# Introduction to Computational Advertising

MS&E 239

Stanford University

Autumn 2011

Instructors: Dr. Andrei Broder and Dr. Vanja Josifovski

Yahoo! Research

1

# General course info

- Course Website: http://www.stanford.edu/class/msande239/
- Instructors
  - **Dr. Andrei Broder**, Yahoo! Research, broder@yahoo-inc.com
  - **Dr. Vanja Josifovski**, Yahoo! Research, vanjaj@yahoo-inc.com
- TA: **Krishnamurthy Iyer**
  - Office hours: Tuesdays 6:00pm-7:30pm, Huang
- Course email lists
  - Staff: msande239-aut1112-staff
  - All: msande239-aut1112-students
  - Please use the staff list to communicate with the staff
- Lectures: 10am ~ 12:30pm Fridays in HP
- Office Hours:
  - After class and by appointment
  - Andrei and Vanja will be on campus for 2 times each to meet and discuss with students. Feel free to come and chat about even issues that go beyond the class.

# Course Overview (subject to change)

1. 09/30 Overview and Introduction
2. 10/07 Marketplace and Economics
3. 10/14 Textual Advertising 1: Sponsored Search
4. 10/21 Textual Advertising 2: Contextual Advertising
5. 10/28 Display Advertising 1
6. 11/04 Display Advertising 2
7. 11/11 Targeting
8. 11/18 Recommender Systems
9. 12/02 Mobile, Video and other Emerging Formats
10. 12/09 Project Presentations

# Lecture 8:
# Recommender Systems for Display Advertising Targeting

# Disclaimers

- This talk presents the opinions of the authors. It does not necessarily reflect the views of Yahoo! inc or any other entity.

- Algorithms, techniques, features, etc mentioned here might or might not be in use by Yahoo! or any other company.

- This lecture is largely based on slides by **Yehuda Koren**. His help very much appreciated! Other contributors: Amr Ahmed and Alex Smola

# Lecture 8 plan

- Problem definition

- Neighborhood based approaches

- Matrix factorization approaches

- Generative models for factorization: Latent Dirichlet Allocation

- External information in campaign optimization (time permitting)

# Checkpoint - targeting

- Targeting is a key step in differentiation of impressions and extracting value!
- Traditional targeting: demo, geo, BT
  - How to get the data from the user?
  - Infer the data from historical activity
- One of the key step in targeting is user profile generation
  - Generative models to assign probability of a sequence of events
  - Weighting based on time, event type and content
  - Predict the counts of events in certain categories
  - Clustering and other unsupervised techniques useful – more to come in the next lecture

# Recommender systems

- What is actually recommender system technology?
- A set of techniques to recommend items based on explicit (rating) or implicit (page visits, ad clicks)
- It collects the user responses and assumes the items are opaque
- Usually does not take in account the content of the item
  - Opposed to matching of ads we have discussed so far
  - Recent techniques combine the two
- Shown to be effective in many domains, including advertising
- Slides mostly on movie ratings, we will discuss the similarities/differences wit the ad domain

# Collaborative filtering

- Recommend items based on past transactions of many users

- Analyze relations between users and/or items

- Specific data characteristics are irrelevant
  - Domain-free: user/item attributes are not necessary
  - Can identify elusive aspects

amazon.com

**Customers who bought items in your Recent History also bought:**

I Own It  Not interested
x|☆☆☆☆☆ Rate it
Add to Cart  Add to Wish List

I Own It  Not interested
x|☆☆☆☆☆ Rate it
Add to Cart  Add to Wish List

LOOK INSIDE!™

I Own It  Not interested
x|☆☆☆☆☆ Rate it
Add to Cart  Add to Wish List

# Movie rating data

## Training data

| user | movie | date | score |
|------|-------|------|-------|
| 1 | 21 | 5/7/02 | 1 |
| 1 | 213 | 8/2/04 | 5 |
| 2 | 345 | 3/6/01 | 4 |
| 2 | 123 | 5/1/05 | 4 |
| 2 | 768 | 7/15/02 | 3 |
| 3 | 76 | 1/22/01 | 5 |
| 4 | 45 | 8/3/00 | 4 |
| 5 | 568 | 9/10/05 | 1 |
| 5 | 342 | 3/5/03 | 2 |
| 5 | 234 | 12/28/00 | 2 |
| 6 | 76 | 8/11/02 | 5 |
| 6 | 56 | 6/15/03 | 4 |

