

Practical Assignment:
Coursework
Designing and Implementing a Database Structure

[Mohammad Faraj P418072]
Database Design and Implementation
16/4/2024

I confirm that this assignment is my own work. Where I have referred to academic sources, I have provided in-text citations and included the sources in the final reference list.

Introduction and Context

The "Internet Movies Database" website, IMDb, was scraped to gather a list of movies that was then compiled into a CSV file named "movies.csv". The objective was to standardize the data in the file and add it to a database using Python code and SQL statements. This method ensures accurate access to the information by utilizing SQL statements. The report will encompass the normalization process, entity-relationship diagram, code implemented to generate tables and input values, SQL statements utilized to extract information and the corresponding outputs, and a reflective segment that will evaluate the design's strengths, weaknesses, and potential areas of improvement.

Normalization Process

To begin the process, it is important to normalize the CSV file housing an extensive 15-column movie database. These columns comprise of crucial details such as the movie's name, rating, genre, year, country-specific release information, score, votes, director, writer, star, shooting location, budget, gross income, production company, and runtime. Given the diverse nature of this information, it is essential to ensure proper normalization to avoid any issues such as update, deletion, or insertion anomalies.

Normalizing to First Normal Form (1NF):

Ensuring that each row is unique, all values are atomic, and there are no repeating groups in the table is crucial for bringing it to the first normal form. Although there are no repeating groups in the table, the "released" column contains two types of data, date and country, which can still be considered atomic, since it does not affect any other data in the table. However, there are non-unique rows in the table, such as two movies with the same name "The Captain," and there may be one production company that produced multiple movies. As a result, unique IDs must be created, and the table must be split into two.

The table will be split into two distinct tables, the first table, "movie_production", will consist of crucial information such as the movie's name, director, stars, writer, and production company. This data can have an impact on other records, as a single writer may create scripts for multiple films, resulting in a one-to-many relationship with the movie's name. However, the name can be used to identify the production company and other pertinent information only if it is unique. Thus, to maintain uniqueness, unique IDs must be created for these fields, alongside a primary key named "production ID".

The second table however, "movie_general", will include all secondary information that does not depend on or rely on any other information, like the rating, genre, year, released, votes, score, country, budget, gross, and runtime. With a

unique ID named “general ID” to serve as a primary key. It will be linked to the first table using the “production ID” however, it shall be a foreign key in the “movie_general” table referencing the “movie_production” table.

Normalizing to Second Normal Form (2NF):

To convert tables to the second normal form (2NF), all non-key attributes in the tables must rely on the primary key after applying normalization to the first normal form (1NF). The "movie_general" table already meets this requirement, since all attributes depend on the primary key, "general ID". However, in the "movie_production" table, each non-key attribute is dependent on its ID, not the primary key "production ID". As a result, the table should be split into six new tables, each with only one attribute and its corresponding ID, with the ID serving as the primary key. There will also be one "movie_production" table that has the "production ID" as its primary key and includes all the other primary keys from the other five tables, to serve as foreign keys to link all tables to it.

Now, the tables are not just in 2NF, but they are also in 3NF since the 3NF points that there must not be any transitive dependencies in the tables, and there are not.

Entity-relationship Diagram

Here is an entity-relationship diagram, Figure 1, that presents and clarifies the relation between the tables mentioned before.

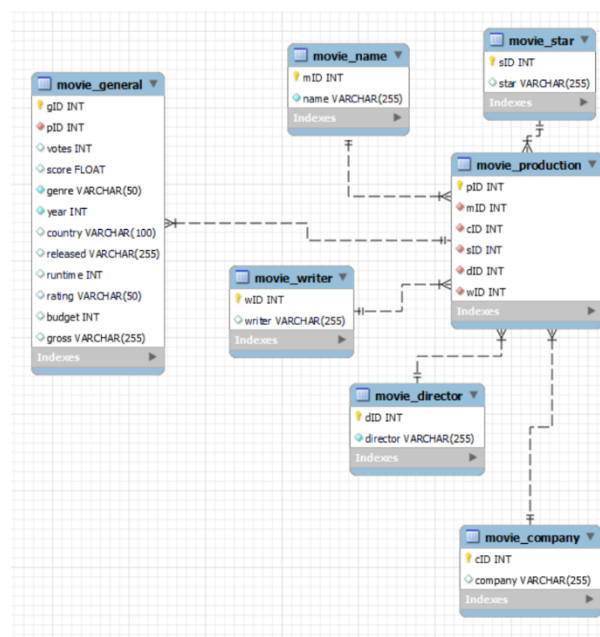


Figure 1: An Entity-relationship Diagram

As how is clear in Figure 1, the relation between the five tables containing the non-key attributes and their IDs, and the “movie_production” tables is a non-identifying one-to-many (1:n) relationship. While the relation between the “movie_production” table and the “movie_general” is also a non-identifying one-to-many relationship.

