

# Building an Implicit Recommendation Engine

Sophie Watson   @sophwats   sophie@redhat.com

#SAISDS12

# Building an Implicit Recommendation Engine

Sophie Watson   @sophwats   sophie@redhat.com

#SAISDS12

**NETFLIX**  
**hulu**

} movies

goodreads

books

**TESCO**

food

 **Expedia**<sup>®</sup> hotels

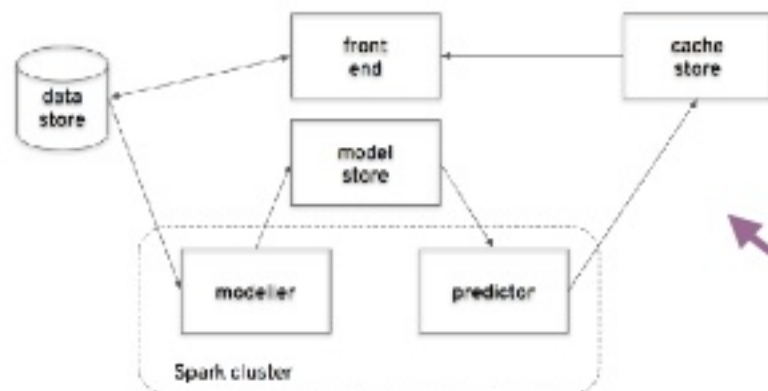
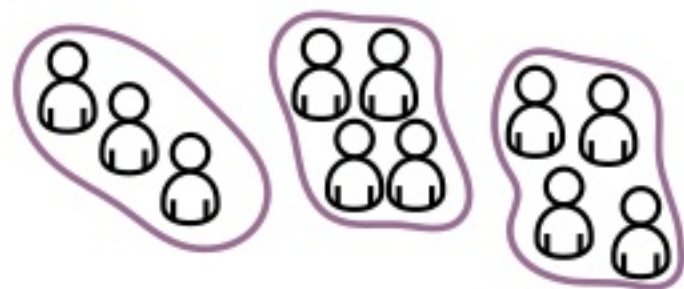
} music



