

# **Factorization Meets the Neighborhood: a Multifaceted Collaborative Filtering Model**

**Yehuda Koren  
AT & T Labs – Research  
2008**

**Present by**

**Hong Ge  
Sheng Qin**

# Info about paper & data-set

## **Factorization Meets the Neighborhood: a Multifaceted Collaborative Filtering Model**

---

- ACM Transactions on Knowledge Discovery from Data (TKDD) archive
- Year of Publication: 2007; cited by 43 times
- Winner of the \$1 Million Netflix Prize (2007)!!!!
  - 9.34% improvement over the original Cinematch accuracy level
- Netflix data:
  - Over 480,000 users, 17,770 movies
  - Over 1 million observed ratings, 1% in total
  - Rating: integer from 1 to 5 (with rating time-stamp)
  - Multivariate, Time-Series



# Title interpretation

## **Factorization Meets the Neighborhood:** **a Multifaceted Collaborative Filtering Model**

---

- Technique about recommender systems
  - Based on: **Collaborative Filtering (CF)**
    - A process often applied to recommender systems
  - Using: **Neighborhood Model** & **Latent Factor Model**
    - Two main disciplines of CF
  - Solution: **Some amazing improvement & integration**
    - Innovative point of this paper
-

# Background

## Collaborative Filtering

Analyze past transactions to establish connections between users and movies.

- Relies on past user behavior
- Does not require explicit profile

Existing  
methods

## Neighborhood

- Computing relationships between movies, or between users
- Not user  $\rightarrow$  movie, but movie  $\rightarrow$  movie

## Latent factor

- Characterize user  $\rightarrow$  movie on factors
- Factors are inferred from user feedback



# The integrated model

❖ **Why integrate?**





# The integrated model-why?

## ❖ Neighborhood Models

- Estimate unknown ratings by using known ratings made by user for similar movies
- Good at capturing localized information
- Intuitive and simple to implement

## ❖ Latent Factor Models

- Estimate unknown ratings by uncover latent features that explain known ratings
  - Efficient at capturing global information
- 

# The integrated model-why?

## ❖ Reasons:

- Neighborhood Model: Good at capture localized information
- Latent Factor Model: Efficient at capturing global information
- Neither is able to capture all information
- Complementary with each other.
- Not account implicit feedback
- It's not tried before, why not?



# The integrated model-how?

## ❖ How ?

- Sum the predications of revised Neighborhood Model(NewNgbr) and revised Latent Model (SVD++)

## ❖ Some details

- I guess you may want take a nap now.
- Just joking!



# Some background before we go further

## ❖ The Netflix data

- Many items in this matrix are missing
- Need find a good estimate for (most of efforts are dealing with this!)

	Ratings			
Users	$r_{11}$	$r_{12}$	$\dots$	$r_{1i}$
	$r_{21}$	$r_{22}$	$\dots$	$r_{2i}$
	$\dots$	$\dots$	$\dots$	$\dots$
	$r_{u1}$	$r_{u2}$	$\dots$	$r_{ui}$
	[Netflix data]			

## ❖ Baseline estimates

- $\mu$  is the **average rating** over all movies
- $b_u, b_j$  indicate the **observed deviations** of user  $u$  and item  $i$ , respectively, from the average

$$b_{ui} = \mu + b_u + b_i$$

[baseline estimator]

# Neighborhood Model

❖ Estimate  $\hat{r}_{ui}$  by using known ratings made by user for similar movies:

$$\hat{r}_{ui} = b_{ui} + \sum_{j \in s^k(i;u)} w_{uj} (r_{uj} - b_{uj})$$

Diagram illustrating the Neighborhood Model equation:

- $\hat{r}_{ui}$ : Estimated rating for user  $u$  on movie  $i$  (represented by a blue bar with a red top).
- $b_{ui}$ : Baseline rating for user  $u$  on movie  $i$  (represented by a blue bar).
- $w_{uj}$ : User specific weights (indicated by a red arrow).
- $j \in s^k(i;u)$ :  $k$  most similar movies rated by  $u$ , also known as Neighbors (indicated by a red arrow).
- $r_{uj}$ : Rating given by user  $u$  to movie  $j$  (represented by a red bar).
- $b_{uj}$ : Baseline rating for user  $u$  on movie  $j$  (represented by a red bar).