SQL Database Implementation using Python

In this section, the Python code, along with the SQL statements, used to create the database “coursework1” and the tables, will be presented in Figures 2 and 3, as in below.

```
3 import mysql.connector
4 from mysql.connector import Error
5
6 # Establish database connection
7 connection = mysql.connector.connect(
8     host="localhost",
9     user="root",
10    password="2542",
11    database='coursework1'
12 )
13
14 cursor = connection.cursor()
15
16 # cursor.execute("CREATE DATABASE IF NOT EXISTS coursework1")
17 # cursor.execute("SHOW DATABASES")
18 # for x in cursor:
19 #     print(x)
20
21 # cursor.execute("CREATE TABLE movie_name (\
22 # mID int NOT NULL,\
23 # name varchar(255) NOT NULL,\
24 # primary key (mID)\
25 # );")
26
27 # cursor.execute("create table movie_company (\
28 # cID int NOT NULL,\
29 # company varchar(255),\
30 # primary key (cID)\
31 # );")
32
33 # cursor.execute("create table movie_star (\
34 # sID int NOT NULL,\
35 # star varchar(255),\
36 # primary key (sID)\
37 # );")
38
39 # cursor.execute("create table movie_director (\
40 # dID int NOT NULL,\
41 # director varchar(255) NOT NULL,\
42 # primary key (dID)\
43 # );")
44
```

Figure 2: Code for creating the database and tables, displayed in Visual Studio

```

45 # cursor.execute("create table movie_writer (\
46 # wID int NOT NULL,\
47 # writer varchar(255),\
48 # primary key (wID)\
49 # );")
50
51 # cursor.execute("create table movie_production (\
52 # pID int NOT NULL,\
53 # mID int NOT NULL,\
54 # cID int NOT NULL,\
55 # sID int NOT NULL,\
56 # dID int NOT NULL,\
57 # wID int NOT NULL,\
58 # primary key (pID),\
59 # foreign key (mID) references movie_name (mID),\
60 # foreign key (cID) references movie_company (cID),\
61 # foreign key (sID) references movie_star (sID),\
62 # foreign key (dID) references movie_director (dID),\
63 # foreign key (wID) references movie_writer (wID)\
64 # );")
65
66 # cursor.execute("create table movie_general (\
67 # gID int NOT NULL,\
68 # pID int NOT NULL,\
69 # votes int,\
70 # score float,\
71 # genre varchar(50) NOT NULL,\
72 # year int NOT NULL,\
73 # country varchar(100),\
74 # released varchar(255),\
75 # runtime int,\
76 # rating varchar(50),\
77 # budget int,\
78 # gross varchar(255),\
79 # primary key (gID),\
80 # foreign key (pID) references movie_production (pID)\
81 # );")
82
83 connection.commit()
84
85 cursor.close()
86 connection.close()