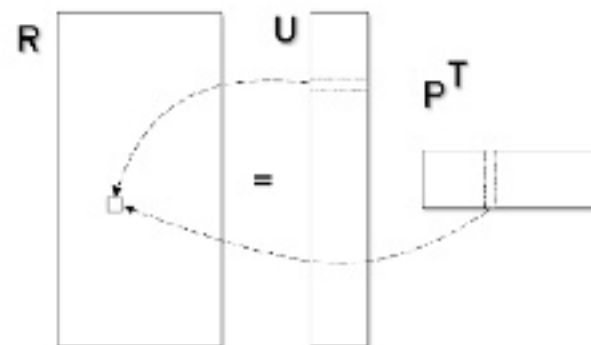
people to follow

# Outline

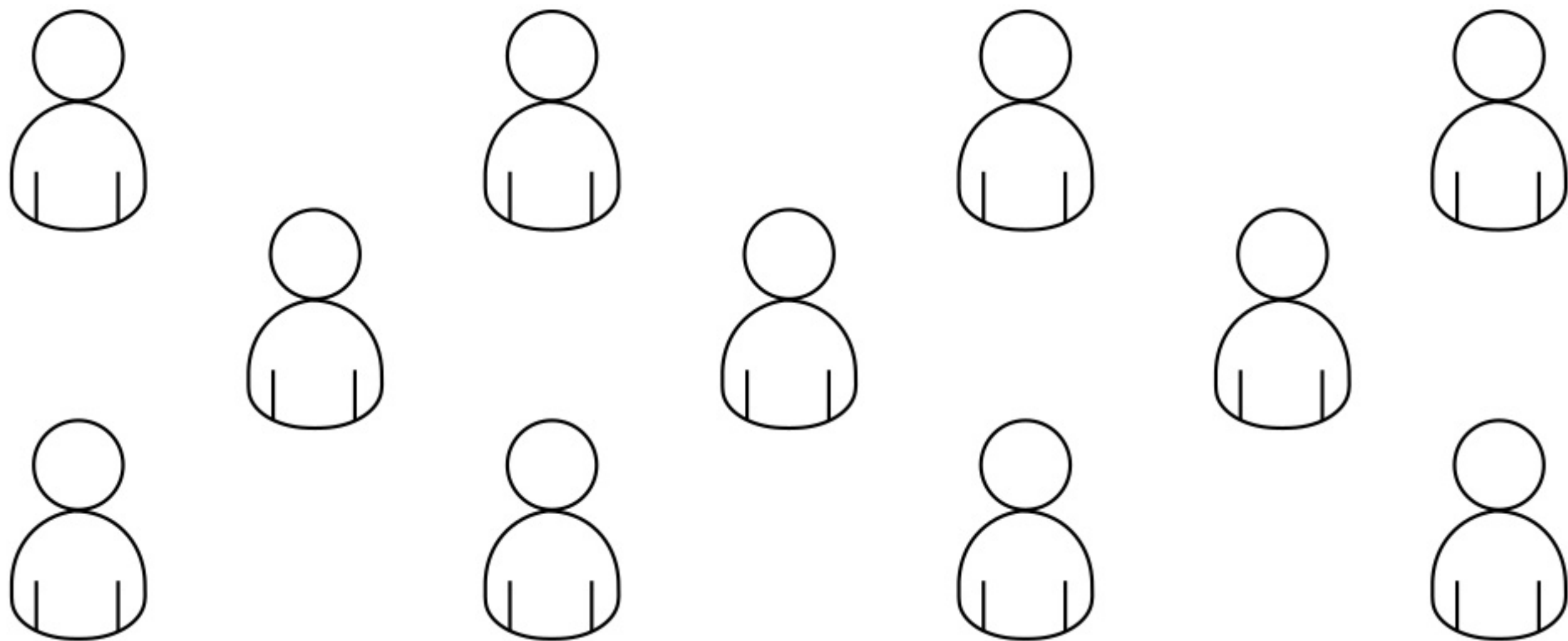
## Collaborative Filtering



## Alternating Least Squares

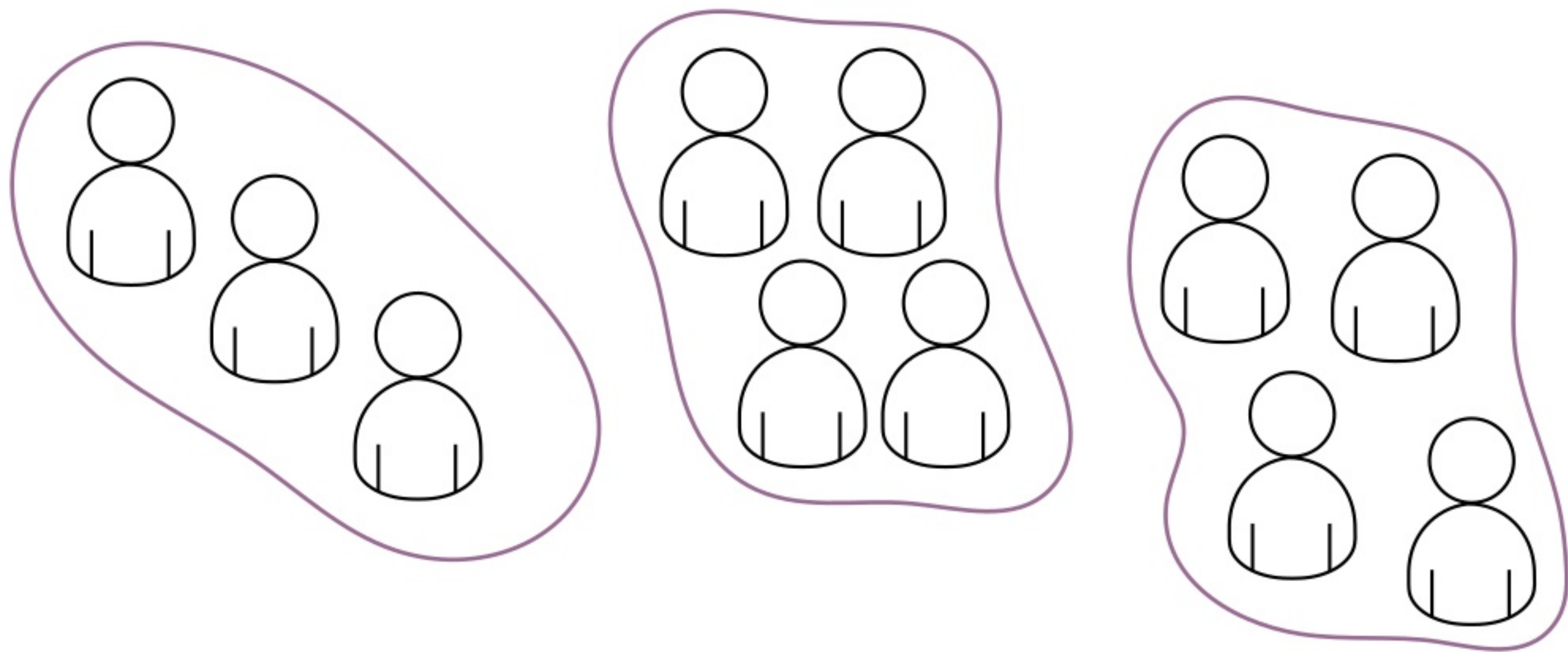


# Collaborative Filtering

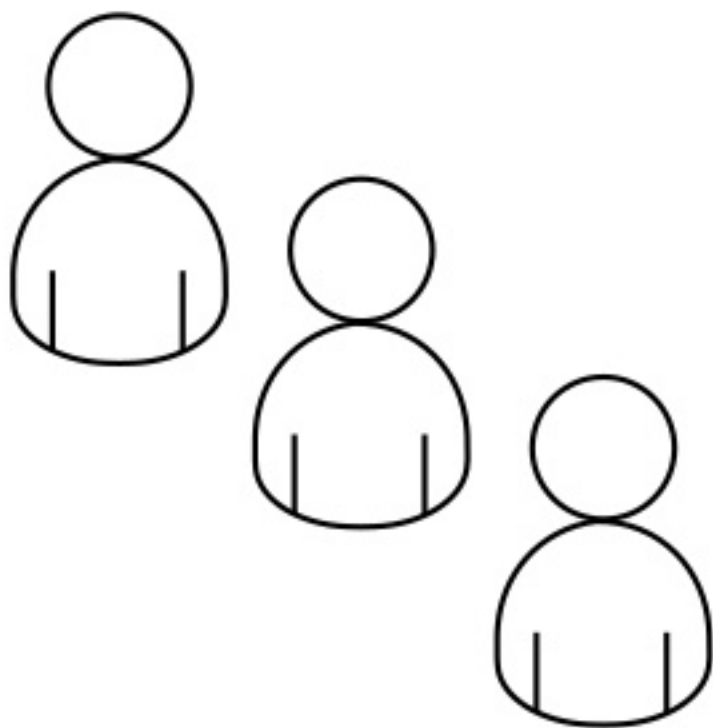




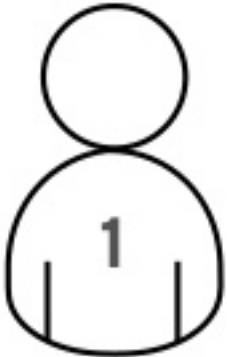
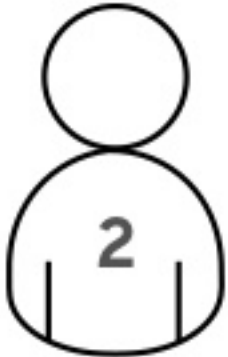
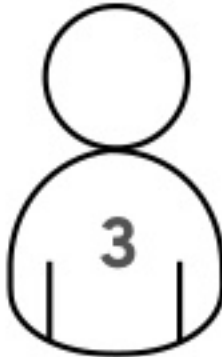
# Collaborative Filtering



# Collaborative Filtering

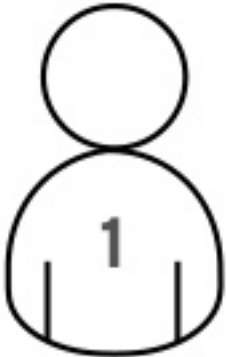
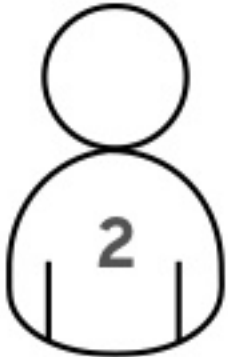
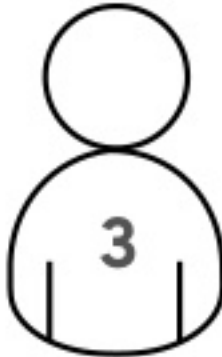


# Collaborative Filtering

|           |  |  |  |
|-----------|--|---|---|
| Product A | 5★   | 5★  | 5★  |
| Product B | 3★   | 3★  | 3★  |
| Product C | 5★   | 3★  | ?★  |

