



**Universidad de Sonora**  
**Introduction to Data Science and its Methodologies**



**Final Project:**  
**Creation and Management of Databases**

Ana Daniela Pérez Romero / A222230117

Professor: Juan Pablo Soto Barrera

## Introduction

Introduction to Data Science and its Methodologies is a course focused on teaching the most common methodologies used in any data science scenario. Throughout the course, it is taught how to ensure that the data used in problem solving is relevant and properly manipulated. Accordingly, we learned how to form a business/research problem, collecting, preparing, and analyzing data, applying the 6 stages of the CRISP-DM methodology, working, and creating databases with SQL, but most importantly, how to make use of these tools and knowledge altogether to develop hands-on experience.

The final project for the second phase of the course consists of developing an application that allows analyzing information contained in a SQL database. This database will be later consulted through a Jupyter notebook and should be analyzed through graphs and tables to solve a problem. Within the requirements of the database created through SQL, is that it contains at least one view, one stored procedure and one function.

## Project Documentation

### I. Database Creation

For the development of this project, I decided to work on one of my greatest passions: watching movies. Perhaps, you are wondering why this topic is so important for me, and the truth is that as a full-time employee and student, I hardly find free time to do those activities that I enjoy the most. I usually watch movies Sunday night, and every time I spend at least 20 minutes just deciding what movie should I watch; there are so many options that I just get overwhelmed. On the other hand, there have been times where friends and family recommend me movies and turns out they do not coincide with my cinematographic taste, meaning I just lost two hours of my short weekend. As a Data Science student, I feel the moral responsibility to apply what I have learned in class to simplify my life; that is one of Data Analytics goals right?

The first stage of any Data Science project is obtaining data, and based on the subject of my project I decided to use a dataset from Internet Movie Database (IMDb), where they display the top 1000 best ranked movies:

<https://www.kaggle.com/datasets/harshitshankhdhar/imdb-dataset-of-top-1000-movies-and-tv-shows>

## IMDB Movies Dataset

Data Code (32) Discussion (2) 242 New Notebook Download (179 kB)

Detail Compact Column 10 of 16 columns imdb\_top\_1000.csv

**About this file**

IMDB Dataset of Top 1000 Movies by IMDB Rating

Poster_Link	Series_Title	Released_Year	Certificate	Runtime	Genre
Poster Link of Movie	Name of the Movie	Released Year of the Movie	Certificate of the Movie	Total Runtime of the Movie	Genre of
<b>1000</b> unique values	<b>999</b> unique values	2014 3% 2004 3% Other (937) 94%	U 3% A 20% Other (569) 57%	23% 100 min 2% Other (954) 95%	Drama Drama, R Other (87)
https://m.media-amazon.com/images/M/MV5BMDFkYTc0MGEtZmNhWC0wZDZlZWFiMTk0MTZmRjY0MwMWFmOTkxYkFjcGde...	The Shawshank Redemption	1994	A	142 min	Drama
https://m.media-amazon.com/images/M/MV5BM2MyNjYxNmUyLTAwNjE0MTYxLWJmNWYyZ2ZlODY3ZTk3OTFlXkEyXkFjcGde...	The Godfather	1972	A	175 min	Crime, I
https://m.media-amazon.com/images/M/MV5BMTRhNTNmODQwNjFhNjE5bWVhbnRkXkF1ZTcwODkyMTkzMw@@_V1_UX67_CR0,0,...	The Dark Knight	2008	UA	152 min	Action,
https://m.media-amazon.com/images/M/MV5BMWwMQzZlLTlY2JlNCBhODZlLWUyODc1NDk2	The Godfather: Part II	1974	A	202 min	Crime, I

