



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

Systems and Methods for Big and Unstructured Data Project

Author(s): **Jiaxiang Yi 10765316**

Yang Hao Mao 10705881

Group Number: **19**

Academic Year: 2024-2025

Contents

Contents	i
1 Introduction	1
1.1 Objective	1
1.2 Problem Specification	1
1.3 DBMS Technology Selection	1
2 Data Wrangling	3
3 Dataset	5
3.1 Global Health Statistics dataset	5
3.2 Fields Overview	5
3.3 Toward to NoSQL	6
4 Queries	9
4.1 Filter documents with a specific Country	9
4.1.1 query	9
4.1.2 output(partial)	9
4.2 Retrieve documents within a certain time range (Year)	11
4.2.1 query	11
4.2.2 output (partial)	11
4.3 Retrieve data for a specific demographic and disease category	13
4.3.1 query	13
4.3.2 output(partial)	13
4.4 Identify the top N countries by document count	15
4.4.1 query	15
4.4.2 output	15
4.5 Retrieve documents matching high mortality rate OR high recovery rate .	17
4.5.1 query	17

4.5.2	output(partial)	17
4.6	Search for disease names allowing minor typos	19
4.6.1	query	19
4.6.2	output(partial)	19
4.7	Calculate global average of doctors per 1000	21
4.7.1	query	21
4.7.2	output	21
4.8	Compare vaccine availability categories side by side	22
4.8.1	query	22
4.8.2	output	23
4.9	Retrieve the latest record for each disease category	24
4.9.1	query	24
4.9.2	output(partial)	25
4.10	Find the costliest diseases by average treatment cost	27
4.10.1	query	27
4.10.2	output	27
4.11	Group year data in intervals and keep high-affected populations only	29
4.11.1	query	29
4.11.2	output(partial)	30
4.12	Get extended statistics for the population affected by a specific disease	31
4.12.1	query	31
4.12.2	output	31
4.13	Prioritize documents from wealthier contexts	33
4.13.1	query	33
4.13.2	output(partial)	33
4.14	Compare age-group mortality side by side	35
4.14.1	query	35
4.14.2	output	36
4.15	Find diseases that stand out in a subpopulation	36
4.15.1	query	36
4.15.2	output	37
4.16	Highlight partial matches in disease names	38
4.16.1	query	39
4.16.2	output(partial)	39
4.17	Break down data by year to compute average mortality	41
4.17.1	query	41
4.17.2	output(partial)	41

4.18	Dynamically compute the gap between healthcare access and education	42
4.18.1	query	43
4.18.2	output(partial)	43
4.19	Smooth data trends over time with a rolling average	44
4.19.1	query	44
4.19.2	output(partial)	45
4.20	Filter results by cost or urbanization, then rank by year	46
4.20.1	query	46
4.20.2	output(partial)	48
5	Dashboard Description	49
5.1	Top 3 Diseases Distribution by Age Group	49
5.2	Top 10 Countries by Alzheimer's Disease Distribution	50
5.3	2000–2024 Count of Records Distribution by Age Group	50
5.4	Top 5 Disease Categories in 5 Countries	51
5.5	2014–2024 Changes in Alzheimer's Disease Counts	51
5.6	Top 10 Countries by Minimum Average Treatment Cost (USD)	52

1 | Introduction

1.1. Objective

The project involves using a publicly accessible dataset and analyzing it with a NoSQL technology studied in the course, specifically **Neo4j**, **MongoDB**, or **Elasticsearch**. The primary goal is to demonstrate how to ingest, store, and query real-world data with twenty specific queries, extracting meaningful insights while showcasing the database's querying capabilities. By implementing and applying these queries, the project highlights how modern NoSQL technologies can support data-driven analysis in the health domain.

1.2. Problem Specification

The dataset, titled “**Global Health Statistics**”, covers diverse health metrics such as disease rates, healthcare expenditures, and life expectancies across various countries and timelines. The problem lies in efficiently accessing and searching these metrics in a way that accommodates complex queries, for example, filtering by specific diseases, regions, or time ranges to reveal patterns and trends. Storing and querying the data in a NoSQL database allows for flexible schema design, fast data retrieval, and robust scalability, all of which are vital for timely health analytics and decision-making.

Reference link of dataset: www.kaggle.com/datasets/malaiarasugraj/global-health-statistics

1.3. DBMS Technology Selection

Elasticsearch is selected for this project because of its powerful full-text search capabilities, scalability, and near real-time data processing. It excels at handling both structured and unstructured data, making it well-suited for an interdisciplinary dataset that includes numerical metrics (e.g., mortality rates, healthcare costs) and textual fields (e.g., disease names, country names). Elasticsearch also provides efficient indexing for complex queries and aggregations, enabling interactive, real-time exploration of large health datasets.

Kibana, a companion tool for Elasticsearch, will be used for data visualization and dashboard creation. Kibana's visual interface allows users to quickly build charts, maps, and tables, offering interactive insights into global health trends and facilitating the discovery of patterns and correlations across the dataset.

By leveraging **Elasticsearch** and **Kibana**, we aim to create a workable proof of concept that demonstrates how global health metrics can be quickly ingested, queried, and visualized. This project will focus on practical queries and basic dashboards, serving as a learning exercise in applying NoSQL tools to real-world health data.

2 | Data Wrangling

We loaded the CSV (Global Health Statistics.csv) into a Python notebook (Data Wrangling.ipynb) environment. After inspecting the dataset (displaying the first few rows, checking column data types, and verifying missing values), we found that:

- No missing values existed in any of the columns.
- All fields had appropriate data types (e.g., numeric or object/string).
- The columns were already in a suitable format for analysis, requiring no further cleaning or wrangling steps.

Since the dataset was already clean, no data wrangling (such as imputations, dropping duplicates, or transforming columns) was performed. No additional synthetic data was generated. We simply proceeded with exploratory data analysis on the dataset as-is.

3 | Dataset

3.1. Global Health Statistics dataset

The Dataset used for the project is named "Global Health Statistics.csv" has 1'000'000 rows, 22 columns and it has been downloaded from Kaggle in following link:
www.kaggle.com/datasets/malaiarasugraj/global-health-statistics

3.2. Fields Overview

The dataset is structured with following fields, each contributing to a comprehensive overview of the channels:

0. **Field (data type):** Description
 1. **Country (string):** The name of the country where the health data was recorded.
 2. **Year (int):** The year in which the data was collected.
 3. **Disease Name (string):** The name of the disease or health condition tracked.
 4. **Disease Category (string):** The category of the disease.
 5. **Prevalence Rate(%) (float):** The percentage of the population affected by the disease.
 6. **Incidence Rate(%) (float):** The percentage of new or newly diagnosed cases.
 7. **Morality Rate(%) (float):** The percentage of the affected population that dies from the disease.
 8. **Age Group (string):** The age range most affected by the disease.
 9. **Gender (string):** The gender(s) affected by the disease (Male, Female, Both).
 10. **Population Affected (int):** The total number of individuals affected by the disease.

