# COMP 4434 Final Presenstation

JIANG Ruixiang

19079662D

# Contents

# Background

- PolyTube
- Large dataset

# Project Objective

- Task 1 build a regression model
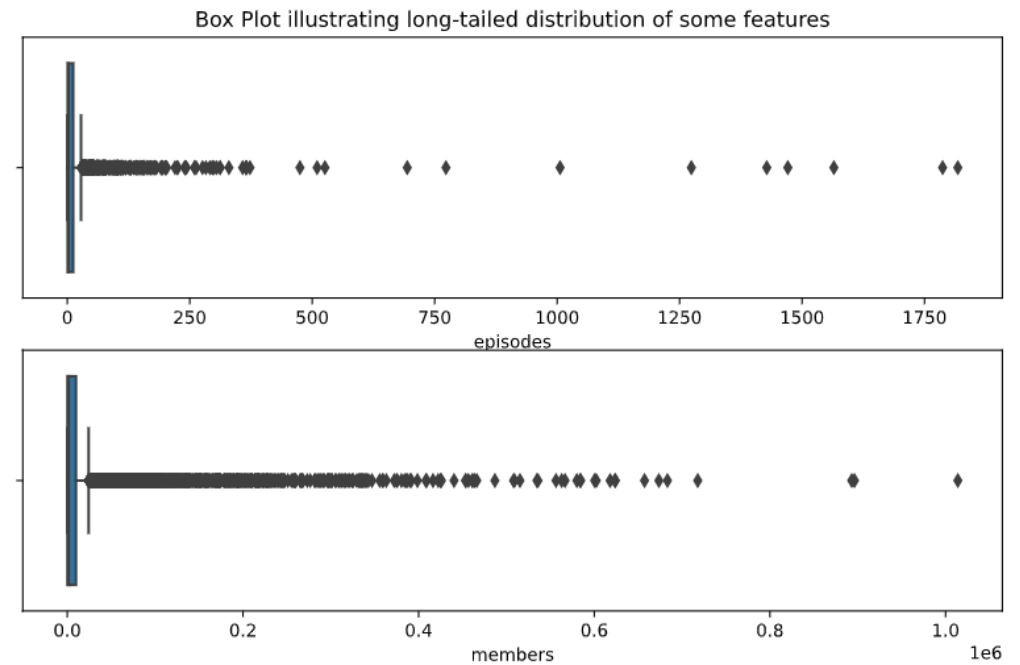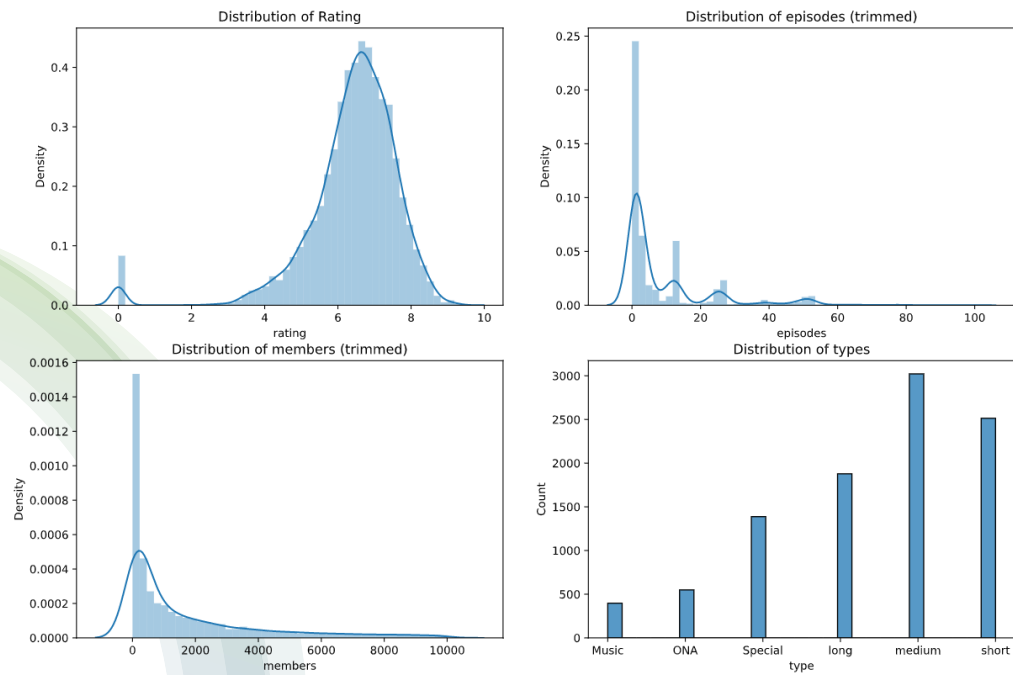- Task 2 build a recommender system

