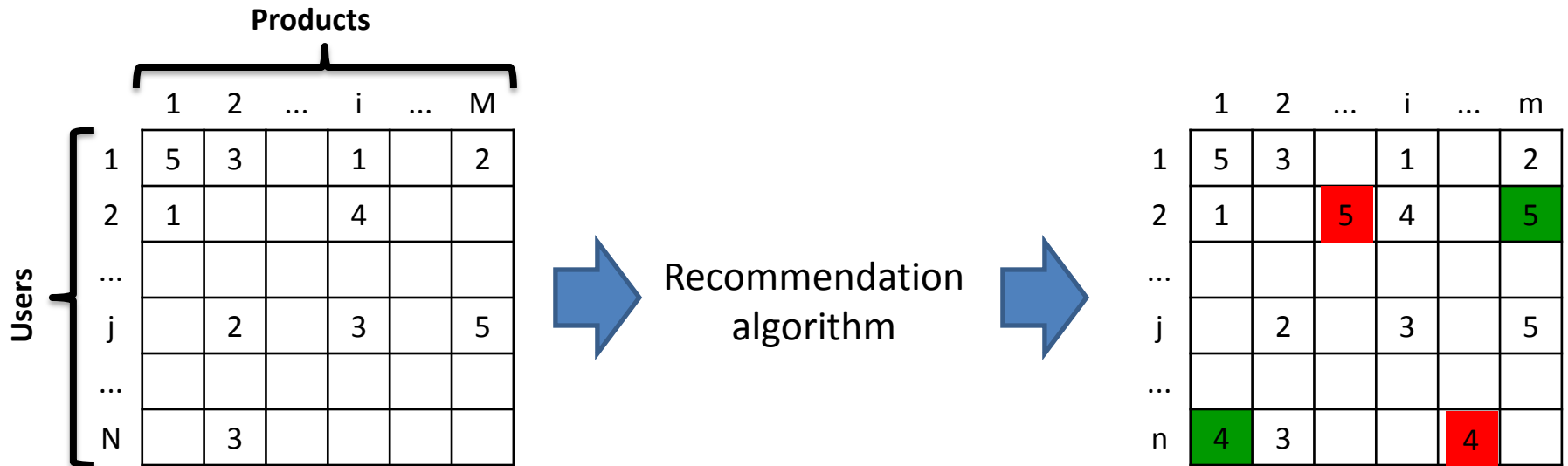# A Regularized Recommendation Algorithm with Probabilistic Sentiment-Ratings