11. **Healthcare Access(%) (float)**: The percentage of the population with access to healthcare.
12. **Doctors per 1000 (float)**: The number of doctors per 1000 people.
13. **Hospital Beds per 1000 (float)**: The number of hospital beds available per 1000 people.
14. **Treatment Type (string)**: The primary treatment method for the disease (e.g., Medication, Surgery).
15. **Average Treatment Cost (USD) (int)**: The average cost of treating the disease in USD.
16. **Availability of Vaccines/Treatment (string)**: Whether vaccines or treatments are available.
17. **Recovery Rate(%) (float)**: The percentage of people who recover from the disease.
18. **DALYs (int)**: Disability-Adjusted Life Years, a measure of disease burden.
19. **Improvement in 5 Years(%) (float)**: The improvement in disease outcomes over the last five years.
20. **Per Capita Income (USD) (int)**: The average income per person in the country.
21. **Education Index (float)**: The average level of education in the country.
22. **Urbanization Rate(%) (float)**: The percentage of the population living in urban areas.

3.3. Toward to NoSQL

Unlike relational databases, NoSQL databases as Elasticsearch do not require a fixed schema. This flexibility allows for the addition of new data types or fields as the dataset evolves. The data was indexed in Elasticsearch with customized mappings for each field to optimize both search and analysis capabilities.

Fields representing categorical data, like **Country**, **Disease Name**, **Disease Category**, **Age Group**, **Gender**, **Treatment Type** and **Availability of Vaccines/Treatment** were indexed as keywords to enable precise, exact match searches. Numerical values, critical for statistical computations, were stored as integers or doubles (float64) depending on their

need for precision. This strategic mapping ensures not only efficient query performance but also delivers highly relevant search results.

For visual reference, a screenshot of the Elasticsearch index mapping is provided in the delivery folder, named “mapping.txt.” This file offers a clear depiction of the data schema utilized in this project, highlighting the detailed organization and indexing strategies employed. Below, a small excerpt from the mapping is displayed to illustrate these points:

```
1  {
2    "globalhealth19": {
3      "aliases": {},
4      "mappings": {
5        "_meta": {
6          "created_by": "file-data-visualizer"
7        },
8        "properties": {
9          "Age Group": {
10            "type": "keyword"
11          },
12          "Availability of Vaccines/Treatment": {
13            "type": "keyword"
14          },
15          "Average Treatment Cost (USD)": {
16            "type": "long"
17          },
18          "Country": {
19            "type": "keyword"
20          },
21          "DALYs": {
22            "type": "long"
23          },
24          "Disease Category": {
25            "type": "keyword"
26          },
27          "Disease Name": {
28            "type": "keyword"
29          }

```

Figure 3.1: small portion of mapping.txt

4 | Queries

4.1. Filter documents with a specific Country

In real-world scenarios, analysts often need to focus on data from a particular country to study localized health outcomes or disease prevalence. For instance, retrieving documents where Country is "Italy" helps you analyze patterns specific to Italy, such as the Mortality Rate (%) in that region.

4.1.1. query

```
5  GET globalhealth19/_search
6  {
7    "query": {
8      "match": {
9        "Country": "Italy"
10       }
11     }
12 }
```

Figure 4.1: query 4.1

4.1.2. output(partial)

```

1  {
2      "took": 3,
3      "timed_out": false,
4      "_shards": {
5          "total": 5,
6          "successful": 5,
7          "skipped": 0,
8          "failed": 0
9      },
10     "hits": {
11         "total": {
12             "value": 10000,
13             "relation": "gte"
14         },
15         "max_score": 3.008238,
16         "hits": [
17             {
18                 "_index": "globalhealth19",
19                 "_id": "nJoxsZMBhWjrNWL3UHVR",
20                 "_score": 3.008238,
21                 "_source": {
22                     "Disease Category": "Autoimmune",
23                     "Incidence Rate (%)": 4.86,
24                     "Mortality Rate (%)": 0.72,
25                     "Doctors per 1000": 4.1,
26                     "Education Index": 0.46,
27                     "Healthcare Access (%)": 51.72,
28                     "Recovery Rate (%)": 84.52,
29                     "Treatment Type": "Medication",
30                     "Improvement in 5 Years (%)": 3.97,
31                     "Gender": "Other",
32                     "Population Affected": 616894,
33                     "Per Capita Income (USD)": 80067,
34                     "Hospital Beds per 1000": 3.2,
35                     "Urbanization Rate (%)": 52.81,
36                     "Year": 2023,
37                     "Prevalence Rate (%)": 0.43,
38                     "Disease Name": "Cancer",
39                     "DALYs": 4909,
40                     "Availability of Vaccines/Treatment": "Yes",
41                     "Country": "Italy",
42                     "Age Group": "0-18",
43                     "Average Treatment Cost (USD)": 42869
44                 }
45             }
46         }
47     }
48 }
```

Figure 4.2: output(partial) of query 4.1

4.2. Retrieve documents within a certain time range (Year)

Health analysts frequently investigate a specific time window (e.g., 2015–2020) to understand disease trends within that interval. This query uses a range filter on the numeric Year field.

4.2.1. query

```
15  GET globalhealth19/_search
16  {
17  "query": {
18    "range": {
19      "Year": {
20        "gte": 2020,
21        "lte": 2022
22      }
23    }
24  }
25 }
```

Figure 4.3: query 4.2

4.2.2. output (partial)

```

1   {
2     "took": 652,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": 1,
16      "hits": [
17        {
18          "_index": "globalhealth19",
19          "_id": "0ZEvsZMBhWjrNWL3RpuN",
20          "_score": 1,
21          "_source": {
22            "Disease Category": "Parasitic",
23            "Incidence Rate (%)": 10.49,
24            "Mortality Rate (%)": 9.51,
25            "Doctors per 1000": 3.5,
26            "Education Index": 0.73,
27            "Healthcare Access (%)": 96.94,
28            "Recovery Rate (%)": 59.59,
29            "Treatment Type": "Vaccination",
30            "Improvement in 5 Years (%)": 1.87,
31            "Gender": "Male",
32            "Population Affected": 640330,
33            "Per Capita Income (USD)": 22233,
34            "Hospital Beds per 1000": 1.96,
35            "Urbanization Rate (%)": 69.02,
36            "Year": 2020,
37            "Prevalence Rate (%)": 16.49,
38            "Disease Name": "Measles",
39            "DALYs": 3237,
40            "Availability of Vaccines/Treatment": "No",
41            "Country": "Nigeria",
42            "Age Group": "61+",
43            "Average Treatment Cost (USD)": 28406
44          }
45        },

```

Figure 4.4: output(partial) of query 4.2

4.3. Retrieve data for a specific demographic and disease category

This query shows how you might isolate data for a specific demographic (e.g., Gender = "Female") combined with a disease category (e.g., "Infectious"). Real-world usage includes analyzing how certain diseases disproportionately affect certain genders within a category like "Infectious."