```

Figure 3: Completion of the Python code

Data Importing using Python

In Figures 4 through 8, there will be the code for importing information from the CSV file “movies” and inserting them into the tables created before.

```

database.py > ...
1 # Importing the random module to generate random numbers
2 import random
3 # Importing the pandas module and aliasing it as pd for ease of use
4 import pandas as pd
5 # Importing the mysql.connector module to connect to the MySQL database
6 import mysql.connector
7 # Importing the Error class from mysql.connector module to handle errors
8 from mysql.connector import Error
9
10 # Function to generate a unique mID
11 def generate_mID(cursor):
12     # Infinite loop to generate a unique mID
13     while True:
14         # Generate a random mID between 10000 and 99999
15         mID = random.randint(10000, 99999)
16         # Check if the generated mID already exists in the movie_name table
17         cursor.execute("SELECT mID FROM movie_name WHERE mID = %s", (mID,))
18         # If mID does not exist, it's unique, so return it
19         if not cursor.fetchone():
20             return mID
21
22 # Function to generate a unique cID
23 def generate_cID(cursor):
24     # Infinite loop to generate a unique cID
25     while True:
26         # Generate a random cID between 10000 and 99999
27         cID = random.randint(10000, 99999)
28         # Check if the generated cID already exists in the movie_company table
29         cursor.execute("SELECT cID FROM movie_company WHERE cID = %s", (cID,))
30         # If cID does not exist, it's unique, so return it
31         if not cursor.fetchone():
32             return cID
33
34 # Function to generate a unique sID
35 def generate_sID(cursor):
36     # Infinite loop to generate a unique sID
37     while True:
38         # Generate a random sID between 10000 and 99999
39         sID = random.randint(10000, 99999)
40         # Check if the generated sID already exists in the movie_star table
41         cursor.execute("SELECT sID FROM movie_star WHERE sID = %s", (sID,))
42         # If sID does not exist, it's unique, so return it
43         if not cursor.fetchone():
44             return sID
45

```

Figure 4: Python code for importing information displayed in Visual Studio

```

46 # Function to generate a unique dID
47 def generate_dID(cursor):
48     # Infinite loop to generate a unique dID
49     while True:
50         # Generate a random dID between 10000 and 99999
51         dID = random.randint(10000, 99999)
52         # Check if the generated dID already exists in the movie_director table
53         cursor.execute("SELECT dID FROM movie_director WHERE dID = %s", (dID,))
54         # If dID does not exist, it's unique, so return it
55         if not cursor.fetchone():
56             return dID
57
58 # Function to generate a unique wID
59 def generate_wID(cursor):
60     # Infinite loop to generate a unique wID
61     while True:
62         # Generate a random wID between 10000 and 99999
63         wID = random.randint(10000, 99999)
64         # Check if the generated wID already exists in the movie_writer table
65         cursor.execute("SELECT wID FROM movie_writer WHERE wID = %s", (wID,))
66         # If wID does not exist, it's unique, so return it
67         if not cursor.fetchone():
68             return wID
69
70 # Function to generate a unique pID
71 def generate_pID(cursor):
72     # Infinite loop to generate a unique pID
73     while True:
74         # Generate a random pID between 10000 and 99999
75         pID = random.randint(10000, 99999)
76         # Check if the generated pID already exists in the movie_production table
77         cursor.execute("SELECT pID FROM movie_production WHERE pID = %s", (pID,))
78         # If pID does not exist, it's unique, so return it
79         if not cursor.fetchone():
80             return pID
81

```

Figure 5: The rest of the code

```

82 # Function to generate a unique gID
83 def generate_gID(cursor):
84     # Infinite loop to generate a unique gID
85     while True:
86         # Generate a random gID between 10000 and 99999
87         gID = random.randint(10000, 99999)
88         # Check if the generated gID already exists in the movie_general table
89         cursor.execute("SELECT gID FROM movie_general WHERE gID = %s", (gID,))
90         # If gID does not exist, it's unique, so return it
91         if not cursor.fetchone():
92             return gID
93
94 # Try block to handle potential errors
95 try:
96     # Establish database connection
97     connection = mysql.connector.connect(
98         host="localhost",
99         user="root",
100         password="2542",
101         database="coursework1"
102     )
103
104     # Creating a cursor object to interact with the database
105     cursor = connection.cursor()
106
107     # Read movie data from CSV file using pandas
108     movieData = pd.read_csv("movies.csv")
109
110     # Iterate over each row in the CSV
111     for index, row in movieData.iterrows():
112         # Generate unique IDs for each table
113         mID = generate_mID(cursor)
114         cID = generate_cID(cursor)
115         sID = generate_sID(cursor)
116         dID = generate_dID(cursor)
117         wID = generate_wID(cursor)
118         pID = generate_pID(cursor)
119         gID = generate_gID(cursor)
120
121         # Replace 'nan' values with None
122         row = row.where(pd.notnull(row), None)
123

```

Figure 6: The rest of the code

As in Figures 4, 5, and 6, there are some functions to generate the unique IDs needed for the tables, where the “random” library was used to generate the random number for each ID, and it only generates a unique ID for each row, by checking if it is present already in the table or not.