Filipa Peleja, Pedro Dias and João Magalhães

Department of Computer Science
Faculdade de Ciências e Tecnologia
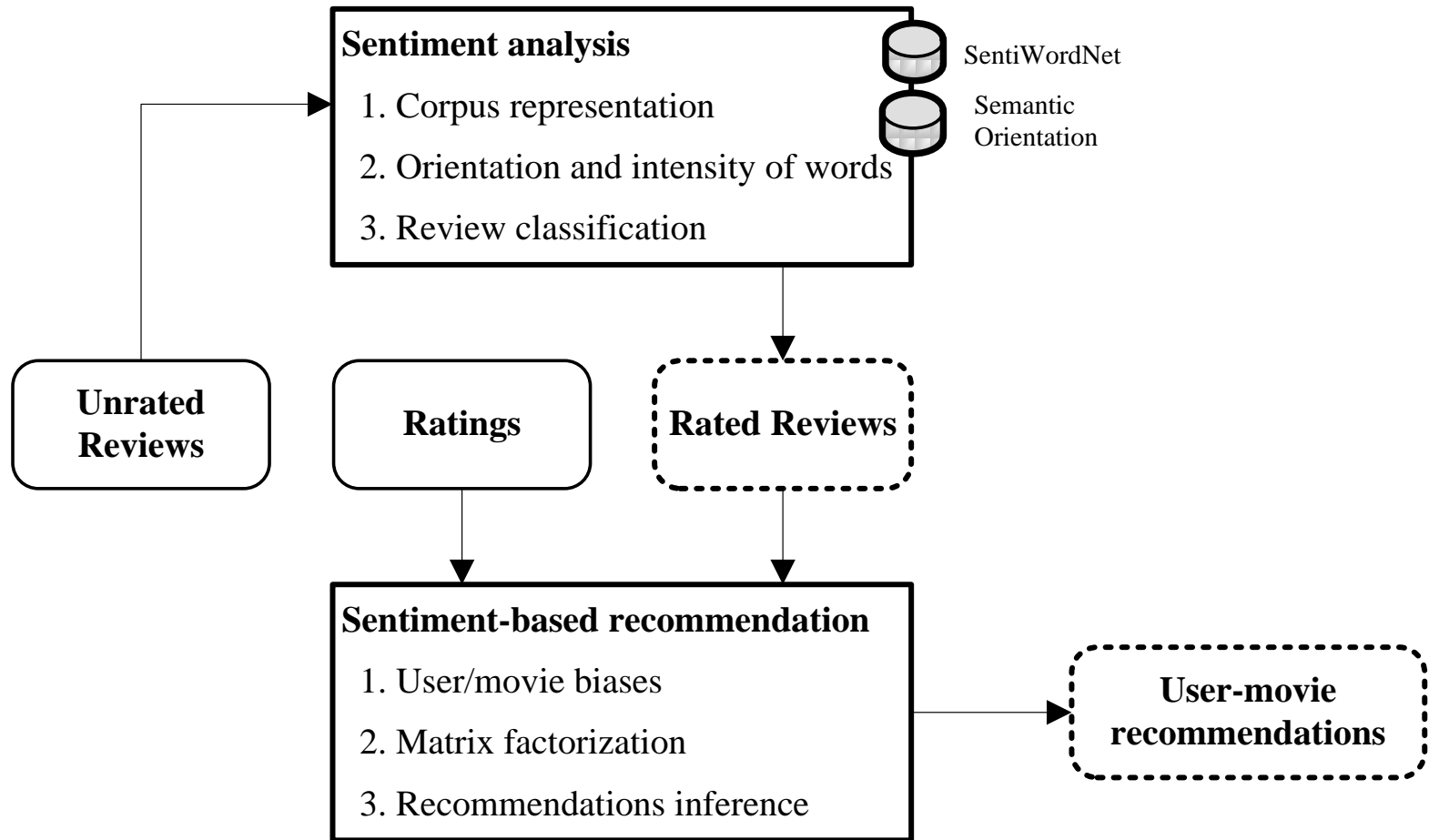Universidade Nova de Lisboa

# OBJECTIVE: HOW TO IMPROVE RECOMMENDATIONS WITH USER COMMENTS AND REVIEWS?



**Rating**: ★★★★
**Review**: Love it or hate it!

...

**Rating**: ★★
**Review**: This is a miserable film.

# PROPOSED SOLUTION

# SENTIMENT ANALYSIS CHALLENGES

- Opinions are written in natural language which implies :

  - subjectivity; - sarcasm; - irony; - idiomatic expressions; misspelling; etc.

- The same **opinion word** may be used in a positive or negative context

- Negative, Conditional and Comparative expressions

# OPINION WORD ORIENTATION AND INTENSITY

- Semantic Orientation[1]:

$$SO(word) = \log_2\left(\frac{hits(word, "excellent")hits("poor")}{hits(word, "poor")hits("excellent")}\right)$$

- How positive or negative is an *opinion word*?
  - SentiWordNet[2]

(1) TURNEY, P. 2002, THUMBS UP OR THUMBS DOWN? SEMANTIC ORIENTATION APPLIED TO UNSUPERVISED CLASSIFICATION OF REVIEWS
(2) ESULI, A. AND SEBASTIANI, F., 2006, SENTIWORDNET: A PUBLICY AVAILABLE LEXICAL RESOURCE FOR OPINION MINING

# Example

"Love it or hate it."