4.3.1. query

```
28  GET globalhealth19/_search
29  {
30    "query": {
31      "bool": {
32        "must": [
33          { "term": { "Gender": "Female" } },
34          { "term": { "Disease Category": "Infectious" } }
35        ]
36      }
37    }
38 }
```

Figure 4.5: query 4.3

4.3.2. output(partial)

```

1   {
2     "took": 164,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": 3.5129046,
16      "hits": [
17        {
18          "_index": "globalhealth19",
19          "_id": "yJYwsZMBhWjrNWL3aH_1",
20          "_score": 3.5129046,
21          "_source": {
22            "Disease Category": "Infectious",
23            "Incidence Rate (%)": 4.56,
24            "Mortality Rate (%)": 9.7,
25            "Doctors per 1000": 4.94,
26            "Education Index": 0.84,
27            "Healthcare Access (%)": 55.64,
28            "Recovery Rate (%)": 80.37,
29            "Treatment Type": "Surgery",
30            "Improvement in 5 Years (%)": 0.02,
31            "Gender": "Female",
32            "Population Affected": 691815,
33            "Per Capita Income (USD)": 72989,
34            "Hospital Beds per 1000": 3.82,
35            "Urbanization Rate (%)": 64.8,
36            "Year": 2010,
37            "Prevalence Rate (%)": 6.79,
38            "Disease Name": "Dengue",
39            "DALYs": 1586,
40            "Availability of Vaccines/Treatment": "Yes",
41            "Country": "Germany",
42            "Age Group": "36-60",
43            "Average Treatment Cost (USD)": 24856
44          }
45        },

```

Figure 4.6: output(partial) of query 4.3

4.4. Identify the top N countries by document count

This aggregation helps identify which countries dominate the dataset. Often used to quickly see the top 5 or 10 countries based on the number of records, possibly indicating higher incidence or broader coverage in your data.

4.4.1. query

```
41  GET globalhealth19/_search
42  {
43    "size": 0,
44    "aggs": {
45      "top_countries": {
46        "terms": {
47          "field": "Country",
48          "size": 5
49        }
50      }
51    }
52 }
```

Figure 4.7: query 4.4

4.4.2. output

```
1  {
2      "took": 212,
3      "timed_out": false,
4      "_shards": {
5          "total": 5,
6          "successful": 5,
7          "skipped": 0,
8          "failed": 0
9      },
10     "hits": {
11         "total": {
12             "value": 10000,
13             "relation": "gte"
14         },
15         "max_score": null,
16         "hits": []
17     },
18     "aggregations": {
19         "top_countries": {
20             "doc_count_error_upper_bound": 49531,
21             "sum_other_doc_count": 748578,
22             "buckets": [
23                 {
24                     "key": "Russia",
25                     "doc_count": 50532
26                 },
27                 {
28                     "key": "South Africa",
29                     "doc_count": 50408
30                 },
31                 {
32                     "key": "South Korea",
33                     "doc_count": 50181
34                 },
35                 {
36                     "key": "Germany",
37                     "doc_count": 50176
38                 },
39                 {
40                     "key": "UK",
41                     "doc_count": 50125
42                 }
43             ]
44         }
45     }
}
```

Figure 4.8: output of query 4.4

4.5. Retrieve documents matching high mortality rate OR high recovery rate

In practice, analysts may want to flag potential outliers: diseases with a mortality rate $\geq 5\%$ or a recovery rate $\geq 90\%$. This query captures documents meeting either condition to investigate them further (e.g., exploring contributing factors or success stories).

4.5.1. query

```
55  GET globalhealth19/_search
56  {
57  "query": {
58    "bool": {
59      "should": [
60        { "range": { "Mortality Rate (%)": { "gte": 5 } } },
61        { "range": { "Recovery Rate (%)": { "gte": 90 } } }
62      ],
63      "minimum_should_match": 1
64    }
65  }
```

Figure 4.9: query 4.5

The `minimum_should_match` ensure must meet at least 1 of the should conditions to be included.

4.5.2. output(partial)

```

1   {
2     "took": 35,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": 2,
16      "hits": [
17        {
18          "_index": "globalhealth19",
19          "_id": "tpEvsZMBhWjrNWL3RpuN",
20          "_score": 2,
21          "_source": {
22            "Disease Category": "Autoimmune",
23            "Incidence Rate (%)": 12.16,
24            "Mortality Rate (%)": 8.84,
25            "Doctors per 1000": 1.02,
26            "Education Index": 0.84,
27            "Healthcare Access (%)": 54.99,
28            "Recovery Rate (%)": 98.48,
29            "Treatment Type": "Surgery",
30            "Improvement in 5 Years (%)": 3.02,
31            "Gender": "Male",
32            "Population Affected": 417787,
33            "Per Capita Income (USD)": 50634,
34            "Hospital Beds per 1000": 9.85,
35            "Urbanization Rate (%)": 31.33,
36            "Year": 2000,
37            "Prevalence Rate (%)": 3.32,
38            "Disease Name": "Malaria",
39            "DALYs": 3052,
40            "Availability of Vaccines/Treatment": "No",
41            "Country": "Nigeria",
42            "Age Group": "36-60",
43            "Average Treatment Cost (USD)": 49373
44          }
45        },

```

Figure 4.10: output(partial) of query 4.5

4.6. Search for disease names allowing minor typos

Analysts searching for “Tuberculosis” might misspell it as “tuberculosi.” Fuzzy matching accommodates typographical errors, increasing recall for textual disease names.

4.6.1. query

```
69  GET globalhealth19/_search
70  {
71    "query": {
72      "match": {
73        "Disease Name": {
74          "query": "tuberculosi",
75          "fuzziness": "AUTO"
76        }
77      }
78    }
79 }
```

Figure 4.11: query 4.6

The `fuzziness:AUTO` helps compensate for potential misspellings

4.6.2. output(partial)

```

1   {
2     "took": 121,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": 2.4573126,
16      "hits": [
17        {
18          "_index": "globalhealth19",
19          "_id": "8ZIvsZMBhWjrNWL3ga3j",
20          "_score": 2.4573126,
21          "_source": {
22            "Disease Category": "Metabolic",
23            "Incidence Rate (%)": 14.07,
24            "Mortality Rate (%)": 8.15,
25            "Doctors per 1000": 2.02,
26            "Education Index": 0.77,
27            "Healthcare Access (%)": 95.73,
28            "Recovery Rate (%)": 77.03,
29            "Treatment Type": "Vaccination",
30            "Improvement in 5 Years (%)": 7.9,
31            "Gender": "Other",
32            "Population Affected": 367660,
33            "Per Capita Income (USD)": 68200,
34            "Hospital Beds per 1000": 5.44,
35            "Urbanization Rate (%)": 62.15,
36            "Year": 2001,
37            "Prevalence Rate (%)": 13.03,
38            "Disease Name": "Tuberculosis",
39            "DALYs": 3869,
40            "Availability of Vaccines/Treatment": "Yes",
41            "Country": "Saudi Arabia",
42            "Age Group": "61+",
43            "Average Treatment Cost (USD)": 36017
44          }
45        }
,
```

Figure 4.12: output(partial) of query 4.6

4.7. Calculate global average of doctors per 1000

Healthcare capacity metrics, like doctors per 1000 people, are crucial. This aggregation calculates an overall average across all documents to give a sense of how well-staffed healthcare systems are in the dataset.

4.7.1. query

```
82  GET globalhealth19/_search
83  {
84    "size": 0,
85    "aggs": {
86      "avg_doctors": {
87        "avg": {
88          "field": "Doctors per 1000"
89        }
90      }
91    }
92 }
```

Figure 4.13: query 4.7

4.7.2. output

```

1   {
2     "took": 68,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": null,
16      "hits": []
17    },
18    "aggregations": {
19      "avg_doctors": {
20        "value": 2.74792922
21      }
22    }
23  }

```

Figure 4.14: output of query 4.7

4.8. Compare vaccine availability categories side by side

It's common to compare how many documents indicate Availability of Vaccines/Treatment = "Yes" vs. "No". This parallel filters aggregation approach helps visualize coverage gaps or success rates in vaccine availability.

