# APPLICATION OF DATA MINING IN MOVIE RECOMMENDATIONS

## TAN YIN YEN
## WQD180108

1

# ANALYSIS GOAL

A movie streaming company (Netflix) seeks to maximize customer's retention by recommending highly rated movies with DVD or streaming options available to their users. A movie can be predicted as a good or bad movie based on the movie information, box office data and sentiment score of user reviews on the movie.

By predicting whether a movie is good or bad based on its reviews and box office achievement, the movie streaming company can filter out latest movies with DVD or streaming options available that are highly rated and recommend them to its users. Customers who are satisfied with the movie recommendations are more likely to subscribe to the movie streaming service in the next month.

# OBJECTIVES

To come out with a model to predict whether a movie is good or bad based on its box office performance and user reviews.

To recommend users the top and latest movies with DVD or streaming options available based on the prediction.

# PROCESS FLOW

Data Acquisition

Data Management

Processing of Data

Descriptive Analysis

Diagnostic Analysis

Predictive Analysis

Prescriptive Analysis

# DATA ACQUISITION

Scope: Scrape data of latest 5000 movies with DVD or streaming options available.

Data acquired:

- Structured data

  - Movie information scraped from _rottentomatoes.com_

  - Movie box office data scraped from _boxofficemojo.com_

- Unstructured data
  - Movie reviews scraped from *rottentomatoes.com*



# DATA MANAGEMENT

- Store the data scraped into hive data warehouse.

# PROCESSING OF DATA

- **SEM**MA
  - ○ Sample: Import the data into csv from hive data warehouse
  - ○ Explore: Identify missing values or patterns in the data



*Histograms of Variables*

- Missing data is identified.

**Proportion of target variable**



*Pie Chart of Target Variable*

- Target variable has a proportion of approximately 50% True (good movie) and 50% False (bad movie)
- Target class is balance

- Modify: Create, transform or remove variables
  - Covert the reviews data to sentiment score
  - Impute missing values
  - Combine multiple genres into a single 'genre' cluster
  - Drop features that are not useful

# DESCRIPTIVE ANALYSIS

- Hierarchical clustering is used for descriptive analysis as the number of clusters no need to be specified.



*Segment Size Plot*

o Two clusters are formed.
o Cluster 1 contains 151 cases.
o Cluster 2 contains 2287 cases.

Good movie

Bad movie

sEMMA

*Profile plot of significant features*

Segment: 2
Count: 2287
Percent: 93.81

Segment: 1
Count: 151
Percent: 6.19

Clustered data distribution

Original distribution

- For segment 1, the overall box office distribution and number of reviews are higher than average.
- Contains very few cases (151 cases)

∴ **Good and recommended movies**

- For segment 2, the overall box office distribution and number of reviews are lower than average.
- Contains many cases (2287 cases)

∴ **More towards bad movies and requires further investigation**

# DIAGNOSTIC ANALYSIS

- Sentiment Analysis using Natural Language Processing

| Positive Sentiment | Negative Sentiment |
|---|---|
| Excellent | Worst |
| Perfect | Awful |
| Great | Boring |
| Wonderful | Waste |
| Amazing | Bad |
| Superb | Poor |
| Enjoyable | Terrible |
| Best | Dull |
| Today | Poorly |
| Fun | Disappointment |
| Enjoyed | Disappointing |
| Brilliant | Unfortunately |
| Must see | Worse |
| Loved | Stupid |
| Fantastic | Horrible |
| Liked | Mess |
| Incredible | Nothing |
| Funniest | Lame |
| Wonderfully | Lacks |
| Better than | Save |

*Top 20 positive features and negative features extracted*

# PREDICTIVE ANALYSIS

- SEM**M**A
  - Model: Apply a model to the data
- A decision tree is created for predictive analysis as the model is simple and easier to interpret.



The nodes are split based on
- Information gain: Higher information gain.
- Interpretability: Easier to understand

