# Assignment 2 Part 1 Report

Ilaria Cattaneo

23rd April 2023

## 1 Exercise 1

### 1.1 Point 1 - Indexing

First thing first, we converted the csv file into a json file. So, the step-by-step procedure consists as follows. We deleted the first column of the given csv, because we didn't need the indices at the beginning of each row, and them we saved data in a new csv file to avoid confusion; the code is written below:

```
1  import csv
2          with open('restaurants.csv', 'r') as csvfile:
3              reader = csv.reader(csvfile)
4              your_list = list(reader)
5              for row in your_list:
6                  del row[0]
7              with open('restaurants2.csv', 'w') as f:
8                  writer = csv.writer(f)
9                  writer.writerows(your_list)
```

We wrote the code to convert the csv into a json file, specifying the data types we would like to have in the mapping when importing it into ElasticSearch, using the code below:

```
1  import json
2  import ast
3  import sys
4  jsonfile = open('restaurants.json', 'w')
5  with open('restaurants2.csv', 'r') as csvfile:
6      reader = csv.DictReader(csvfile)
7      for row in reader:
8          document = {}
9          for column, value in row.items():
10             if column == 'RestaurantName':
11                 document[column] = value
12             elif column == 'City':
13                 document[column] = value
14             elif column == 'AverageCostForTwo':
15                 document[column] = float(value)
16             elif column == 'AggregateRating':
17                 document[column] = float(value)
18             elif column == 'RatingText':
19                 document[column] = value
20             elif column == 'Votes':
21                 document[column] = float(value)
22             elif column == 'Date':
23                 document[column] = value
24             elif column == 'Coordinates':
25                 res = ast.literal_eval(value)
26                 document[column] = res
27
28          jsonfile.write(json.dumps(document))
29          jsonfile.write('\n')
```

Then, using the terminal, we used the commands the professor provided us on iCorsi, in order to make a bulk import of the json file into ElasticSearch into an index named 'restaurants1'; the code is written below:

```
$ sed -i .bak -E 's/.*/{"index":{}}\&/g' restaurants.json
$ split -l 10000 restaurants.json restaurants_split_
$ for f in `ls restaurants_split_*` ; do curl -k --user
"elastic:${ES_PASS}" -X POST https://localhost:9200/restaurants1/
    _bulk
-H 'Content-type:_application/x-ndjson' --data-binary @$f; done
```

After that, on ElasticSearch we checked the mapping of our imported file, using the following command.

```
1       GET restaurants1/_mapping
```

We noticed that, the imported file has a slightly different mapping with respect to the one we've written in our python code; for instance, 'AverageCostForTwo' in our code was written as a float, but it was mapped as an integer, so after looking at the data we decided to leave it as an integer. To be sure everything was working smoothly, we created a new index called 'restaurantsindex' and its mapping, and imported the json again with the new index.

```
1  PUT restaurantsindex
2
3  PUT restaurantsindex/_mapping
4  {
5    "properties": {
6      "AggregateRating": {
7        "type": "float"
8      },
9      "AverageCostForTwo": {
10       "type": "integer"
11     },
12     "City": {
13       "type": "text",
14       "fields": {
15            "keyword": {
16               "type": "keyword",
17               "ignore_above": 256
18          }
19        }
20      },
21      "Coordinates": {
22        "type": "geo_point"
23      },
24      "Date": {
25        "type": "date"
26      },
27      "RatingText": {
28        "type": "text",
```

```
29        "fields": {
30            "keyword": {
31                "type": "keyword",
32                "ignore_above": 256
33            }
34        }
35    },
36    "RestaurantName": {
37        "type": "text",
38        "fields": {
39            "keyword": {
40                "type": "keyword",
41                "ignore_above": 256
42            }
43        }
44    },
45    "Votes": {
46        "type": "float"
47    }
48  }
49 }
```

```
$ for f in `ls restaurants_split_*` ; do curl -k --user
"elastic:${ES_PASS}" -X POST https://localhost:9200/restaurantsindex
    /_bulk
-H 'Content-type: application/x-ndjson' --data-binary @$f; done
```

