# DW&DM

## MS. NUSRAT SHAHEEN

## LAB PROJECT

**Submitted By:**
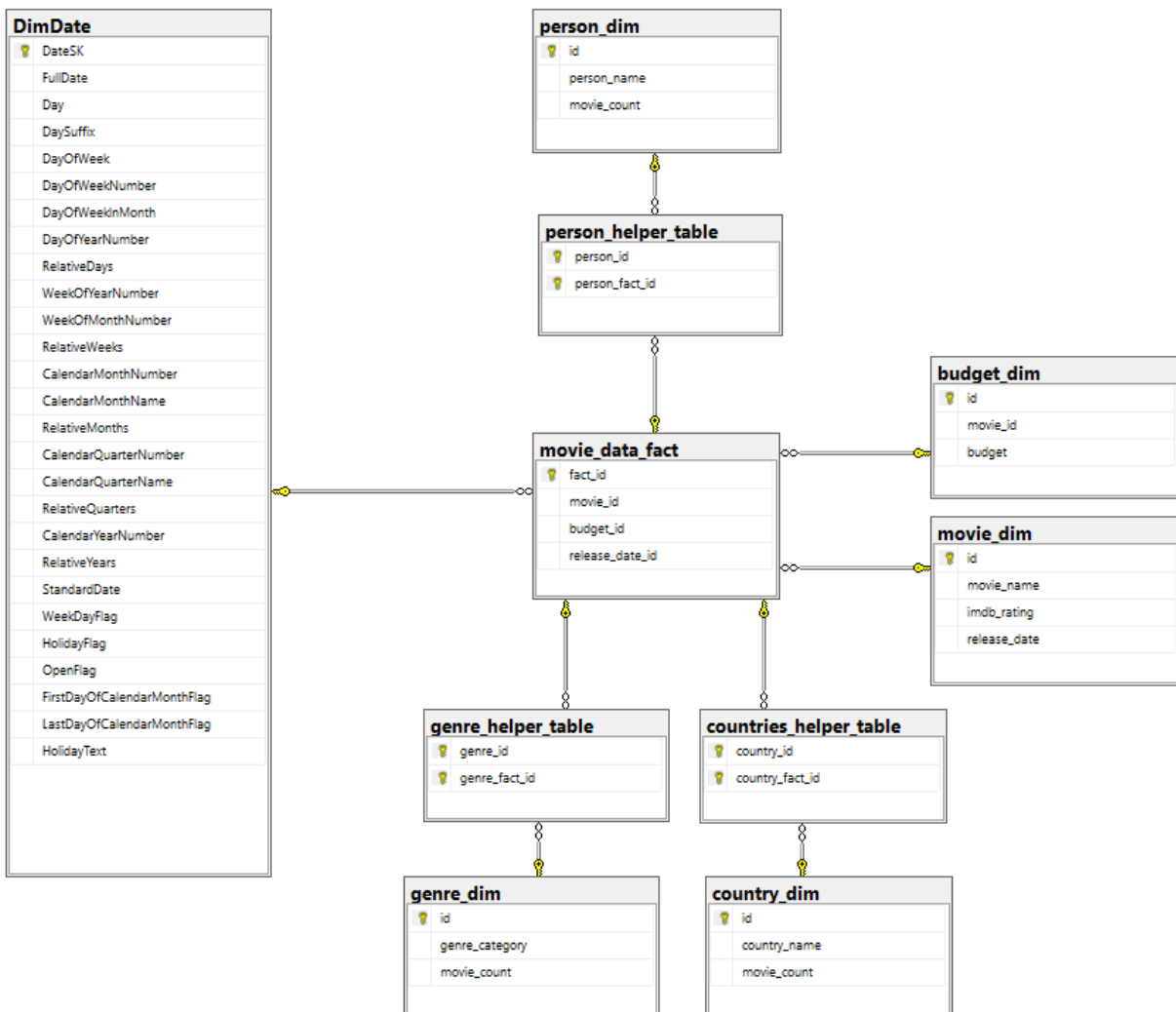
*Faiq Shahzad* (*FA19-BCS-021*)

*Muhammad Ahmed* (*FA19-BCS-041*)