4.8.1. query

```
95   GET globalhealth19/_search
96   {
97     "size": 0,
98     "aggs": {
99       "vaccine_comparison": {
100         "filters": {
101           "filters": {
102             "has_vaccine": {
103               "term": {
104                 "Availability of Vaccines/Treatment": "Yes"
105               }
106             },
107             "no_vaccine": {
108               "term": {
109                 "Availability of Vaccines/Treatment": "No"
110               }
111             }
112           }
113         }
114       }
115     }
116   }
```

Figure 4.15: query 4.8

4.8.2. output

```
1  {
2    "took": 1,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
11     "total": {
12       "value": 10000,
13       "relation": "gte"
14     },
15     "max_score": null,
16     "hits": []
17   },
18   "aggregations": {
19     "vaccine_comparison": {
20       "buckets": {
21         "has_vaccine": {
22           "doc_count": 500354
23         },
24         "no_vaccine": {
25           "doc_count": 499646
26         }
27       }
28     }
29   }
30 }
```

Figure 4.16: output of query 4.8

4.9. Retrieve the latest record for each disease category

Each disease category may span multiple years. Grouping by Disease Category and returning only the doc from the latest Year is a common approach for summarizing the most recent data across categories.

4.9.1. query

```
120  {
121    "size": 0,
122    "aggs": {
123      "by_category": {
124        "terms": {
125          "field": "Disease Category",
126          "size": 10
127        },
128        "aggs": {
129          "latest_year": {
130            "top_hits": {
131              "sort": [
132                { "Year": { "order": "desc" } }
133              ],
134              "size": 1
135            }
136          }
137        }
138      }
139    }
140  }
141
```

Figure 4.17: query 4.9

4.9.2. output(partial)

```

18 "aggregations": {
19   "by_category": {
20     "doc_count_error_upper_bound": 0,
21     "sum_other_doc_count": 90445,
22     "buckets": [
23       {
24         "key": "Metabolic",
25         "doc_count": 91332,
26         "latest_year": {
27           "hits": {
28             "total": {
29               "value": 91332,
30               "relation": "eq"
31             },
32             "max_score": null,
33             "hits": [
34               {
35                 "_index": "globalhealth19",
36                 "_id": "P5oxsZMBhWjrNWL3UIZT",
37                 "_score": null,
38                 "_source": {
39                   "Disease Category": "Metabolic",
40                   "Incidence Rate (%)": 6.28,
41                   "Mortality Rate (%)": 4.95,
42                   "Doctors per 1000": 3.2,
43                   "Education Index": 0.61,
44                   "Healthcare Access (%)": 81.48,
45                   "Recovery Rate (%)": 71.49,
46                   "Treatment Type": "Vaccination",
47                   "Improvement in 5 Years (%)": 3.91,
48                   "Gender": "Male",
49                   "Population Affected": 642914,
50                   "Per Capita Income (USD)": 75213,
51                   "Hospital Beds per 1000": 0.88,
52                   "Urbanization Rate (%)": 48.67,
53                   "Year": 2024,
54                   "Prevalence Rate (%)": 6.02,
55                   "Disease Name": "Hepatitis",
56                   "DALYs": 4104,
57                   "Availability of Vaccines/Treatment": "No",
58                   "Country": "France",
59                   "Age Group": "61+",
60                   "Average Treatment Cost (USD)": 44697
61                 },
62                 "sort": [
63                   2024
64                 ]
65               }
66             ]
67           }
68         }
69       }
    ],
    "size": 0
  }
}

```

Figure 4.18: output(partial) of query 4.9

4.10. Find the costliest diseases by average treatment cost

Some diseases can be very expensive to treat. This query groups documents by disease name and calculates the average Average Treatment Cost (USD), sorting them in descending order to highlight the top costliest diseases in the dataset.

4.10.1. query

```
143   GET globalhealth19/_search
144   {
145     "size": 0,
146     "aggs": {
147       "costliest_diseases": {
148         "terms": {
149           "field": "Disease Name",
150           "size": 5,
151           "order": {
152             "avg_cost": "desc"
153           }
154         },
155         "aggs": {
156           "avg_cost": {
157             "avg": {
158               "field": "Average Treatment Cost (USD)"
159             }
160           }
161         }
162       }
163     }
164   }
```

Figure 4.19: query 4.10

4.10.2. output

```
18 "aggregations": {
19   "costliest_diseases": {
20     "doc_count_error_upper_bound": -1,
21     "sum_other_doc_count": 780796,
22     "buckets": [
23       {
24         "key": "Measles",
25         "doc_count": 39734,
26         "avg_cost": {
27           "value": 25195.426863643228
28         }
29       },
30       {
31         "key": "Influenza",
32         "doc_count": 49919,
33         "avg_cost": {
34           "value": 25138.26358701096
35         }
36       },
37       {
38         "key": "Cancer",
39         "doc_count": 40163,
40         "avg_cost": {
41           "value": 25129.76767173767
42         }
43       },
44       {
45         "key": "Hypertension",
46         "doc_count": 39418,
47         "avg_cost": {
48           "value": 25091.27550865087
49         }
50       },
51       {
52         "key": "Hepatitis",
53         "doc_count": 49970,
54         "avg_cost": {
55           "value": 25083.49457674605
56         }
57       }
58     ]
59   }
60 }
61 }
```

Figure 4.20: output of query 4.10

4.11. Group year data in intervals and keep high-affected populations only

This approach divides the Year field into multi-year buckets (e.g., every 5 years) and computes the average Population Affected. It then filters to retain only the buckets that exceed a certain threshold (e.g., >500 000). Researchers can spot time intervals with particularly large affected populations.

4.11.1. query

```
168 GET globalhealth19/_search
169 {
170   "size": 0,
171   "query": {
172     "bool": {
173       "must": [
174         { "term": { "Disease Category": "Chronic" } }
175       ]
176     }
177   },
178   "aggs": {
179     "year_hist": {
180       "histogram": {
181         "field": "Year",
182         "interval": 5
183       },
184       "aggs": {
185         "avg_population": {
186           "avg": {
187             "field": "Population Affected"
188           }
189         },
190         "population_filter": {
191           "bucket_selector": {
192             "buckets_path": {
193               "avgPop": "avg_population"
194             },
195             "script": "params.avgPop > 1000000"
196           }
197         }
198       }
199     }
200   }
201 }
```

Figure 4.21: query 4.11

4.11.2. output(partial)

```

1  {
2    "took": 2,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
11     "total": {
12       "value": 10000,
13       "relation": "gte"
14     },
15     "max_score": null,
16     "hits": []
17   },
18   "aggregations": {
19     "year_hist": {
20       "buckets": [
21         {
22           "key": 2005,
23           "doc_count": 18026,
24           "avg_population": {
25             "value": 501765.6664817486
26           }
27         },
28         {
29           "key": 2015,
30           "doc_count": 18195,
31           "avg_population": {
32             "value": 500894.8512228634
33           }
34         }
35       ]
36     }
37   }
38 }
```

Figure 4.22: output(partial) of query 4.11

4.12. Get extended statistics for the population affected by a specific disease

In many analyses, you want more than just an average or sum of the affected population. Extended statistics (e.g., min, max, variance, std dev) provide a deeper look at the distribution of the Population Affected field for a specific disease like "Diabetes". This helps epidemiologists or health data scientists see not only central tendencies but also variability within the affected population.