# Neighborhood models- Revised

## ❖ New Neighborhood model:

- introduce implicit feedback effect
- use global rather than user-specific weights

## ❖ New predicting rule:

The diagram illustrates the new predicting rule equation. It features a vertical stack of a red square on top of a dark blue rectangle on the left, followed by an equals sign, another dark blue rectangle, a plus sign, a red-bordered box containing the formula  $R^{-1/2} \sum_{j \in R^k(i;u)} w_{ij}$ , a red square, a plus sign, another red-bordered box containing the formula  $N^{-1/2} \sum_{j \in N^k(i;u)} c_{ij}$ , and a yellow square. Red arrows point from the text 'introduce implicit feedback effect' to the red square in the second term, and from 'use global rather than user-specific weights' to the yellow square in the third term.

$$\hat{r}_{ui} = b_{ui} + R^{-1/2} \sum_{j \in R^k(i;u)} w_{ij} (r_{uj} - b_{uj}) + N^{-1/2} \sum_{j \in N^k(i;u)} c_{ij}$$

# Latent Models

- ❖ Estimate  $\hat{r}_{ui}$  by uncover latent features that explain observed ratings:

$$\begin{array}{c} \text{[Red bar]} \\ \text{[Blue bar]} \end{array} = \begin{array}{c} \text{[Blue bar]} \end{array} + \begin{array}{c} ( \quad \cdot \quad \cdot ) \end{array} \begin{array}{c} \left( \begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array} \right) \end{array}$$

$\hat{r}_{ui}$     $b_{ui}$     $p_u^T$     $q_i$

- $p_u, q_i$  are user-factors vector and item-factors vector respectively

# Latent Model- Revised

## ❖ Introduce implicit feedback information

- Asymmetric-SVD

$$\hat{r}_{ui} = \underbrace{b_{ui}}_{\text{baseline estimate}} + q_i^T \left( R^{-1/2} \sum_{j \in R(u)} (r_{uj} - b_{uj}) + \underbrace{N^{-1/2} \sum_{j \in N(u)} y_j}_{\text{Implicit feedback effect}} \right)$$

## ❖ SVD+ +

- No theoretical explanation, it just works!

$$\hat{r}_{ui} = b_{ui} + q_i^T \left( p_u + N^{-1/2} \sum_{j \in N(u)} y_j \right)$$

- This model will be integrated with Neighborhood Model



# The integrated model

## ❖ How well does it work?

- Here is the result.