# Contents

# Exploratory data analysis
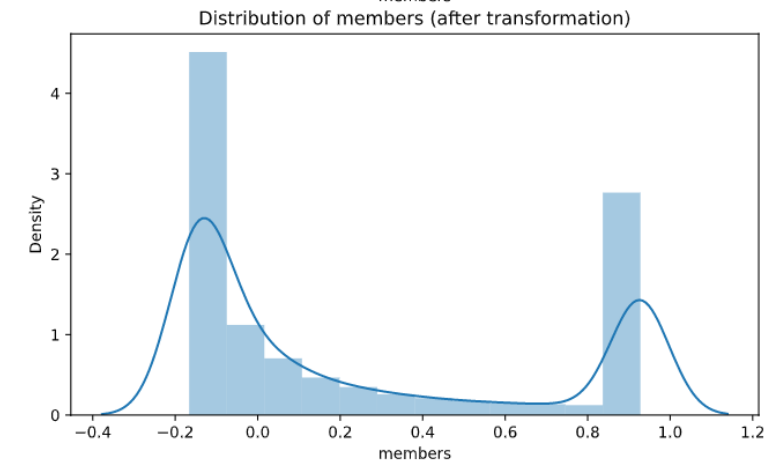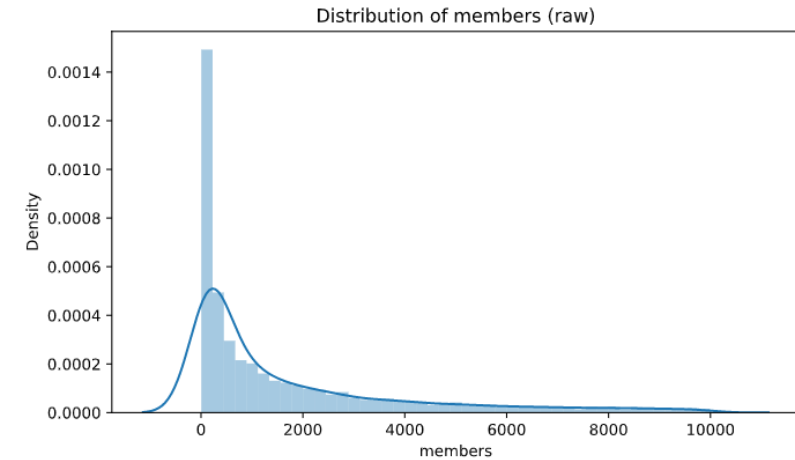
Long-tail distribution
Outliers