## 1.2 Point 2 - Queries

### 1.2.1 Question a

In the first task we have to get restaurants that have 'pizza' in the name and not 'pasta', that have reviews greater than or equal to 'Very Good'. The code is as follows.

```
1 GET restaurantsindex/_search
2 {
3   "size": 52,
4   "query": {
5     "bool": {
6       "must": [
7         {
8           "match": {
9           "RestaurantName": "pizza"
10        }},
11        {
12          "bool": {
13            "should": [
14              {
```

```
15            "match": {
16              "RatingText": "Very Good"
17            }
18          },
19          {
20            "match": {
21              "RatingText": "Excellent"
22            }
23          }
24        ]
25      }
26    },
27    {
28      "bool": {
29        "must_not": [
30          {
31            "match": {
32              "RestaurantName": "pasta"
33            }
34          }
35        ]
36      }
37    }
38  ]
39  }
40  }
41 }
```

From the running of this query, we found 52 documents returned; in particular, there are 52 restaurants that have a rating of 'Very Good' or 'Excellent' that have 'pizza' and not 'pasta' in their names.

### 1.2.2 Question b

In the second task we have to retrieve the 5 most expensive restaurants, which have reviews from 2018 and are located within 20km from Athens, using the following code.

```
1 GET restaurantsindex/_search
2 {
3   "size": 5,
4   "query": {
5     "bool": {
6       "must": [
7         {
8           "geo_distance": {
9             "distance": "20km",
10            "Coordinates": {
11              "lat": 33.9259,
12              "lon": -83.3389
```

```
13                    }
14
15                }},
16                {
17                   "range": {
18                      "Date": {
19                         "gte": "2018||/y",
20                         "lte": "2018||/y",
21                         "format": "yyyy"
22                      }
23                   }
24                }
25             ]
26          }
27       },
28    "sort": [
29          {
30          "AverageCostForTwo": {
31             "order": "desc"
32          }
33       }
34    ]
35 }
```

The query returned 8 documents, this means that 8 restaurants satisfy the conditions written in the code to answer the task.

### 1.2.3   Question c

In the third task we have to retrieve all restaurants which contain the sub-string 'pizz' in their name but that do not contain neither 'pizza' nor 'pizzeria'. To do so, we used the code below.

```
1 GET restaurantsindex/_search
2 {
3    "query": {
4       "bool": {
5          "must": [
6             {
7                "wildcard": {
8                   "RestaurantName": {
9                      "value": "*pizz*"
10                   }
11                }
12             },
13             {
14                "bool": {
15                   "must_not": [
16                      {
17                         "wildcard": {
```

```
18              "RestaurantName": {
19                "value": "*pizza*"
20              }
21            }
22          },
23          {
24            "wildcard": {
25              "RestaurantName": {
26                "value": "*pizzeria*"
27              }
28            }
29          }
30        ]
31      }
32    }
33    ]
34  }
35  }
36 }
```

The query returned 1 document and this means that there is only 1 restaurant, whose name contains the sub-string 'pizz' but does not contain neither 'pizza' nor 'pizzeria'; as a matter of fact, the restaurant returned by the query is called 'Pizzoccheri'.