4.12.1. query

```
203  GET globalhealth19/_search
204  {
205    "size": 0,
206    "query": {
207      "term": {
208        "Disease Name": "Diabetes"
209      }
210    },
211    "aggs": {
212      "population_stats": {
213        "extended_stats": {
214          "field": "Population Affected"
215        }
216      }
217    }
218 }
```

Figure 4.23: query 4.12

4.12.2. output

```
1  {
2      "took": 1,
3      "timed_out": false,
4      "_shards": {
5          "total": 5,
6          "successful": 5,
7          "skipped": 0,
8          "failed": 0
9      },
10     "hits": {
11         "total": {
12             "value": 10000,
13             "relation": "gte"
14         },
15         "max_score": null,
16         "hits": []
17     },
18     "aggregations": {
19         "population_stats": {
20             "count": 50020,
21             "min": 1004,
22             "max": 999981,
23             "avg": 501187.25407836866,
24             "sum": 25069386449,
25             "sum_of_squares": 16734459843636860,
26             "variance": 83366711072.23279,
27             "variance_population": 83366711072.23279,
28             "variance_sampling": 83368377773.1079,
29             "std_deviation": 288732.94074669207,
30             "std_deviation_population": 288732.94074669207,
31             "std_deviation_sampling": 288735.8269649056,
32             "std_deviation_bounds": {
33                 "upper": 1078653.1355717527,
34                 "lower": -76278.62741501548,
35                 "upper_population": 1078653.1355717527,
36                 "lower_population": -76278.62741501548,
37                 "upper_sampling": 1078658.9080081799,
38                 "lower_sampling": -76284.39985144255
39             }
40         }
41     }
42 }
```

Figure 4.24: output of query 4.12

4.13. Prioritize documents from wealthier contexts

Some health analyses may start by examining diseases in regions with higher resources first. Here, each document's importance is influenced by its Per Capita Income (USD) for a specific disease category (like Respiratory), so that higher-income entries appear more prominently in the results.

4.13.1. query

```
221  GET globalhealth19/_search
222  {
223    "query": {
224      "function_score": {
225        "query": {
226          "term": {
227            "Disease Category": "Non-communicable"
228          }
229        },
230        "script_score": {
231          "script": {
232            "source": "doc['Per Capita Income (USD)'].value"
233          }
234        },
235        "boost_mode": "replace"
236      }
237    }
238  }
```

Figure 4.25: query 4.13

4.13.2. output(partial)

```

1   {
2     "took": 1343,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "skipped": 0,
8       "failed": 0
9     },
10    "hits": {
11      "total": {
12        "value": 10000,
13        "relation": "gte"
14      },
15      "max_score": 99999,
16      "hits": [
17        {
18          "_index": "globalhealth19",
19          "_id": "NJMvsZMBhWjrNWL3rF9V",
20          "_score": 99999,
21          "_source": {
22            "Disease Category": "Respiratory",
23            "Incidence Rate (%)": 10.33,
24            "Mortality Rate (%)": 6.62,
25            "Doctors per 1000": 2.4,
26            "Education Index": 0.71,
27            "Healthcare Access (%)": 56.34,
28            "Recovery Rate (%)": 67.93,
29            "Treatment Type": "Surgery",
30            "Improvement in 5 Years (%)": 3.99,
31            "Gender": "Other",
32            "Population Affected": 78699,
33            "Per Capita Income (USD)": 99999,
34            "Hospital Beds per 1000": 4.8,
35            "Urbanization Rate (%)": 61.2,
36            "Year": 2022,
37            "Prevalence Rate (%)": 16.91,
38            "Disease Name": "Parkinson's Disease",
39            "DALYs": 2479,
40            "Availability of Vaccines/Treatment": "Yes",
41            "Country": "Australia",
42            "Age Group": "36-60",
43            "Average Treatment Cost (USD)": 23134
44          }
45        }
5

```

Figure 4.26: output(partial) of query 4.13

4.14. Compare age-group mortality side by side

To investigate how mortality differs across age brackets (like "0-18" vs. "19-35" vs. "36-60" vs. "61+"), we set up separate conditions for each bracket. This yields parallel results and can highlight which diseases are more severe in each age range.

4.14.1. query

```
241   GET globalhealth19/_search
242   {
243     "size": 0,
244     "aggs": {
245       "age_comparison": {
246         "filters": {
247           "filters": {
248             "0_to_18": {
249               "term": { "Age Group": "0-18" }
250             },
251             "19_to_35": {
252               "term": { "Age Group": "19-35" }
253             },
254             "36_to_60": {
255               "term": { "Age Group": "36-60" }
256             },
257             "61_plus": {
258               "term": { "Age Group": "61+" }
259             }
260           }
261         },
262         "aggs": {
263           "avg_mortality": {
264             "avg": {
265               "field": "Mortality Rate (%)"
266             }
267           }
268         }
269       }
270     }
271 }
```

Figure 4.27: query 4.14

4.14.2. output

```

18 "aggregations": {
19   "age_comparison": {
20     "buckets": {
21       "0_to_18": {
22         "doc_count": 249605,
23         "avg_mortality": {
24           "value": 5.049840588129244
25         }
26       },
27       "19_to_35": {
28         "doc_count": 251201,
29         "avg_mortality": {
30           "value": 5.050999956210365
31         }
32     },
33     "36_to_60": {
34       "doc_count": 249205,
35       "avg_mortality": {
36         "value": 5.0527357797797
37       }
38     },
39     "61_plus": {
40       "doc_count": 249989,
41       "avg_mortality": {
42         "value": 5.046102588513895
43       }
44     }
45   }
46 }
47 }
48 }
```

Figure 4.28: output of query 4.14

4.15. Find diseases that stand out in a subpopulation

Within a certain subpopulation (e.g., Gender="Male"), some diseases might appear disproportionately often compared to the general dataset. Identifying these outliers provides insight into at-risk groups.

4.15.1. query

```
274   GET globalhealth19/_search
275   {
276     "query": {
277       "term": {
278         "Gender": "Male"
279       }
280     },
281     "size": 0,
282     "aggs": {
283       "significant_diseases": {
284         "significant_terms": {
285           "field": "Disease Name",
286           "size": 5
287         }
288       }
289     }
290   }
```

Figure 4.29: query 4.15

4.15.2. output

```

18 "aggregations": {
19     "significant_diseases": {
20         "doc_count": 333676,
21         "bg_count": 1000000,
22         "buckets": [
23             {
24                 "key": "Leprosy",
25                 "doc_count": 16974,
26                 "score": 0.0008186725862958704,
27                 "bg_count": 50064
28             },
29             {
30                 "key": "Dengue",
31                 "doc_count": 17004,
32                 "score": 0.000679556277939564,
33                 "bg_count": 50289
34             },
35             {
36                 "key": "Cancer",
37                 "doc_count": 16882,
38                 "score": 0.00031088804524608964,
39                 "bg_count": 50285
40             },
41             {
42                 "key": "Cholera",
43                 "doc_count": 16867,
44                 "score": 0.00030182709683728906,
45                 "bg_count": 50249
46             },
47             {
48                 "key": "Diabetes",
49                 "doc_count": 16783,
50                 "score": 0.0002788316099633605,
51                 "bg_count": 50020
52             }
53         ]
54     }
55 }
```

Figure 4.30: output of query 4.15

4.16. Highlight partial matches in disease names

When searching disease names by substring (e.g. “flu”), it can be helpful to highlight where that substring appears. This is particularly relevant for UI or user-facing search

results.

