

DejaVu Sans

Visual Analytics Assignment 2 Report

Zitian

April 28, 2025

1 Section 1: Data Indexing

1.1 Data Cleaning

The raw CSV data was first cleaned using a Python script to ensure all fields are well-typed, missing values are handled, and city names are standardized. The main steps include:

- Converting numeric fields to the correct type (e.g., ratings, votes, cost).
- Filling missing values with defaults.
- Standardizing date and location formats.
- Removing duplicates.
- Outputting a cleaned CSV file for indexing.

Key code snippet:

Listing 1: Key data cleaning code

```
df['AggregateRating'] = pd.to_numeric(df['AggregateRating'],
    errors='coerce')
df['Votes'] = pd.to_numeric(df['Votes'], errors='coerce')
df['AverageCostForTwo'] = pd.to_numeric(df['
    AverageCostForTwo'], errors='coerce')
df['RestaurantName'] = df['RestaurantName'].fillna('Unknown
    Restaurant')
df['Date'] = df['Date'].apply(validate_date)
df.to_csv('Assignment2/restaurants_cleaned.csv', sep=';',
    index=False)
```

1.2 Indexing with Grok Pipeline

To ingest the cleaned data into Elasticsearch, we initially used a Grok pipeline but later switched to a more reliable direct JSON document approach. Here's our journey and lessons learned:

1.2.1 Initial Approach: Grok Pipeline

The initial approach used a Grok pipeline to parse raw CSV lines:

- Converting each row to a single raw string
- Using Grok patterns to parse fields
- Converting field types after parsing

However, this approach had limitations:

- Complex string parsing could lead to data loss
- Field type conversion was error-prone
- Debugging was difficult when fields were missing

1.2.2 Improved Approach: Direct JSON Documents

We improved the indexing process by directly constructing JSON documents:

- Each field is explicitly typed and mapped
- No intermediate string parsing required
- Better control over data integrity

Key code snippet of the improved approach:

Listing 2: Direct JSON document construction

```
def prepare_bulk_data(csv_file):
    df = pd.read_csv(csv_file, sep=';')
    bulk_data = []
    for _, row in df.iterrows():
        doc = {
            "SerialNumber": int(row['SerialNumber']),
            "RestaurantName": str(row['RestaurantName']),
            "AverageCostForTwo": int(row['AverageCostForTwo']),
            "AggregateRating": float(row['AggregateRating']),
            "RatingText": str(row['RatingText']),
            "Votes": float(row['Votes']),
            "Date": str(row['Date']),
            "Coordinates": str(row['Coordinates']),
            "City": str(row['City']),
            "Country": str(row['Country']),
            "Continent": str(row['Continent']),
            "City/Country/Continent": f"{str(row['City'])}/{str(row['Country'])}/{str(row['Continent'])}"
        }
```

```

        bulk_data.append(json.dumps({"index": {"_index":
            INDEX_NAME}}))
        bulk_data.append(json.dumps(doc))
    return '\n'.join(bulk_data) + '\n'

```

1.2.3 Key Learnings

Through this process, we learned several important lessons:

- Direct data mapping is more reliable than string parsing
- Explicit type conversion prevents data loss
- Simpler processing pipelines are easier to debug
- Maintaining data integrity requires careful handling at each step

Result: All cleaned restaurant data is now correctly indexed in Elasticsearch as structured JSON documents, with all fields properly preserved and typed.

1.3 Query 1: Find restaurants with 'bar' in the name but not 'barbecue(s)' or 'barbeque(s)'

The goal is to retrieve all restaurants whose names contain the substring `bar`, but do not contain `barbecue`, `barbeque`, `barbecues`, or `barbeques`. Only the restaurant name, city/country/continent, and number of votes are returned.

Listing 3: Elasticsearch query for 'bar' but not 'barbecue(s)' or 'barbeque(s)'

```

1 GET restaurants/_search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         {
7           "wildcard": {
8             "RestaurantName": "*bar*"
9           }
10        }
11      ],
12      "must_not": [
13        {
14          "wildcard": {
15            "RestaurantName": "*barbecue*"
16          }
17        },
18        {
19          "wildcard": {
20            "RestaurantName": "*barbeque*"
21          }

```

```

22     },
23     {
24         "wildcard": {
25             "RestaurantName": "*barbecues*"
26         }
27     },
28     {
29         "wildcard": {
30             "RestaurantName": "*barbeques*"
31         }
32     },
33     {
34         "match": {
35             "RestaurantName": "bar"
36         }
37     }
38 ]
39 }
40 },
41 "_source": ["RestaurantName", "City/Country/Continent", "
    Votes"]
42 }

```

Query Result Summary

The query returned **89** restaurants whose names contain “bar” but do not contain “barbecue(s)” or “barbeque(s)”.