## 1.3   Point 3 - Aggregations

### 1.3.1   Question a

In the first activity of this section, we have to aggregate restaurants that have 'Good' as review by number of votes, by creating 4 buckets and for each of them we have to return the maximum and the minimum value of 'AverageCostForTwo'.

```
1  GET restaurantsindex/_search
2  {
3    "size": 0,
4    "query": {
5      "bool": {
6        "must": [
7          {
8            "match": {
9              "RatingText": "Good"
10           }
11         }
12       ]
13     }
14   },
15   "aggs": {
16     "rating_buckets": {
17       "range": {
```

```
18          "field": "Votes",
19          "ranges": [
20             {
21                "from": 0,
22                "to": 250
23             },
24             {
25                "from": 250,
26                "to": 500
27             },
28             {
29                "from": 500,
30                "to": 750
31             },
32             {
33                "from": 750,
34                "to": 1000
35             }
36          ]
37       },
38       "aggs": {
39          "min_cost": {
40             "min": {
41                "field": "AverageCostForTwo"
42             }
43          },
44          "max_cost":{
45             "max": {
46                "field": "AverageCostForTwo"
47             }
48          }
49       }
50    }
51  }
52 }
```

The aggregation query returned a total of 3162 documents; in this case we are not facing the shard size problem. The result is reported below.

```
1  {
2    "took": 1,
3    "timed_out": false,
4    "_shards": {
5      "total": 1,
6      "successful": 1,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
```

```
11      "total": {
12        "value": 3162,
13        "relation": "eq"
14      },
15      "max_score": null,
16      "hits": []
17    },
18    "aggregations": {
19      "rating_buckets": {
20        "buckets": [
21          {
22            "key": "0.0-250.0",
23            "from": 0,
24            "to": 250,
25            "doc_count": 2060,
26            "max_cost": {
27              "value": 350000
28            },
29            "min_cost": {
30              "value": 0
31            }
32          },
33          {
34            "key": "250.0-500.0",
35            "from": 250,
36            "to": 500,
37            "doc_count": 583,
38            "max_cost": {
39              "value": 450000
40            },
41            "min_cost": {
42              "value": 10
43            }
44          },
45          {
46            "key": "500.0-750.0",
47            "from": 500,
48            "to": 750,
49            "doc_count": 217,
50            "max_cost": {
51              "value": 5000
52            },
53            "min_cost": {
54              "value": 10
55            }
56          },
57          {
58            "key": "750.0-1000.0",
```

```
59      "from": 750,
60      "to": 1000,
61      "doc_count": 99,
62      "max_cost": {
63        "value": 200000
64      },
65      "min_cost": {
66        "value": 10
67      }
68    }
69   ]
70  }
71  }
72 }
```

### 1.3.2 Question b

The second activity consists of finding cities that have at least 10 restaurants and restaurants that have at least 100 votes, and returning the 7 cities with the highest 'AverageCostForTwo'.

```
1 GET restaurantsindex/_search
2 {
3   "size": 0,
4   "query": {
5     "range": {
6       "Votes": {
7         "gte": 100
8       }
9     }
10  },
11  "aggs": {
12    "frequent_cities": {
13      "terms": {
14        "field": "City.keyword",
15        "min_doc_count": 10,
16        "order": {
17          "average": "desc"
18        },
19        "shard_size": 1000,
20        "size": 7
21      },
22      "aggs": {
23        "average": {
24          "max": {
25            "field": "AverageCostForTwo"
26          }
27        }
28      }
```

9

```
29          }
30        }
31  }
```

The aggregation query returned 2799 documents; here we are facing the shard problem, as we find "doc_count_error_upper_bound": -1, so to have the right result we must specify the shard size to 1000. The result is reported below.

```
1   {
2     "took": 1,
3     "timed_out": false,
4     "_shards": {
5       "total": 1,
6       "successful": 1,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 2799,
13        "relation": "eq"
14      },
15      "max_score": null,
16      "hits": []
17    },
18    "aggregations": {
19      "frequent_cities": {
20        "doc_count_error_upper_bound": 0,
21        "sum_other_doc_count": 970,
22        "buckets": [
23          {
24            "key": "Jakarta",
25            "doc_count": 16,
26            "average": {
27              "value": 800000
28            }
29          },
30          {
31            "key": "New Delhi",
32            "doc_count": 1174,
33            "average": {
34              "value": 8000
35            }
36          },
37          {
38            "key": "Gurgaon",
39            "doc_count": 274,
40            "average": {
41              "value": 5000
```

```
42                 }
43               },
44               {
45                 "key": "Colombo",
46                 "doc_count": 14,
47                 "average": {
48                   "value": 4500
49                 }
50               },
51               {
52                 "key": "Noida",
53                 "doc_count": 192,
54                 "average": {
55                   "value": 3200
56                 }
57               },
58               {
59                 "key": "Jaipur",
60                 "doc_count": 18,
61                 "average": {
62                   "value": 2500
63                 }
64               },
65               {
66                 "key": "Bangalore",
67                 "doc_count": 20,
68                 "average": {
69                   "value": 2400
70                 }
71               }
72             ]
73           }
74         }
75 }
```