4.16.1. query

```
292 #16.Highlight partial matches in disease names
293 GET globalhealth19/_search
294 {
295   "query": {
296     "wildcard": {
297       "Disease Name": "*flu*"
298     }
299   },
300   "highlight": {
301     "fields": {
302       "Disease Name": {}
303     }
304   }
305 }
```

Figure 4.31: query 4.16

4.16.2. output(partial)

```

10    "hits": {
11        "total": {
12            "value": 10000,
13            "relation": "gte"
14        },
15        "max_score": 1,
16        "hits": [
17            {
18                "_index": "globalhealth19",
19                "_id": "xJEvsZMBhWjrNWL3RpyN",
20                "_score": 1,
21                "_source": {
22                    "Disease Category": "Autoimmune",
23                    "Incidence Rate (%)": 9.79,
24                    "Mortality Rate (%)": 8.9,
25                    "Doctors per 1000": 2.85,
26                    "Education Index": 0.81,
27                    "Healthcare Access (%)": 99.98,
28                    "Recovery Rate (%)": 86.15,
29                    "Treatment Type": "Medication",
30                    "Improvement in 5 Years (%)": 1.58,
31                    "Gender": "Other",
32                    "Population Affected": 968551,
33                    "Per Capita Income (USD)": 18834,
34                    "Hospital Beds per 1000": 8.7,
35                    "Urbanization Rate (%)": 28.83,
36                    "Year": 2004,
37                    "Prevalence Rate (%)": 5.47,
38                    "Disease Name": "Influenza",
39                    "DALYs": 3808,
40                    "Availability of Vaccines/Treatment": "No",
41                    "Country": "Argentina",
42                    "Age Group": "61+",
43                    "Average Treatment Cost (USD)": 33391
44                },
45                "highlight": {
46                    "Disease Name": [
47                        "<em>Influenza</em>"
48                    ]
49                }
50            },

```

Figure 4.32: output(partial) of query 4.16

4.17. Break down data by year to compute average mortality

Here, the Year field is numeric, so we treat each distinct year as a bucket and then calculate the average Mortality Rate (%). This is useful for year-by-year comparisons without a formal date field.

4.17.1. query

```
308   GET globalhealth19/_search
309   {
310     "size": 0,
311     "aggs": {
312       "yearly_distribution": {
313         "histogram": {
314           "field": "Year",
315           "interval": 1
316         },
317         "aggs": {
318           "avg_mortality": {
319             "avg": {
320               "field": "Mortality Rate (%)"
321             }
322           }
323         }
324       }
325     }
326   }
```

Figure 4.33: query 4.17

4.17.2. output(partial)

```

1  {
2    "took": 59,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
11     "total": {
12       "value": 10000,
13       "relation": "gte"
14     },
15     "max_score": null,
16     "hits": []
17   },
18   "aggregations": {
19     "yearly_distribution": {
20       "buckets": [
21         {
22           "key": 2000,
23           "doc_count": 40268,
24           "avg_mortality": {
25             "value": 5.054627992450581
26           }
27         },
28         {
29           "key": 2001,
30           "doc_count": 39896,
31           "avg_mortality": {
32             "value": 5.038353719671145
33           }
34         }
35       ],
36       "summary": {
37         "min": 2000,
38         "max": 2001,
39         "count": 80154
40       }
41     }
42   }
43 }
```

Figure 4.34: output(partial) of query 4.17

4.18. Dynamically compute the gap between health-care access and education

Policy analysts sometimes compare Healthcare Access (%) and the Education Index. By generating a difference on the fly for each document, it becomes easier to see where healthcare is outpacing or lagging behind education levels.

4.18.1. query

```
329   GET globalhealth19/_search
330   {
331     "_source": [ "Country", "Healthcare Access (%)", "Education
332       Index" ],
333     "script_fields": {
334       "access_education_gap": {
335         "script": {
336           "lang": "painless",
337           "source": "doc['Healthcare Access (%)'].value - doc
338             ['Education Index'].value"
339         }
340       },
341       "query": {
342         "match_all": {}
343     }
344   }
```

Figure 4.35: query 4.18

4.18.2. output(partial)

```

1  {
2    "took": 17,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
11     "total": {
12       "value": 10000,
13       "relation": "gte"
14     },
15     "max_score": 1,
16     "hits": [
17       {
18         "_index": "globalhealth19",
19         "_id": "sJEvsZMBhWjrNWL3RpuN",
20         "_score": 1,
21         "_source": {
22           "Education Index": 0.42,
23           "Healthcare Access (%)": 50.99,
24           "Country": "South Korea"
25         },
26         "fields": {
27           "access_education_gap": [
28             50.57
29           ]
30         }
31       },
32     ],
33     "sort": [
34       {
35         "field": "Country"
36       }
37     ]
38   }
39 }
```

Figure 4.36: output(partial) of query 4.18

4.19. Smooth data trends over time with a rolling average

Yearly metrics can fluctuate significantly. Applying a rolling average (e.g. 3-year window) to Hospital Beds per 1000 helps reveal long-term trends.

4.19.1. query

```
346   GET globalhealth19/_search
347   {
348     "size": 0,
349     "aggs": {
350       "year_buckets": {
351         "histogram": {
352           "field": "Year",
353           "interval": 1
354         },
355         "aggs": {
356           "beds_avg": {
357             "avg": {
358               "field": "Hospital Beds per 1000"
359             }
360           },
361           "moving_avg_beds": {
362             "moving_fn": {
363               "buckets_path": "beds_avg",
364               "window": 3,
365               "script": "MovingFunctions.unweightedAvg(values)"
366             }
367           }
368         }
369       }
370     }
371   }
```

Figure 4.37: query 4.19

4.19.2. output(partial)

```

1  {
2    "took": 107,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "skipped": 0,
8      "failed": 0
9    },
10   "hits": {
11     "total": {
12       "value": 10000,
13       "relation": "gte"
14     },
15     "max_score": null,
16     "hits": []
17   },
18   "aggregations": {
19     "year_buckets": {
20       "buckets": [
21         {
22           "key": 2000,
23           "doc_count": 40268,
24           "beds_avg": {
25             "value": 5.272970597000099
26           },
27           "moving_avg_beds": {
28             "value": null
29           }
30         },
31       ]
32     }
33   }
34 }
```

Figure 4.38: output(partial) of query 4.19

4.20. Filter results by cost or urbanization, then rank by year

Some researchers may want to see diseases where the treatment cost is above a threshold or the urbanization rate is high, while excluding unknown vaccine availability. Once filtered, results are ranked higher if they're from more recent years, making it easier to focus on the latest data.

