



# IMDB and Kaggle Datasets



Group BHNK: Brionna Huynh and  
Neha Kumar



# Datasets

---

- IMDB (Hollywood and other international film industries) and Kaggle (Bollywood only) datasets
- IMDB: Title\_Basics (tconst, title, year, runtime), Name\_Basics (nconst, name, birth year, death year), Title\_Akas (tconst, title, region, language), Title\_Principals (tconst, characters), Title\_Episode (tconst, season, episode), Title\_Ratings (tconst, rating, numVotes), Title\_Crew (tconst, directors, writers)
- Kaggle: Bolly\_actors (name, height), Bolly\_actress (name, height, debut as lead), Bolly\_movies (title, director, cast, genre)

# Staging/ Modeled Tables

---

- IMDB: Combine the Title\_Basics and Title\_Akas tables (both referring to attributes of movies, but have some different attributes)
  - One-to-one: Title\_Basics and Title\_Akas, Title\_Basics and Title\_Ratings
  - One-to-many: Title\_basics tconst to genres column (new Genres table), Name\_basics nconst to primaryProfession column (new Primary\_profession table), Name\_basics nconst to knownForTitles column (new Known\_for table), Title\_crew tconst to directors column (new Directors table), Title\_crew tconst to writers column (new Writers table),
  - Many-to-many: Title\_principal tconst/ nconst to characters column (new Characters table)
- Kaggle:
  - One-to-many: Bolly\_movies to Director column (new Title\_director table), Bolly\_movies to Cast column (new Title\_cast table), Bolly\_movies to Genre column (new Title\_genre table)

# Beam pipelines

---

- Made a new child table that refers to the main parent table
- Split the element by the commas to get rid of the one to many relationship
- Tables: Characters, Directors, Genre, Known\_For, Primary\_Professions, Writers, Title\_cast, Title\_director, Title\_genre

id	values
id0717	12, 34, 87, 43



id	values
id0717	12
id0717	34
id0717	87
id0717	43

# Areas of interest

---

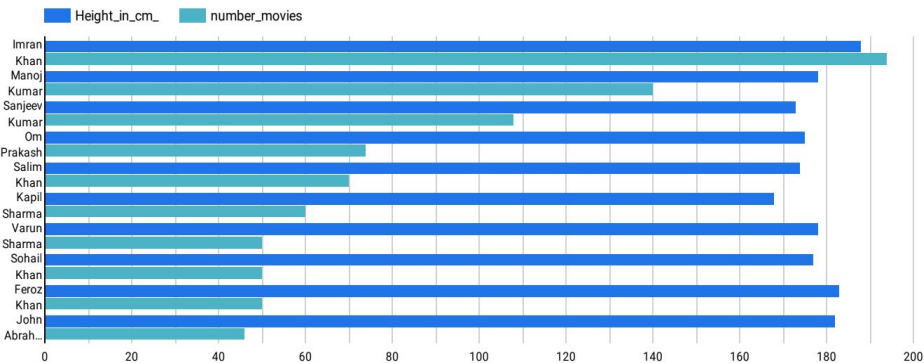
- Generally tried to learn more about Bollywood actors and movies by joining these tables with tables in the imdb\_refined dataset
  - The kaggle dataset did not contain many details about both the movies and the actors, and the imdb\_refined dataset included a lot more details about the same movies and actors
- Wanted to see if height was an important aspect of the Bollywood film industry
- Wanted to further examine the average ratings of Bollywood movies

