# VOD RECOMMENDATION SYSTEM

설재완

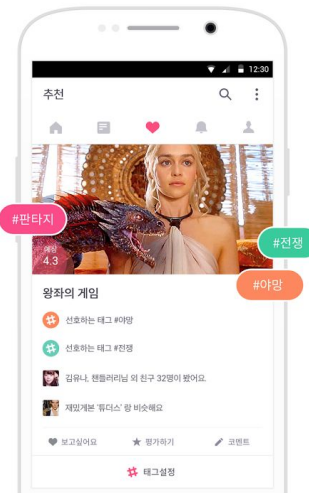# Goal

# How about others







개인화 추천, 드라마까지
내게 맞는 영화뿐 아니라
이제 드라마도 쑥쑥

# No rating but history



| | USER_ID | SERIES_ID | ASSET_ID | DURATION | EVENT_TIME |
|---|---|---|---|---|---|
| 23451632 | 303428 | 1122 | 22022 | 3128 | 2017-09-30 23:58:57 |
| 23451633 | 5059 | 90 | 5233 | 52 | 2017-09-30 23:58:57 |
| 23451634 | 1444 | 22 | 588 | 722 | 2017-09-30 23:58:59 |
| 23451635 | 459223 | 592 | 24531 | 613 | 2017-09-30 23:58:59 |
| 23451636 | 459214 | 150 | 3677 | 1761 | 2017-09-30 23:59:00 |

# Conversion From History To Rating



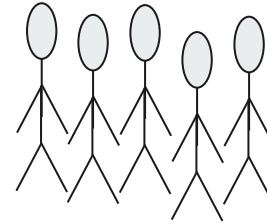**Rating Based Recommendation**

개인화 추천, 드라마까지
내게 맞는 영화뿐 아니라
이제 드라마도 쑥쑥

# How to convert

29:37 / 59:57

45:11 / 59:57

mean watching time = 30

watching time = 60

# Recommendation Overview

Watch History ⟹ Rating ⟹ **Rating Based Recommendation**

ex)
collaborative filtering
matrix factorization

history based recommendation

개인화 추천, 드라마까지
내게 맞는 영화뿐 아니라
이제 드라마도 쓱쓱

# Implement - Rating Conversion

| 사람 | 영상 | 시청시간 |
|---|---|---|
| 갑 | 벡터맨 | 10 |
| 을 | 벡터맨 | 40 |
| 갑 | 벡터맨 | 10 |
| 갑 | 스칼라맨 | 30 |
| 병 | 스칼라맨 | 10 |

absolute time

| 사람 | 영상 | 시청시간 |
|---|---|---|
| 갑 | 벡터맨 | 20 |
| 갑 | 스칼라맨 | 30 |
| 을 | 벡터맨 | 40 |
| 병 | 스칼라맨 | 10 |

relative time

| | 벡터맨 | 스칼라맨 |
|---|---|---|
| 갑 | 20 / 40 = 0.5 | 30 / 30 = 1 |
| 을 | 40 / 40 = 1 | |
| 병 | | 10 / 30 = 0.33 |

user - asset matrix

# Rating based recommendation ex)1



Items

|   | 1 | 2 | ... | i | ... | m |
|---|---|---|-----|---|-----|---|
| 1 | 5 | 3 |     | 1 | 2   |   |
| 2 |   | 2 |     |   |     | 4 |
| : |   |   | 5   |   |     |   |
| u | 3 | 4 |     | 2 | 1   |   |
| : |   |   |     |   | 4   |   |
| n |   |   | 3   | 2 |     |   |

Users

**baseline model**

- **consider each user and asset**
- **rui = mean + bu + bi**
- **also can be applied to another model**

# Rating based recommendation ex)2

**collaborative filtering**



- Similarity
  - between item
  - between user

# Rating based recommendation ex)3

**matrix factorization**



Item

|   | W | X | Y | Z |
|---|---|---|---|---|
| A |   | 4.5 | 2.0 |   |
| B | 4.0 |   | 3.5 |   |
| C |   | 5.0 |   | 2.0 |
| D |   | 3.5 | 4.0 | 1.0 |

Rating Matrix

=

latent factor

|   |   |   |
|---|---|---|
| A | 1.2 | 0.8 |
| B | 1.4 | 0.9 |
| C | 1.5 | 1.0 |
| D | 1.2 | 0.8 |

User Matrix

X

|   | W | X | Y | Z |
|---|---|---|---|---|
|   | 1.5 | 1.2 | 1.0 | 0.8 |
|   | 1.7 | 0.6 | 1.1 | 0.4 |

Item Matrix

parameters

# Experiment

- 21,000,000 records
- 555,000 users
- 55,000 assets
- Compare and select model between
    - baseline
    - collaborative filtering only(CFonly)
    - collaborative filtering with baseline(CFwithBase)
    - matrix factorization(MF)
- Tain set: watch history of July, August
- Test set: watch history of September

# Results



Final accuracy using CFwithBase: 47.24%

# Limitation and future work

- If there were validation set, results will be better
- If there were good optimization technique(optimizer in ML/DL framework, fine tunning, etc), results will be better(It might not be trained well)
- If we can consider watching pattern(interval between each asset, consecutive assets, etc), results will be better
- If we use 'genre' and 'series' data well, results will be better
- If we pre-processed raw data well(there were so many histories of '뽀로로'), results will be better

# Conclusion

- sutdy and implement basic ML concepts
    - use only numpy
    - gradient descnt, early stopping
- study and implement recommendation algorithm
    - collaborative filtering
    - matrix factorization
- read paper and implement
- problem solving life cycle
    - recognition -> background study -> design and implement -> evaluation -> result analysis

# Thanks!

# BACKUP SLIDES

# Rating Conversion

- Let arr[i, j] be rating of user i at asset j (after converting history to rating)
- for jdx in range(len(asset)):
-     ratings = arr[:, j] # all ratings for asset j
-     median = sorted(ratings)[len(ratings) / 2] # median
-     arr[:, jdx] /= median # divide by median

# Rating Conversion

- **Relative time**
    - **t = (total watching time) / (running time)**
    - **can over-estimate for cases that short running time**
- **Absolute time**
    - **t = (total wathcing time)**
    - **can over-estimate for cases that long running time**
- **So, let's use (total watching time) / (mean or median watching time)**
    - **can consider other people(relative time) and pure wathcing time(absolute time) at once**
    - **median is better when extreme value exists**

19