Also to deal with ‘nan’ values, this line of code: “row = row.where(pd.notnull(row), None)”, replaces ‘nan’ values in a Pandas DataFrame row with ‘None’, using the “where()” function from the Pandas library. (Geeks for Geeks, 2023)

```

124 # Insert values into movie_name table
125 cursor.execute("INSERT INTO movie_name (mID, name) VALUES (%s, %s)", (mID, row['name']))
126
127 # Insert values into movie_company table
128 cursor.execute("INSERT INTO movie_company (cID, company) VALUES (%s, %s)", (cID, row['company']))
129
130 # Insert values into movie_star table
131 cursor.execute("INSERT INTO movie_star (sID, star) VALUES (%s, %s)", (sID, row['star']))
132
133 # Insert values into movie_director table
134 cursor.execute("INSERT INTO movie_director (dID, director) VALUES (%s, %s)", (dID, row['director']))
135
136 # Insert values into movie_writer table
137 cursor.execute("INSERT INTO movie_writer (wID, writer) VALUES (%s, %s)", (wID, row['writer']))
138
139 # Insert values into movie_production table
140 cursor.execute("INSERT INTO movie_production (pID, mID, cID, sID, dID, wID) VALUES (%s, %s, %s, %s, %s, %s)",
141               (pID, mID, cID, sID, dID, wID))
142
143 # Insert values into movie_general table
144 cursor.execute("INSERT INTO movie_general (gID, pID, votes, score, genre, year, country, released, runtime, rating, budget, gross) \
145               VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s)",
146               (gID, pID, row['votes'], row['score'], row['genre'], row['year'], row['country'],
147               row['released'], row['runtime'], row['rating'], row['budget'], row['gross']))
148
149 # Commit the changes for each row
150 connection.commit()
151
152 # Print a success message if no errors occurred
153 print("Data inserted successfully.")
154
155 # Exception block to handle any errors that occur during execution
156 except Error as e:
157     # Print the error message
158     print(f"Error: {e}")
159

```

Figure 7: Rest of the code

```

160 # Finally block to ensure cursor and connection are closed regardless of errors
161 finally:
162     # Check if the connection is still open
163     if connection.is_connected():
164         # Close the cursor object
165         cursor.close()
166         # Close the database connection
167         connection.close()
168         # Print a message indicating that the database connection is closed
169         print("Database connection closed.")

```

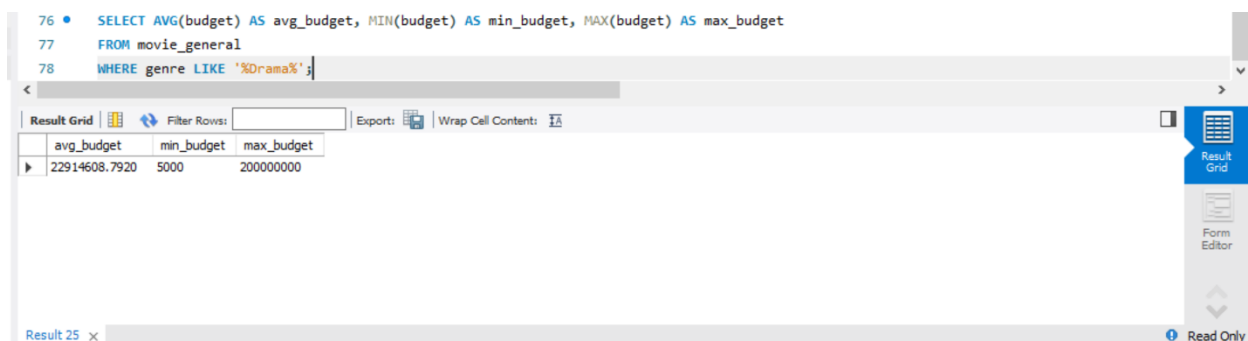
Figure 8: The end of the code

SQL Statements to Show Correct Implementation

There are 5 types of information needed to be obtained about the data in this database, and they are:

1. The average, minimum and maximum budget of drama movies.
2. The minimum budget for each genre.
3. The average score of each rating (ignore not rated movies).
4. The average runtime of all movies between the years 2014 and 2016.
5. The number of movies released before 2000.

