



EMSE 6586

DATABASE CREATION AND ANALYSIS OF NOBEL PRIZE WINNERS

Anbu Ezhilmathi Nambi
G33350186

INTRODUCTION

- The Nobel Prize is an esteemed international accolade awarded yearly across multiple fields for exceptional achievements.
- The dataset encompasses information on more than 900 laureates dating back to 1901.
- The goal of this project is to construct a well-organized SQL database from the Nobel Prize dataset using Python scripts.
- This database facilitates the analysis of the dataset, allowing for the identification of trends and patterns among laureates across different prize categories and eras.
- The aim of this project is to provide a comprehensive understanding of the dataset and reveal valuable insights.

DATA TRANSLATION

01 FETCH JSON DATA

Fetching Data from Nobel Prize APIs

02 PARSE JSON DATA

After fetching the JSON data, it is parsed to extract the necessary information.

03 CREATE SQLITE DATABASE SCHEMA

Set up a SQLite database to store the parsed data

04 CREATE DATABASE TABLES

Design and create database tables to efficiently store and query the Nobel Prize data.

ABOUT THE DATA

- The dataset contains detailed information on Nobel laureates from 1901 to 2023.
- It includes data from all six Nobel Prize categories: Peace, Literature, Chemistry, Physics, Medicine, and Economic Sciences.
- Each record includes details about the laureates such as names, birthdates, birthplaces, and affiliations.
- The dataset provides specifics on the prize, including the year of the award, the motivations for each prize, and information on instances where multiple laureates shared a prize.
- The data is in JSON format, which simplifies the process of data handling and analysis.

1. Fetch the JSON Data from API

```
1 # Step 1: Fetch JSON data
2 # prize
3 response1 = requests.get("https://api.nobelprize.org/v1/prize.json")
4 data1 = response1.json()
5 # laureate
6 response2 = requests.get("https://api.nobelprize.org/v1/laureate.json")
7 data2 = response2.json()
```