### 1.3.3 Question c

In the third activity, we have to filter for restaurants that are located within 9000km of New Dehli and show the highest number of votes for different rating types in descending order.

```
1 GET restaurantsindex/_search
2 {
3   "size": 0,
4   "query": {
5     "bool": {
6       "must": [
7         {
8           "geo_distance": {
```

```
 9              "distance": "9000km",
10              "Coordinates": {
11                "lat": 28.642449499999998,
12                "lon": 77.10684570000001
13              }
14            }
15          }
16        ]
17      }
18    },
19    "aggs": {
20      "ratings": {
21        "terms": {
22          "field": "RatingText.keyword",
23          "size": 10
24        },
25        "aggs": {
26          "max_rate": {
27            "max": {
28              "field": "AggregateRating"
29            }
30          },
31          "max_votes": {
32            "max": {
33              "field": "Votes"
34            }
35          },
36          "sorting_votes": {
37           "bucket_sort": {
38             "sort": [
39               {
40                 "max_votes": {
41                   "order": "desc"
42                 }
43               }
44             ]
45           }
46          }
47        }
48      }
49    }
50 }
```

The aggregation query returned 8930 documents. Here we are not facing the shard problem. The result is reported below.

```
1 {
2    "took": 0,
3    "timed_out": false,
```

```
 4    "_shards": {
 5      "total": 1,
 6      "successful": 1,
 7      "skipped": 0,
 8      "failed": 0
 9    },
10    "hits": {
11      "total": {
12        "value": 8930,
13        "relation": "eq"
14      },
15      "max_score": null,
16      "hits": []
17    },
18    "aggregations": {
19      "ratings": {
20        "doc_count_error_upper_bound": 0,
21        "sum_other_doc_count": 0,
22        "buckets": [
23          {
24            "key": "Excellent",
25            "doc_count": 196,
26            "max_rate": {
27              "value": 4.900000095367432
28            },
29            "max_votes": {
30              "value": 10934
31            }
32          },
33          {
34            "key": "Very Good",
35            "doc_count": 834,
36            "max_rate": {
37              "value": 4.400000095367432
38            },
39            "max_votes": {
40              "value": 7931
41            }
42          },
43          {
44            "key": "Good",
45            "doc_count": 1902,
46            "max_rate": {
47              "value": 3.9000000953674316
48            },
49            "max_votes": {
50              "value": 4914
51            }
```

```
52          },
53          {
54            "key": "Average",
55            "doc_count": 3683,
56            "max_rate": {
57              "value": 3.4000000953674316
58            },
59            "max_votes": {
60              "value": 2460
61            }
62          },
63          {
64            "key": "Poor",
65            "doc_count": 180,
66            "max_rate": {
67              "value": 2.4000000953674316
68            },
69            "max_votes": {
70              "value": 2412
71            }
72          },
73          {
74            "key": "Not rated",
75            "doc_count": 2135,
76            "max_rate": {
77              "value": 0
78            },
79            "max_votes": {
80              "value": 3
81            }
82          }
83        ]
84      }
85    }
86 }
```