4.20.1. query

```
374  GET globalhealth19/_search
375  {
376    "query": {
377      "function_score": {
378        "query": {
379          "bool": {
380            "should": [
381              {
382                "range": { "Average Treatment Cost (USD)": {
383                  "gt": 2000 } }
384                },
385                {
386                  "range": { "Urbanization Rate (%)": { "gt": 60 } }
387                }
388            ],
389            "must_not": [
390              {
391                "term": { "Availability of Vaccines/ Treatment": "Unknown" }
392              }
393            ],
394            "minimum_should_match": 1
395          }
396        },
397        "script_score": {
398          "script": {
399            "source": "doc['Year'].value"
400          }
401        },
402        "boost_mode": "multiply"
403      }
404    }
```

Figure 4.39: query 4.20

4.20.2. output(partial)

```

1  {
2      "took": 1295,
3      "timed_out": false,
4      "_shards": {
5          "total": 5,
6          "successful": 5,
7          "skipped": 0,
8          "failed": 0
9      },
10     "hits": {
11         "total": {
12             "value": 10000,
13             "relation": "gte"
14         },
15         "max_score": 4048,
16         "hits": [
17             {
18                 "_index": "globalhealth19",
19                 "_id": "PZEvsZMBhWjrNWL3RqCN",
20                 "_score": 4048,
21                 "_source": {
22                     "Disease Category": "Respiratory",
23                     "Incidence Rate (%)": 7.77,
24                     "Mortality Rate (%)": 7,
25                     "Doctors per 1000": 1.04,
26                     "Education Index": 0.44,
27                     "Healthcare Access (%)": 75.12,
28                     "Recovery Rate (%)": 66.06,
29                     "Treatment Type": "Vaccination",
30                     "Improvement in 5 Years (%)": 6.77,
31                     "Gender": "Female",
32                     "Population Affected": 526173,
33                     "Per Capita Income (USD)": 17993,
34                     "Hospital Beds per 1000": 9.02,
35                     "Urbanization Rate (%)": 80.42,
36                     "Year": 2024,
37                     "Prevalence Rate (%)": 9.29,
38                     "Disease Name": "COVID-19",
39                     "DALYs": 904,
40                     "Availability of Vaccines/Treatment": "No",
41                     "Country": "China",
42                     "Age Group": "61+",
43                     "Average Treatment Cost (USD)": 26968
44                 }
45             }
46         ],
47     }
48 }
```

Figure 4.40: output(partial) of query 4.20

5 | Dashboard Description

<https://calcio-c70f89.kb.us-east-1.aws.elastic.cloud/app/r/s/ESwdn>

5.1. Top 3 Diseases Distribution by Age Group

This chart displays the distribution of the top 3 disease categories for each age group, segmented by 0–18, 19–35, 36–60, and 61+. The inner segments show the proportion of individuals affected, while the outer segments break down the disease categories.

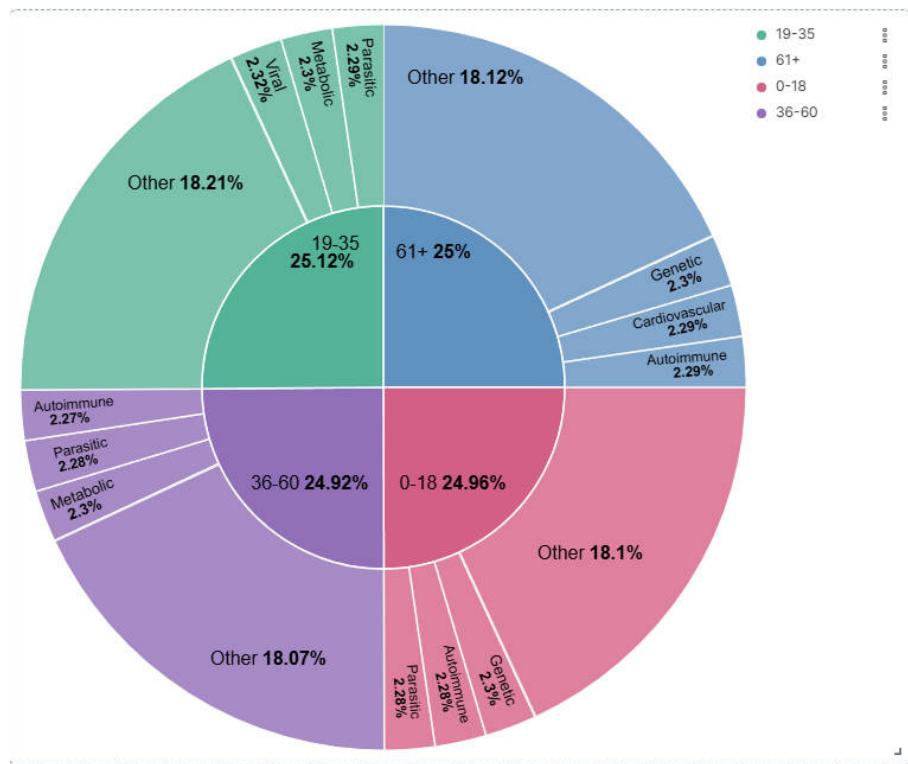


Figure 5.1: Top 3 Diseases Distribution by Age Group

5.2. Top 10 Countries by Alzheimer's Disease Distribution

A treemap chart highlighting the top 10 countries based on the proportion of Alzheimer's disease cases. Each segment represents a country, with the percentage denoting its share.

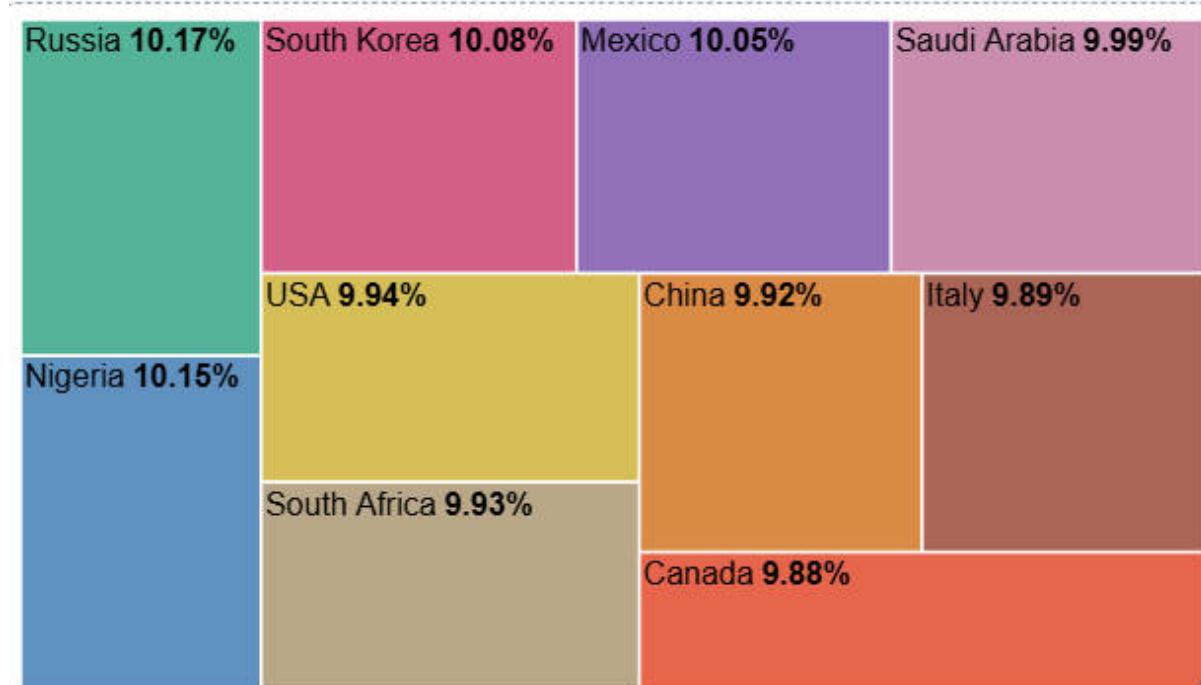


Figure 5.2: Top 10 Countries by Alzheimer's Disease Distribution

5.3. 2000–2024 Count of Records Distribution by Age Group

A line chart showing the trends in record counts across different age groups (0–18, 19–35, 36–60, and 61+) over the years 2000–2024, indicating variations in population affected.