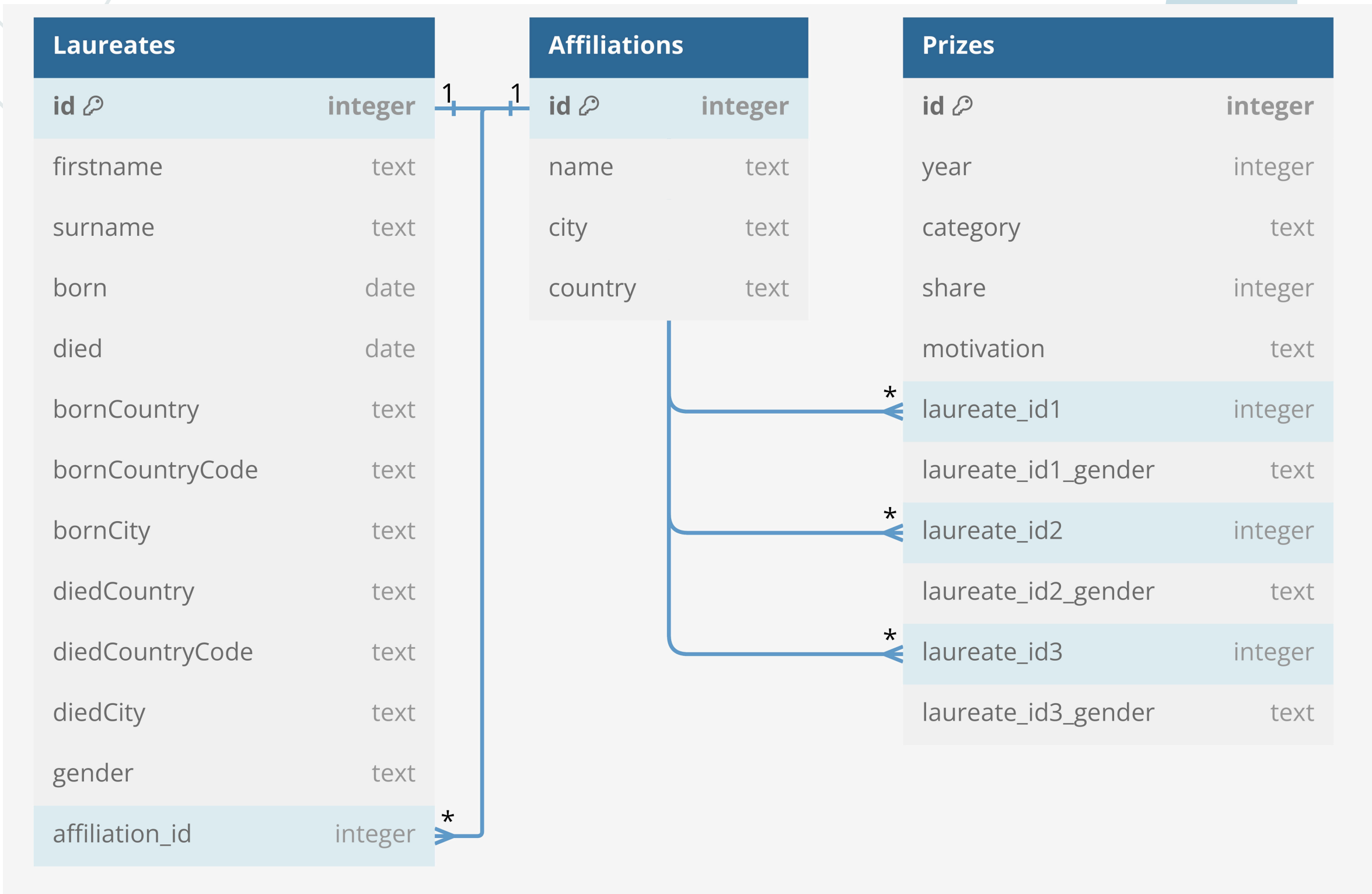
2. Sample of Prize data.json

```
1 {
2   "id": "1",
3   "firstname": "Wilhelm Conrad",
4   "surname": "R\u00f6ntgen",
5   "born": "1845-03-27",
6   "died": "1923-02-10",
7   "bornCountry": "Prussia (now Germany)",
8   "bornCountryCode": "DE",
9   "bornCity": "Lennep (now Remscheid)",
10  "diedCountry": "Germany",
11  "diedCountryCode": "DE",
12  "diedCity": "Munich",
13  "gender": "male",
14  "prizes": [
15    {
16      "year": "1901",
17      "category": "physics",
18      "share": "1",
19      "motivation": "\"in recognition of the extraordinary services he has rendered by the discovery of the remarkable rays subsequently named after him\"",
20      "affiliations": [
21        {
22          "name": "Munich University",
23          "city": "Munich",
24          "country": "Germany"
25        }
26      ]
27    }
28  ]
29 }
```

3. Sample of Laureate data.json

```
1 {
2   "year": "2023",
3   "category": "chemistry",
4   "laureates": [
5     {
6       "id": "1029",
7       "firstname": "Moungi",
8       "surname": "Bawendi",
9       "motivation": "\"for the discovery and synthesis of quantum dots\"",
10      "share": "3"
11    },
12    {
13      "id": "1030",
14      "firstname": "Louis",
15      "surname": "Brus",
16      "motivation": "\"for the discovery and synthesis of quantum dots\"",
17      "share": "3"
18    },
19    {
20      "id": "1031",
21      "firstname": "Aleksey",
22      "surname": "Yekimov",
23      "motivation": "\"for the discovery and synthesis of quantum dots\"",
24      "share": "3"
25    }
26  ]
27 }
```

ENTITY RELATIONSHIP DIAGRAM



1. Creating SQLite Database

```
1 # Step 3: Create SQLite database schema
2 conn = sqlite3.connect('nobel_prizes.db')
3 cursor = conn.cursor()
```