The results for these are obtained in MySQL Workbench by the following SQL statements in Figures 9 through 13.



The screenshot shows the MySQL Workbench interface. The SQL editor at the top contains the following query:

```

76 • SELECT AVG(budget) AS avg_budget, MIN(budget) AS min_budget, MAX(budget) AS max_budget
77 FROM movie_general
78 WHERE genre LIKE '%Drama%';

```

Below the editor, the 'Result Grid' tab is active, displaying the results of the query in a table format. The table has three columns: avg_budget, min_budget, and max_budget. The first row shows the values 22914608.7920, 5000, and 2000000000 respectively.

avg_budget	min_budget	max_budget
22914608.7920	5000	2000000000

Figure 9: Average, minimum, and maximum budget


```

80 • SELECT genre, MIN(budget) AS min_budget
81 FROM movie_general
82 GROUP BY genre;

```

genre	min_budget
Action	7000
Comedy	3000
Horror	15000
Drama	5000
Crime	6000
Biography	111000
Adventure	500000
Animation	70000
Fantasy	1000000
Mystery	6900000
Thriller	2500000
Family	10500000
Sci-Fi	370000
Musical	NO DATA
Western	10000000
Romance	10000000
Music	NO DATA
History	323562
Sport	NO DATA

Result 26 x Read Only

Figure 10: Minimum budget

```

84 • SELECT rating, AVG(score) AS avg_score
85 FROM movie_general
86 WHERE rating != 'Not Rated'
87 GROUP BY rating;

```

rating	avg_score
R	6.446347403255376
PG-13	6.287215909497305
PG	6.223562303156898
TV-PG	6.939998664855955
G	6.590196081236297
TV-MA	7.02222220102946
NC-17	6.547826041346011
Unrated	6.617307699643648
X	6.900000095367432
Approved	3.4000000953674316
TV-14	6.300000190734863

Result 27 x Read Only

Figure 11: Average score

```

89 • SELECT AVG(runtime) AS avg_runtime
90 FROM movie_general
91 WHERE year BETWEEN 2014 AND 2016;

```

avg_runtime
110.6150

Result 28 x Read Only

Figure 12: Average runtime

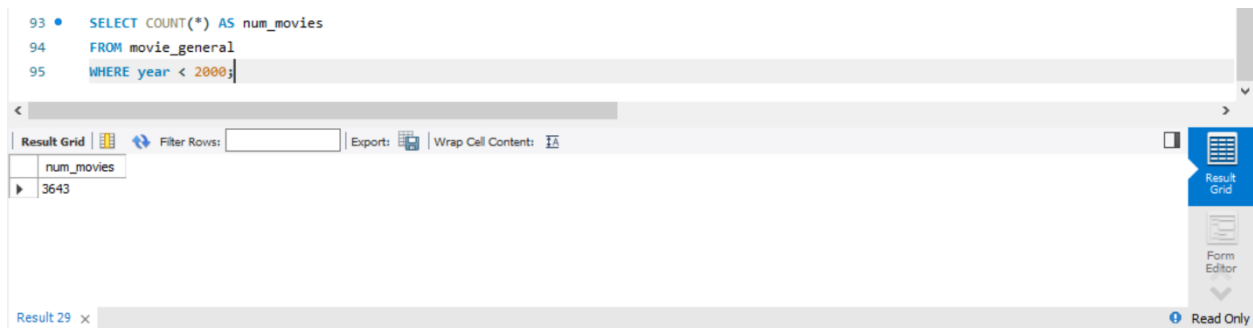


Figure 13: Number of movies before the year 2000

Also, the results obtained by these statements in Visual Studio, using Python, will be presented in Figures 14 through 18.