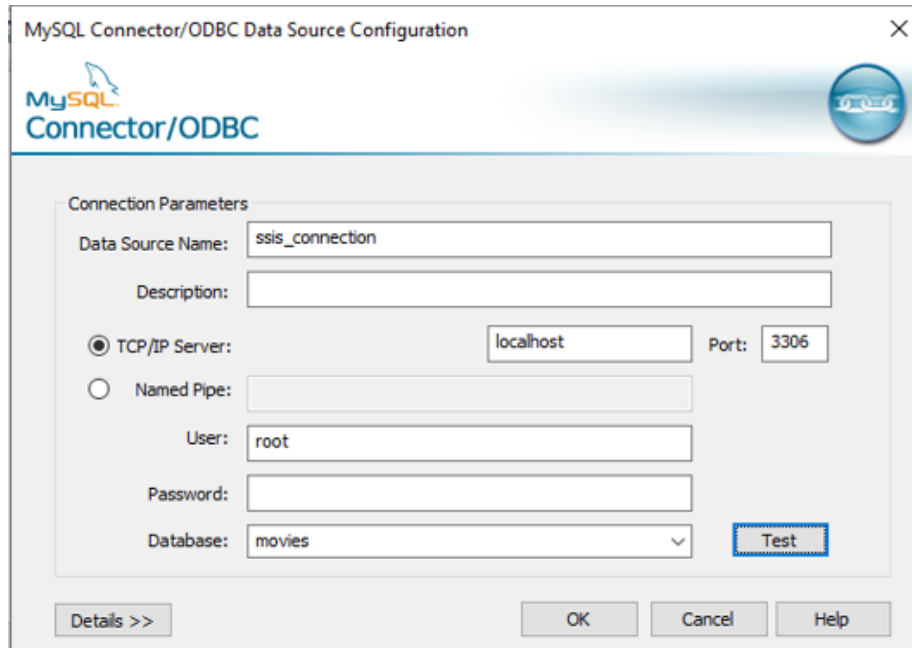*Muhammad Mubir Shami* (*FA19-BCS-051*)

# Business Questions:

- What are the total number of movies in the specific genre?
- What are the total number of Movies produced in a specific country?
- What are the total number of Movies a person has participated in?
- What is the list of movies released in a specific year?
- What is the list of movies under a specific budget?

## Dimension Model:



**DimDate**
- DateSK
- FullDate
- Day
- DaySuffix
- DayOfWeek
- DayOfWeekNumber
- DayOfWeekInMonth
- DayOfYearNumber
- RelativeDays
- WeekOfYearNumber
- WeekOfMonthNumber
- RelativeWeeks
- CalendarMonthNumber
- CalendarMonthName
- RelativeMonths
- CalendarQuarterNumber
- CalendarQuarterName
- RelativeQuarters
- CalendarYearNumber
- RelativeYears
- StandardDate
- WeekDayFlag
- HolidayFlag
- OpenFlag
- FirstDayOfCalendarMonthFlag
- LastDayOfCalendarMonthFlag
- HolidayText

**person_dim**
- id
- person_name
- movie_count

**person_helper_table**
- person_id
- person_fact_id

**budget_dim**
- id
- movie_id
- budget

**movie_data_fact**
- fact_id
- movie_id
- budget_id
- release_date_id

**movie_dim**
- id
- movie_name
- imdb_rating
- release_date

**genre_helper_table**
- genre_id
- genre_fact_id

**countries_helper_table**
- country_id
- country_fact_id

**genre_dim**
- id
- genre_category
- movie_count

**country_dim**
- id
- country_name
- movie_count

# ETL:

## Step 1: Creating Operational Database Connection:

## Step 2: Extraction:

1. Create Connections:



2. Extracting Table from Operation Database:

## Step 3: Transformations:

1. Choose a Transformation to be performed on the Extracted Data such as:
   - **Aggregate (Group By)**

- **Sort**



- **Merge Join**

- **Data Conversion**



- **Row Sampling**

- **Conditional Split:**



2. Perform the Transformations:
   - **Person_Dim:**

- **Countries_dim:**



- **Genre_dim**

- **Movies_dim**



- **Budget_dim**

- **Movie_Data_Fact**



- **Person_Helper_Table**

- **Genre_Helper_Table**



- **Country_Helper_Table**

# Step 4: Loading:

1. Select the OLE DB from the destination panel
2. Connect to the desired table in the database

# Business Queries:

- **No of Movies in a Genre:**

```sql
USE [movies_dw]
GO

SELECT [genre_category]
      ,[movie_count]
  FROM [dbo].[genre_dim]
  WHERE [genre_category] = 'Comedy'

GO
```