2. Create Laureates table

```
1 # Step 4: Create database tables
2
3 # Create Laureates table
4 cursor.execute('''CREATE TABLE IF NOT EXISTS Laureates (
5     id INTEGER PRIMARY KEY,
6     firstname TEXT,
7     surname TEXT,
8     born DATE,
9     died DATE,
10    bornCountry TEXT,
11    bornCountryCode TEXT,
12    bornCity TEXT,
13    diedCountry TEXT,
14    diedCountryCode TEXT,
15    diedCity TEXT,
16    gender TEXT,
17    affiliation_id INTEGER
18    )''')
```

3. Create Affiliations table

```
1 # Create Affiliations table
2 cursor.execute('''CREATE TABLE Affiliations (
3     id INTEGER PRIMARY KEY AUTOINCREMENT,
4     name TEXT,
5     city TEXT,
6     country TEXT
7     )''')
```

4. Create Prize table

```
1 # Create Prizes table
2 cursor.execute('''CREATE TABLE Prizes (
3     id INTEGER PRIMARY KEY AUTOINCREMENT,
4     year INTEGER,
5     category TEXT,
6     share INTEGER,
7     motivation TEXT,
8     laureate_id1 INTEGER,
9     laureate_id1_gender TEXT,
10    laureate_id2 INTEGER,
11    laureate_id2_gender TEXT,
12    laureate_id3 INTEGER,
13    laureate_id3_gender TEXT
14    )''')
```

1. Inserting Values into Laureates Table

```
1 for laureate in laureate_data:
2     cursor.execute('''INSERT INTO Laureates VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?)''',
3                     (laureate['id'],
4                       laureate.get('firstname', None),
5                       laureate.get('surname', None),
6                       laureate.get('born', None),
7                       laureate.get('died', None),
8                       laureate.get('bornCountry', None),
9                       laureate.get('bornCountryCode', None),
10                      laureate.get('bornCity', None),
11                      laureate.get('diedCountry', None),
12                      laureate.get('diedCountryCode', None),
13                      laureate.get('diedCity', None),
14                      laureate.get('gender', None),
15                      None))
```

2. Inserting Values into Affiliations Table

```
1 cursor.execute("INSERT INTO Affiliations (name, city, country) VALUES (?, ?, ?)", (name, city, country))
```

3. Inserting Values into Prizes Table

```
1 cursor.execute("INSERT INTO Prizes (year, category, share, motivation) VALUES (?, ?, ?, ?)", (year, category, share, motivation))
```


SAMPLE OUTPUT

1 Laureates who won more than one Nobel Prize:

	firstname	surname	prize_count
0	Marie	Curie	2
1	John	Bardeen	2
2	Linus	Pauling	2
3	Frederick	Sanger	2
4	International Committee of the Red Cross		3
5	Office of the United Nations High Commissioner...		2
6	Barry	Sharpless	2

2 First Females to win Nobel Prize in each category

	Year	Full Name	Category	Birth Country
0	1903	Marie Curie	physics	Russian Empire (now Poland)
1	1905	Bertha von Suttner	peace	Austrian Empire (now Czech Republic)
2	1909	Selma Lagerlöf	literature	Sweden
3	1911	Marie Curie	chemistry	Russian Empire (now Poland)
4	1947	Gerty Cori	medicine	Austria-Hungary (now Czech Republic)
5	2009	Elinor Ostrom	economics	USA

3 Youngest Nobel Laureates

	Age	Full Name	Year	Category
0	17	Malala Yousafzai	2014	peace
1	25	Lawrence Bragg	1915	physics
2	31	Carl D. Anderson	1936	physics
3	31	Paul A.M. Dirac	1933	physics
4	31	Tsung-Dao Lee	1957	physics
5	31	Werner Heisenberg	1932	physics