**Summary**

- 1 file
- 16 columns

Through MySQL a new schema was created and called “movies”; firstly, it was supposed that the table was going to be created from Table Data Import Wizard option, however it resulted in missing records. For this reason, the table was created through a Jupyter Notebook, creating a data frame with pandas, and using sqlalchemy library to upload the data frame to MySQL.

```
import pandas as pd
from sqlalchemy import create_engine
import pymysql

# Creation of dataframe from a .csv file
movies = pd.read_csv("top_imdb.csv")

# We store the necessary information to access MySQL
myUser = 'root'
myPass = '1'
myEndpoint = 'localhost'
myPort = 3306
myPort2 = ':' + str(myPort) + '/'
myDb = 'movies'

myDataConnect = [myUser, myPass, myEndpoint, myPort, myDb]

db_data = 'mysql+mysqldb://' + myUser + ':' + myPass + '@' + myEndpoint + myPort2 \
+ myDb
engine = create_engine(db_data)

# Execute the to_sql for writing DF into MySQL
movies.to_sql('top_movies', engine, if_exists='replace', index=False)
```

Once the program is executed, the table “top\_movies” is uploaded to MySQL’s schema “movies”. With *SELECT \* FROM movies.top\_movies* query, 1000 records should be displayed.

Query 1 movies top\_movies dany's\_favorite\_genre - View dany's\_favorite\_genre

1. SELECT \* FROM movies.top\_movies;

Result Grid

Movie_title	Released_Year	Runtime_min	Genre	IMDB_Rating	Overview	Director	Star 1	Star 2	Star 3	Star 4
Saving Private Ryan	1998	169	Drama, War	8.6	Following the Normandy Landings, a group of U...	Steven Spielberg	Tom Hanks	Matt Damon	Tom Sizemore	Edward E...
The Green Mile	1999	189	Crime, Drama, Fantasy	8.6	The lives of guards on Death Row are affected ...	Frank Darabont	Tom Hanks	Michael Clarke Duncan	David Morse	Bonnie H...
La vita è bella	1997	116	Comedy, Drama, Romance	8.6	When an open-minded Jewish Italian and his s...	Roberto Benigni	Roberto Benigni	Nicoletta Braschi	Gorgio Cantarini	Giustino I...
Se7en	1995	127	Crime, Drama, Mystery	8.6	Two detectives, a rookie and a veteran, hunt a ...	David Fincher	Morgan Freeman	Brad Pitt	Kevin Spacey	Andrew J...
The Silence of the Lambs	1991	118	Crime, Drama, Thriller	8.6	A young F.B.I. cadet must receive the help of a ...	Jonathan Demme	Jodie Foster	Anthony Hopkins	Lawrence A. Bon...	Kasi Lem...
Star Wars	1977	121	Action, Adventure, Fantasy	8.6	Luke Skywalker joins forces with a Jedi Knight, ...	George Lucas	Mark Hamill	Harrison Ford	Carrie Fisher	Alec Guin...
Seppuku	1962	133	Action, Drama, Mystery	8.6	When a ronin requesting seppuku at a feudal lo...	Masaki Kobayashi	Tatsuya Nakadai	Akira Ishihama	Shima Iwashita	Tetsuzō Ō...
Shichinin no samurai	1954	207	Action, Adventure, Drama	8.6	A poor village under attack by bandits recruits s...	Akira Kurosawa	Toshiro Ō...	Takashi Shimura	Kikko Tachibana	Yukiko Shi...

Data is uploaded, however, there are still a thousand movies that need to be filtered based on my cinematographic taste. A new view called “dany's\_favorite\_genre” is created selecting my favorite movie genres: romance, history, and comedy. With this view 636 records are left behind, leaving only 364 movies to be furthered analyzed.

Navigator: dany's\_favorite\_genre

Query 1 movies top\_movies dany's\_favorite\_genre - View dany's\_favorite\_genre

Name: dany's\_favorite\_genre The name of the view is parsed automatically from the DDL statement. The DDL is parsed automatically while you type.

DDL:

```

20 FROM
21 'movies'.top_movies
22 WHERE
23 (('movies'.top_movies.Genre LIKE 'Romance')
24 OR ('movies'.top_movies.Genre LIKE 'History'))
25 OR ('movies'.top_movies.Genre LIKE 'Comedy'))

```

**Rows Processed:**

Rows affected: 0

Rows sent to client: 364

Rows examined: 1000

Another important factor when selecting a film is its length; there are days with less time availability, where shorter movies must be selected. A stored procedure “length” was created and allows to filter movies based on the maximum time in minutes that can be consigned to watch them. This procedure allows the user to type a maximum number of minutes long for a movie and return movies equal or less than the number typed.