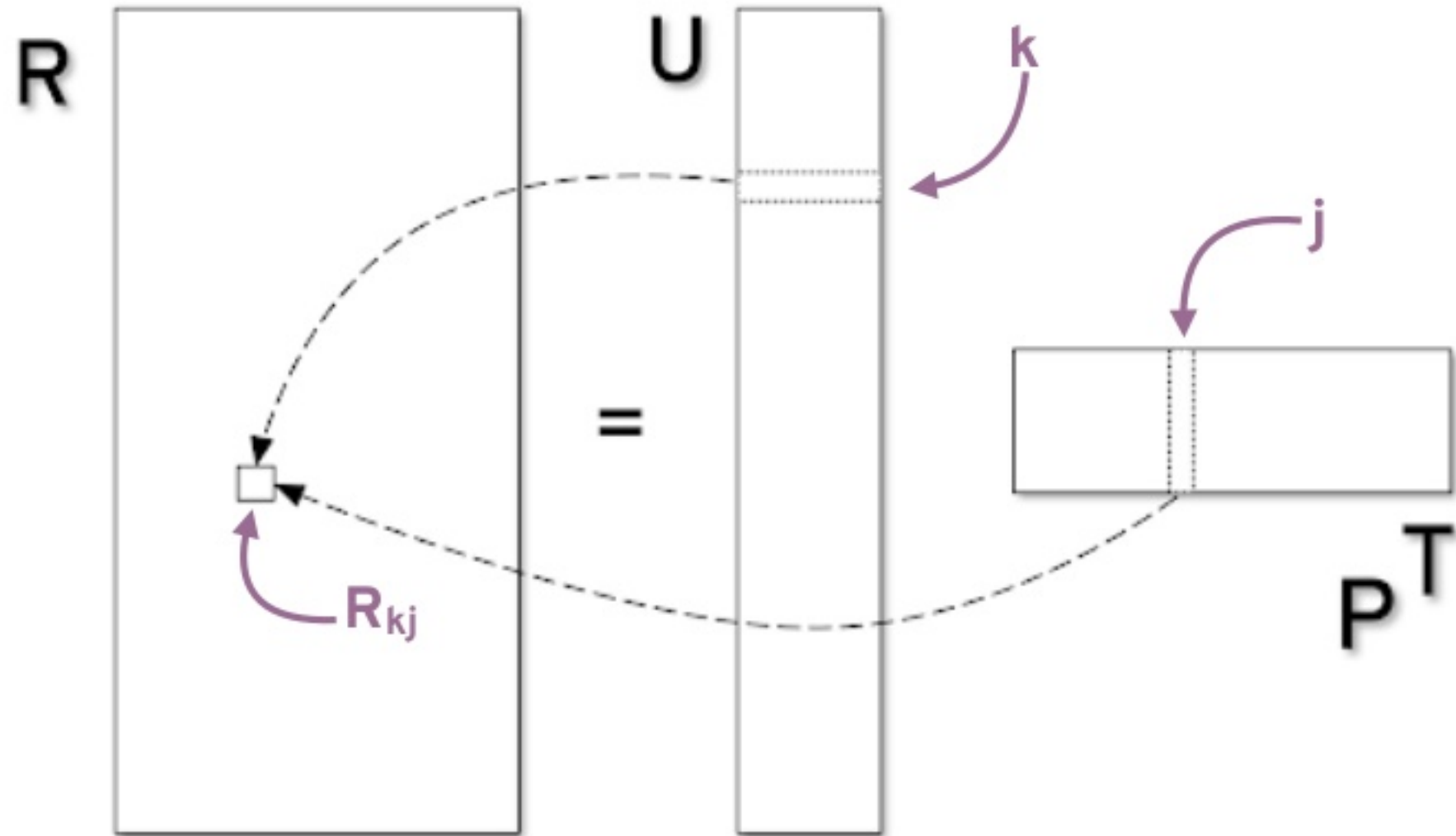


# Collaborative Filtering

|           |  |  |  |
|-----------|--|---|---|
| Product A | 5★   | 5★  | 5★  |
| Product B | 3★   | 3★  | 3★  |
| Product C | 5★   | 3★  | 4★  |

# Alternating Least Squares

|     | user 1 | user 2 | user 3 | ... | user N |           |
|-----|--------|--------|--------|-----|--------|-----------|
| R = | 1      | 4.5    | ?      | ... | 3      | product 1 |
|     | ?      | 3      | 3      | ... | 4      | product 2 |
|     | 5      | 3      | ?      | ... | ?      | product 3 |
|     | ⋮      | ⋮      | ⋮      | ⋱   | ⋮      | ⋮         |
|     | 2      | 4      | 1      | ... | ?      | product M |



# Alternating Least Squares

|     | user 1 | user 2 | user 3 | ... | user N |           |
|-----|--------|--------|--------|-----|--------|-----------|
| R = | 1      | 4.5    | 3.8    | ... | 3      | product 1 |
|     | 3.2    | 3      | 3      | ... | 4      | product 2 |
|     | 5      | 3      | 3.4    | ... | 3.1    | product 3 |
|     | ⋮      | ⋮      | ⋮      | ⋱   | ⋮      | ⋮         |
|     | 2      | 4      | 1      | ... | 2.7    | product M |

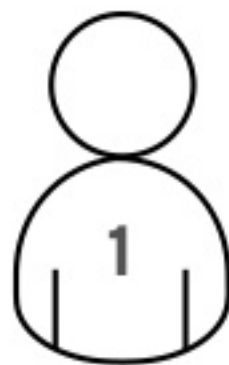
# Implicit Data



**Song A**

**1 play**

# Implicit Data



Song A

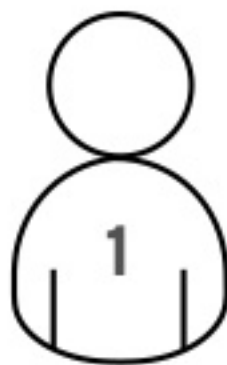
1 play

**Song B**

**0 plays**



# Implicit Data



Song A

1 play

Song B

0 plays

**Song C**

**100 plays**

# Implicit Alternating Least Squares

## Collaborative Filtering for Implicit Feedback Datasets

Yifan Hu

AT&T Labs – Research  
Florham Park, NJ 07932

Yehuda Koren\*

Yahoo! Research  
Haifa 31905, Israel

Chris Volinsky

AT&T Labs – Research  
Florham Park, NJ 07932



# Implicit Alternating Least Squares

The aim:

$$p_{ui} \in (0, 1) \quad \text{preference}$$

# Implicit Alternating Least Squares

The recorded data:

$p_{ui} \in (0,1)$  preference

$r_{ui} \in \mathbb{R}$  recording

# Implicit Alternating Least Squares

The recorded data:

$$p_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0 \\ 0 & \text{if } r_{ui} = 0 \end{cases}$$

$p_{ui} \in (0,1)$  preference

$r_{ui} \in \mathbb{R}$  recording

# Implicit Alternating Least Squares

Confidence:

$$C_{ui} = 1 + \alpha r_{ui}$$

$p_{ui} \in (0, 1)$  preference

$r_{ui} \in \mathbb{R}$  recording

$$p_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0 \\ 0 & \text{if } r_{ui} = 0 \end{cases}$$



# Implicit Alternating Least Squares

Confidence:

$$C_{ui} = 1 + \alpha r_{ui}$$

tuning parameter



$p_{ui} \in (0, 1)$  preference

$r_{ui} \in \mathbb{R}$  recording

$$p_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0 \\ 0 & \text{if } r_{ui} = 0 \end{cases}$$

# Implicit Alternating Least Squares

Minimisation:

$p_{ui} \in (0,1)$  preference

$r_{ui} \in \mathbb{R}$  recording

$$C_{ui} (p_{ui} - u_u X_i^T)$$

$$p_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0 \\ 0 & \text{if } r_{ui} = 0 \end{cases}$$

user vector

Item vector

# What does Spark offer?

```
1 from pyspark.mllib.recommendation import ALS
```

# Data

## Lastfm dataset

17 million recordings / 360k users / 200k artists

```
[ '00000c289a1829a808ac09c00daf10bc3c4e223b',  
  '8bfac288-ccc5-448d-9573-c33ea2aa5c30',  
  'red',  
  'hot',  
  'chili',  
  'peppers',  
  '691' ],
```

The diagram illustrates the structure of a recording in the Last.fm dataset. It shows a JSON array with seven elements. Annotations with arrows point to specific parts of the array:

- user id**: Points to the first element, a long hexadecimal string.
- artist id**: Points to the second element, a hexadecimal string with hyphens.
- artist name**: A bracket groups the next three elements ('red', 'hot', 'chili') and points to them.
- number of plays**: Points to the last element, '691'.