# Preprocessing

- Null value imputation –mean values
- Outliers -clipping
- Feature scaling –transformation
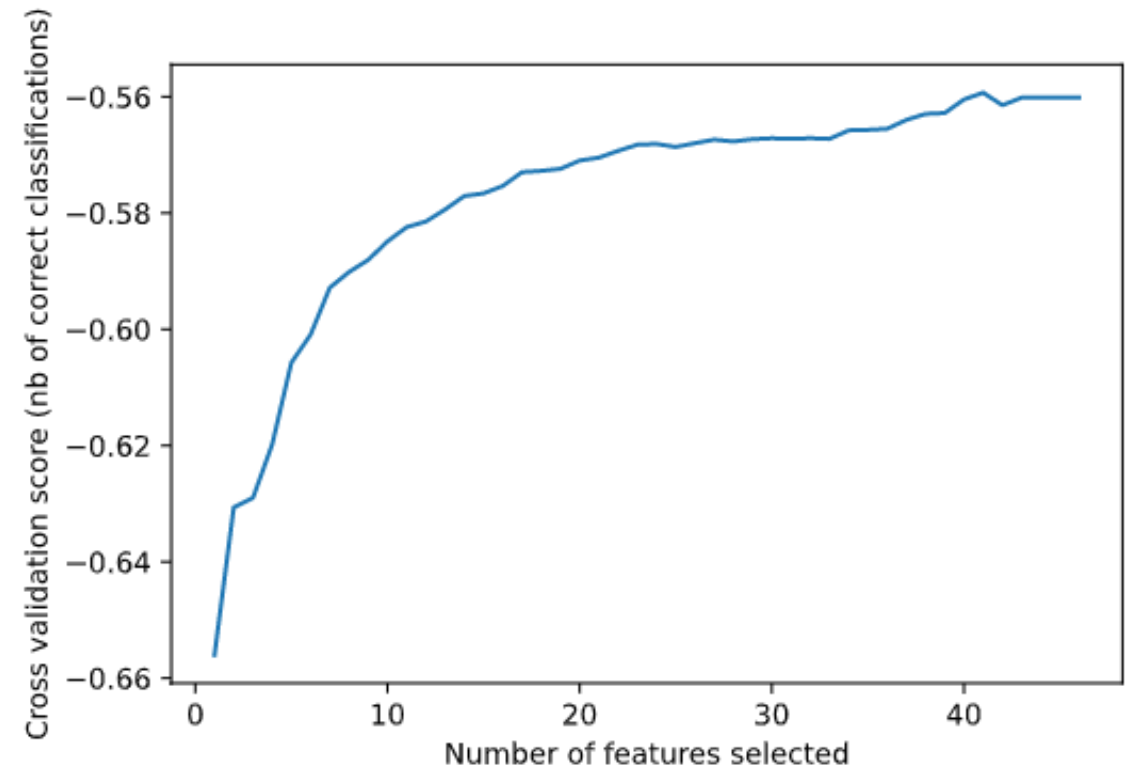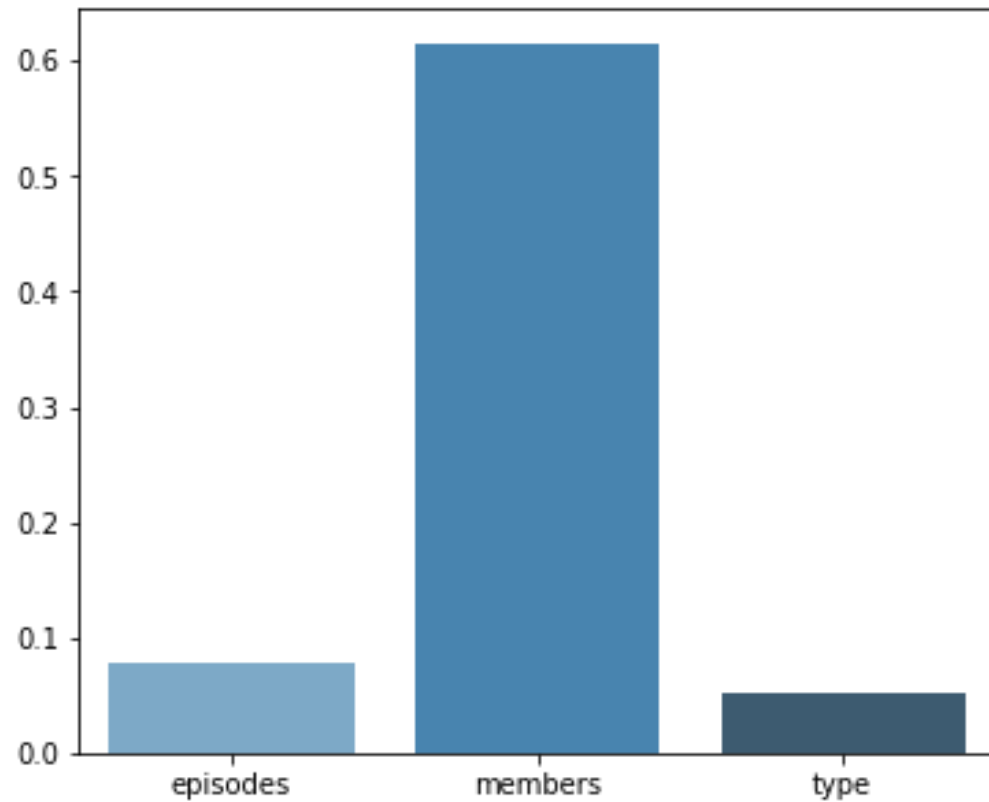- Feature encoding -vectorize

