

# IMDb and Bollywood



Jacob Helwig and Ian Chang

# Raw IMDb Tables

imdb_staging.NameBasics		
	nconst	STRING
	primaryName	STRING
	birthYear	STRING
	deathYear	STRING
	primaryProfession	STRING
	knownForTitles	STRING

imdb_staging.TitleBasics		
	tconst	STRING
	titleType	STRING
	primaryTitle	STRING
	originalTitle	STRING
	isAdult	INTEGER
	startYear	INTEGER
	endYear	STRING
	runtimeMinutes	STRING
	genres	STRING

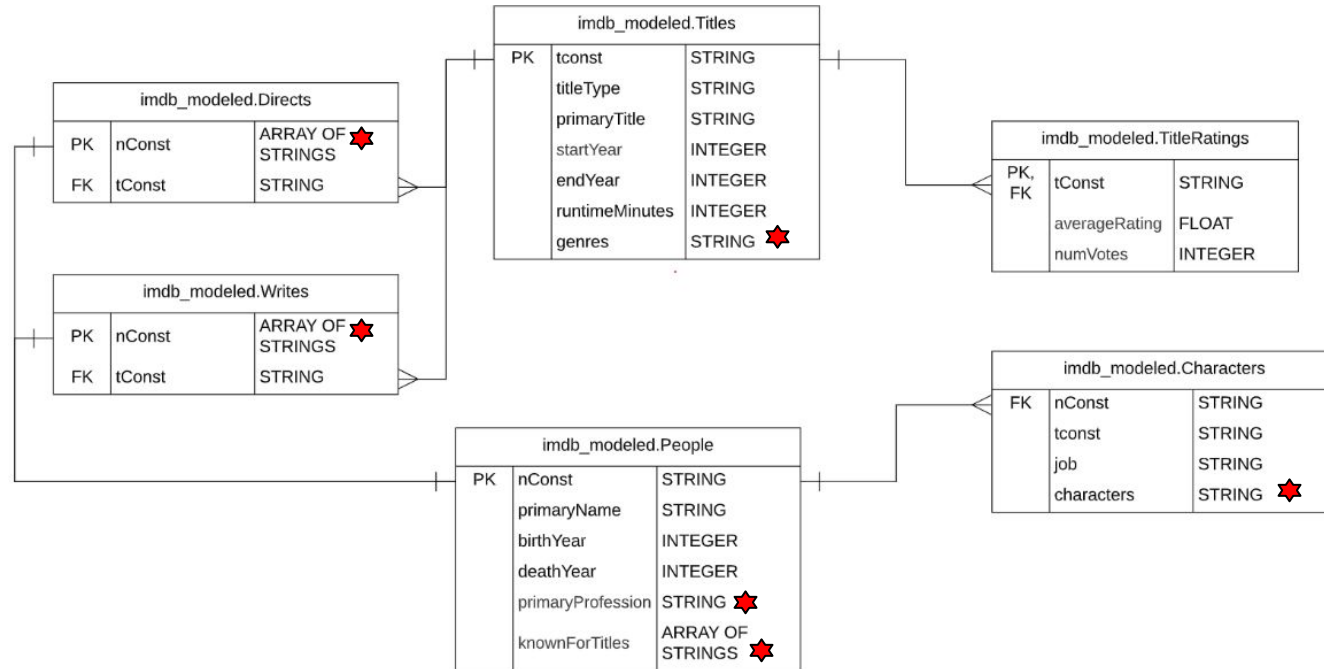
imdb_staging.TitleCrew		
	tconst	STRING
	directors	STRING
	writers	STRING

imdb_staging.TitleRatings		
	tconst	STRING
	averageRating	FLOAT
	numVotes	INTEGER

imdb_staging.TitlePrincipals		
	tconst	STRING
	ordering	INTEGER
	nconst	STRING
	category	STRING
	job	STRING
	characters	STRING

- No primary / foreign keys
- Issues with data types (years, runtime)
- Multiple entities in TitleCrew