Apply SQL Script to Database

Review SQL Script

Apply SQL Script

Online DDL

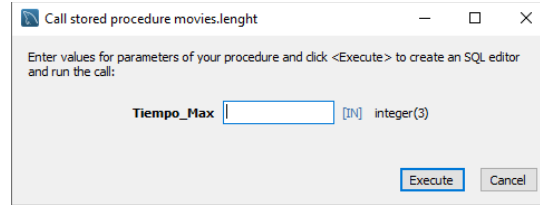
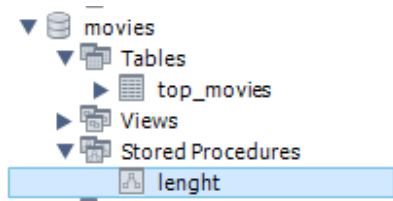
Algorithm: Default Lock Type: Default

```

1  USE `movies`;
2  DROP procedure IF EXISTS `length`;
3
4  USE `movies`;
5  DROP procedure IF EXISTS `movies`.`length`;
6  ;
7
8  DELIMITER $$
9  USE `movies` $$
10 CREATE DEFINER=`root`@`localhost` PROCEDURE `length` (Tiempo_Max integer(3))
11 BEGIN
12 SELECT * FROM top_movies WHERE Runtime_min <= Tiempo_Max;
13 END$$
14
15 DELIMITER ;
16 ;
17