4 Top 10 Affiliations that won Nobel Prizes:

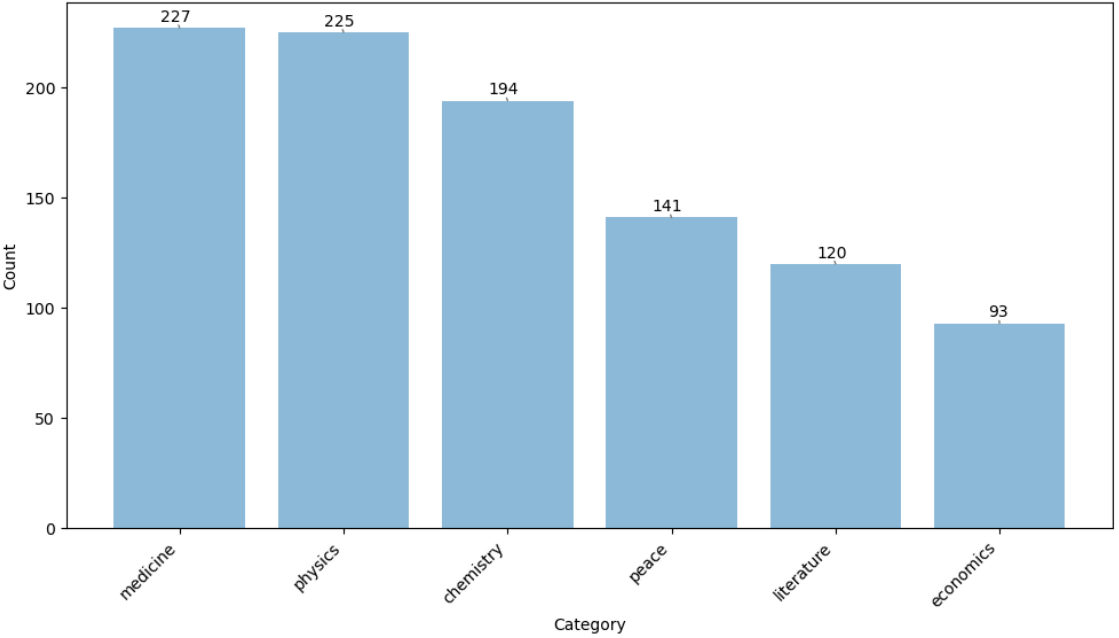
	Affiliation Name	Prize Count
0	University of California	36
1	Harvard University	28
2	Massachusetts Institute of Technology (MIT)	23
3	Stanford University	22
4	University of Chicago	19
5	California Institute of Technology (Caltech)	19
6	University of Cambridge	18
7	Columbia University	18
8	Princeton University	17
9	Rockefeller University	13