Full results are provided in the file: section1QueryResult/query1_results.json (see submission folder).

Below are some sample results:

```

1  {
2    "RestaurantName": "Paribar",
3    "City/Country/Continent": "Sao Paulo/Brazil/",
4    "Votes": 46
5  },
6  {
7    "RestaurantName": "Bardenay",
8    "City/Country/Continent": "Boise/United States/",
9    "Votes": 879
10 },
11 {
12    "RestaurantName": "Barbacoa Restaurant",
13    "City/Country/Continent": "Boise/United States/",
14    "Votes": 538
15 }

```

A Full Query 1 Results

The complete results for Query 1 are provided below:

Listing 4: All results for Query 1

```
1 {
2   "took": 34,
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": 89,
13      "relation": "eq"
14    },
15    "max_score": 1,
16    "hits": [
17      {
18        "_index": "restaurants",
19        "_id": "QE5hfZYBQKwotppyNThx",
20        "_score": 1,
21        "_source": {
22          "RestaurantName": "Paribar",
23          "City/Country/Continent": "S    O Paulo/Brazil/"
24        },
25        "Votes": 46
26      },
27      {
28        "_index": "restaurants",
29        "_id": "hk5hfZYBQKwotppyNThx",
30        "_score": 1,
31        "_source": {
32          "RestaurantName": "Bardenay",
33          "City/Country/Continent": "Boise/United States/"
34        },
35        "Votes": 879
36      },
37      {
38        "_index": "restaurants",
39        "_id": "i05hfZYBQKwotppyNThx",
40        "_score": 1,
41        "_source": {
42          "RestaurantName": "Barbacoa Restaurant",
```

```

43         "City/Country/Continent": "Boise/United States/"
44     },
45     "Votes": 538
46 },
47 {
48     "_index": "restaurants",
49     "_id": "ik5hfZYBQKwotppyNTly",
50     "_score": 1,
51     "_source": {
52         "RestaurantName": "Barrett Junction Cafe",
53         "City/Country/Continent": "Potrero/United States",
54         "Votes": 9
55     }
56 },
57 {
58     "_index": "restaurants",
59     "_id": "uU5hfZYBQKwotppyNTly",
60     "_score": 1,
61     "_source": {
62         "RestaurantName": "Rhubarb Le Restaurant",
63         "City/Country/Continent": "Singapore/Singapore/"
64     },
65     "Votes": 33
66 },
67 {
68     "_index": "restaurants",
69     "_id": "f05hfZYBQKwotppyNTpy",
70     "_score": 1,
71     "_source": {
72         "RestaurantName": "The Cafe Baraco",
73         "City/Country/Continent": "Ahmedabad/India/",
74         "Votes": 317
75     }
76 },
77 {
78     "_index": "restaurants",
79     "_id": "Rk5hfZYBQKwotppyNTty",
80     "_score": 1,
81     "_source": {
82         "RestaurantName": "Doon Darbar",
83         "City/Country/Continent": "Dehradun/India/",
84         "Votes": 121
85     }
86 },
87 {
88     "_index": "restaurants",
89     "_id": "J05hfZYBQKwotppyNTzo",

```

```

90     "_score": 1,
91     "_source": {
92         "RestaurantName": "The Grill Darbar",
93         "City/Country/Continent": "Faridabad/India/",
94         "Votes": 17
95     }
96 },
97 {
98     "_index": "restaurants",
99     "_id": "TU5hfZYBQKwotppyNTzo",
100     "_score": 1,
101     "_source": {
102         "RestaurantName": "Barista",
103         "City/Country/Continent": "Ghaziabad/India/",
104         "Votes": 33
105     }
106 },
107 {
108     "_index": "restaurants",
109     "_id": "dk5hfZYBQKwotppyNTzo",
110     "_score": 1,
111     "_source": {
112         "RestaurantName": "Zambar",
113         "City/Country/Continent": "Gurgaon/India/",
114         "Votes": 802
115     }
116 }
117 ]
118 }
119 }

```