# Predicting Success of Cable Sitcoms using NLP on Script Data

Emily Pfeifer
Capstone Project
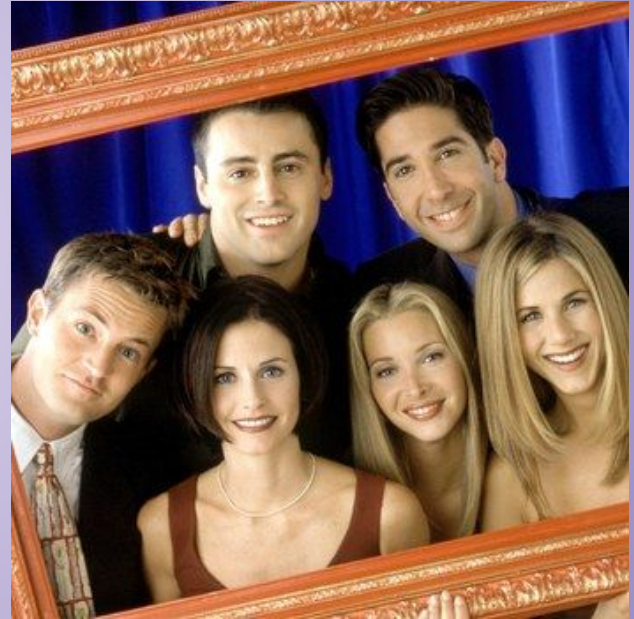
# Obtaining the data

- 50 sitcom pilots – 25 "failed" and 25 "successful"
- 14 columns of descriptive data – year aired, network, etc.
- Collect script for each pilot using web-scraping API "Selenium"
- Generate dataframe

# CRITERIA FOR CHOSEN PROGRAMS



- All "Situational Comedies"
- Aired on American Cable Network
- Made in last 30 years
- Successful Sitcoms:
  - Rotten Tomatoes Score > 60%
  - IMDB Score > 7
  - Minimum of 3 Seasons
- Failed Sitcoms:
  - Rotten Tomatoes Score < 45%
  - IMDB Score < 6.5
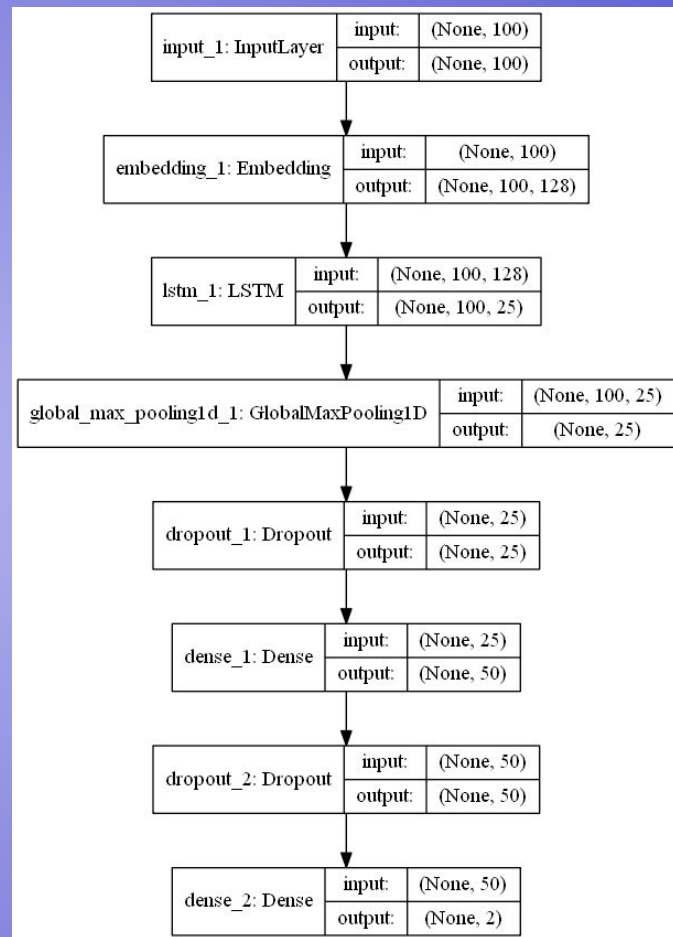  - Canceled after 1 season

# A Closer look

- Average RT Score – Success: **88.93**
- Average RT Score – Failed: **15.33**

- Average IMDB Score – Success: **8.18**
- Average IMDB Score – Failed: **5.17**

- Average # of Episodes – Success: **168.3**
- Average # of Episodes – Failed: **9.0**

# Classification

- TF-IDF
  - Highest Test Score – 62%
  - Highest Train Score – 94%
- TF-IDF with LSA
- TF-IDF with Text Descriptive Features
- Text Descriptive Features Alone

# Deep LEarning

- LSTM Recurrent Neural Network
- Word embeddings
- K-folds cross evaluation
- Average accuracy of 80%

# LIMItations

- ❏ Small sample size
- ❏ Uneven Data
- ❏ Availability of scripts
- ❏ Time

# Future Work

➢ Data

- ○ Increase sample size
- ○ Expand criteria (include HBO/Netflix, older shows)
- ○ Add more descriptive features

➢ Technical

- ○ Data exploration with unused variables (actors, viewer data, etc.)
- ○ Additional feature engineering
- ○ Hyperparameter tuning classifiers
- ○ Hyperparameter tuning RNN