5 Gender Distribution: A Trend Analysis from 2009 to 2023

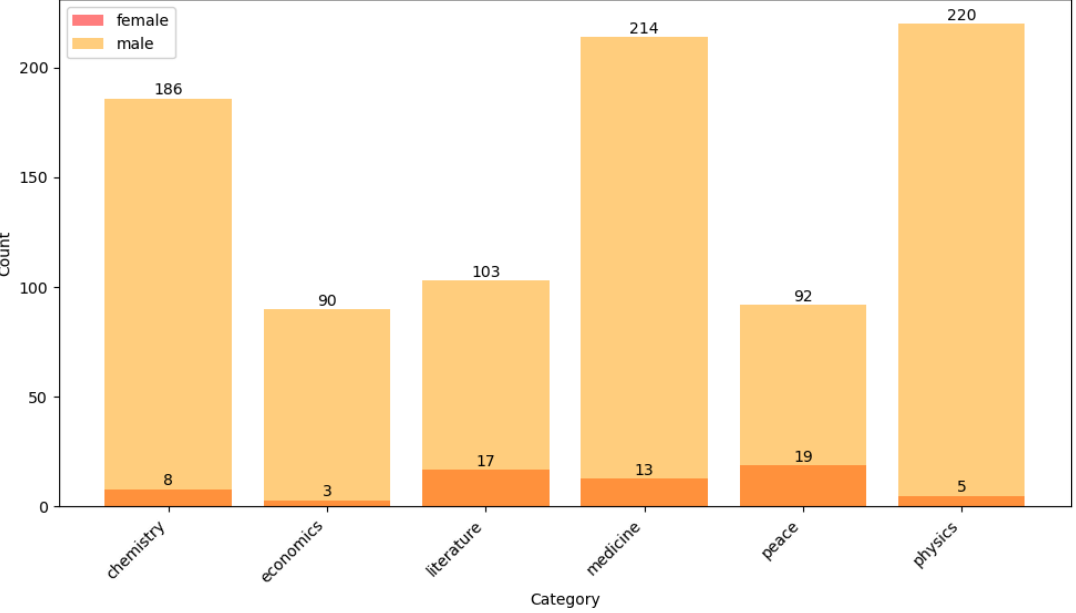
gender	female	male	org	Total	Male %	Female %
year						
2009	5	8	0	13	61.538462	38.461538
2010	0	11	0	11	100.000000	0.000000
2011	3	10	0	13	76.923077	23.076923
2012	0	9	1	10	90.000000	0.000000
2013	1	11	1	13	84.615385	7.692308
2014	2	11	0	13	84.615385	15.384615
2015	2	8	1	11	72.727273	18.181818
2016	0	11	0	11	100.000000	0.000000
2017	0	11	1	12	91.666667	0.000000
2018	4	9	0	13	69.230769	30.769231
2019	1	13	0	14	92.857143	7.142857
2020	4	7	1	12	58.333333	33.333333
2021	1	12	0	13	92.307692	7.692308
2022	2	10	2	14	71.428571	14.285714
2023	4	7	0	11	63.636364	36.363636