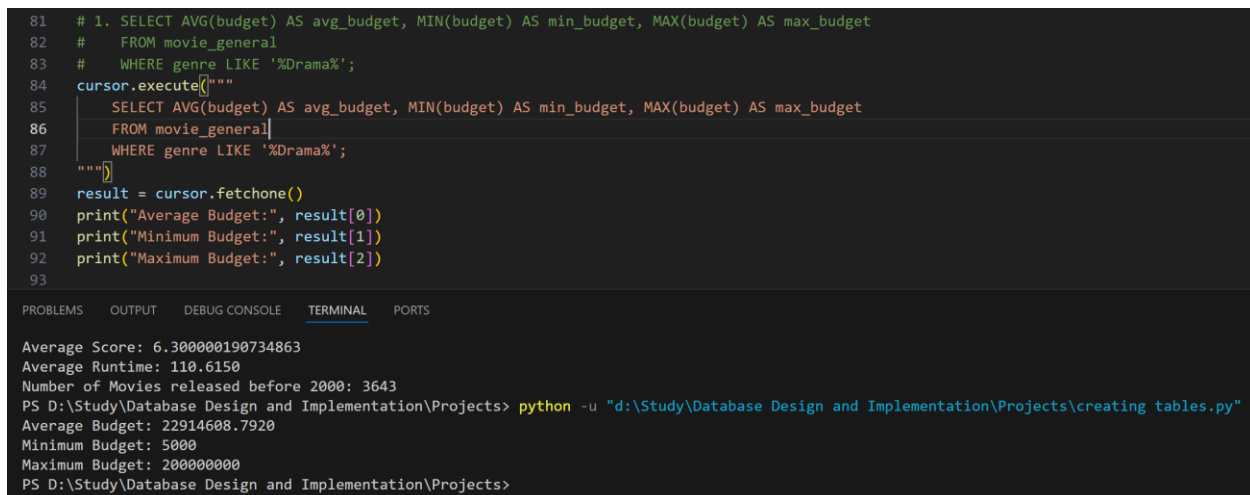


Figure 14: Average, minimum, and maximum budget

```

94 # 2. SELECT genre, MIN(budget) AS min_budget
95 # FROM movie_general
96 # GROUP BY genre;
97 cursor.execute("""
98     SELECT genre, MIN(budget) AS min_budget
99     FROM movie_general
100     GROUP BY genre;
101 """)
102 results = cursor.fetchall()
103 for row in results:
104     print("Genre:", row[0], ", Minimum Budget:", row[1])
105
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

PS D:\Study\Database Design and Implementation\Projects> python -u "d:\Study\Database Design and Implementation\Projects\creating tables.py"
Genre: Action , Minimum Budget: 7000
Genre: Comedy , Minimum Budget: 3000
Genre: Horror , Minimum Budget: 15000
Genre: Drama , Minimum Budget: 5000
Genre: Crime , Minimum Budget: 6000
Genre: Biography , Minimum Budget: 111000
Genre: Adventure , Minimum Budget: 500000
Genre: Animation , Minimum Budget: 70000
Genre: Fantasy , Minimum Budget: 1000000
Genre: Mystery , Minimum Budget: 6900000
Genre: Thriller , Minimum Budget: 2500000
Genre: Family , Minimum Budget: 10500000
Genre: Sci-Fi , Minimum Budget: 370000
Genre: Musical , Minimum Budget: None
Genre: Western , Minimum Budget: 10000000
Genre: Romance , Minimum Budget: 10000000
Genre: Music , Minimum Budget: None
Genre: History , Minimum Budget: 323562
Genre: Sport , Minimum Budget: None
PS D:\Study\Database Design and Implementation\Projects>

```

Figure 15: Minimum budget

```

106 # 3. SELECT rating, AVG(score) AS avg_score
107 # FROM movie_general
108 # WHERE rating != 'Not Rated'
109 # GROUP BY rating;
110 cursor.execute("""
111     SELECT rating, AVG(score) AS avg_score
112     FROM movie_general
113     WHERE rating != 'Not Rated'
114     GROUP BY rating;
115 """)
116 results = cursor.fetchall()
117 for row in results:
118     print("Rating:", row[0], ", Average Score:", row[1])
119
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

PS D:\Study\Database Design and Implementation\Projects> python -u "d:\Study\Database Design and Implementation\Projects\creating tables.py"
Rating: R , Average Score: 6.446347403255376
Rating: PG-13 , Average Score: 6.287215909497305
Rating: PG , Average Score: 6.223562303156898
Rating: TV-PG , Average Score: 6.939998664855955
Rating: G , Average Score: 6.590196081236297
Rating: TV-MA , Average Score: 7.0222220102946
Rating: NC-17 , Average Score: 6.547826041346011
Rating: Unrated , Average Score: 6.617307699643648
Rating: X , Average Score: 6.900000095367432
Rating: Approved , Average Score: 3.4000000953674316
Rating: TV-14 , Average Score: 6.300000190734863
PS D:\Study\Database Design and Implementation\Projects>

