

INFO 7374 : Algorithmic Digital Marketing

Group 6 : Omkar Kalange and Sanghamitra Shanmugam

Claa:

https://codelabs-preview.appspot.com/?file_id=1XOuJ7S7FMm9Nfqhbg6RiSjJA26sIJ1vlo1iktMOm1D4#0

Google doc:

<https://docs.google.com/document/d/1XOuJ7S7FMm9Nfqhbg6RiSjJA26sIJ1vlo1iktMOm1D4/edit?ts=5f774829#>

MovieLens 25M : Insights and Analytics

Summary	Uncovering marketing insights and performing analysis on the 25M reviews for movies from 1960 to 2019
Category	Web

[About the Dataset](#)

[Working on Dataset using XSV](#)

[Trifacta Wrangler : Data Cleaning and creating flows](#)

[Created recipe for regular expressions for deducing new columns](#)

[Recipe for Date Time objects](#)

[Snowflake : Staging into Data Warehouse](#)

[Connecting Snowflake to Salesforce](#)

[Analytics Dashboard using Salesforce Einstein Analytics](#)

[Movie Lens Dashboard](#)

[Promotions](#)

[Recommendations](#)

[Insights in all genres, ratings and tags](#)

[Strength and Weaknesses of Wrangling tools](#)

[XSV](#)

[Trifacta](#)

[Answered Questions related to Dashboards](#)

About the Dataset

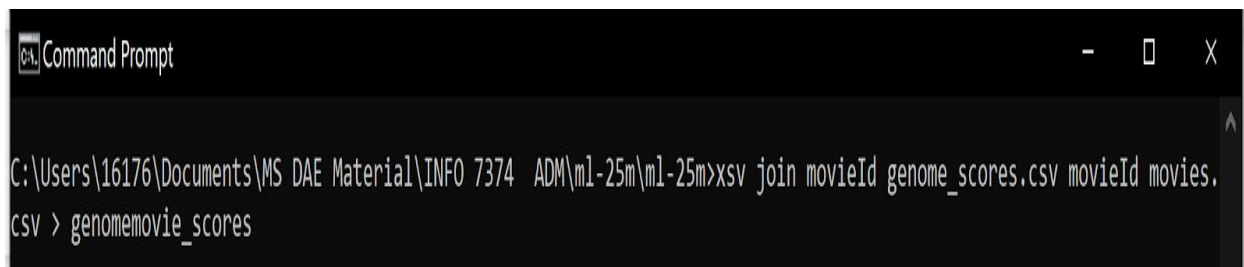
This dataset describes 5-star rating and free-text tagging activity from [MovieLens] (<http://movielens.org>), a movie recommendation service. These ratings were given for the movies during 1995 to 2019. Users were selected at random for inclusion. All selected users had rated at least 20 movies. The dataset also contains the relevance of the tags given by users.

The files contained in dataset were :

- Ratings.csv
- GenomeTags.csv
- Movies.csv
- GenomeScores.csv

Working on Dataset using XSV

1. Using XSV Join command - Performed join operations using the MovieID to get the title for relevance scores of different tags



```
Command Prompt
C:\Users\16176\Documents\MS DAE Material\INFO 7374 ADM\m1-25m\m1-25m>xsv join movieId genome_scores.csv movieId movies.csv > genomemovie_scores
```

```
Command Prompt
C:\Users\16176\Documents\MS DAE Material\INFO 7374 ADM\ml-25m\ml-25m>xsv join tagId genomemovie_scores.csv tagId genome-tags.csv > genome_movies.csv
```

2. Taking samples from genome scores using XSV Sample : For loading into Snowflake Datawarehouse

```
Command Prompt
C:\Users\16176\Documents\MS DAE Material\INFO 7374 ADM\ml-25m\ml-25m>xsv sample 900000 genome-scores.csv > genome_score_s.csv
```

3. Calculated the Stats of the Sampled data for ratings using XSV stats command

```
C:\Users\16176\Documents\MS DAE Material\INFO 7374 ADM\ml-25m\ml-25m>xsv stats movies_ratings.csv
field,type,sum,min,max,min_length,max_length,mean,stddev
movieId,Integer,3794322266,1,209133,1,6,32586.072363448922,48975.42959183392
title,Unicode,, Days of Summer (2009),Я хуею ,1,153,,
release_year,Integer,231945957,1408,2150,4,4,1991.978332188251,20.823757079725702
genres,Unicode,,(no genres listed),Western,3,77,,
,NULL,,,0,0,,
year_of_rating,Integer,234067971,1996,2019,4,4,2010.2024304362765,6.396004684215687
Avg_rating,Float,384777.19999992487,0.5,5,1,3,3.304510477499171,0.8757024519860734
```

