on the hdfs repo. I copied the orders file to the /users/maria\_dev folder with the following command:

```
hdfs dfs -put /home/maria_dev/orders.csv /user/maria_dev/
```

### Script

```
1 | ordersCSV = LOAD '/user/maria_dev/orders.csv' USING PigStorage(',') AS
2 |   (game_id:chararray,
3 |    unit_id:chararray,
4 |    unit_order:chararray,
5 |    location:chararray,
6 |    target:chararray,
7 |    target_dest:chararray,
8 |    success:chararray,
9 |    reason:chararray,
10 |    turn_num:chararray);
11 |
12 | filter_data = FILTER ordersCSV BY (target == "Holland");
13 | grp_by_target = GROUP filter_data BY location;
14 | count_holland = FOREACH grp_by_target GENERATE group as location, COUNT($1);
15 | order_result = ORDER count_holland BY location;
16 |
17 | DUMP order_result;
```

First, I filter the data, so that the only data that is left is the rows with the target "Holland". After having filtered the data, I grouped the data by the location. This would result in a dataset with the following structure: {location: {rows}, location2: {rows}, etc.}

After having grouped the data, I counted the times Holland was the target from the location. To close it out I ordered the data by the location to get a list on alphabetical order.

Result

assign2 - COMPLETED

Job ID

job\_1629002342303\_0047

Started

2021-08-19 08:48

▼ Results

Download

("Adriatic Sea",1)

("Aegron Sea",5)

("Albania",1)

("Armenia",1)

("Baltic Sea",326)

("Barents Sea",18)

("Belgium",25124)

("Berlin",1232)

("Black Sea",3)

("Bohemia",5)

("Brest",22)

("Budepest",1)

("Bulgaria",2)

("Burgundy",1155)

("CJyle",18)

("Constantinople",4)

("Denmark",4051)

("Eastern Mediterranean",4)

("Edinburgh",2023)

("English Channel",1232)

("Finland",1)

("Galicia",1)

("Gascony",4)

("Greece",3)

("Gulf of Bothnia",2)

("Gulf of Lyons",3)

("Hailand Sign",9187)

("Holland",325)

("Ionian Sea",11)

▼ Results

("Ionian Sea",11)

("Kiel",44658)

("Liverpool",90)

("Livonia",2)

("London",3519)

("Marseilles",2)

("Mid-Atlantic Ocean",30)

("Moscow",2)

("Munich",2297)

("Naples",1)

("North Africa",1)

("North Sea",16250)

("Norway",740)

("Norwegian Sea",233)

("Paris",10)

("Picardy",501)

("Portugal",5)

("Prussia",8)

("Rome",3)

("Ruhr",22142)

("Rumania",2)

("Serbia",4)

("Silesia",4)

("Skagerrack",164)

("Smyrna",1)

("Spain (South Coast)",1)

("Spain",5)

("St. Petersburg (North Coast)",10)

("St. Petersburg",24)

("Sweden",73)

("St. Petersburg",24)

("Sweden",73)

("Syria",1)

("Tunis",2)

("Tyrolia",4)

("Tyrrhenian Sea",11)

("Venice",2)

("Vienna",2)

("Wales",37)

("Warsaw",1)

("Western Mediterranean",13)

("Yorkshire",2882)