# Queries and Data Visualization

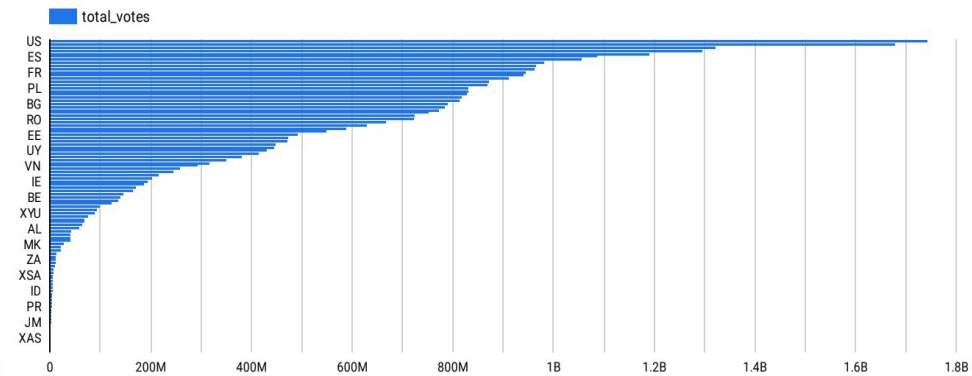
---

- Find all the movies (both from imdb and kaggle datasets) that fit into the genre War
  - Used a join to get all of the collective output of all the movies that are War based
- Find the ratings of a specific director's movies
  - Use the inner query to find the Titles of this director's movies and in the outer query, get the ratings of these Titles
- Are present-day Bollywood actresses taller than past actresses?
- Find the titles that have higher than average ratings
- Find the regions in which more than 1 million ratings (votes) for Titles are cast
- Find other professions of Bollywood actors
- What are the most common genres in Bollywood and Hollywood?
- Are the taller Bollywood actors more popular than the shorter ones?
  - Used knownForTitles to interpret popularity → if known for more Titles, then must be more popular

Height correlation to popularity(number of movies they've acted in)

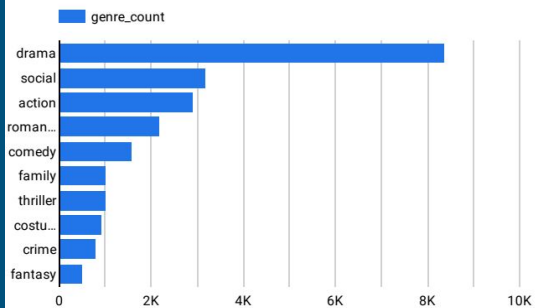


Regions with over 1 million total votes from films

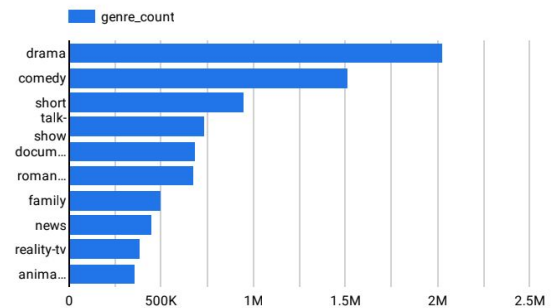


Most common genres in:

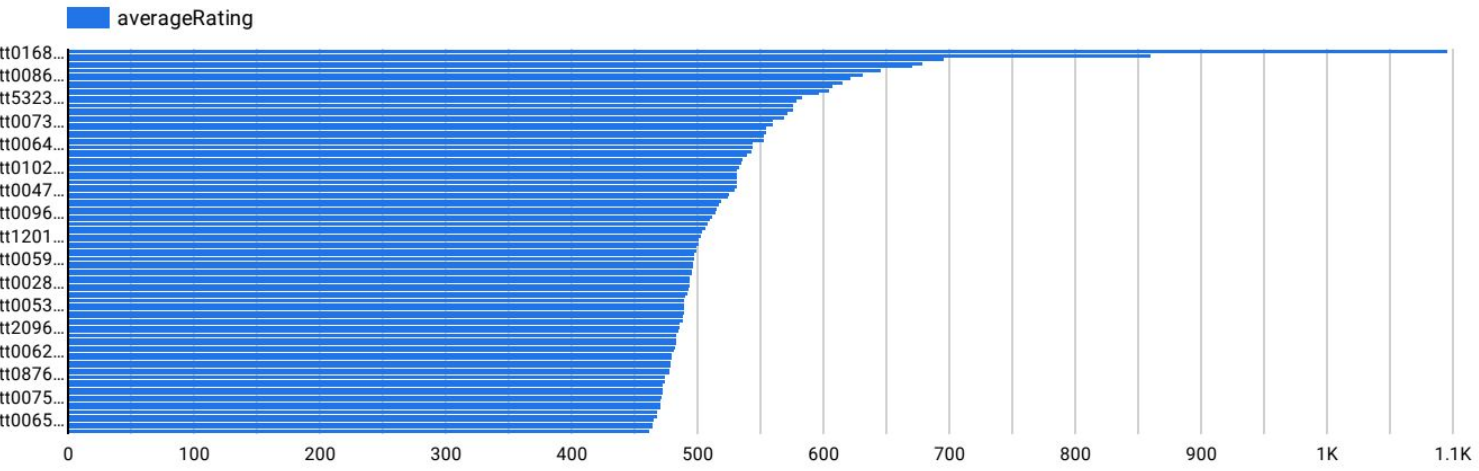
Bollywood



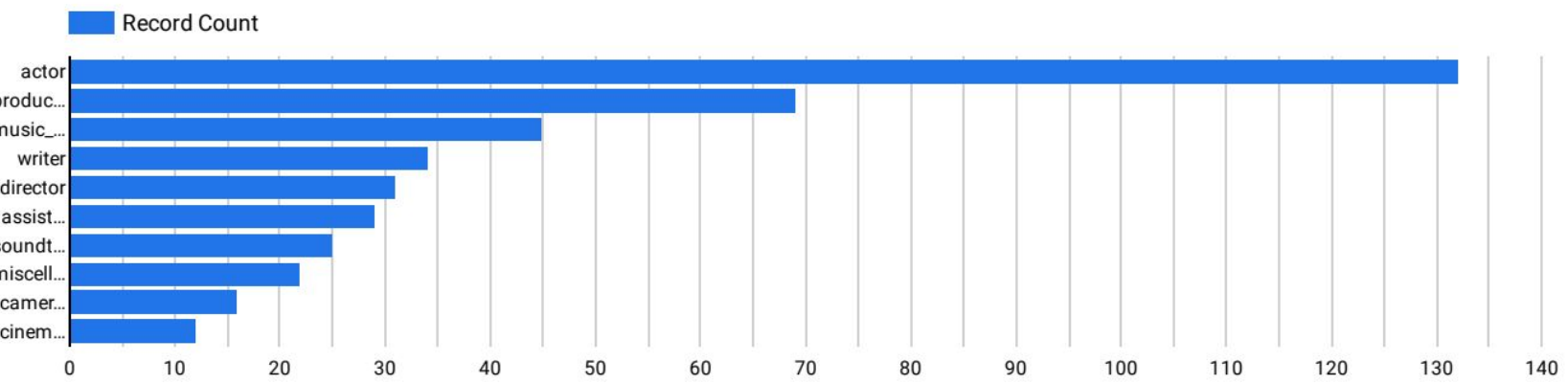
Hollywood



Movies with higher than average ratings



Primary professions of Bollywood "actors"





# Challenges and Accomplishments

---

- The hardest part was deciding how to remodel the data
  - Had to think of the possible questions/ queries that would come up and account for this when remodeling the data
- Inconsistencies across datasets
  - Had to account for differences in formatting between tables and datasets and can to update this in order to run queries
  - “War” vs “war” in the 2 genre tables
- Were able to overcome the errors we encountered and ultimately learned a lot from the interesting queries we were able to run

# Future improvements

---

- Look at the genres of both the datasets a little closer
  - Have multiple names for the same genre, so would try to make this more consistent between the datasets
- Remove duplicate titles from the Title\_Basics table in the imdb dataset
  - Director's movie ratings query shows that there are repeat Title names with different attributes for each (different ratings for the same movie)
- Combine the Title\_basics and Bolly\_movies table to make it easier to access titles and information surrounding these titles
  - Can then combine the 2 cast tables, 2 directors tables and 2 genres tables