Trifacta Wrangler : Data Cleaning and creating flows

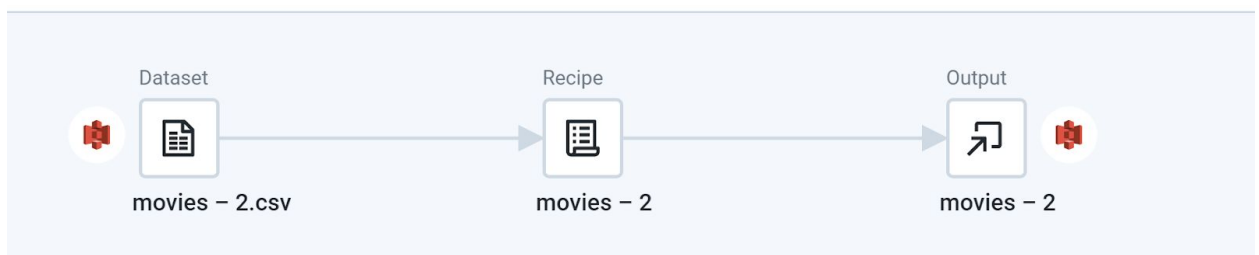
Wrangled the datasets using Trifacta

1. Joined the datasets using XSV and imported them to Trifacta for further analysis

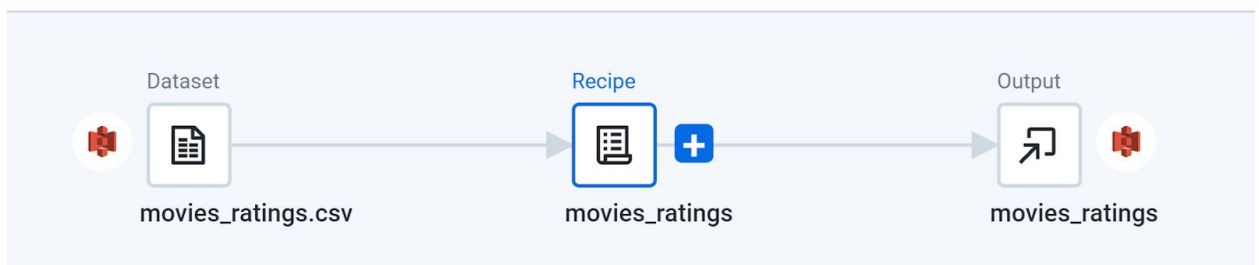
 Rating-dates




 movie rating cleaning



 Missing Values for Datasets



2. Created recipe for regular expressions for deducing new columns

 movies – 2

Edit Recipe

Add

...

Recipe

Data

Steps Preview

- 1 Replace matches of "\" from title with "
- 2 Extract matches of `/[\d]{4}/` from title
- 3 Rename title1 to 'release_year'
- 4 Delete rows where `ISMISMATCHED(release_year, [Datetime,'yy',yyyy])`
- 5 Delete rows where `ISMISSING([release_year])`
- 6 Replace matches of `/[\d]{4}*/` from title with "
- 7 Sort rows by title
- 8 Sort rows by title
- 9 Split title on delimiters matching '(' into 2 columns
- 10 Delete title2
- 11 Rename title1 to 'title'

3. Recipe for Date Time objects

movies_ratings

Edit Recipe

Add

...

Recipe

Data

Steps Preview

1 Create removesymbols_title from REMOVESYMBOLS(title)

2 Delete title

3 Delete rows where ISMISMATCHED(release_year, [Datetime;'yy','yyyy'])

4 Change removesymbols_title type to String

5 Delete rows where ISMISSING([removesymbols_title])

Steps

5

Updated

Last Tuesday at 11:17 AM

Created

Last Tuesday at 11:14 AM

ratings5l

Edit Recipe

Add

...

Recipe

Data

Steps Preview

1 Delete userId

2 Change timestamp type to Datetime, mm-dd-yy, mm*dd*yy

3 Change timestamp type to Datetime, mm-dd-yy, mm*dd*yy

4 Create unixtimeformat_timestamp from UNIXTIMEFORMAT(timestamp, 'yyyy-MM-dd')

5 Delete unixtimeformat_timestamp

6 Change timestamp type to Integer

7 Create Unixtimestamp from timestamp * 1000

8 Delete timestamp

9 Create Date from UNIXTIMEFORMAT(Unixtimestamp, 'yyyy-MM-dd')

10 Delete Unixtimestamp

Recipe Data

8 Delete timestamp

9 Create Date from
UNIXTIMEFORMAT(Unixtimestamp, 'yyyy-MM-dd')

10 Delete Unixtimestamp

11 Change Date type to Datetime, yy-mm-dd, yyyy*mm*dd

12 Create year_Date from YEAR(Date)

13 Create month_Date from MONTH(Date)

14 Delete month_Date

15 Create average_rating from AVERAGE(rating) grouped by movieId, year_Date

16 Create Avg_rating from
ROUND(average_rating, 1)

17 Delete average_rating

18 Sort rows by year_Date

19 Sort rows by movieId

20 Rename year_Date to 'year'

4. Jobs ran successfully eliminating the null values and the data is ready to be staged in Snowflake

Completed stages

✓

Transform with profile

Completed Last Tuesday at 11:19 AM, started Last Tuesday at 11:17 AM • Ran for 1 min