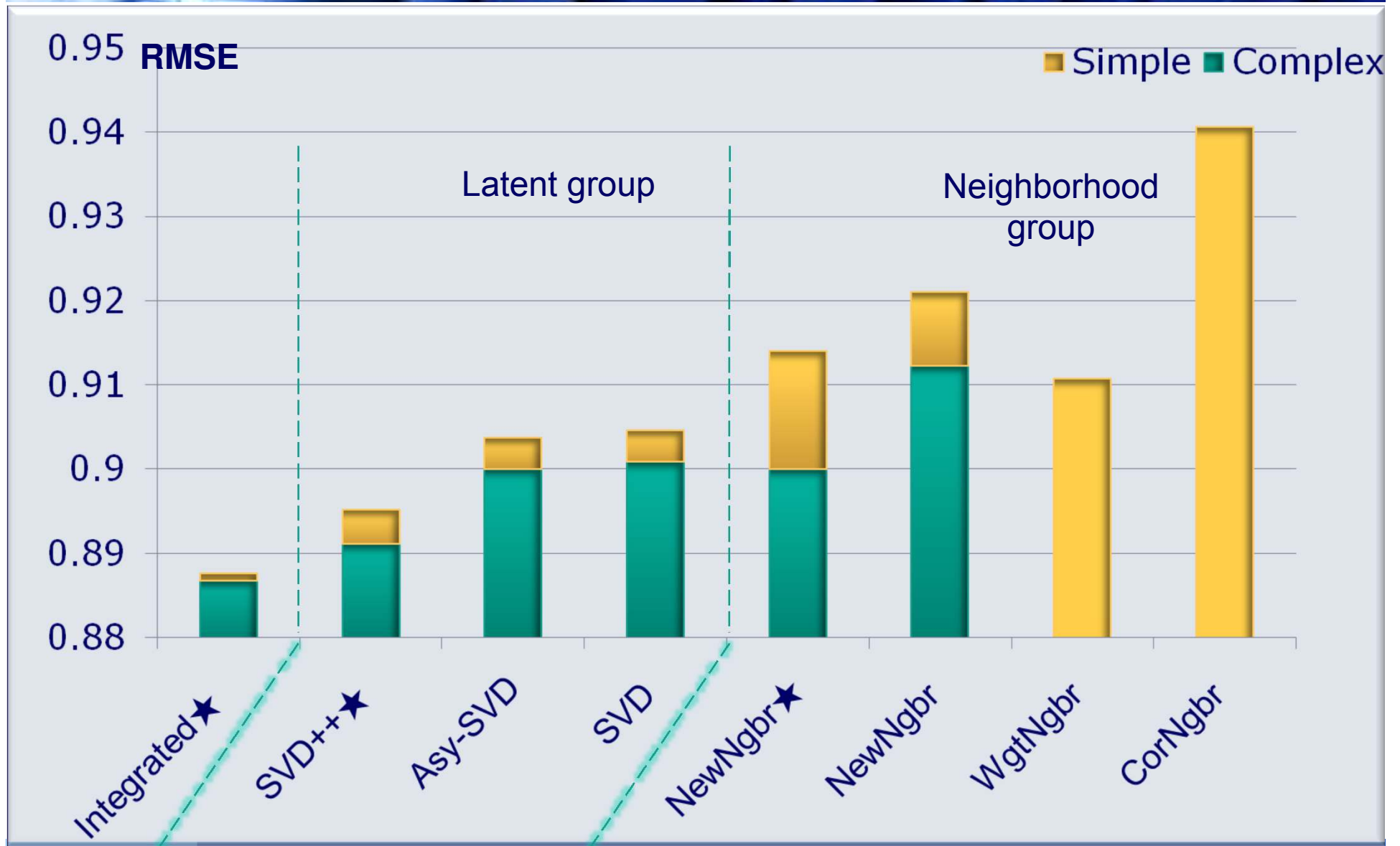
# Test (Instructions)

Measured by Root Mean Square Error (RMSE)

$$\sqrt{\sum_{(u,i) \in TestSet} (r_{ui} - \hat{r}_{ui})^2 / |TestSet|}$$

Abbreviation instructions	
Integrated★	Proposed Integrated Model
SVD+ + ★	Proposed improved Latent Factor
SVD	Common Latent Factor
New Ngbr★	Proposed neighborhood, with implicit feedback
New Ngbr	Proposed neighborhood, without implicit feedback
WgtNgbr	improved neighborhood of the same user
CorNgbr	Popular neighborhood method