# Building the model

(user id, product id, recording)



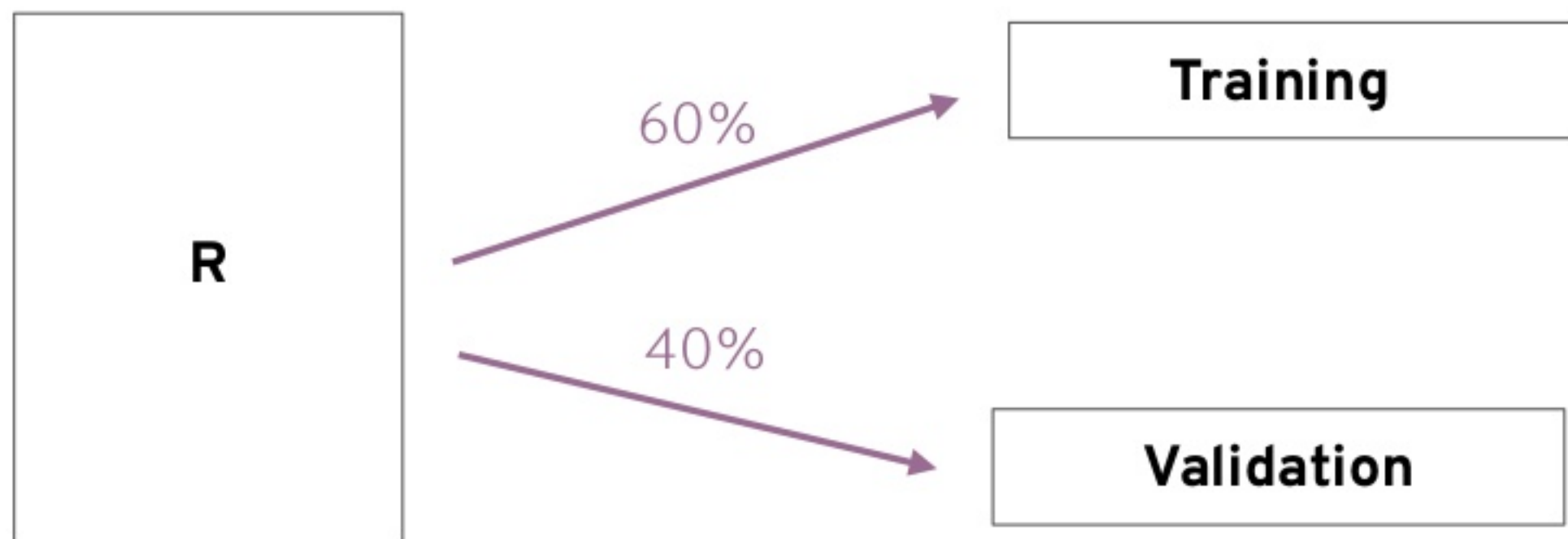
```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```



tuning parameters

# Tuning Parameters

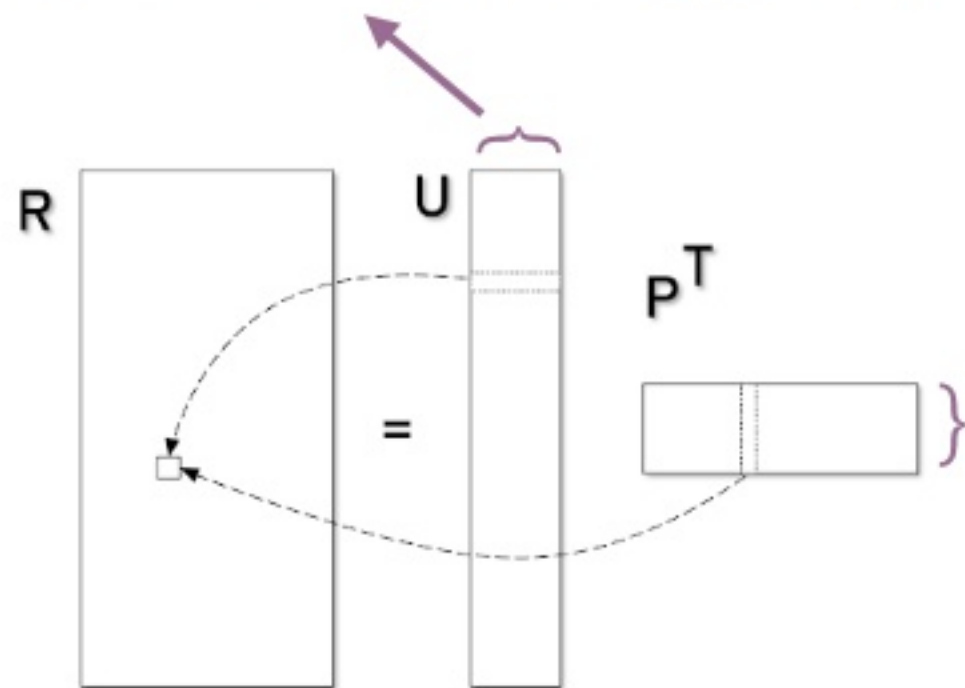
```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```





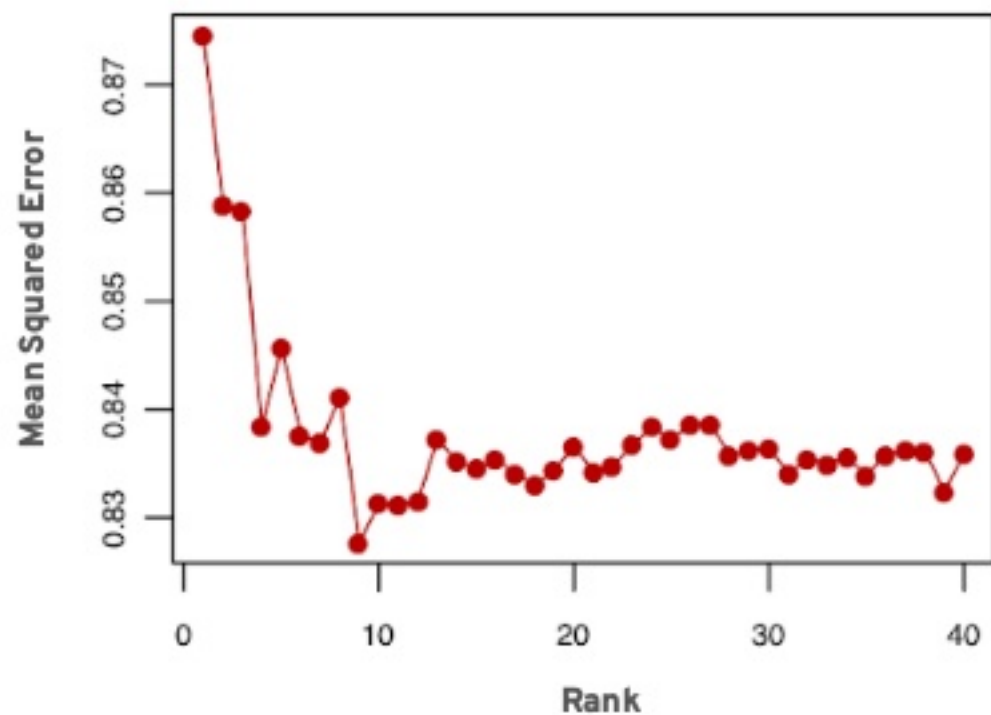
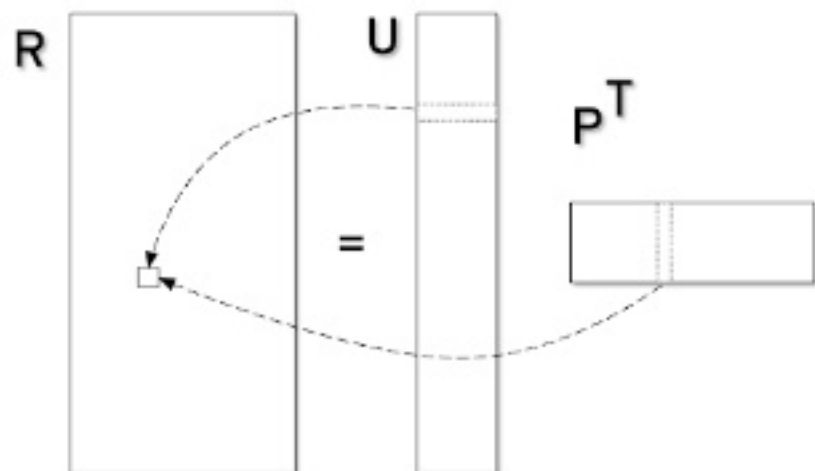
# Rank