100 %

Results | Messages

| | genre_category | movie_count |
|---|---|---|
| 1 | Comedy | 1722 |

- **No of Movies Produced in a Country:**

```sql
USE [movies_dw]
GO

SELECT [country_name]
      ,[movie_count]
  FROM [dbo].[country_dim]
  WHERE [country_name] = 'United States of America'

GO
```

100 %

Results | Messages

| | country_name | movie_count |
|---|---|---|
| 1 | United States of America | 3956 |

- **No of Movies respective of Cast:**

```sql
USE [movies_dw]
GO

SELECT [person_name]
      ,[movie_count]
  FROM [dbo].[person_dim]
  WHERE [person_name] = 'Arnon Milchan'

GO
```

100 %

Results | Messages

| | person_name | movie_count |
|---|---|---|
| 1 | Arnon Milchan | 54 |

- **Movies in a Specific Range of Budget:**

```sql
USE [movies_dw]
GO

SELECT md.id
       ,md.movie_name
       ,md.imdb_rating
       ,bd.budget
    FROM [dbo].[movie_dim] md INNER JOIN [dbo].[movie_data_fact] mdf on md.id = mdf.movie_id
                              INNER JOIN [dbo].[budget_dim] bd on mdf.budget_id = bd.id
    WHERE bd.budget > 200000000 and bd.budget < 350000000
```

100 %

▦ Results | ⯃ Messages

| | id | movie_name | imdb_rating | budget |
|---|---|---|---|---|
| 1 | 254 | King Kong | 6 | 207000000 |
| 2 | 285 | Pirates of the Caribbean: At World's End | 6 | 300000000 |
| 3 | 767 | Harry Potter and the Half-Blood Prince | 7 | 250000000 |
| 4 | 19995 | Avatar | 7 | 237000000 |
| 5 | 49026 | The Dark Knight Rises | 7 | 250000000 |
| 6 | 49051 | The Hobbit: An Unexpected Journey | 7 | 250000000 |
| 7 | 49521 | Man of Steel | 6 | 225000000 |
| 8 | 57158 | The Hobbit: The Desolation of Smaug | 7 | 250000000 |
| 9 | 122917 | The Hobbit: The Battle of the Five Armies | 7 | 250000000 |
| 10 | 127585 | X-Men: Days of Future Past | 7 | 250000000 |
| 11 | 209112 | Batman v Superman: Dawn of Justice | 5 | 250000000 |
| 12 | 271110 | Captain America: Civil War | 7 | 250000000 |

- **Movies Released in a Specific Year:**

```sql
USE [movies_dw]
GO

SELECT md.id
       ,md.movie_name
       ,md.release_date
       ,md.imdb_rating
    FROM [dbo].[movie_dim] md INNER JOIN [dbo].[movie_data_fact] mdf on md.id = mdf.movie_id
                              INNER JOIN [dbo].[DimDate] dt on mdf.release_date_id = dt.DateSK
    WHERE dt.RelativeYears = 2010-2022

GO
```