Environment Spark

100% valid values

0% mismatching values

0% missing values

[View steps and dependencies](#)
[View profile](#)

✓

Publish

Completed Last Tuesday at 11:19 AM, started Last Tuesday at 11:19 AM • Ran for <1 sec

Activity

movies_ratings.csv

Completed

[View all](#)

Completed stages

✓

Transform with profile
Completed Last Monday at 9:46 PM, started Last Monday at 9:43 PM • Ran for 2 min
Environment Spark

100% valid values


0% mismatching values

0% missing values

[View steps and dependencies](#) [View profile](#)

✓

Publish
Completed Last Monday at 9:46 PM, started Last Monday at 9:46 PM • Ran for <1 sec
Activity

 ratings5l.csv	<div>✓ Completed</div>
---	------------------------

[View all](#)

Snowflake : Staging into Data Warehouse

1. After wrangling data into Trifacta, the data has been loaded into Snowflake. The first thing to be done was to create two tables : Movie Ratings and Genome Relevance

Databases > MOVIES > MOVIE_RATINGS (PUBLIC)

Tables

Views


Schemas

Stages

File Formats

Sequences

Pipes

 Load Table

Column Name	Ordinal ▲	Type	Nullable
MOVIEID	1	NUMBER(38,0)	false
RELEASE_YEAR	2	NUMBER(38,0)	true
GENRES	Table Name: MOVIE_RATINGS		true
YEAR_OF_RATING	4	NUMBER(38,0)	true
AVG_RATING	5	FLOAT	true
TITLE	6	VARCHAR(200)	true

Databases > MOVIES > GENOME_RELEVANCE (PUBLIC)

Tables

Views

Schemas

Stages

File Formats

Sequences

Pipes

 Load Table

Column Name	Ordinal ▲	Type	Nullable
MOVIEID	1	NUMBER(38,0)	false
TITLE	Table Name: GENOME_RELEVANCE		false
RELEASE_YEAR	3	NUMBER(38,0)	false
TAG	4	VARCHAR(100)	false
RELEVANCE	5	FLOAT	false
GENRE	6	VARCHAR(100)	true