# Preprocessing Result

# Feature selection

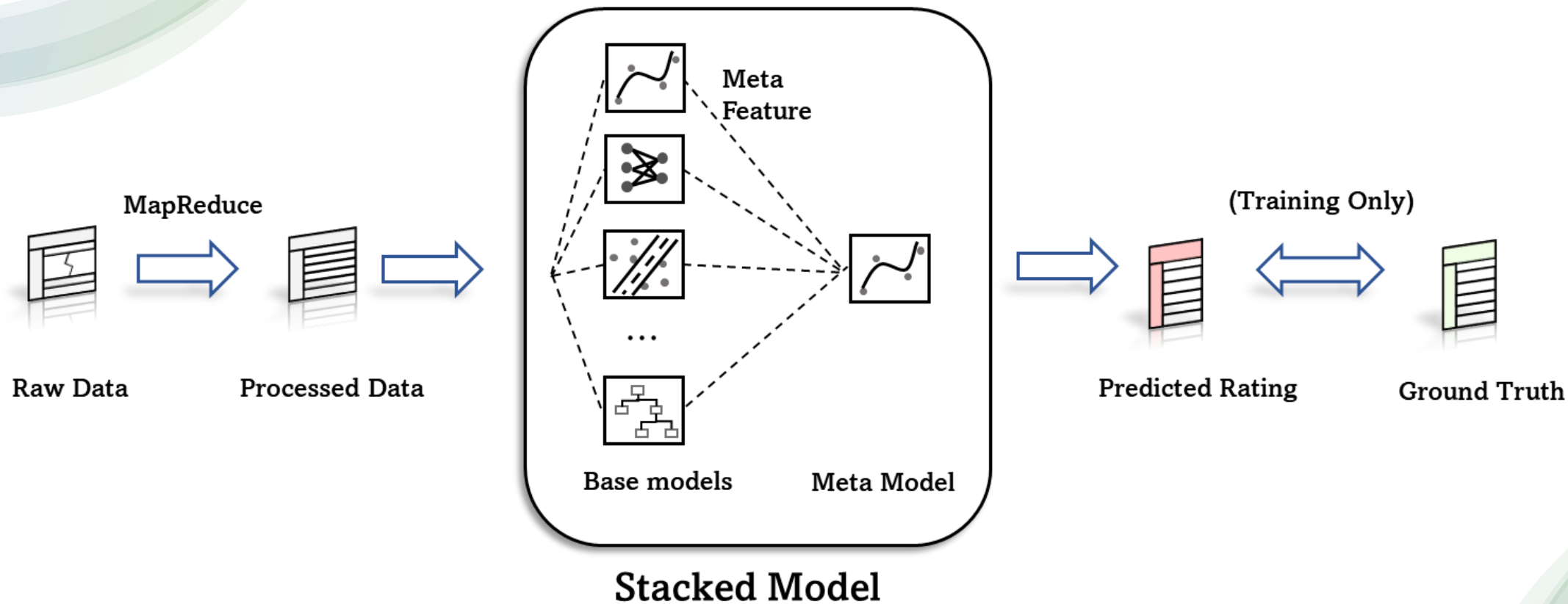- Recursive Feature Elimination
- Feature importance

# Contents

# Model Design (Task 1 )

- Baseline -linear regression, neural networks
- Ours –Stacked model (ensembling)