Figure 5.3: 2000–2024 Count of Records Distribution by Age Group

5.4. Top 5 Disease Categories in 5 Countries

A stacked bar chart showcasing the distribution of the top 5 disease categories in five countries. Each bar represents a country, with colored segments denoting the disease categories.

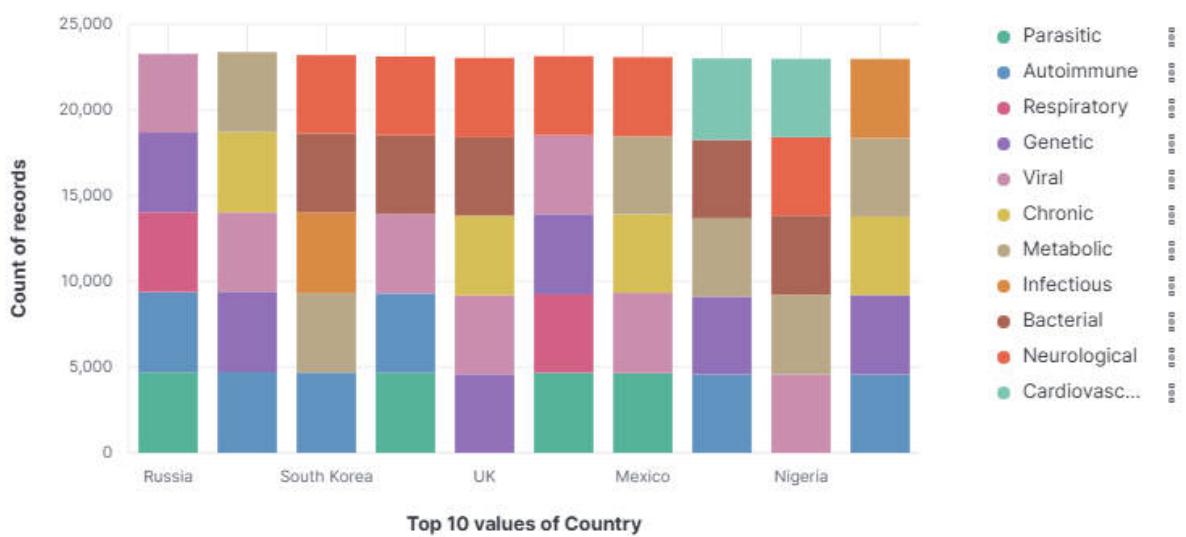


Figure 5.4: Top 5 Disease Categories in 5 Countries

5.5. 2014–2024 Changes in Alzheimer's Disease Counts

A line graph showing the trends in Alzheimer's disease case counts from 2014 to 2024, highlighting changes over time.

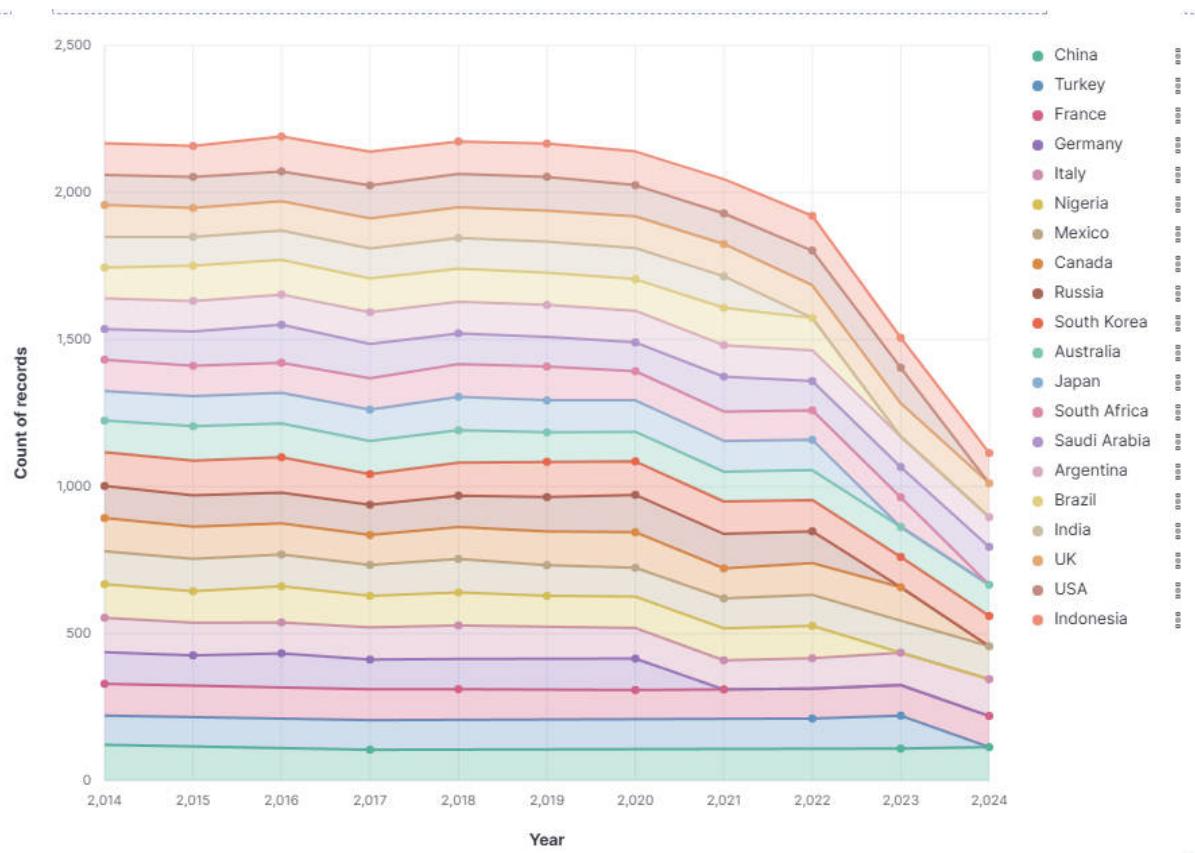


Figure 5.5: 2014–2024 Changes in Alzheimer’s Disease Counts

5.6. Top 10 Countries by Minimum Average Treatment Cost (USD)

A word cloud visualizing the top 10 countries with the lowest average treatment costs in USD. The size of the country name correlates with its rank, with larger text indicating lower costs.



Top 10 values of Country - Minimum of Average Treatment Cost (USD)

Figure 5.6: Top 10 Countries by Minimum Average Treatment Cost (USD)