SAMPLE OUTPUT

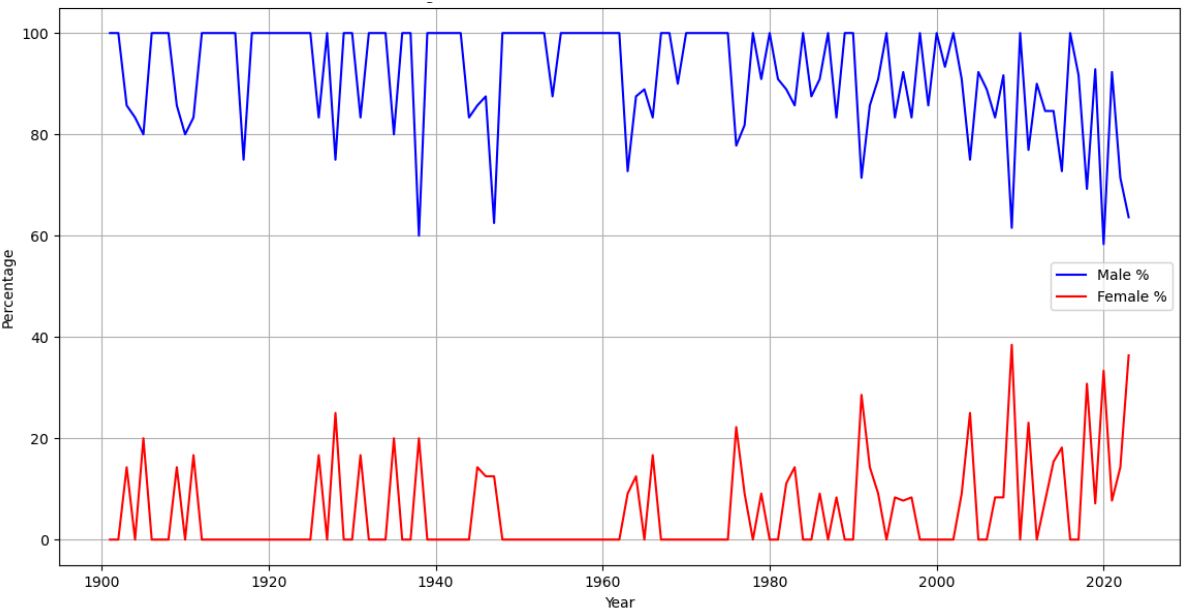
1. Prize by Category



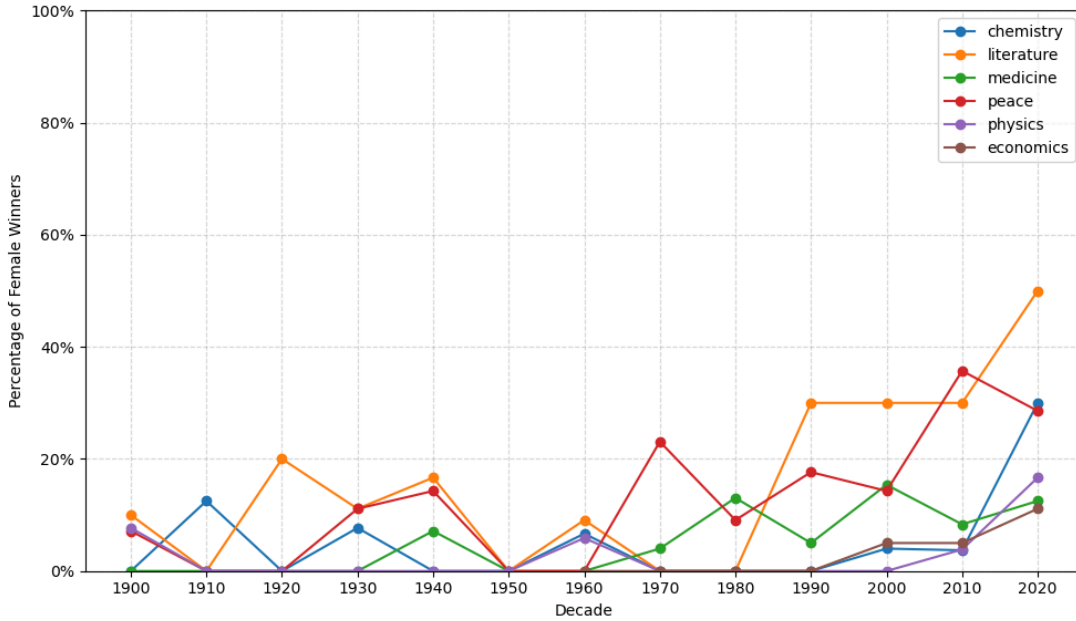
2. Prize won by Category and Gender



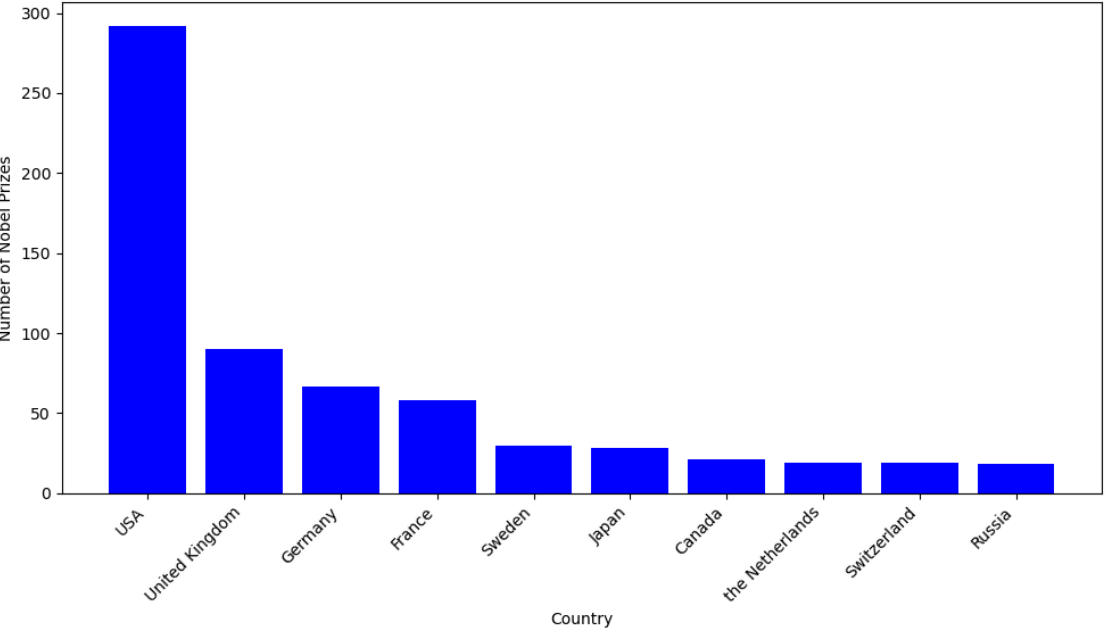
3. Percentage of Male vs Female Nobel Laureates Over the Years



4. Percentage Female Winners by Decade and Category



5. Top 10 Countries with Most Nobel Prizes by Birth Country of Laureates



CHALLENGES

- When working with Nobel Prize data from different databases, it's common to face inconsistencies and missing values. This makes it crucial to ensure data consistency, handle missing values, and clean up inconsistencies such as empty strings, 'None', or 'Unknown' values, as seen in SQL queries. It can be a challenging task, but it's essential for data cleaning, integration, and integrity.
- Integrating data from various sources like laureates, affiliations, and prizes is complex. Effective table joining and handling different data formats or schemas require careful attention to ensure successful data integration.
- Maintaining data quality requires ensuring the accuracy and integrity of data. To achieve this, it is important to validate data entries, verify relationships between entities, and prevent duplicate records. In particular, avoiding duplicated data entries is crucial as data is often collected from multiple columns that may have overlapping content.
- It is crucial to collect complete data that includes laureate affiliation, country, and prize category to avoid any impact on the analysis results due to missing data. Taking necessary measures to prevent data loss is important.

OVERCOME CHALLENGES

- To ensure consistency across datasets, use data cleaning techniques to standardize data formats. This involves trimming whitespace, converting text to a consistent case, and handling special characters.