# SQL Transforms

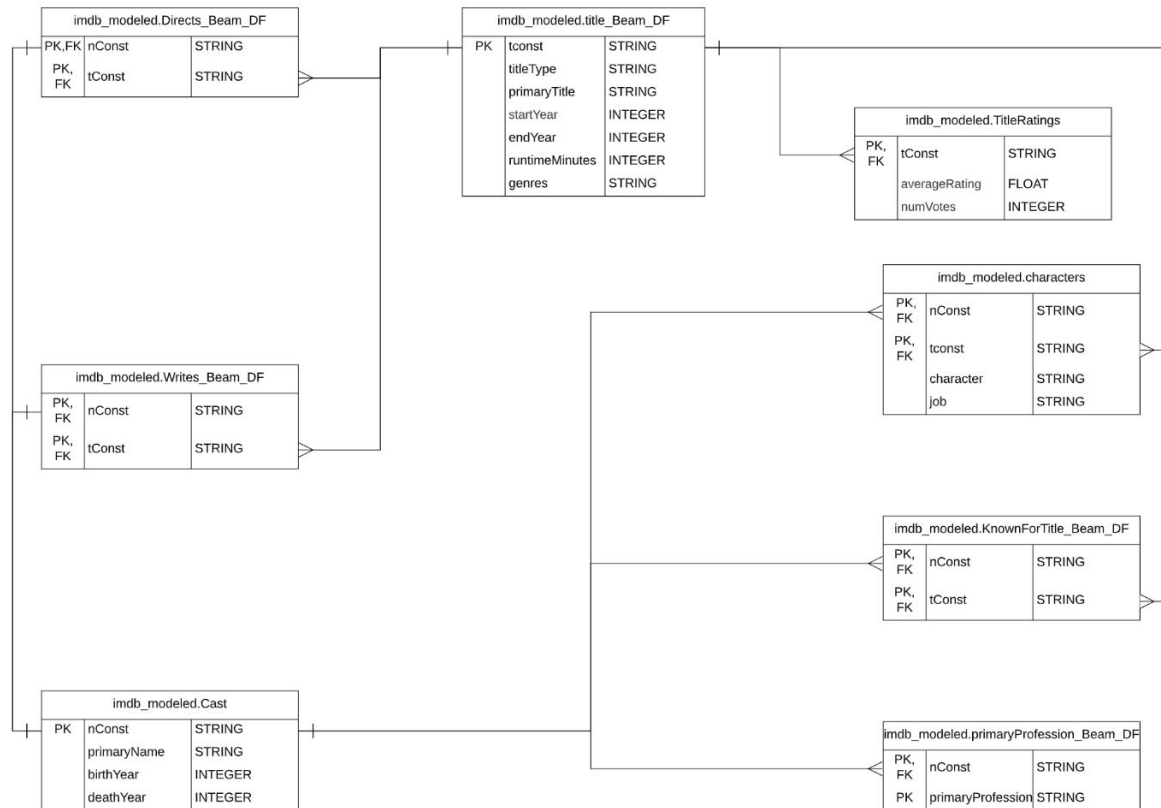


- SQL Transforms
- Remove duplicates to establish primary keys
- Break TitleCrew into two tables, only add directors with titles in Titles
- Fix Data Type Issues

## Remaining Issues

- 1NF violations (6)
- 2NF violations (primaryProfessions, knownForTitles)
- No primary key in Characters table