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```



# Rank

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

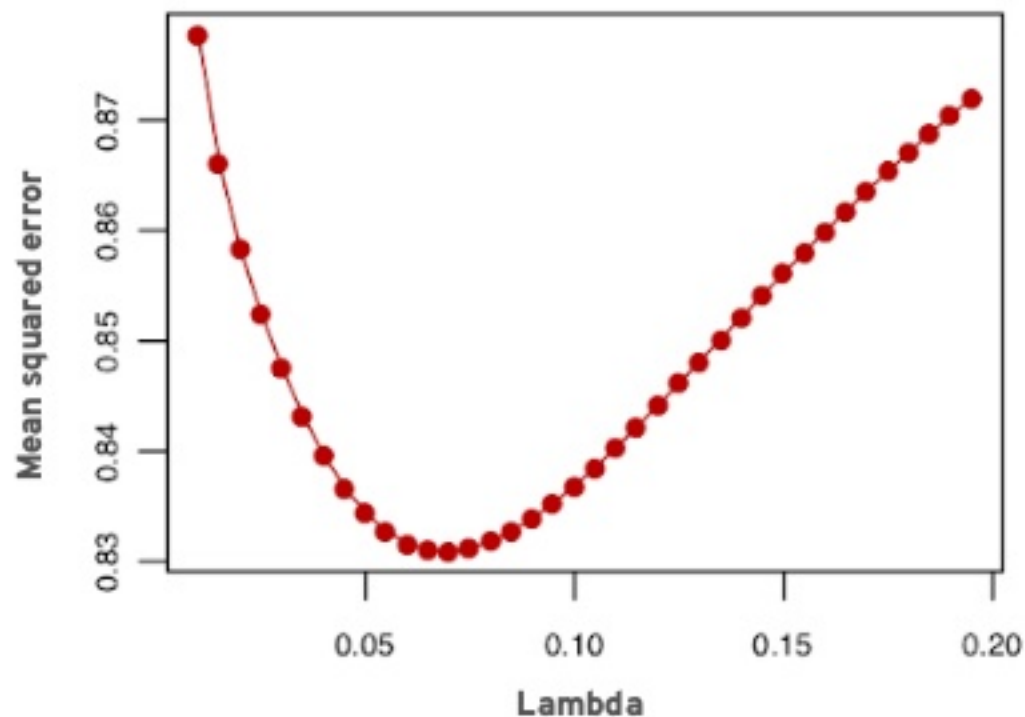


# Lambda

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda=0.01, alpha = 1.0, iterations=5)
```

Minimisation:

$$C_{ui} (p_{ui} - U_u X_i^T) + \lambda (\|U\|^2 + \|X\|^2)$$



# Alpha

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

$$C_{ui} = 1 + \alpha r_{ui}$$

$r_{ui} \in \mathbb{R}$  recording

# Alpha

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

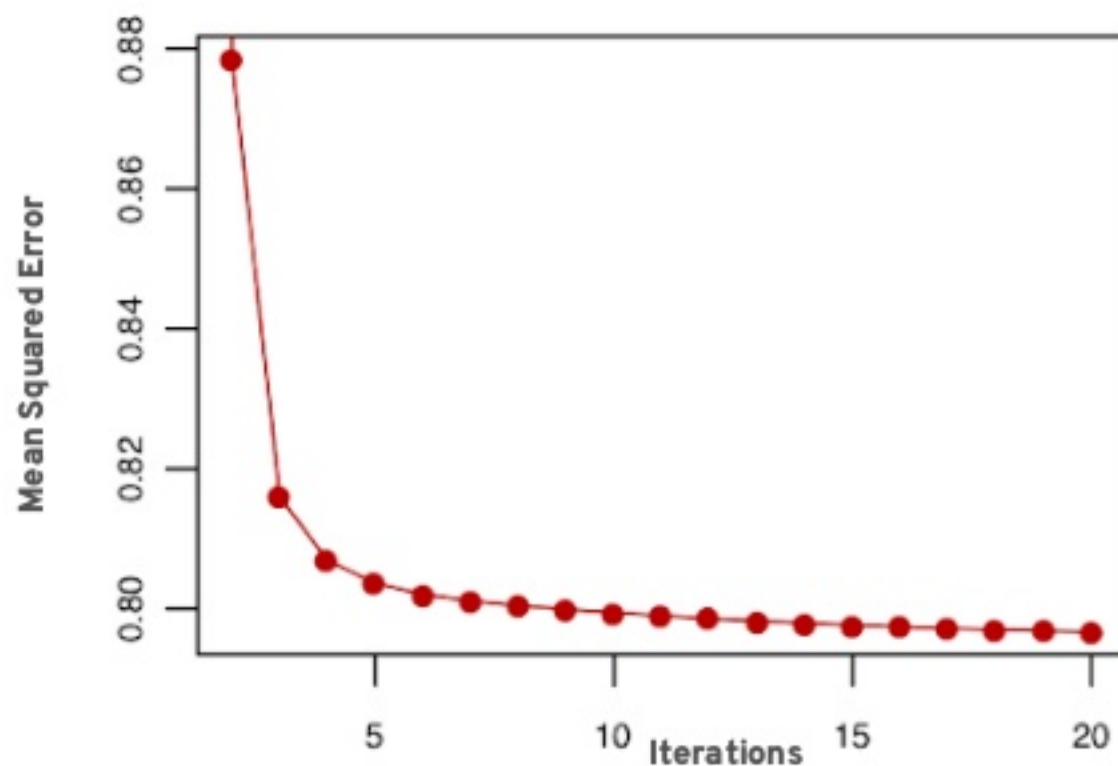
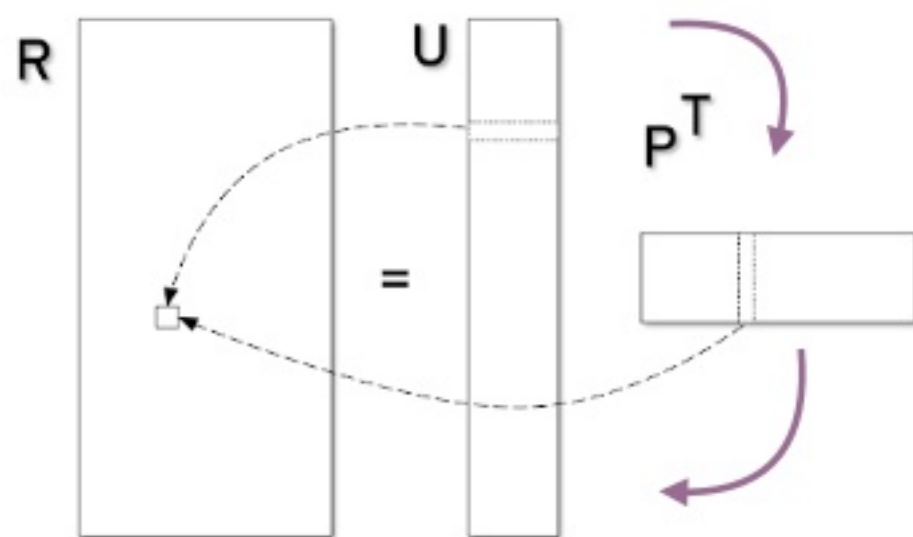
$$C_{ui} = 1 + \alpha r_{ui}$$

$r_{ui} \in \mathbb{R}$  recording

relates to scale of recording

# Iterations

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```



# Making Predictions