2. Creating file formats for csvs to be loaded into database tables :

Edit File Format

Name CSV

Compression Method Auto

Column separator Comma

Row separator New Line

Header lines to skip 1

Field optionally enclosed by None

Null String \\N

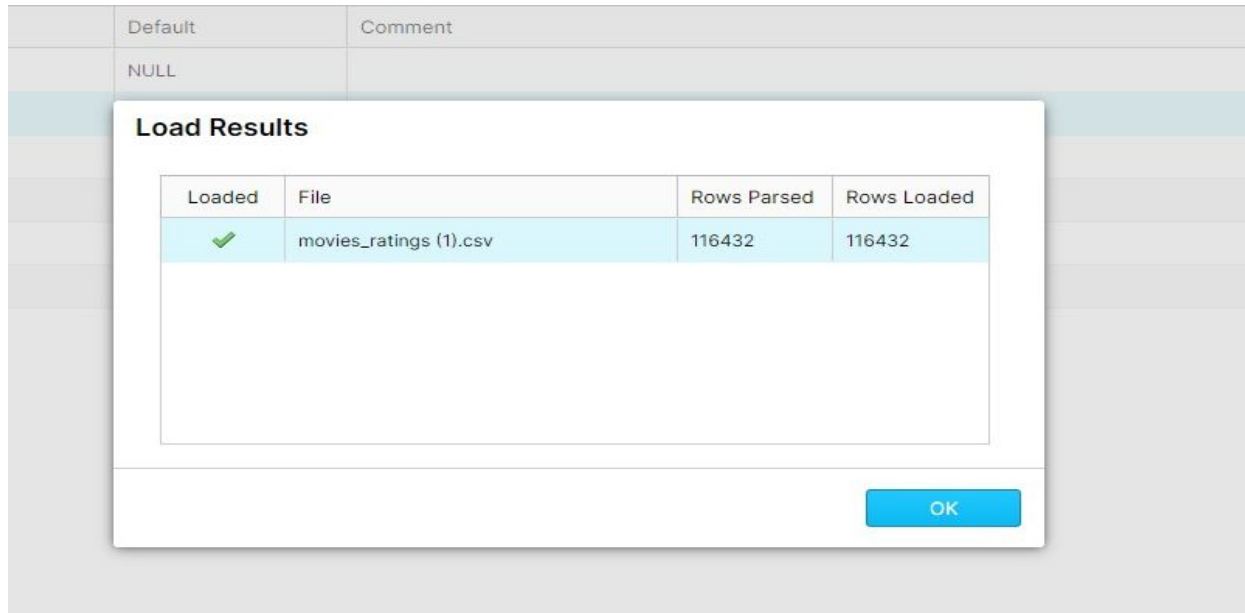
☐ Trim space before and after

☒ Error on Column Count Mismatch

Escape Character None

[Show SQL](#) Cancel Finish

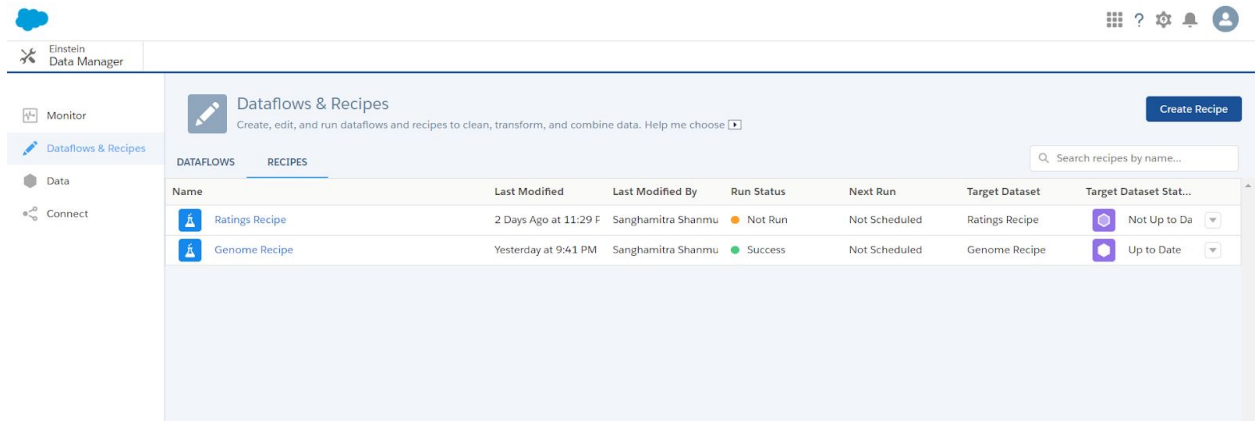
3. Loading the data into tables, ensuring the same count of rows and columns from csv to database table :



4. Verifying the data using SQL queries on the data warehouse :

1 SELECT * FROM MOVIE_RATINGS;							
Results Data Preview Open History							
✓ Query ID SQL 93ms 116,432 rows							
Filter result... Download Copy Columns							
Row	MOVIEID	RELEASE_YEAR	GENRES	YEAR_OF_RATING	AVG_RATING	TITLE	
1	1	1995	Adventure Animation Childre...	2002	4.1	Toy Story	
2	1	1995	Adventure Animation Childre...	2006	3.7	Toy Story	
3	1	1995	Adventure Animation Childre...	2015	3.7	Toy Story	
4	1	1995	Adventure Animation Childre...	2008	3.7	Toy Story	
5	1	1995	Adventure Animation Childre...	2001	4.1	Toy Story	
6	1	1995	Adventure Animation Childre...	1998	4	Toy Story	

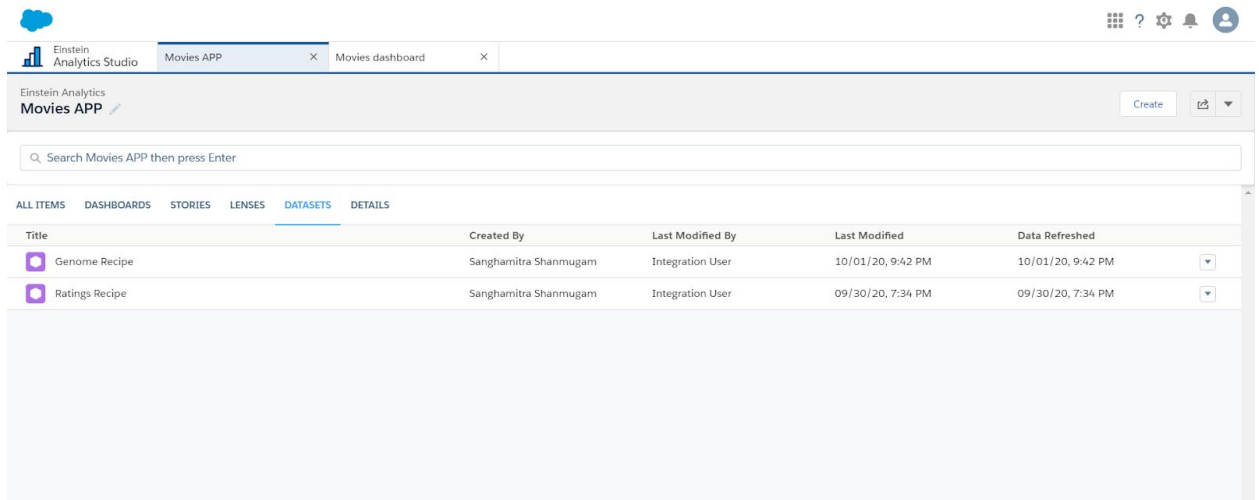
2) Creating Recipes out of connected data



The screenshot shows the Einstein Data Manager interface. On the left is a sidebar with navigation options: Monitor, Dataflows & Recipes (selected), Data, and Connect. The main area is titled 'Dataflows & Recipes' with a subtitle 'Create, edit, and run dataflows and recipes to clean, transform, and combine data. Help me choose'. A 'Create Recipe' button is in the top right. Below the title is a search bar 'Search recipes by name...'. A table lists the recipes:

Name	Last Modified	Last Modified By	Run Status	Next Run	Target Dataset	Target Dataset Stat...
Ratings Recipe	2 Days Ago at 11:29 F	Sanghamitra Shanmu	Not Run	Not Scheduled	Ratings Recipe	Not Up to Da
Genome Recipe	Yesterday at 9:41 PM	Sanghamitra Shanmu	Success	Not Scheduled	Genome Recipe	Up to Date

3) Connecting the App to the recipe



The screenshot shows the Einstein Analytics Studio interface. At the top, there are tabs for 'Einstein Analytics Studio', 'Movies APP' (selected), and 'Movies dashboard'. Below the tabs is a search bar 'Search Movies APP then press Enter'. The main area is titled 'Einstein Analytics Movies APP'. Below the title is a navigation bar with tabs: ALL ITEMS, DASHBOARDS, STORIES, LENSES, DATASETS (selected), and DETAILS. A table lists the datasets:

Title	Created By	Last Modified By	Last Modified	Data Refreshed
Genome Recipe	Sanghamitra Shanmugam	Integration User	10/01/20, 9:42 PM	10/01/20, 9:42 PM
Ratings Recipe	Sanghamitra Shanmugam	Integration User	09/30/20, 7:34 PM	09/30/20, 7:34 PM