# Model Architecture (Task 1 )



Raw Data → MapReduce → Processed Data → Stacked Model (Meta Feature, Base models, Meta Model) → Predicted Rating ← (Training Only) → Ground Truth
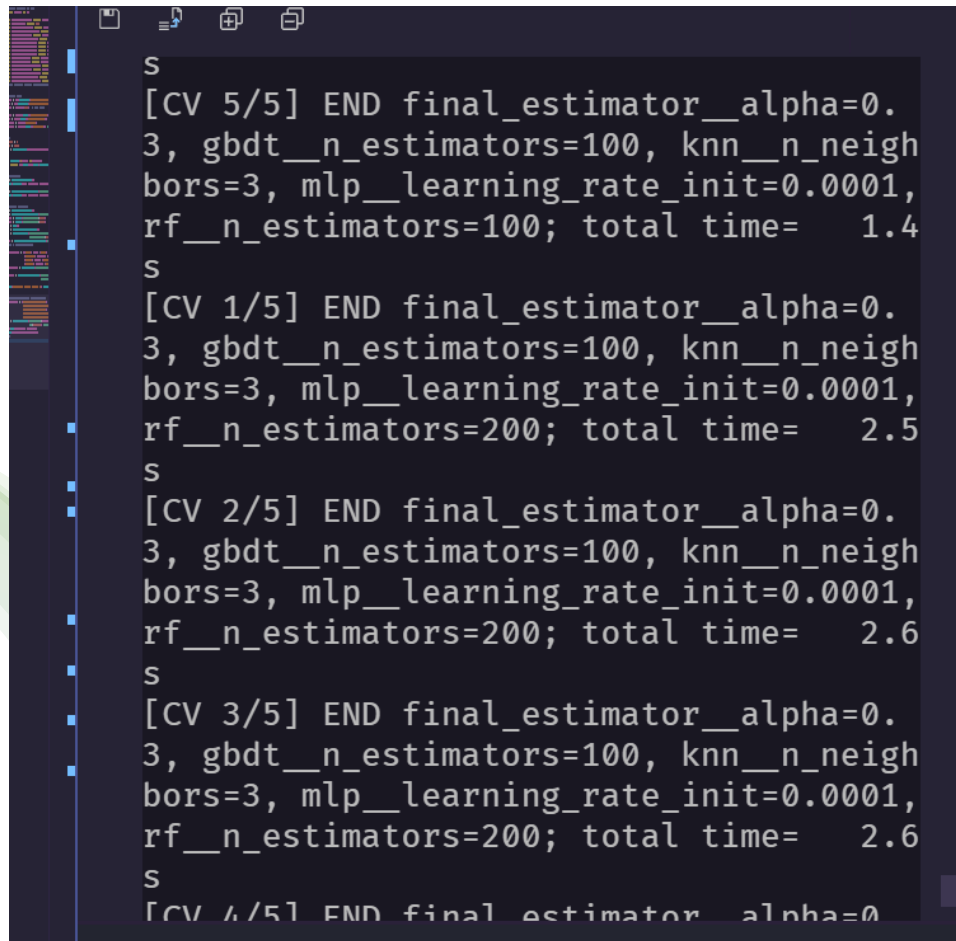
# Model Design (Task 2)

- Content-based recommender system

- Many baselines

- Final solution: linear regression

# Regularizations

- L1, L2 norms for LR, MLPs
- Dropout
- Batchnormalization
- Earlystopping
- reduceLRonPlatau
- Adaptive learning rate scheduler
- Gradient Clipping

# Model Training



- Grid Search to optimization hyper-param
- 9-12 hours on a server to optimize
- Singel training round ~10 min

# Contents

Introduction

Preprocessing

Model

**Result**

# Model Merics

- Baselines : plain linear regression, linear regression with polynomial feature expansion, Lasso regression ($\lambda = 0.01$), multilayer perceptron

- Metrics: 5-fold RMSE

# Quantitative Results

| Models | LR | PolyLR(n = 2) | Lasso | MLP | **Stacked** |
|---|---|---|---|---|---|
| 5-fold CV RMSE | 0.748 | 0.762 | 0.764 | 0.688 | **0.649** |

Table 1: Quantitative result of cross-validation task 1 model performance

| Models | LR | PolyLR(n = 2) | Lasso | **MLP** |
|---|---|---|---|---|
| 5-fold CV RMSE | 1.273 | 1.611 | 1.266 | **1.242** |

Table 2: Quantitative result of cross-validation model performance for task 2

# Discussion on models

- Non-linear model outperform others, but training is hard
- Linear model is good for task 2
- Conclusion

# Thanks