"However, can someone tell me what on earth the last page…"

| word | family | SO (Google) | + SentiWordNet | - SentiWordNet |
|------|--------|-------------|----------------|----------------|
| love | n | -0.0824 | 1.375 | 0.0 |
| it | nointerest | na | na | na |
| or | nointerest | na | na | na |
| hate | v | -0.8399 | 0.0 | 0.75 |
| it | nointerest | na | na | na |
| however | r | -0.34153 | 0.5 | 0.5 |
| someone | N | -0.65935 | 0.0 | 0.0 |
| tell | V | -0.3956 | 0.875 | 0.625 |
| me | nointerest | na | na | na |
| what | nointerest | na | na | na |
| on | nointerest | na | na | na |
| earth | n | -0.4041 | 0.0 | 0.625 |

# Multiple Bernoulli Classification

$$p(ra_i = r \mid re_i) = \frac{f^r(ra_i, re_i)}{\sum_{L=1}^{10} f^L(ra_i, re_i)}$$

A classifier is learned for every rating value. Thus, for each rating value there is a prediction for each review

Rating range IMDb dataset is 1 to 10

The prediction is normalized accordingly to the predictions of all ratings

# DATASET
# REVIEWS FROM IMDB

## A TOTAL OF 1,729,293 REVIEWS WERE COLLECTED

| Split | #Reviews | Description |
|:-----:|---------|-------------|
| A | 335,975 | Only to train SA |
| B | 335,975 | Test SA/Train RS |
| C | 417,147 | Train RS (**no explicit ratings**) |
| D | 335,976 | Train RS |
| E | 201,586 | Test RS |
| F | 102,634 | Validate RS |

# PERFORMANCE - SENTIMENT ANALYSIS



F-score on the IMDb corpus

# PERFORMANCE - SENTIMENT ANALYSIS

# Inferred Ratings in Recommendation Algorithm

# RATING MATRIX

$$R_{ra} = \begin{bmatrix} r_{11} & \cdots & r_{1m} \\ \vdots & \ddots & \vdots \\ r_{n1} & \cdots & r_{nm} \end{bmatrix}$$

$\longleftarrow$ HIGHLY INCOMPLETE SINCE MOST ELEMENTS ARE EMPTY

# PREDICT AN UNKNOWN RATING:

$$\hat{r}_{ui} = p_u . q_i$$

USERS AND PRODUCTS REPRESENTED IN THE SAME LATENT FACTOR SPACE

# WITH A SVD DECOMPOSITION THE RATING MATRIX

$$R_{ra} = \begin{bmatrix} u_{11} & \cdots & u_{1m} \\ \vdots & \ddots & \vdots \\ u_{n1} & \cdots & u_{nm} \end{bmatrix} . \begin{bmatrix} p_{11} & \cdots & p_{1m} \\ \vdots & \ddots & \vdots \\ p_{n1} & \cdots & p_{nm} \end{bmatrix}^T = P.Q^T$$

MATRIX FACTORIZATION ENABLES THE ASSESSMENT OF USERS PREFERENCES REGARDING THE PRODUCTS BY CALCULATING THEIR FACTOR REPRESENTATIONS

# RATINGS MATRIX $R_{ra}$
## *FACTORIZATION WITH BIASES*

GOAL: MINIMIZE THE
PREDICTION ERROR

$$[P,Q] = \arg\min_{p_u, q_i} \sum_{r_{ui} \in R_{ra}} (r_{ui} - \hat{r}_{ui})^2 + \lambda(\|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2)$$

# RATINGS MATRIX $R_{ra}$ FACTORIZATION WITH BIASES

$$[P,Q] = \underset{p_u,q_i}{\arg\min} \sum_{r_{ui} \in R_{ra}} (r_{ui} - \hat{r}_{ui})^2 + \sum_{r_{ui} \in \hat{R}_{rev}} (\hat{c}_{ui} - \hat{r}_{ui})^2 + \lambda(\| p_u \|^2 + \| q_i \|^2 + b_u^2 + b_i^2)$$