- First level of decision tree:
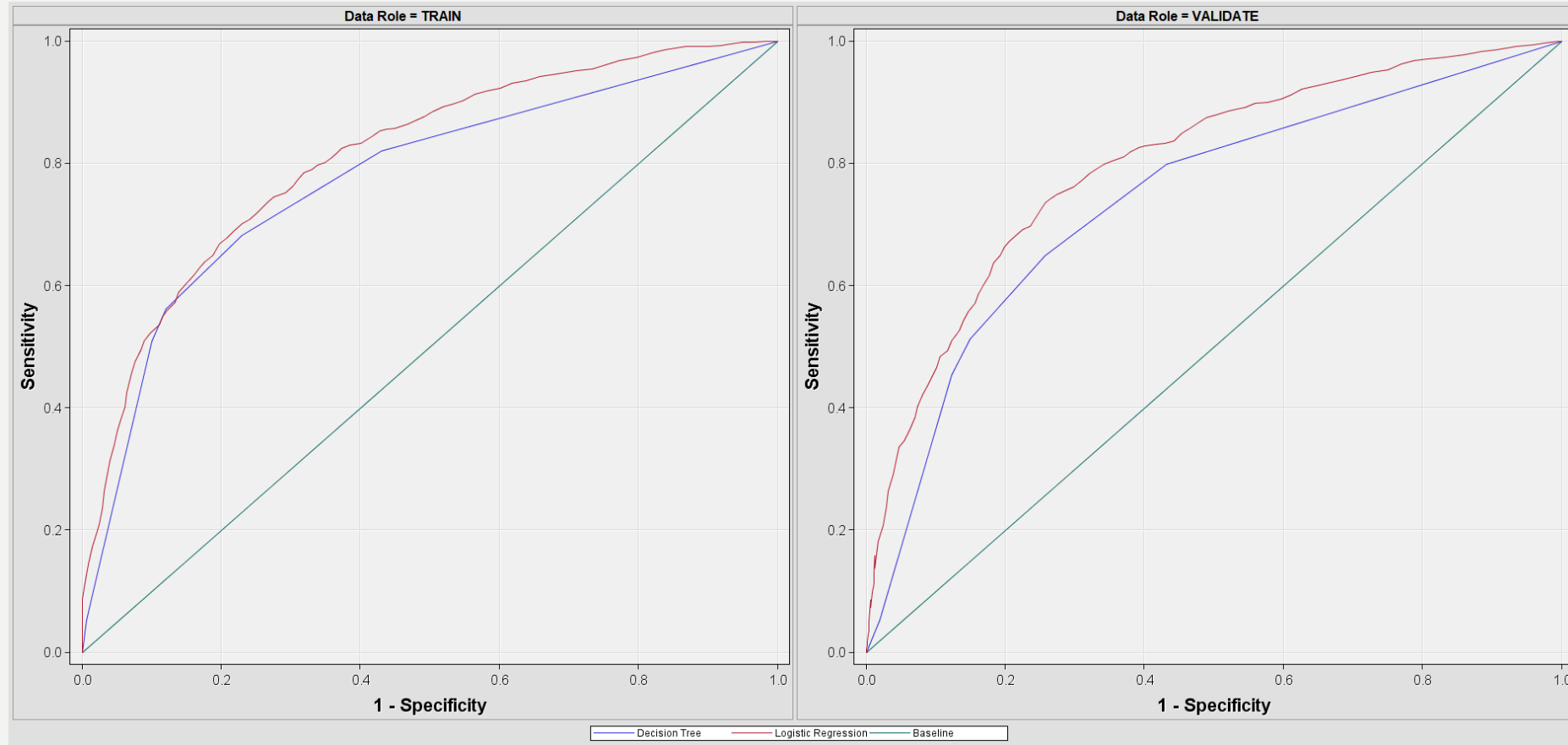  - Movie with average sentiment score **higher** than or equal to 0.77 is likely to be a **good** movie
  - Movie with average sentiment score **lower** than or equal to 0.77 is likely to be a **bad** movie

12

- A logistic regression is developed as the model is simple to implement and efficient to train.

| Parameter | | DF | Estimate | Standard Error | Wald Chi-Square | Pr > ChiSq | Standardized Estimate | Exp(Est) |
|---|---|---|---|---|---|---|---|---|
| Intercept | | 1 | -4.5011 | 0.5828 | 59.65 | <.0001 | | 0.011 |
| IMP_runtime | | 1 | 0.0154 | 0.00328 | 22.13 | <.0001 | 0.1593 | 1.016 |
| IMP_ss_mean | | 1 | 5.1110 | 0.3489 | 214.62 | <.0001 | 0.4918 | 165.837 |
| R | 0 | 1 | 0.2007 | 0.0562 | 12.74 | 0.0004 | | 1.222 |
| genre_AnimationManga | 0 | 1 | -0.3999 | 0.1367 | 8.56 | 0.0034 | | 0.670 |
| genre_DramaTele | 0 | 1 | -0.2091 | 0.0565 | 13.68 | 0.0002 | | 0.811 |
| genre_FitnessSports | 0 | 1 | -0.9172 | 0.3881 | 5.59 | 0.0181 | | 0.400 |
| genre_HistDocument | 0 | 1 | -0.5957 | 0.0797 | 55.93 | <.0001 | | 0.551 |
| genre_Horror | 0 | 1 | 0.2259 | 0.0943 | 5.74 | 0.0166 | | 1.253 |
| genre_Sci_fi | 0 | 1 | 0.1687 | 0.0936 | 3.25 | 0.0716 | | 1.184 |
| genre_ThrillMysSusp | 0 | 1 | 0.1377 | 0.0673 | 4.18 | 0.0409 | | 1.148 |
| tomatometer_count | | 1 | 0.00581 | 0.000726 | 64.03 | <.0001 | 0.2692 | 1.006 |

- A positive regression coefficient indicates that the mean of the dependent variable (good or bad movie) increases with the value of independent variable.
- Higher regression coefficient indicates higher influence.
- Mean sentiment score has highest positive influence
  - Higher mean sentiment score is likely to be good movie.
- For genre, horror movie has highest positive influence, followed by science fiction movie.
  - Horror movie is often combined with sci-fi movie (e.g. Earth is threatened by Aliens)
- Sports movie has highest negative influence.
  - The 'plot' of the game is the same every time, hence is boring.
  - Waste of time watching people win their game.

- Compare decision tree and logistic regression.
- Accuracy (decision tree): 68.11%
- Accuracy (logistic regression): 73.36%
- Logistic regression has better performance compared to decision tree, hence is selected for prescriptive analysis.

# PRESCRIPTIVE ANALYSIS

- SEMM**A**

  o Assess: Determine whether the result is useful

- Recommend movies with **high sentiment score** to users

  o Movie with a positive review are likely to be a good movie

- Recommend movies with **genre of horror and science fiction** to users

  o Movie with both of these genres are likely to be a good movie

# CONCLUSION

- Box office and user reviews can be used to predict whether a movie is good or bad.

- Words like excellent, perfect and great normally appear in a good movie review.

- Words like worst, awful and boring normally appear in a bad movie review.

- Movies with positive reviews and genres of horror and science fiction are likely to be good.