Analytics Dashboard using Salesforce Einstein Analytics

Movie Lens Dashboard

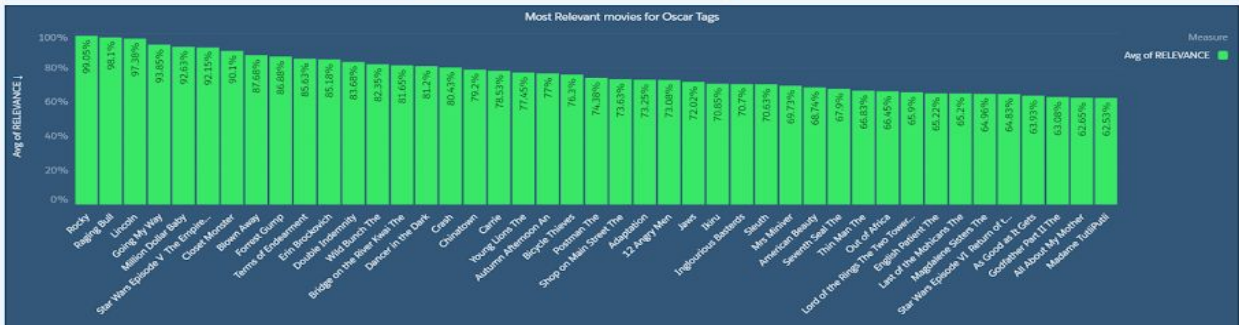
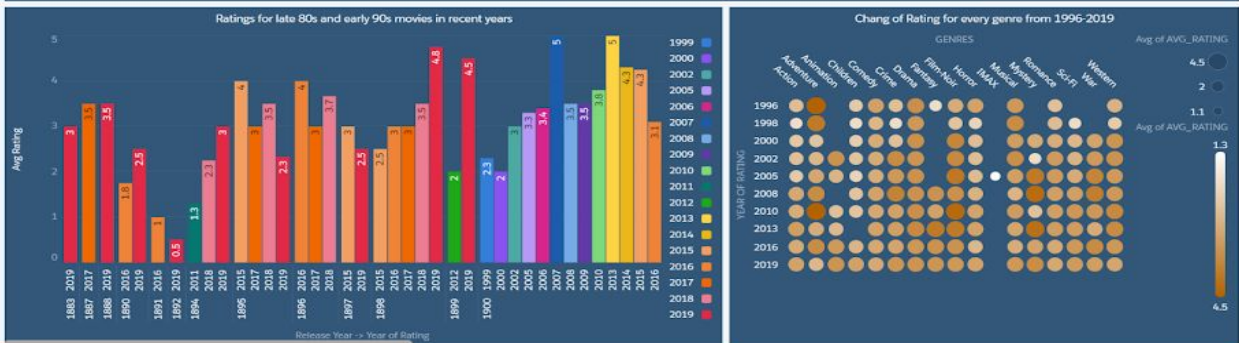
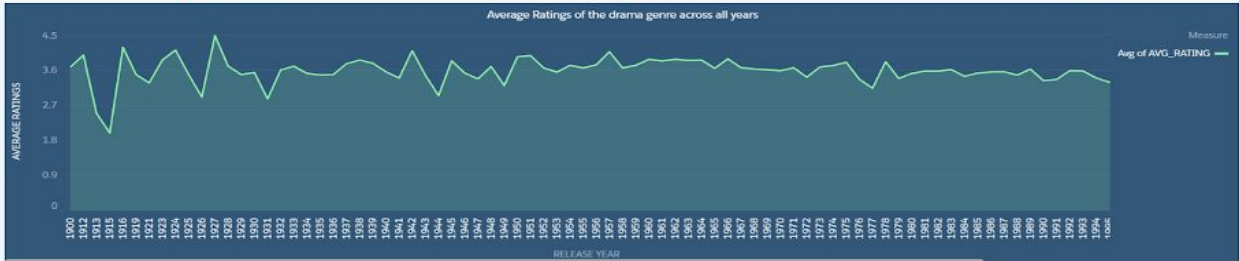
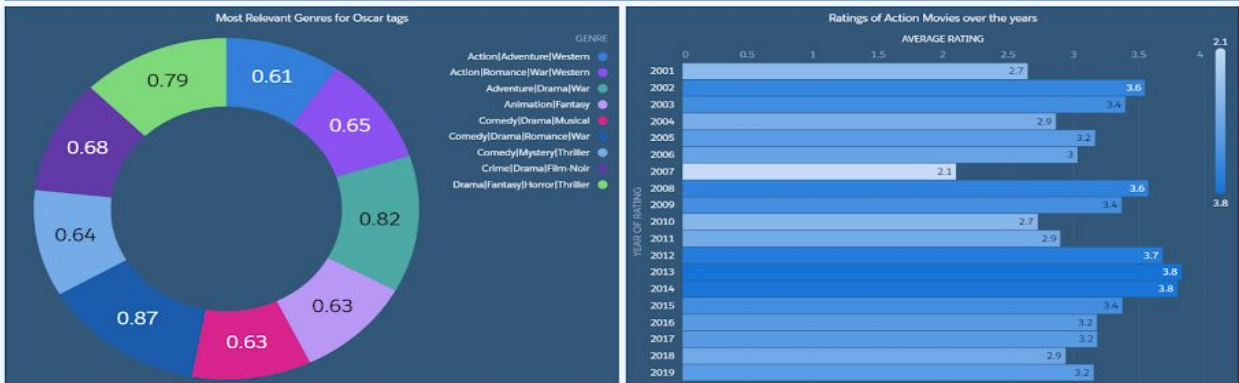
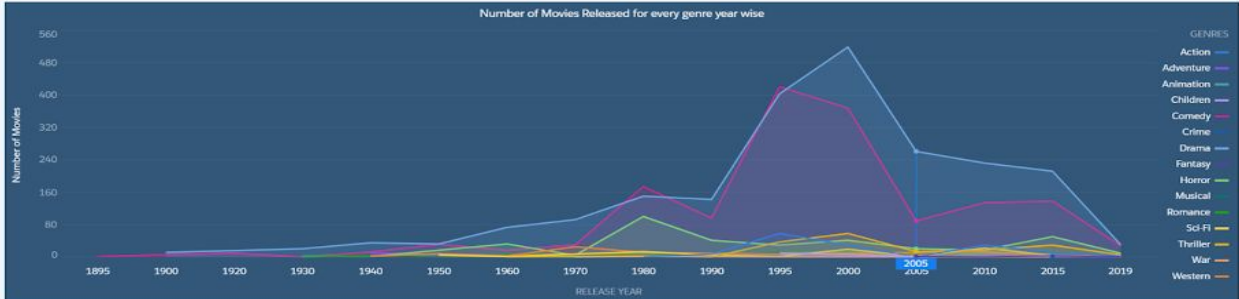
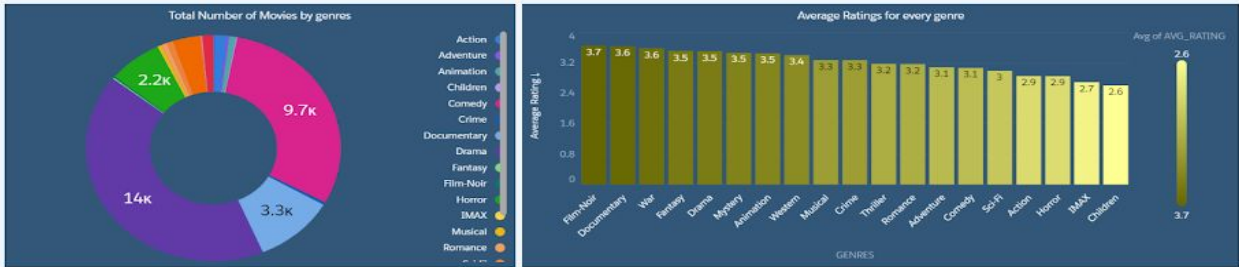
Total Genres and their categories:- 1,279