```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

```
1 predictions = model.predictAll(zero_listens)
```

(user id, item id)



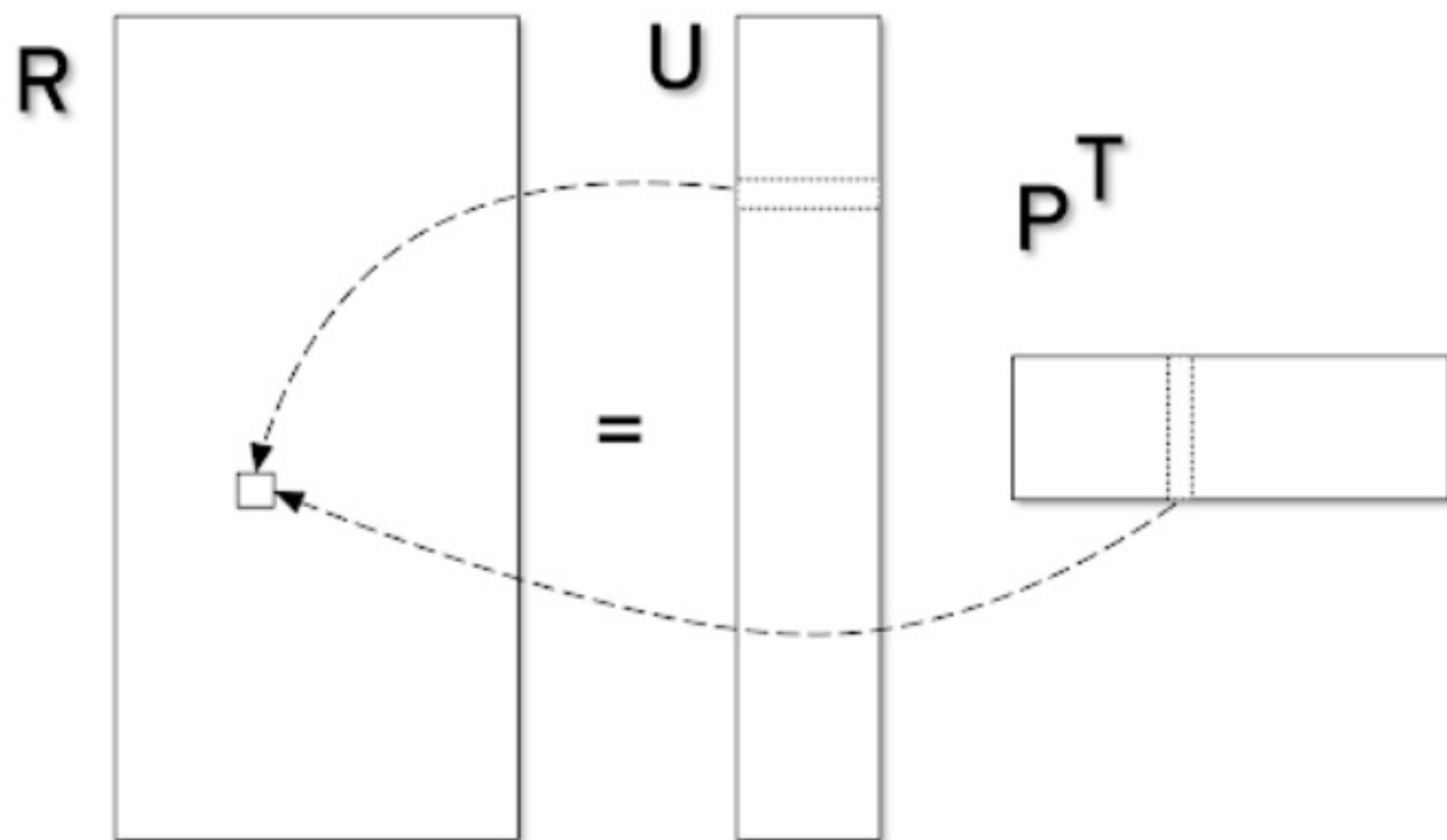
```
user1_listened.take(10)  
user1_listened.map(lambda x:(x[1][1][1:])).take(10)
```

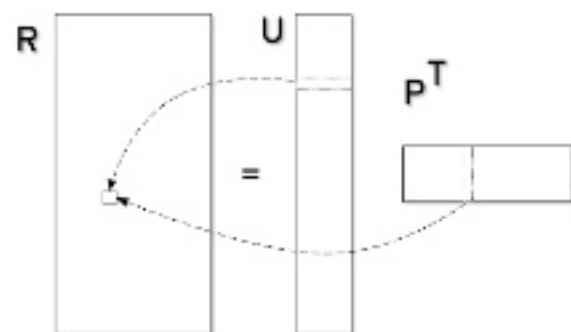
```
[['beirut', '609'],  
 ['dredg', '605'],  
 ['calexico', '562'],  
 ['led', 'zeppelin', '456'],  
 ['laura', 'marling', '401'],  
 ['minus', 'the', 'bear', '377'],  
 ['zion', 'i', '352'],  
 ['bon', 'iver', '313'],  
 ['xavier', 'rudd', '306'],  
 ['passion', 'pit', '273']]
```



```
user1_pred=model.predictAll(user1_unlistened)
```

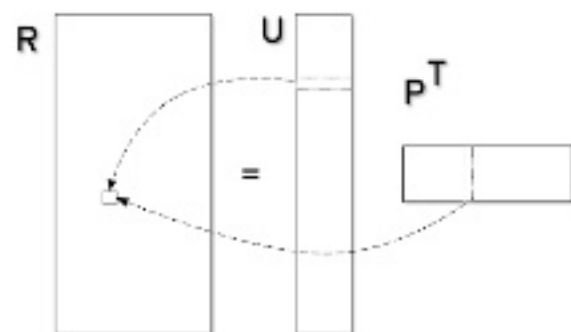
```
[['john', 'frusciante', '2140'],  
 ['red', 'hot', 'chili', 'peppers', '1614'],  
 ['waglewski', 'fisz', 'emade', '566'],  
 ['coldplay', '461'],  
 ['the', 'mars', 'volta', '402'],  
 ['pj', 'harvey', '398'],  
 ['muchy', '397'],  
 ['maria', 'peszek', '373'],  
 ['fisz', '305'],  
 ['ataxia', '301']]
```