# Experimental results — RMSE



# Time cost

## NewNeighborhood

Time*(min)	10	27	58
Neighbors	250	500	Infinity
Precision	0.9014	-0.0010	-0.0004

## SVD++

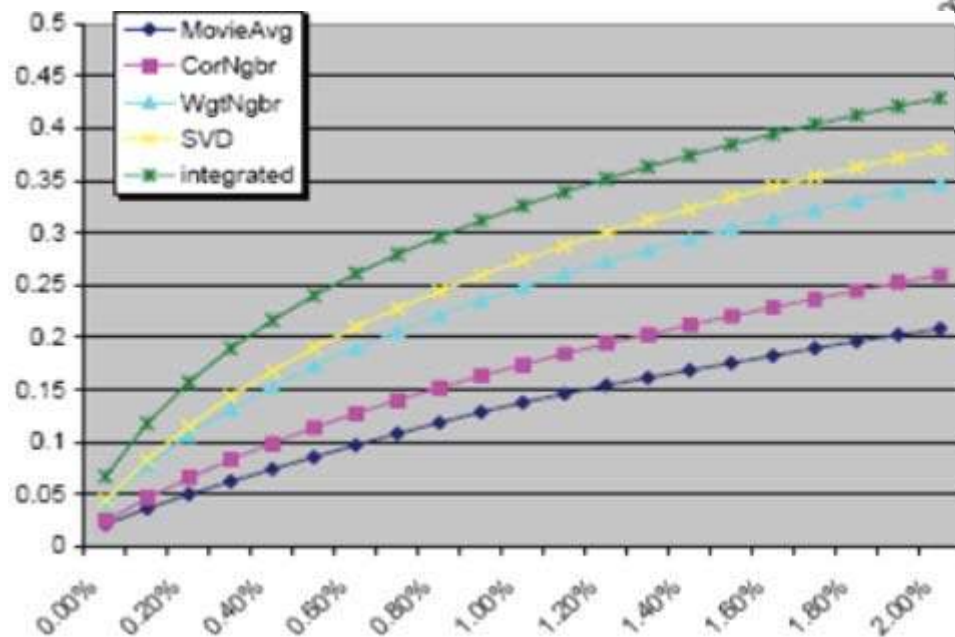
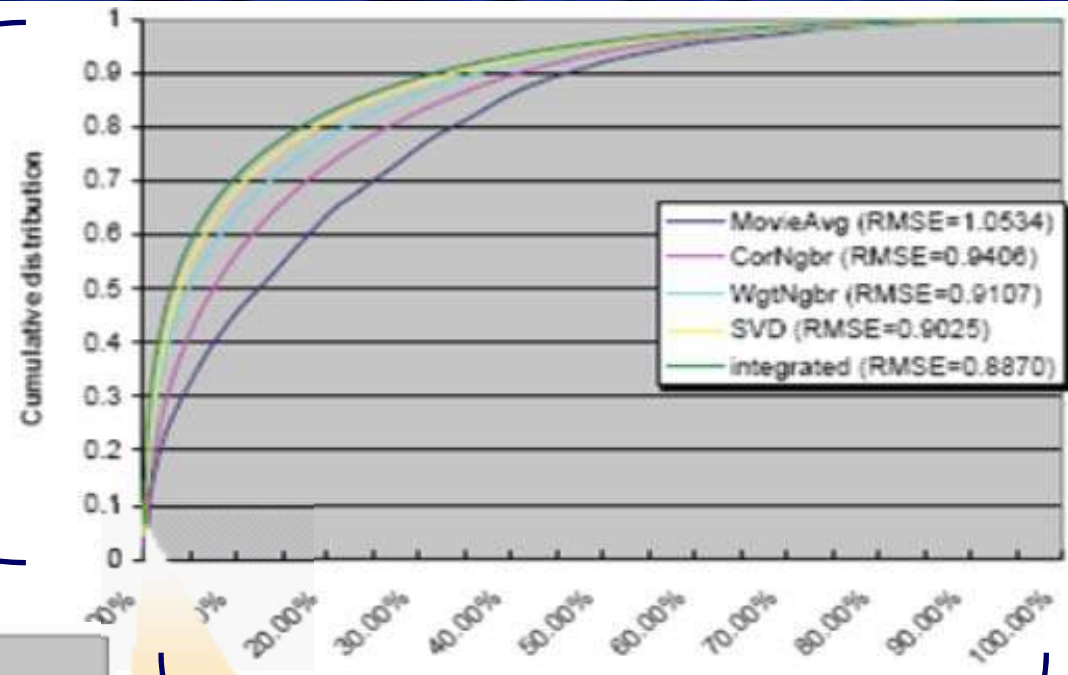
Time*(min)	--	--	--
Factors	50	100	200
Precision	0.8952	-0.0028	-0.0013

## Integrated

Time(min)	17	20	25
Neighbors	300	300	300
Factors	50	100	200
Precision	0.8877	-0.0007	-0.0002