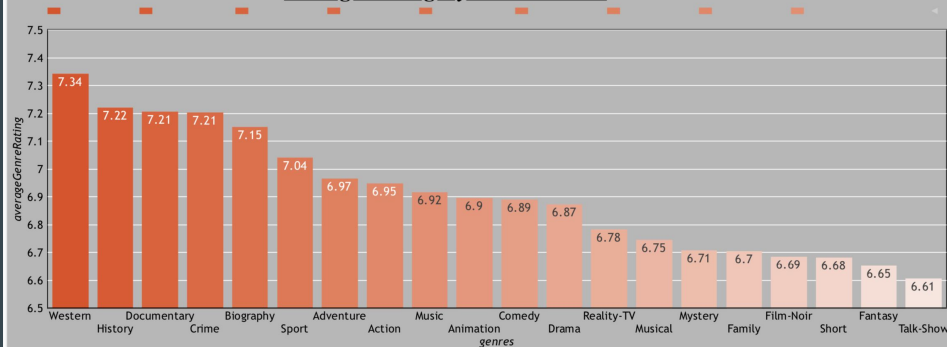
# Beam Transforms



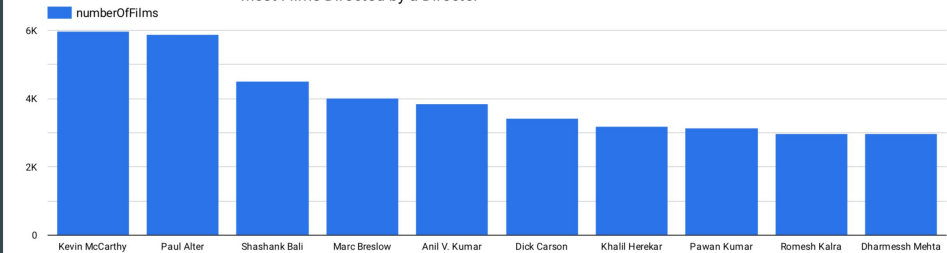
- Resolved 6 1NF issues
- Established Composite PK relationship in Directs and Writes
- Using the SQL query in beam, FK restriction on nConst in Writes and Directs, and on tConst in knownForTitle
- 2NF violations avoided by breaking up People
- Established composite PK in knownForTitle and primaryProfession

# Exploratory Data Analysis

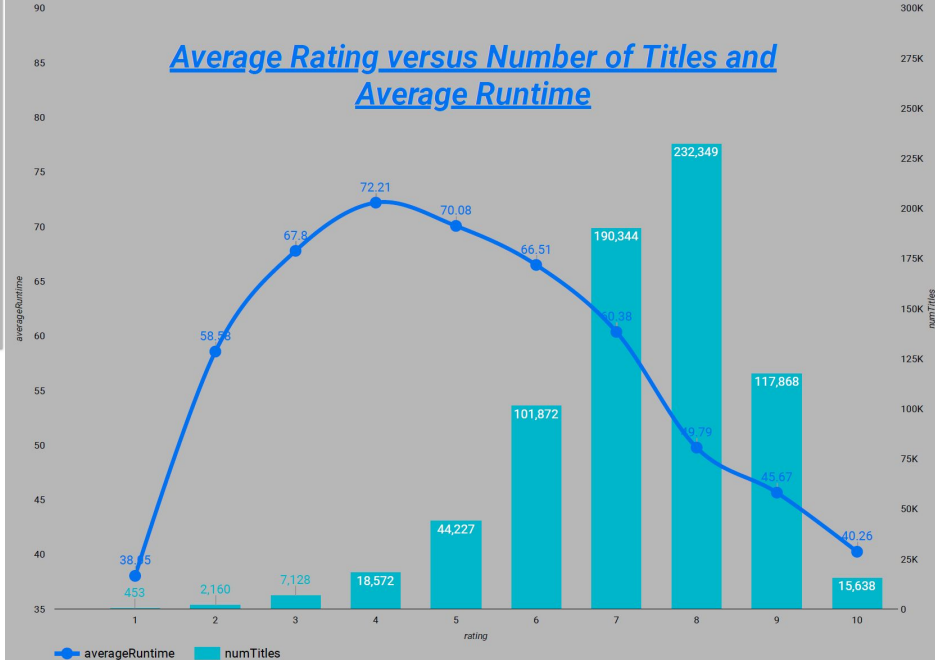
Average Rating By Genre-IMDb



Most Films Directed by a Director



Average Rating versus Number of Titles and Average Runtime



# Secondary Dataset: Bollywood DB, Raw Data

## Bollywood

	Year	Title	Director	Cast	Genre	Release_Month	Release_Date	Highest_Grosser_By_Year_in_crores_
0	1920	Shakuntala	Shree Nath Patankar	unknown	mythology,drama	None	None	None
1	1921	Belgian Emperor's Visit To India	Nitin Bose	unknown	documentary	None	None	None
2	1921	Bilet Pherat a.k.a. England Returned	Dhirendranath Ganguly	dhirendranath,chakraborty,sushilabala,kunjalal...	social,comedy	None	None	None
3	1921	King Gopichand a.k.a. Gopichand	Vishnupant Divekar	mama,koregaonkar,hira,bhatt	legend	None	None	None
4	1922	Bhartrahari a.k.a. King Bhartrahari	S. N. Patankar	tara,koregaonkar,thatte	legend	None	None	None