```

Back Apply Cancel



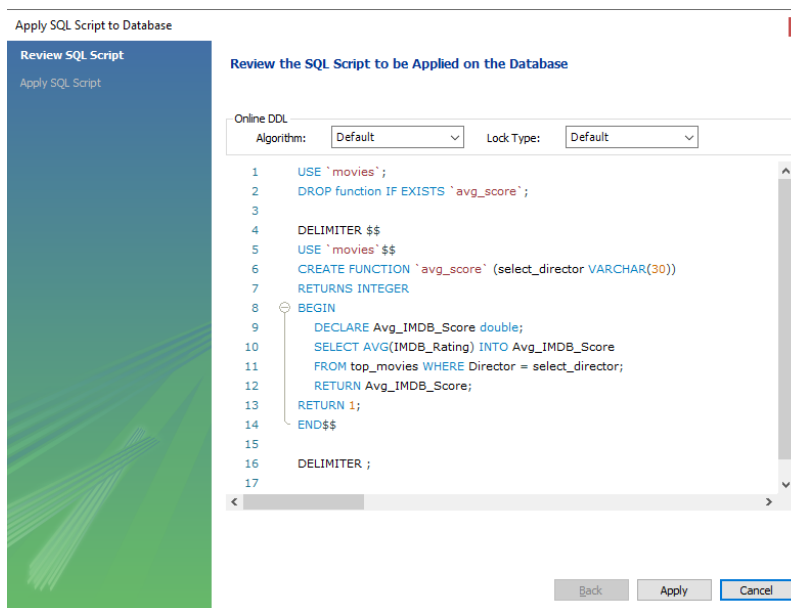
```
1 • call movies.lenght(10);
2
```

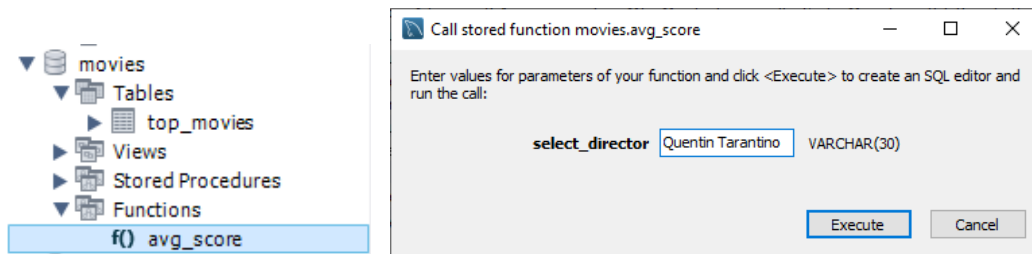
Movie_title	Released_Year	Runtime_min	Genre	IMDB_Rating	Overview	Director	Star1	Star2	Star3	Star4
The Lion King	1994	88	Animation, Adventure, Drama	8.5	Lion prince Simba and his father are targeted b...	Roger Allers	Rob Minkoff	Matthew Broderick	Jeremy Irons	Jai
Hotaru no haka	1988	89	Animation, Drama, War	8.5	A young boy and his little sister struggle to surv...	Isao Takahata	Tsutomu Tatsumi	Ayano Shirashi	Akemi Yamaguchi	Yo
Modern Times	1936	87	Comedy, Drama, Family	8.5	The Tramp struggles to live in modern industrial ...	Charles Chaplin	Charles Chaplin	Paulette Goddard	Henry Bergman	Tir
City Lights	1931	87	Comedy, Drama, Romance	8.5	With the aid of a wealthy erratic tippler, a dewy...	Charles Chaplin	Charles Chaplin	Virginia Cherrill	Florence Lee	Ha
Dr. Strangelove or: How I Learned to Stop Wor...	1964	95	Comedy	8.4	An insane general triggers a path to nuclear hol...	Stanley Kubrick	Peter Sellers	George C. Scott	Sterling Hayden	Ke
Paths of Glory	1957	88	Drama, War	8.4	After refusing to attack an enemy position, a g...	Stanley Kubrick	Kirk Douglas	Ralph Meeker	Adolphe Menjou	Ge
Bacheba-Ye aseman	1997	89	Drama, Family, Sport	8.3	After a boy loses his sister's pair of shoes, he g...	Majid Majidi	Mohammad Amir Naj	Amir Farrokh Hashemian	Bahare Seddiqi	Na
Toy Story	1995	81	Animation, Adventure, Comedy	8.3	A cowboy doll is profoundly threatened and jeal...	John Lasseter	Tom Hanks	Tim Allen	Don Rickles	Jin
Ladri di biciclette	1948	89	Drama	8.3	In post-war Italy, a working-class man's bicycle ...	Vittorio De Sica	Lamberto Maggiorani	Enzo Staiola	Lianella Carell	Ele
The Kid	1921	58	Comedy, Drama, Family	8.3	The Tramp cares for an abandoned child, but e...	Charles Chaplin	Charles Chaplin	Edna Purviance	Jackie Coogan	Ca

For the creation of functions, a global variable must be defined: *log\_bin\_trust\_function\_creators* = 1; by default, it comes to 0. Once the global variable is defined, we can create functions

```
1 • SET GLOBAL log_bin_trust_function_creators = 1;
```

An important aspect while selecting movies is its director; how are their films ranked? To answer this question, a function that returns the average rating per director was created.





“Quentin Tarantino” was selected as a director, once the function was runned, it returned an average score value of his movies of 8.175.

	movies.avg_score('Quentin Tarantino')
▶	8.175

## II. Data Analysis through Python

Through a Jupyter Notebook, a connection is created to call the database “movies” from MySQL, which is set as a data frame.

```
connection=pymysql.connect(host= myEndpoint,
port= myPort,
user=myUser,
password= myPass,
db= myDb)
```

```
dataframe=pd.read_sql_query("SELECT * FROM movies.top_movies", connection)
dataframe.head(8)
```