## Test data

| user | movie | date | score |
|------|-------|------|-------|
| 1 | 62 | 1/6/05 | ? |
| 1 | 96 | 9/13/04 | ? |
| 2 | 7 | 8/18/05 | ? |
| 2 | 3 | 11/22/05 | ? |
| 3 | 47 | 6/13/02 | ? |
| 3 | 15 | 8/12/01 | ? |
| 4 | 41 | 9/1/00 | ? |
| 4 | 28 | 8/27/05 | ? |
| 5 | 93 | 4/4/05 | ? |
| 5 | 74 | 7/16/03 | ? |
| 6 | 69 | 2/14/04 | ? |
| 6 | 83 | 10/3/03 | ? |

# Alternate view of the data: matrix

users

|     | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|-----|---|---|---|---|---|---|---|---|---|----|----|----|
| 1   | 1 |   | 3 |   | ? | 5 |   |   | 5 |    | 4  |    |
| 2   |   |   | 5 | 4 |   |   | 4 |   |   | 2  | 1  | 3  |
| 3   | 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3  | 5  |    |
| 4   |   | 2 | 4 |   | 5 |   |   | 4 |   |    | 2  |    |
| 5   |   |   | 4 | 3 | 4 | 2 |   |   |   |    | 2  | 5  |
| 6   | 1 |   | 3 |   | 3 |   |   | 2 |   |    | 4  |    |

items

# Major Challenges

1. Countless factors may affect preferences
   - Genre, movie/TV series/other
   - Style of action, dialogue, plot, music et al.
   - Director, actors

2. Large imbalances
   - Most user-item preferences are unknown
   - Number of ratings per user or item may vary by several orders of magnitude
   - Information to estimate individual parameters varies widely

3. Scalability
   - Some datasets contain millions of users/items

# Conventions

- $r_{ui}$ - rating by user u to item i

- $\hat{r}_{ui}$ - predicted rating by user u to item i

- Error function:

$$rmse(S) = \sqrt{\frac{\sum_{(u,i) \in S} \left( \hat{r}_{ui} - r_{ui} \right)^2}{|S|}}$$

# How does this map to the ad world

- Ad matrix a lot sparser
- As with movies, no info does not mean negative response
  - We could determine negative responses by analysis of user history
- Ranking metrics might be better option
  - AUC of ROC curve
- Need to limit to the top-k items
  - We cannot show every ad to every user
- In practice – combine rec sys methods with predictive modeling for best performance

# Neighborhood methods



Joe

#3

#2

#1

#4

# Neighborhood-based CF

- Earliest and most popular collaborative filtering method
- Derive unknown ratings from those of "similar" items (item-item variant)
- A parallel user-user flavor: rely on ratings of like-minded users

# Neighborhood-based CF

users

| items | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |  | 3 |  |  | 5 |  |  | 5 |  | 4 |  |
| 2 |  |  | 5 | 4 |  |  | 4 |  |  | 2 | 1 | 3 |
| 3 | 2 | 4 |  | 1 | 2 |  | 3 |  | 4 | 3 | 5 |  |
| 4 |  | 2 | 4 |  | 5 |  |  | 4 |  |  | 2 |  |
| 5 |  |  | 4 | 3 | 4 | 2 |  |  |  |  | 2 | 5 |
| 6 | 1 |  | 3 |  | 3 |  |  | 2 |  |  | 4 |  |

- unknown rating    - rating between 1 to 5

# Neighborhood-based CF

users

| items | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |   | 3 |   | ? | 5 |   |   | 5 |   | 4 |   |
| 2 |   |   | 5 | 4 |   |   | 4 |   |   | 2 | 1 | 3 |
| 3 | 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |   |
| 4 |   | 2 | 4 |   | 5 |   |   | 4 |   |   | 2 |   |
| 5 |   |   | 4 | 3 | 4 | 2 |   |   |   |   | 2 | 5 |
| 6 | 1 |   | 3 |   | 3 |   |   | 2 |   |   | 4 |   |

- estimate rating of item 1 by user 5

# Neighborhood-based CF

users

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |  | 3 |  | ? | 5 |  |  | 5 |  | 4 |  |
| 2 |  |  | 5 | 4 |  |  | 4 |  |  | 2 | 1 | 3 |
| **3** | 2 | 4 |  | 1 | 2 |  | 3 |  | 4 | 3 | 5 |  |
| 4 |  | 2 | 4 |  | 5 |  |  | 4 |  |  | 2 |  |
| 5 |  |  | 4 | 3 | 4 | 2 |  |  |  |  | 2 | 5 |
| **6** | 1 |  | 3 |  | 3 |  |  | 2 |  |  | 4 |  |

items

**Neighbor selection:**
Identify items similar to 1, rated by user 5

# Neighborhood-based CF

users

| items | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 1 | | 3 | | ? | 5 | | | 5 | | 4 | |
| **2** | | | 5 | 4 | | | 4 | | | 2 | 1 | 3 |
| **3** | 2 | 4 | | 1 | 2 | | 3 | | 4 | 3 | 5 | |
| **4** | | 2 | 4 | | 5 | | | 4 | | | 2 | |
| **5** | | | 4 | 3 | 4 | 2 | | | | | 2 | 5 |
| **6** | 1 | | 3 | | 3 | | | 2 | | | 4 | |

