

## Movie Recommendation System

### [dataset](#)

The dataset contained two databases which were merged based on movie\_id. Majority of the work was done in data per-processing. After merging the table contained

Data columns (total 23 columns):

#	Column	Non-Null Count	Dtype
0	budget	4809 non-null	int64
1	genres	4809 non-null	object
2	homepage	1713 non-null	object
3	id	4809 non-null	int64
4	keywords	4809 non-null	object
5	original_language	4809 non-null	object
6	original_title	4809 non-null	object
7	overview	4806 non-null	object
8	popularity	4809 non-null	float64
9	production_companies	4809 non-null	object
10	production_countries	4809 non-null	object
11	release_date	4808 non-null	object
12	revenue	4809 non-null	int64
13	runtime	4807 non-null	float64
14	spoken_languages	4809 non-null	object
15	status	4809 non-null	object
16	tagline	3965 non-null	object
17	title	4809 non-null	object
18	vote_average	4809 non-null	float64
19	vote_count	4809 non-null	int64
20	movie_id	4809 non-null	int64
21	cast	4809 non-null	object
22	crew	4809 non-null	object

Here, genres, keywords, overview, cast (first 3), crew (director only) were merged together in tag field. And new data field contained title, movie\_id and tag only.

Word vectorization was done after selecting top 5000 words from entire tags. Tags were stemmed to remove duplicates. Cosine similarity was calculated against every other vector. Based upon their similarity top 5 movies were selected for recommendation.