100 %

▦ Results | ⯃ Messages

| | id | movie_name | release_date | imdb_rating |
|---|---|---|---|---|
| 1 | 10138 | Iron Man 2 | 2010-04-28 00:00:00.000 | 6 |
| 2 | 10140 | The Chronicles of Narnia: The Voyage of the Dawn... | 2010-08-13 00:00:00.000 | 6 |
| 3 | 10191 | How to Train Your Dragon | 2010-03-05 00:00:00.000 | 7 |
| 4 | 10192 | Shrek Forever After | 2010-05-16 00:00:00.000 | 6 |
| 5 | 11324 | Shutter Island | 2010-02-18 00:00:00.000 | 7 |
| 6 | 11439 | The Ghost Writer | 2010-02-12 00:00:00.000 | 6 |
| 7 | 12155 | Alice in Wonderland | 2010-03-03 00:00:00.000 | 6 |
| 8 | 12819 | Alpha and Omega | 2010-09-17 00:00:00.000 | 5 |
| 9 | 16290 | Jackass 3D | 2010-10-15 00:00:00.000 | 6 |
| 10 | 170480 | The Deported | 2010-06-15 00:00:00.000 | 0 |
| 11 | 20504 | The Book of Eli | 2010-01-14 00:00:00.000 | 6 |
| 12 | 20526 | TRON: Legacy | 2010-12-10 00:00:00.000 | 6 |
| 13 | 20662 | Robin Hood | 2010-05-12 00:00:00.000 | 6 |
| 14 | 22538 | Scott Pilgrim vs. the World | 2010-07-27 00:00:00.000 | 7 |
| 15 | 22894 | Legion | 2010-01-21 00:00:00.000 | 5 |
| 16 | 22907 | Takers | 2010-08-26 00:00:00.000 | 6 |
| 17 | 22972 | Green Zone | 2010-03-11 00:00:00.000 | 6 |
| 18 | 23168 | The Town | 2010-09-15 00:00:00.000 | 7 |
| 19 | 23483 | Kick-Ass | 2010-03-22 00:00:00.000 | 7 |
| 20 | 23631 | Machete | 2010-09-01 00:00:00.000 | 6 |
| 21 | 24021 | The Twilight Saga: Eclipse | 2010-06-23 00:00:00.000 | 5 |
| 22 | 26022 | My Name Is Khan | 2010-02-12 00:00:00.000 | 7 |
| 23 | 26388 | Buried | 2010-09-24 00:00:00.000 | 6 |
| 24 | 26389 | From Paris with Love | 2010-02-05 00:00:00.000 | 6 |
| 25 | 27022 | The Sorcerer's Apprentice | 2010-07-13 00:00:00.000 | 5 |
| 26 | 27205 | Inception | 2010-07-14 00:00:00.000 | 8 |
| 27 | 27569 | Extraordinary Measures | 2010-01-21 00:00:00.000 | 6 |
| 28 | 27573 | The Bounty Hunter | 2010-03-16 00:00:00.000 | 5 |
| 29 | 27578 | The Expendables | 2010-08-03 00:00:00.000 | 6 |
| 30 | 27585 | Rabbit Hole | 2010-12-16 00:00:00.000 | 6 |

# NAÏVE BAYES:

```
In [17]:   1  import pandas as pd
           2  from sklearn.naive_bayes import GaussianNB
           3  from sklearn.model_selection import train_test_split
           4  from sklearn.preprocessing import LabelEncoder
           5
           6  file = 'genre_prediction.csv'
           7  col_names = ['cast_1', 'cast_2', 'cast_3', 'cast_4', 'genre']
           8  feature_cols = ['cast_1', 'cast_2', 'cast_3', 'cast_4']
           9
          10  dataset = pd.read_csv(file, names=col_names)
          11  print(dataset.head())
          12
          13  le = LabelEncoder()
          14  dataset[feature_cols] = dataset[feature_cols].apply(LabelEncoder().fit_transform)
          15
          16  X = dataset[feature_cols]
          17  Y = dataset.genre
          18
          19
          20  x_train, x_test, y_train, y_test = train_test_split(X, Y, test_size = 0.05)
          21
          22  naiveBase = GaussianNB()
          23  naiveBase.fit(x_train, y_train)
          24
          25  y_predict = naiveBase.predict(x_test)
          26  print("\n\nPrediction: ", y_predict)
          27
          28  print("\n\nNaive Bayes score: ",naiveBase.score(x_test, y_test))
          29
```

```
                 cast_1              cast_2            cast_3  \
0                cast_1              cast_2            cast_3
1         Lucais Fullom   Ilyssa Matityahu   Kiri Fulleylove
2       Ester Ransfield      Shani Salasar     Kirstin Awty
3  Mellisent Matyushenko   Gillie Docharty     Becka Belloch
4       Bernadine Figgs          Wait Nund   Cary MacHostie

              cast_4                      genre
0             cast_4                      genre
1     Donetta Bocken  Adventure|Documentary|Drama
2     Padriac Broadis                     Comedy
3  Perkin Besnardeau                     Comedy
4    Shandie Drohane             Drama|Thriller


Prediction:  ['Drama' 'Adventure|Romance' 'Drama' 'Drama' 'Drama' 'Drama' 'Drama'
 'Drama' 'Drama' 'Drama' 'Drama' 'Drama' 'Drama' 'Drama' 'Drama' 'Comedy'
 'Drama' 'Adventure|Drama|Romance' 'Drama' 'Drama' 'Drama' 'Drama' 'Drama'
 'Drama' 'Drama' 'Drama' 'Drama' 'Comedy' 'Drama' 'Drama' 'Drama' 'Drama'
 'Drama' 'Comedy' 'Comedy' 'Drama' 'Comedy' 'Drama' 'Drama' 'Drama'
 'Drama' 'Drama' 'Drama' 'Drama' 'Drama' 'Comedy' 'Drama' 'Drama' 'Drama'
 'Drama' 'Drama']


Naive Bayes score:  0.1568627450980392
```