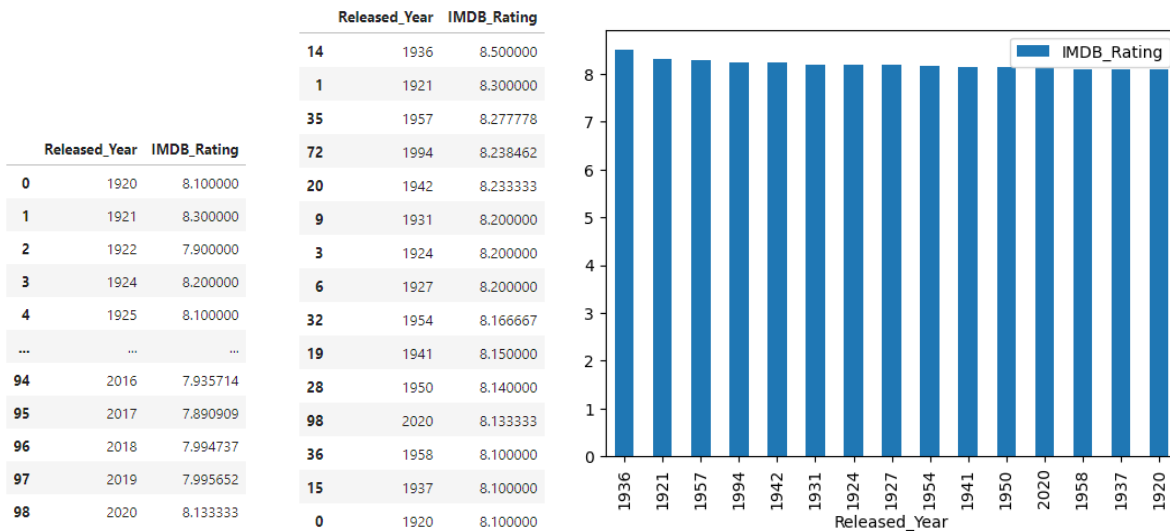
	Movie_title	Released_Year	Runtime_min	Genre	IMDB_Rating	Overview	Director	Star1	Star2	Star3	Star4	No_of_Votes	Gross
0	The Shawshank Redemption	1994	142	Drama	9.3	Two imprisoned men bond over a number of years...	Frank Darabont	Tim Robbins	Morgan Freeman	Bob Gunton	William Sadler	2343110	28,341,469
1	The Godfather	1972	175	Crime, Drama	9.2	An organized crime dynasty's aging patriarch L...	Francis Ford Coppola	Marlon Brando	Al Pacino	James Caan	Diane Keaton	1620367	134,966,411
2	The Dark Knight	2008	152	Action, Crime, Drama	9.0	When the menace known as the Joker wreaks havoc...	Christopher Nolan	Christian Bale	Heath Ledger	Aaron Eckhart	Michael Caine	2303232	534,858,444

Firstly, I would like to know if cinematography has been evolving positively or negatively over the years; for this matter, a table and a top 15 values bar graph showing the average IMDB Score Rate was created.

```
score_year= dataframe.groupby("Released_Year").mean()["IMDB_Rating"].reset_index()
score_year
```

```
scoreyear_pbar = score_year.sort_values('IMDB_Rating', ascending=False).head(15)
scoreyear_pbar
```

```
scoreyear_pbar.plot.bar(x="Released_Year",y="IMDB_Rating")
```



Apparently, it seems like the 21st century has not been the best cinematographic era; 14 out of the 15 best positions based on average score rate are taken by the 20th century. However, this information seems intriguing, so I decide to count how many films were evaluated per year. My suspicions where true: the number of films reviewed from the 20th century were smaller compared to those of the 21st century.