- Multiple entities in Bollywood table
- Date in 3 columns; data type issues (crores as String)
- 1NF in cast, genres
- Similar entities in Actors and Actresses; height was a string

## Actors

	Name	Height_in_cm_
0	Aamir Khan	163
1	Himesh Reshamiya	163
2	Kamal Haasan	165
3	Kunal Khemu	166
4	Shahid Kapoor	167

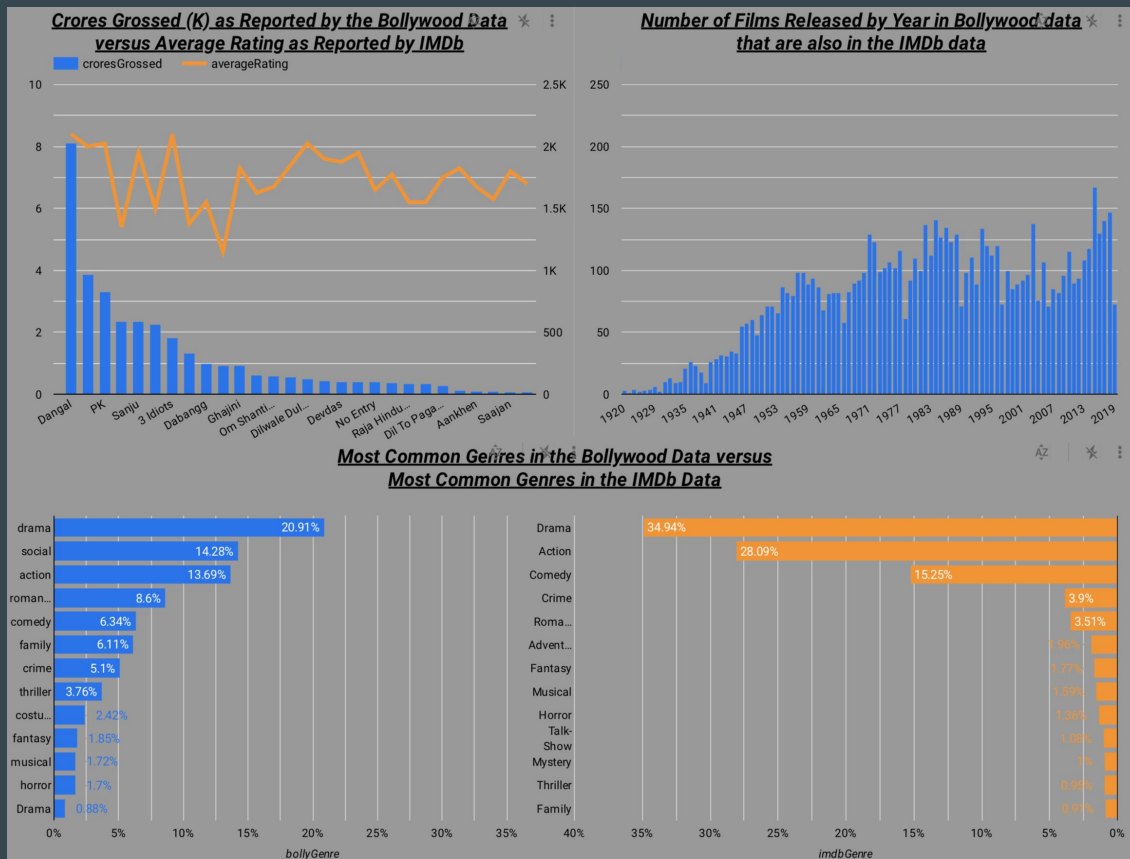
## Actresses

	Name	Height_in_cm_	Debut_aslead_role	_____
0	Ameesha Patel	152	Kaho Naa... Pyaar Hai	None
1	Soha Ali Khan	152	Rang De Basanti	None
2	Tanuja Samarth	152	Hamari Yaad Aayegi	None
3	Alia Bhatt	155	Student of the Year	None
4	Farida Jalal	155	Laal Patthar	None



# Cross-dataset visualizations

- Look at the croresGrossed reported by the bollywood db versus the average rating reported by the imdb db
- Shared movies between imdb\_modeled and bollywood\_modeled
- Compare counts of genres in each data set





Q&A