RATINGS INFERRED FROM THE SENTIMENT ANALYSIS
FRAMEWORK ARE GIVEN TO THE RS
THE REVIEWS ACTUAL RATING ARE KNOWN

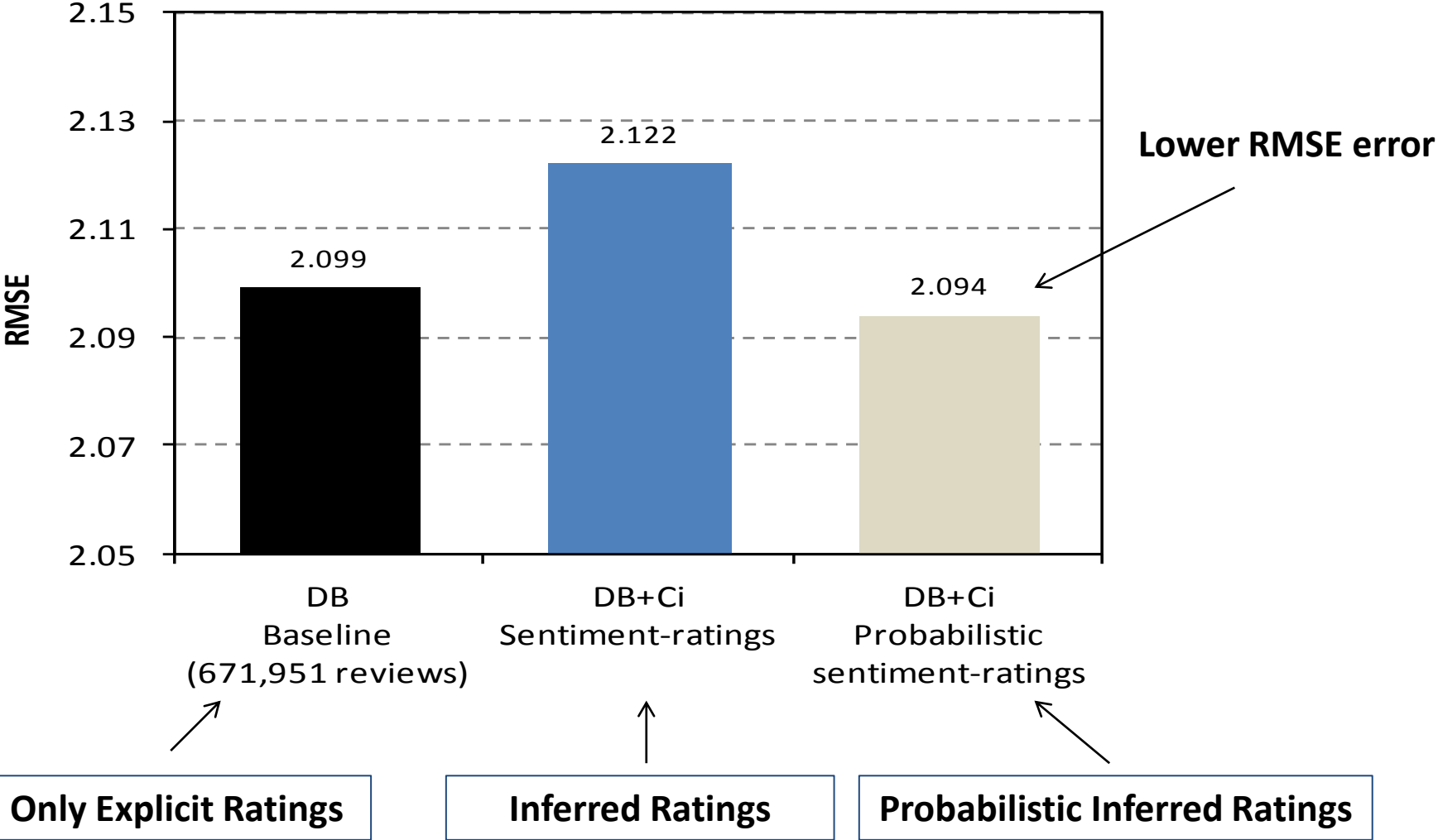# RATINGS MATRIX $R_{ra}$ FACTORIZATION WITH SENTIMENT-BASED REGULARIZATION

$$\mathbb{R} = R_{ra}, \mathbb{R}_{rev}$$

ENRICH THE MATRIX R WITH RATINGS INFERRED FROM REVIEWS WITH **KNOWN** AND **UNKNOWN EXPLICIT RATINGS** $\mathbf{R_{REV}}$
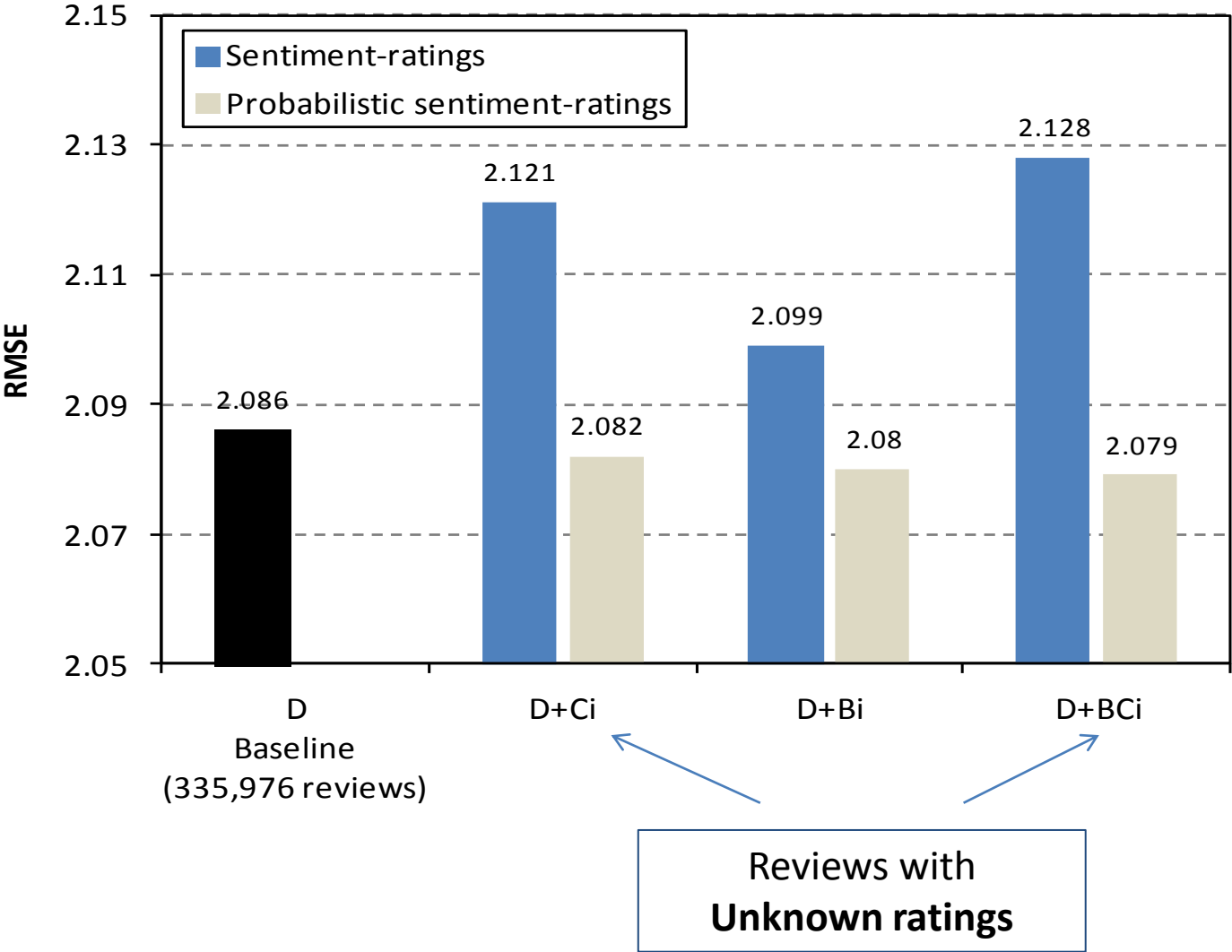
$$\hat{R}_{ra} = \arg\min_{\hat{r}_{ui}} \sum_{r_{ui} \in R_{ra}} (r_{ui} - \hat{r}_{ui})^2 + \sum_{c_{ui} \in \hat{R}_{rev}} \theta_{ui} \cdot (\hat{c}_{ui} - \hat{r}_{ui})^2$$
$$+ \lambda(\| p_u \|^2 + \| q_i \|^2 + b_u{}^2 + b_i{}^2),$$

The confidence level is given by de Sentiment Analysis framework

# RECOMMENDATIONS: IMDB DATASET



15

# RECOMMENDATIONS: IMDB DATASET

# Summary

- **Achievements**:
  - Extraction and sentiment analysis of users reviews
  - Introduced sentiment-based ratings in a recommendation algorithm

- **Next step**:
  - alternatives to SentiWordNet
  - semantic orientation metric
  - improve algorithm with opinion targets information

# Thank you for your attention

# Questions?