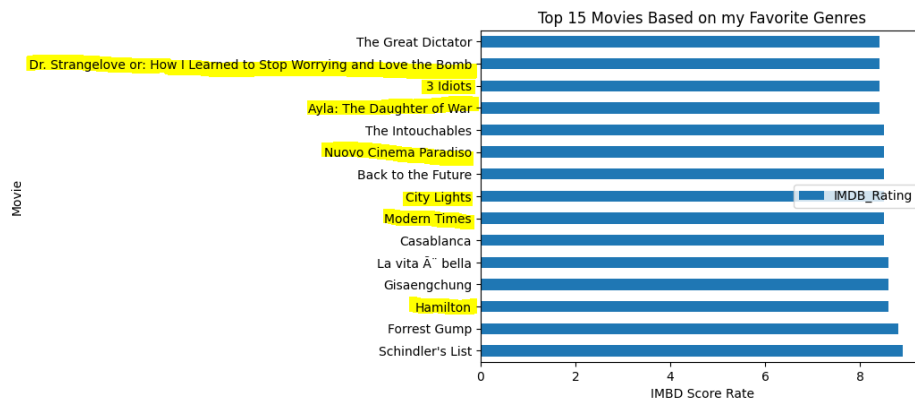
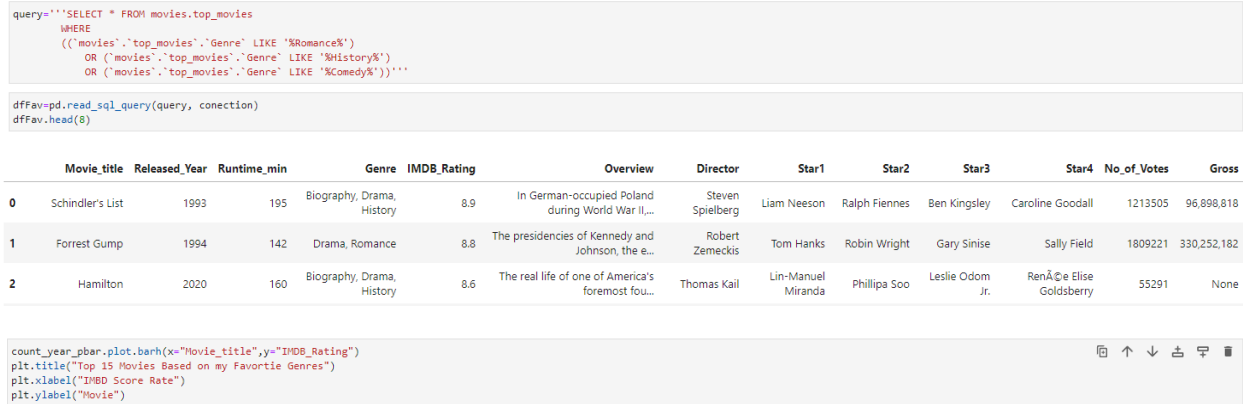
Released_Year IMDB_Rating		
92	2014	32
82	2004	31
87	2009	29
94	2016	28
91	2013	28
79	2001	27
84	2006	26
85	2007	26
93	2015	25
90	2012	24
88	2010	23
71	1993	23
97	2019	23
81	2003	22
95	2017	22

Released_Year IMDB_Rating		
0	1920	1
1	1921	1
2	1922	1
3	1924	1
4	1925	2
...	...	...
94	2016	28
95	2017	22
96	2018	19
97	2019	23
98	2020	6

Thanks to this analysis it is known that selecting a movie based on the average IMDB Score Rate from its released year is not the best method.

For a second analysis, I'd like to know which movies has the best IMDB score rate. A query was executed to filter the movies based on my favorite genres: Romance, History and Comedy. These values where saved in a data frame, that will later be used to create a plot bar that shows top 15 best ranked movies.



From the results above, and removing the movies that I have already seen in the past, I know which seven movies I'll see during my winter vacations:

- Hamilton
- Modern Times
- City Lights
- Cinema Paradiso
- Ayla: The Daughter of War
- 3 Idiots
- Dr. Strangelove or: How I Learned to Stop Worrying and Love

## Conclusion

Python is an incredibly useful tool with a huge range of open-source libraries for the development of Data Science project, however, without data, it is not that useful for analytics, so that's where datasets enter the game. Relational databases are an extremely efficient, powerful, and widely used way to create, read, update, and delete data of all kinds. The most common relational database management systems are all based on Structured Query Language (SQL); this means that Data Scientists needs a strong understanding on using Python and SQL together to have an advantage when it comes to working with data analytics.



With the completion of this project, we were able to put into practice the knowledge acquired over the last few months. By establishing a statement problem and using tools like SQL and Python we were able to build and analyze a dataset that provided a data-driven solution. Without a doubt, it was a great challenge whose conclusion brings satisfaction, but above all, the desire to continue learning and improving our knowledge in the world of Data Science.