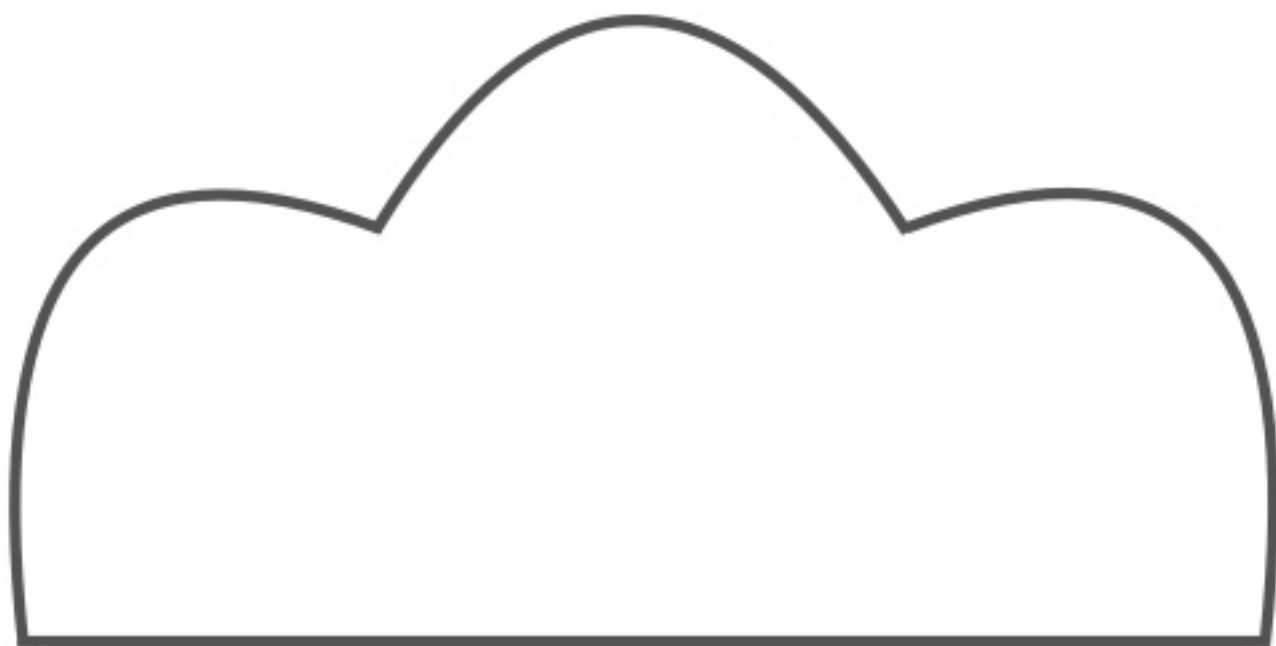
```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

```
1 predictions = model.predictAll(zero_listens)
```

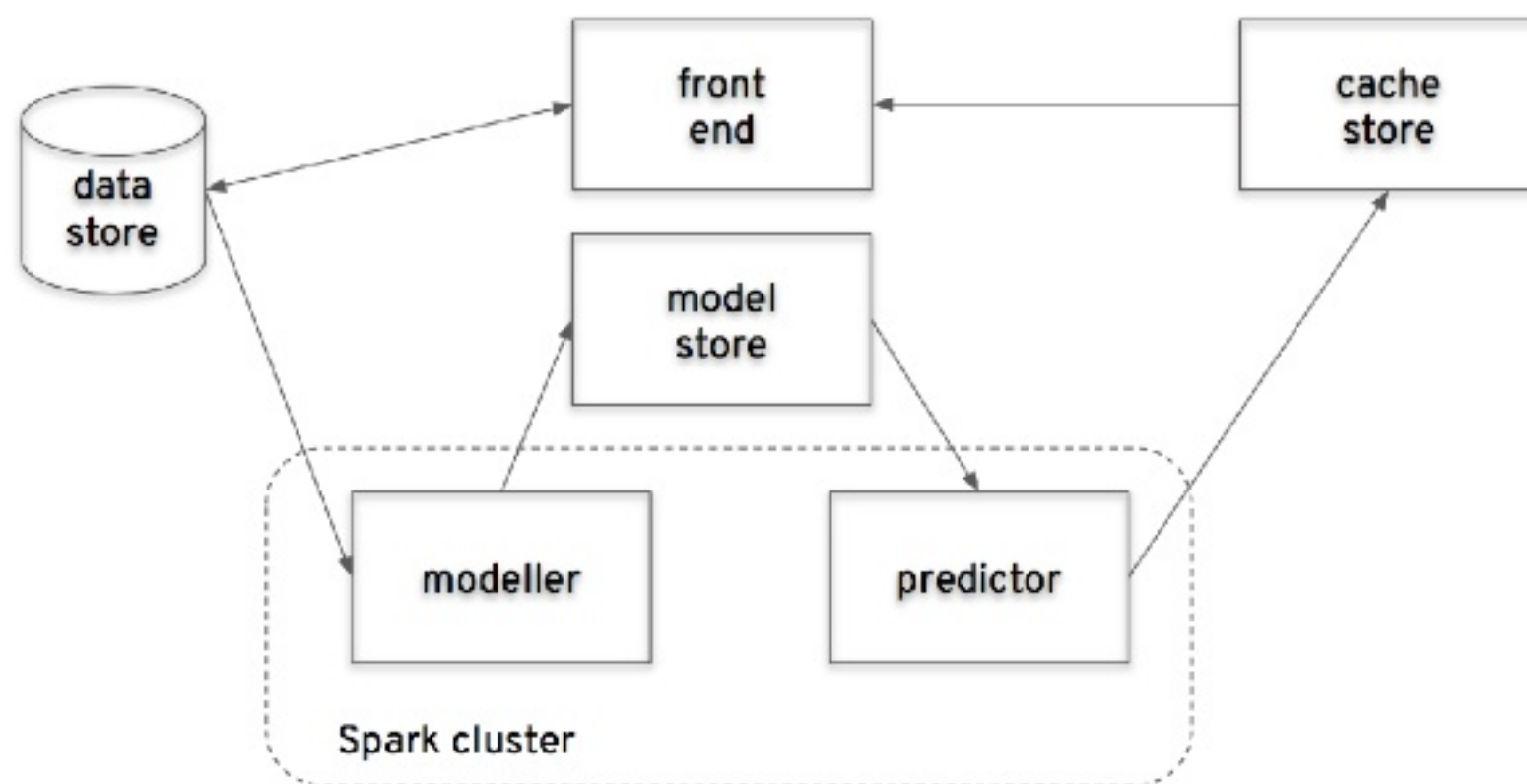