Total Tags:- 1,128

Promotions

We created these below graphs to get insights-

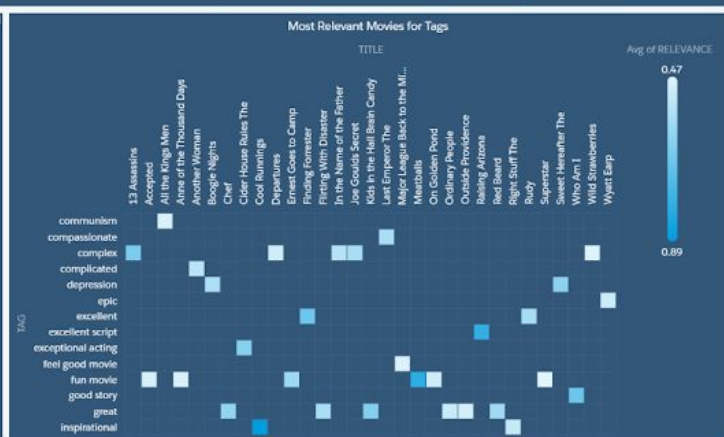
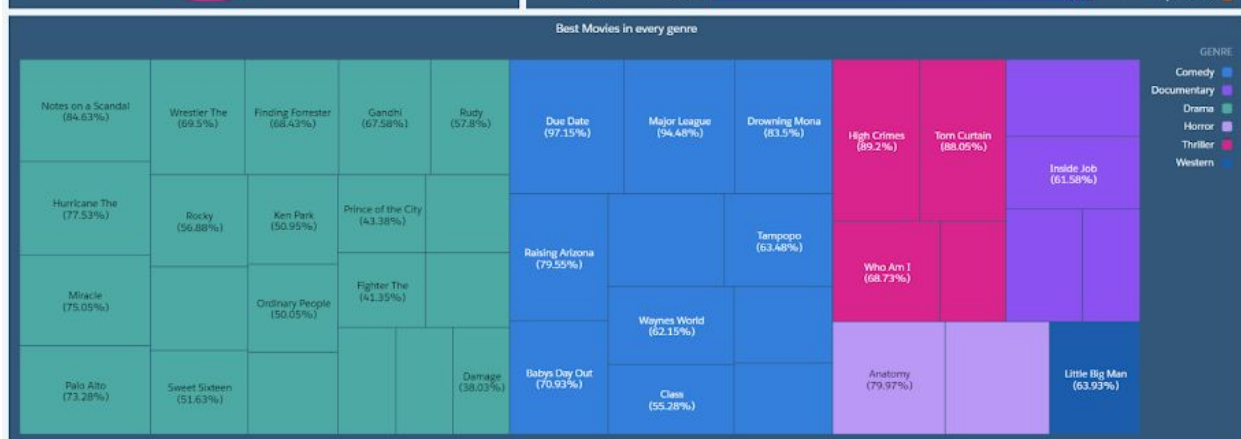
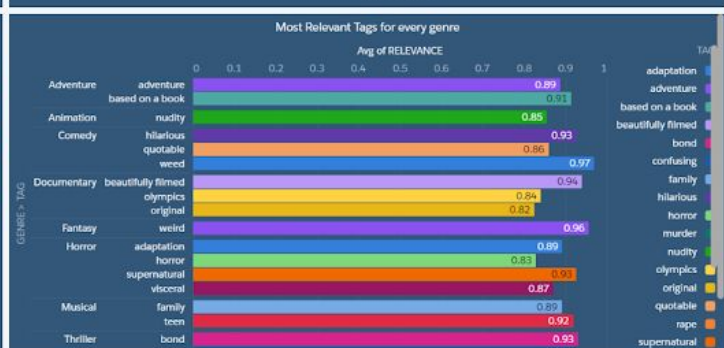
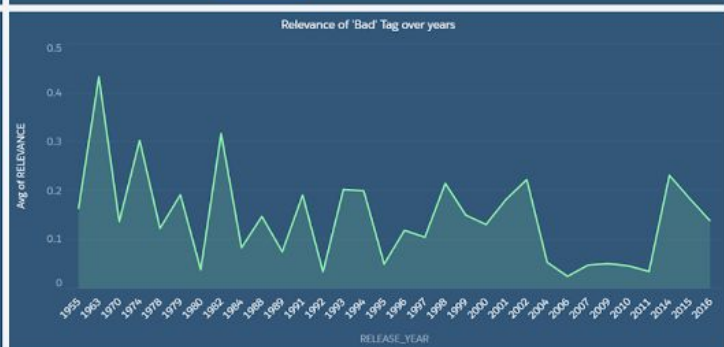
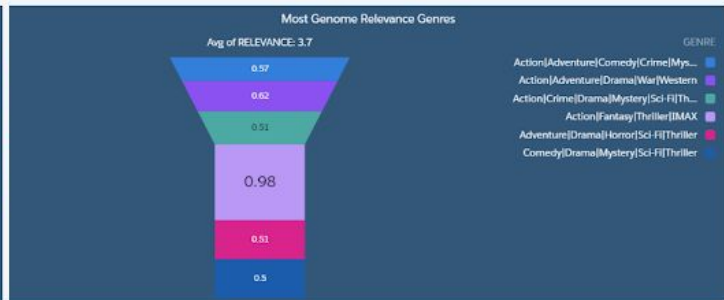
1. **Number of movies in every genre and their ratings-** To get insights on the range of movies released every year and their reach to the public
2. **Track the genre of movies released across all years-** To see the effect of ratings make a significant change in the number of movies released in genres
3. **Analyse users ratings for old movies in recent years -** To see the demand for old movies in the current generation and compare it against the genre of movies released recently
4. **Track the movies with the tag oscar with its ratings and genre:-** To get insights on the most relevant genres with oscar tags



Recommendations

We created these below graphs to get insights-

1. **Most Relevant Tags in the 20th century vs 90s vs recent years** - These three graphs combined together can be used to get insights on the type of movies that are being made across all the years
2. **Movie preferences for every genre** - Well rated movies for every genre based on ratings and tags
3. **Study the relevance of 'Bad' Tags over the years** - To get insights on how relevant the tags related to the word 'bad' for all the movies across all years
4. Further we can use collaborative filtering or association rules to recommend movies based on these insights.