- Another important aspect to consider is replacing missing values, which can be done using appropriate placeholders or statistical methods to impute missing values.
- Validating data is crucial to ensure consistent and accurate information in the database.
- To prevent duplicate entries in a database, it's important to establish procedures for detecting and resolving them. You can achieve this by using SQL's GROUP BY and HAVING clauses to find repeated entries or implementing uniqueness constraints within the database design.

CONCLUSION

- Constructed a well-organized SQL database from the Nobel Prize dataset using Python scripts.
- Dataset encompasses information on over 900 laureates spanning from 1901 to 2023 across six Nobel Prize categories.
- Facilitates analysis to identify trends and patterns among laureates, enabling a comprehensive understanding of the dataset and revealing valuable insights.
- Managed inconsistencies, missing values, and data integration complexities through data cleaning, validation, and de-duplication techniques.
- Ensured accuracy and integrity by standardizing formats, replacing missing values, and implementing procedures for detecting and resolving duplicates.

The background features four decorative geometric patterns in the corners. The top-left corner has a series of parallel diagonal lines in a light blue-grey color. The top-right corner contains a cluster of overlapping semi-circles in yellow, red, teal, and dark blue. The bottom-left corner also features a cluster of overlapping semi-circles in red, teal, and dark blue. The bottom-right corner has a series of parallel diagonal lines in a light blue-grey color, mirroring the top-left pattern.

THANK YOU