```
1 model=ALS.trainImplicit(data_set, rank=5, lambda_=0.01, alpha = 1.0, iterations=5)
```

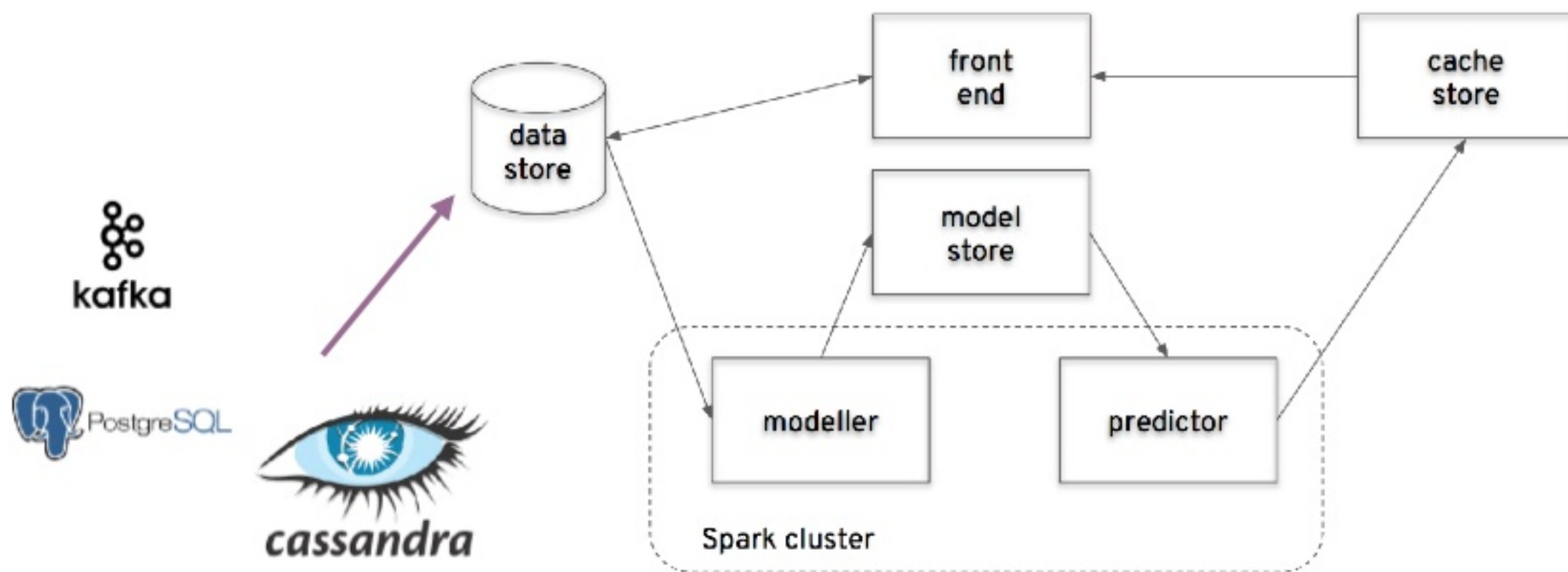
```
1 predictions = model.predictAll(user listens)
```



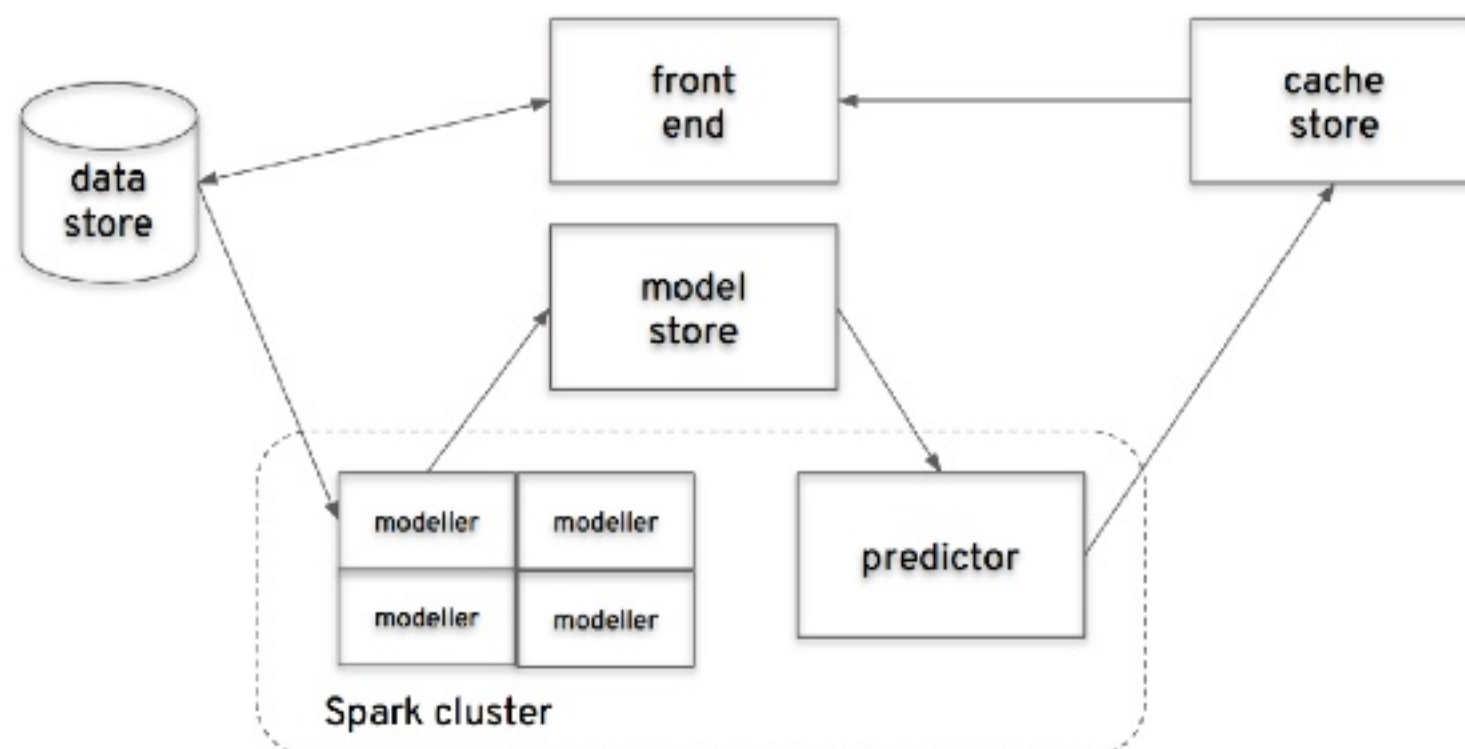
# Microservices



# Microservices



# Microservices



# radanalytics.io



## Contents

[Introduction](#)  
[Architecture](#)  
[Installation](#)  
[Usage](#)  
[Expansion](#)  
[Videos](#)

## Recommendation engine service with Apache Spark

### Introduction

Project Jiminy is a service based application that implements a simple [recommendation system](#) using [collaborative filtering](#) based on an alternating least squares methodology. That may sound complicated but through the source repositories and these instructions you will find that creating a recommendation engine is more straightforward than expected.

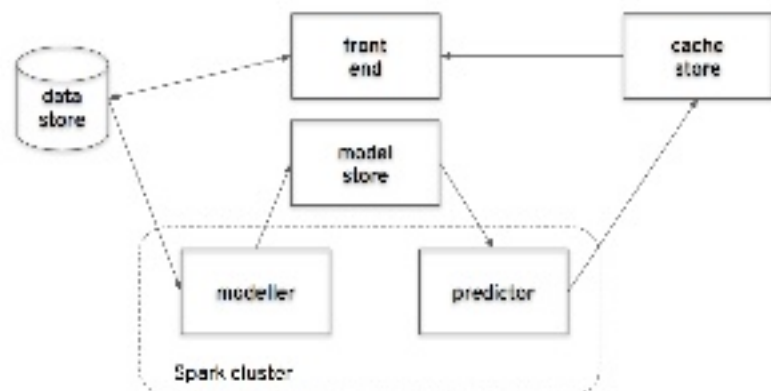
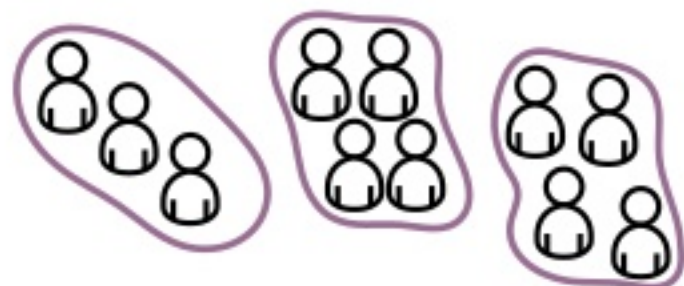
With these instructions you will learn how to deploy jiminy with the [MovieLens dataset](#) by the GroupLens Research organization. This dataset represents a set of movies, users and their ratings of the movies. Although Jiminy uses this dataset as the starting point, you will see how easily the services can be modified to utilize your own datasets.



@sophwats

sophie@redhat.com

## Collaborative Filtering



## Alternating Least Squares