Insights in all genres, ratings and tags

Metrics:

Movie Rating

Relevance scores for tags

1. **Rating Fluctuations & Effect on movie releases-** Even though the Average Ratings for Film Noir and Documentary is higher compared to other genres, more movies are being released in Drama and Comedy genre compared to the rest of the genres.
2. **Tracking the genres of movies being released-** In the 19th century, Comedy seems to be the most popular genre, however since the 90s there has been a significant surge in the genre Drama followed by comedy which is not far behind
3. **Oscar Tag vs Genre-** Majority of the Oscar tags are associated with the genre which is a mix of Comedy, Drama, Romance and War genre followed by Adventure, Drama and War. Almost all the subcategories of the genres in this category have war associated with them.
4. **Tracking ratings of the genre Drama-** Since the overall drama genre has the maximum number of movies released, we decided to analyse its ratings across all years. The ratings of drama movies has been very consistent throughout though its facing a downfall in the recent years compared to the early days
5. **Ratings of old movies in the recent years-** The movies released in 1900 to 1903 seem to have good average ratings compared to the previous years. However, 80s movies don't have great ratings in recent years.
6. **Change of Ratings for every genre in the recent years-** Even though Adventure movies had good ratings in the late 90s, the ratings for the same movies seem to go down in the recent years which shows that the change in the kind of movies preferred by the public. Likewise, even though Film-Noir was not much appreciated in the late 90s, they seem to have gathered attention in the recent years.
7. **Most Relevant Tags in the early 18th century, 90s and recent years-** In the early 20th century, 'Melancholy' tag seems to be the most associated with the movies and in the late 90s, tags like 'original', 'corruption' and 'vengeance' seems to be more related whereas in the years of 2000-2019 'powerful ending', 'brutality', 'visually appealing' have better ratings.
8. **Most Relevant Tags for each Genre-** For the adventures Genre, the most relevant tag is 'based on a book' and for thriller 'bond' seems to have gained popularity the most
9. **Recommending best movies in every popular genre -** Filtering based on the tags focusing on good scripts, good direction and the tags relating to it and the genre, and by also filtering out the movies with the maximum relevance for the above tags, we recommend movies.

10. **Recommending movies based on appreciated tags** - We created a heat map based on the tags that focuses on good script, acting etc and movies that have a high relevance ratings with the above tags
11. **Recommending old movies** - All the best rated movies in the 20th century mainly come from three genres which are Drama, Musical, Comedy and Film-Noir among which Drama takes a huge percentage of the movies. Based on high relevance of tags that mean well, we can recommend movies.

Strength and Weaknesses of Wrangling tools

XSV

Xsv is a program for indexing, slicing, evaluating, separating and joining CSV files on the command line.

- In contrast to other resources, it is much quicker.
- But there is only a small number of operations in terms of Wrangling that can be done as of now as per the source github.
- It is mainly used for handling 10s of GBs of very wide csv files.

Trifacta

For data discovery and self-service data preparation for study, Trifacta develops data wrangling tools.

- Trifacta Wrangler is a connected framework for downstream analytics and visualization to transform knowledge.
 - It enables analyst teams in an enterprise to explore and turn data with centralized security, governance and operationalization management using self-service.
 - Particularly for individuals who do not like coding and want to prepare their data for forecasts / modeling, it is an easy to use wrangling tool.
 - It can handle even some complicated operations such as data imputation, quickly creating new columns from existing columns.
 - It also provides you with a facility for using different cloud platforms to export your data.
-

Answered Questions related to Dashboards

Which columns are dimensions, which columns are measures?

Dimensions -

1. MovieId
2. Title
3. Genre
4. Tag
5. Release Year
6. Year of Rating

Measure

1. Average Rating
2. Relevance

How would you generate new dimensions? What will you do to summarize measures?

1. We can create new columns/dimensions in wrangling tools mentioned above as well as Salesforce.
2. SAQL is used to create new dimensions from existing dimensions

Who would use this dashboard?

1. Business Analyst
2. Data Scientist
3. Data Analyst
4. Business Intelligence
5. Data Scientist - Recommendations

What value would be generated using this dashboard ?

Insights can be generated that could be useful for recommending movies to the user and also the kind of movies that are doing well in the current market. This dashboard can also be used to track how the customers response changes periodically over the years.