**Compute similarity weights:**

$s_{13}=0.2$, $s_{16}=0.3$

# Neighborhood-based CF

users

| items | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |  | 3 |  | 2.6 | 5 |  |  | 5 |  | 4 |  |
| 2 |  |  | 5 | 4 |  |  | 4 |  |  | 2 | 1 | 3 |
| **3** | 2 | 4 |  | 1 | 2 |  | 3 |  | 4 | 3 | 5 |  |
| 4 |  | 2 | 4 |  | 5 |  |  | 4 |  |  | 2 |  |
| 5 |  |  | 4 | 3 | 4 | 2 |  |  |  |  | 2 | 5 |
| **6** | 1 |  | 3 |  | 3 |  |  | 2 |  |  | 4 |  |

**Predict by taking a weighted average:**

(0.2*2+0.3*3)/(0.2+0.3)=2.6

# Properties of neighborhood-based CF

- Intuitive

- Easy to explain reasoning behind a recommendation

- Handles new ratings/users seamlessly

- No substantial preprocessing is required (?)

- Accurate (enough?)

# Data normalization

- Need to identify relations and mix ratings across items/ users

- However:

- User and item-specific variability masks fundamental relationships

- Examples:
  - Some items are systematically rated higher
  - Some items were rated by users that tend to rate low
  - Ratings change along time

- Normalization is critical to the success of a k-NN approach

# Data normalization

- Remove data characteristics that are unlikely to be explained by k-NN

- Common practice is to use <span style="color:red">centering</span>:
  Remove user- and item-means

- A more comprehensive approach eliminates additional interfering variability such as time effects
  See "global effects" @ "Scalable Collaborative Filtering with Jointly Derived Neighborhood Interpolation Weights", ICDM'07

- Here, we normalize by removing the <span style="color:red">baseline predictors</span>:

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

| Global mean | User bias | Item bias |

# Baseline predictors (biases)

- Mean rating: 3.7 stars ( $\mu$ )
- *The Sixth Sense* is 0.5 stars above avg ($b_i$)
- *Joe* rates 0.2 stars below avg ($b_u$)
- ➔ Baseline estimation:
  *Joe* will rate *The Sixth Sense* 4 stars ( $\mu + b_i + b_u$)

# Estimation of biases

- Try to explain each $r_{ui}$ in the train set as $\qquad \mu + b_u + b_i$

- Solve the **regularized** least squares problem:

$$\min_{b_*} \sum_{(u,i) \in K} \underbrace{(r_{ui} - \mu - b_u - b_i)^2}_{\substack{\text{Error for a} \\ \text{training case}}} + \lambda_1 \underbrace{\left( \sum_u b_u^{\,2} + \sum_i b_i^{\,2} \right)}_{\text{Regularization}}$$

**An alternative:**

- First, estimate item biases by averaging over users that rated the item:

$$b_i = \frac{\sum_{u \in R(i)} (r_{ui} - \mu)}{\lambda_2 + |R(i)|}$$

- Then, estimate user biases by averaging residuals over items rated by the user:

$$b_u = \frac{\sum_{i \in R(u)} (r_{ui} - \mu - b_i)}{\lambda_3 + |R(u)|}$$

# Residual: ratings of similar items by the same user

1. Define a similarity measure between items: $s_{ij}$

2. Use $s_{ij}$ to select neighbors – $s^k(i;u)$:
   k items most similar to i, that were rated by u

3. Estimate unknown rating, $r_{ui}$, as the weighted average rating that u gave to the neighbors:

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

- How to compute item-item similarity?

# Estimating item-item similarities

- Common practice – rely on Pearson correlation coeff
- Challenge – non-uniform user support of item ratings, each item rated by a distinct set of users

User ratings for item i:

| 1 | ? | ? | 5 | 5 | 3 | ? | ? | ? | 4 | 2 | ? | ? | ? | ? | 4 | ? | 5 | 4 | 1 | ? |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

User ratings for item j:

| ? | ? | 4 | 2 | 5 | ? | ? | 1 | 2 | 5 | ? | ? | 2 | ? | ? | 3 | ? | ? | ? | 5 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- Compute correlation over shared support