# Experimental results — top K

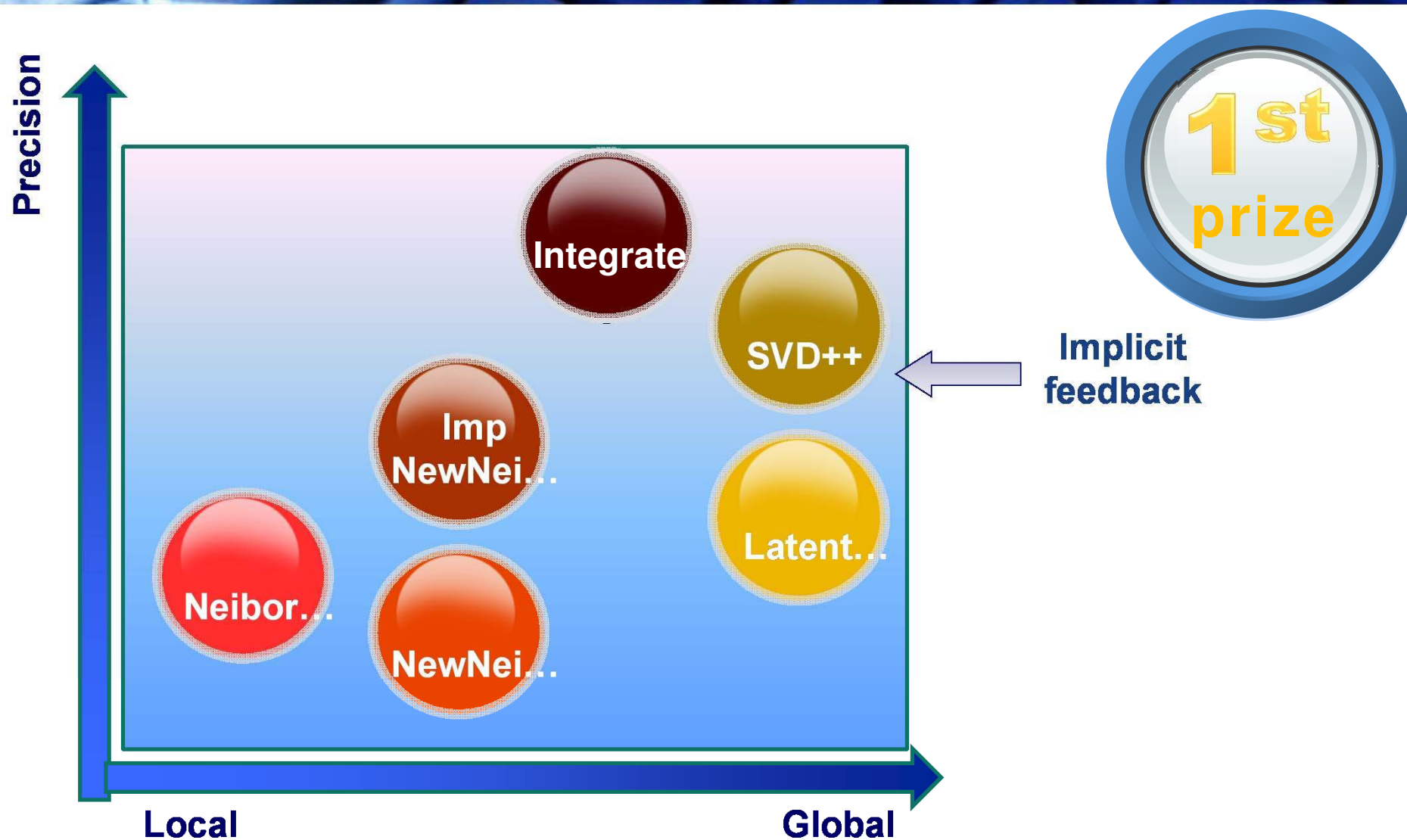
**Y axis:**  
Probability distribution of the  
observed best movie returned



0%~2%

**X axis:**  
Threshold of return in  
percentile

# Conclusion





# Hard to beat, but...



## Ignored time-stamps

- Time-stamps available (from 1998 to 2005)
- Temporal dynamics matters

### Example 1



Action

6 years later...



Romance





# Hard to beat, but...

## ● Ignored time-stamps

- Time-stamps available (from 1998 to 2005)
- Temporal dynamics matters

### Example 2





# Hard to beat, but...



## Temporal dynamics are too personal

- Represented in author's latest publication, with comparison
- May move the model towards local level





# References

- ❖ **Yehuda Koren, Factorization meets the neighborhood: a multifaceted collaborative filtering model, in Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining (Las Vegas, Nevada, USA: ACM, 2008), 426-434**
- ❖ **Yehuda Koren, The BellKor Solution to the Netflix Grand Prize, August 2009**



❖ **Questions?**