```

Figure 16: Average score

```
120 # 4. SELECT AVG(runtime) AS avg_runtime
121 #     FROM movie_general
122 #     WHERE year BETWEEN 2014 AND 2016;
123 cursor.execute("""
124     SELECT AVG(runtime) AS avg_runtime
125     FROM movie_general
126     WHERE year BETWEEN 2014 AND 2016;
127 """)
128 result = cursor.fetchone()
129 print("Average Runtime:", result[0])
130
```

PROBLEMS OUTPUT DEBUG CONSOLE **TERMINAL** PORTS

```
PS D:\Study\Database Design and Implementation\Projects> python -u "d:\Study\Database Design and Implementation\Projects\creating tables.py"
Average Runtime: 110.6150
PS D:\Study\Database Design and Implementation\Projects> |
```

Figure 17: Average runtime

```
131 # 5. SELECT COUNT(*) AS num_movies
132 #     FROM movie_general
133 #     WHERE year < 2000;
134 cursor.execute("""
135     SELECT COUNT(*) AS num_movies
136     FROM movie_general
137     WHERE year < 2000;
138 """)
139 result = cursor.fetchone()
140 print("Number of Movies released before 2000:", result[0])
141
```

PROBLEMS OUTPUT DEBUG CONSOLE **TERMINAL** PORTS

```
PS D:\Study\Database Design and Implementation\Projects> python -u "d:\Study\Database Design and Implementation\Projects\creating tables.py"
Number of Movies released before 2000: 3643
PS D:\Study\Database Design and Implementation\Projects> |
```

Figure 18: Number of movies before the year 2000

Reflective Segment

This section aims to provide an in-depth analysis of the challenges encountered by the designer throughout the project. Additionally, any improvements that can be made to the design will be evaluated and their impact on the overall outcome.

Challenges in the design:

Throughout the design process, the initial obstacle was to standardize the database schema. At first, the "movie_production" and "movie_general" tables had a shared primary key, which was the "movie ID" attribute of the "movie_name" table. This allowed users to effortlessly access all the information for each movie in other tables via the "movie ID" attribute as a foreign key. However, this approach had the potential for data anomalies and made maintaining referential integrity difficult. Additionally, it could create complications for database operations, such as updates and deletions. As a result, the design was modified to its current version, where each table has its own primary key that is not a foreign key. This guarantees the same output without any issues.

Another issue arose in the Python code for inserting values into tables. Initially, the logic was to extract all values from a particular column in the CSV file and insert them into the relevant table after assigning IDs if necessary. The code would then repeat the process with the next column. However, this approach was not efficient, as MySQL Workbench automatically rearranges rows in alphabetical order. Consequently, when moving to insert other columns, incorrect values could be assigned from one table to another, given that the ID was randomly generated rather than sequenced. Furthermore, the values would differ every time the code ran. Therefore, the logic was revised to navigate through the CSV file row by row, extract values and store them in variables before inserting them into the tables, and then move to the next row.

Improvements to the design:

An enhancement that could be made to the code involved in utilizing Python functions and the "random" library to generate distinct IDs. Although this method is effective and comprehensible, it may appear slightly complex when compared to an alternative technique that utilizes the "AUTO_INCREMENT" SQL statement during table creation. The latter approach is more straightforward and eliminates the necessity of creating a specific function for each ID.

To handle 'nan' values in a CSV file and ensure their proper storage in a database, we utilized the "where()" function from the Pandas library. This function effectively transformed 'nan' values to "None". While this approach generally works well with most databases, there exists a simpler and more direct method. This method involves assigning a default value of "0" to all 'nan' values. However, this method has its drawbacks as it may lead to biased analysis in certain scenarios by concealing missing data. In our particular database, this method would have sufficed since no cell in any table holds the value "0".

References

Geeks for Geeks (2023) *Python | Pandas DataFrame.where()*. Available from: <https://www.geeksforgeeks.org/python-pandas-dataframe-where/> [Accessed 16 April 2024].