# Estimating item-item similarities

- Empirical Pearson correlation coefficient on shared support of items i and j:

$$\hat{\rho}_{ij} = \frac{\sum\limits_{u \in U(i,j)} (r_{ui} - b_{ui})(r_{uj} - b_{uj})}{\sqrt{\sum\limits_{u \in U(i,j)} (r_{ui} - b_{ui})^2 \cdot \sum\limits_{u \in U(i,j)} (r_{uj} - b_{uj})^2}}$$

U(i,j) contains the users who rated both items i and j

- Estimates with smaller supports are less reliable

- Use shrunk correlation coeff as a similarity measure:

$$s_{ij} = \frac{|U(i,j)| - 1}{|U(i,j)| - 1 + \lambda} \hat{\rho}_{ij}$$

- λ penalizes small supports:

$$|U(i,j)| << \lambda \rightarrow s_{ij} \rightarrow 0$$
$$|U(i,j)| >> \lambda \rightarrow s_{ij} \rightarrow \hat{\rho}_{ij}$$

# Improvements to common practice

- Use transformed similarities as interpolation coeff's

- E.g., by squaring we emphasize stronger relations:

$$\hat{r}_{ui} = b_{ui} + \frac{\sum\limits_{j \in S^k(i;u)} s_{ij}^2 (r_{uj} - b_{uj})}{\sum\limits_{j \in S^k(i;u)} s_{ij}^2}$$

- See A. Toscher, M. Jahrer and R. Legenstein, "Improved Neighborhood-Based Algorithms for Large-Scale Recommender Systems" for sigmodial transformations

- Shrink towards baseline when not enough neighborhood info is available:

$$\hat{r}_{ui} = b_{ui} + \frac{\sum\limits_{j \in S^k(i;u)} s_{ij}^2 (r_{uj} - b_{uj})}{\lambda + \sum\limits_{j \in S^k(i;u)} s_{ij}^2}$$

$$\left( \sum\limits_{j \in S^k(i;u)} s_{ij}^2 \ll \lambda \Rightarrow \hat{r}_{ui} \to b_{ui} \right)$$

# Similarities for binary data

- Often data is not ratings but binary. Example: ad clicks and conversions

- This requires other natural similarity measures

- Notation:

  $m_i$ - #users acting on i

  $m_{ij}$ - #users acting on both i and j

  $m$ - overall #users

(1) Jaccard similarity:

$$s_{ij} = \frac{m_{ij}}{m_i + m_j - m_{ij}}$$

- Shrink estimates:

$$s_{ij} = \frac{m_{ij}}{\alpha + m_i + m_j - m_{ij}}$$

# Similarities for binary data #2

- Under random sampling, expected i-j |intersection|:

$$m_i \cdot m_j / m$$

(expectation of a hypergeometric distribution)

$$\frac{\text{observed}}{\text{expected}} = \frac{m_{ij}}{(m_i \cdot m_j / m)}$$

- As usual, need to shrink:

$$s_{ij} = \frac{m_{ij}}{\alpha + (m_i \cdot m_j / m)}$$

# A user-user approach

- Dual to the so-far described item-item approach (with similar derivation)

- Predict a rating from ratings of similar users on the same item

- Building stones are user-user similarities $s_{uv}$:

$$\hat{r}_{ui} = b_{ui} + \frac{\displaystyle\sum_{v \in S^k(u;i)} s_{uv}(r_{vi} - b_{vi})}{\lambda + \displaystyle\sum_{v \in S^k(u;i)} s_{uv}}$$

- Item-item is commonly considered advantageous over user-user
  - when #items < #users: less item-item relations to store, more stable relations, more reliable estimation
  - Item-item meshes better with new users and explaining rec's

- In some cases user-user becomes more sensible:
  - When users are the more stable anchor of the system (e.g. items are web articles that quickly expire)
  - When #users < #items

# Part II: Matrix factorization techniques

# Latent factor models

# Basic matrix factorization model

users

items

| 1 |   | 3 |   |   | 5 |   |   | 5 |   | 4 |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 5 | 4 |   |   | 4 |   |   | 2 | 1 | 3 |
| 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |   |
|   | 2 | 4 |   | 5 |   |   | 4 |   |   | 2 |   |
|   |   | 4 | 3 | 4 | 2 |   |   |   |   | 2 | 5 |
| 1 |   | 3 |   | 3 |   |   | 2 |   |   | 4 |   |

~

items

| .1 | -.4 | .2 |
|-----|-----|-----|
| -.5 | .6 | .5 |
| -.2 | .3 | .5 |
| 1.1 | 2.1 | .3 |
| -.7 | 2.1 | -2 |
| -1 | .7 | .3 |

●

users

| 1.1 | -.2 | .3 | .5 | -2 | -.5 | .8 | -.4 | .3 | 1.4 | 2.4 | -.9 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| -.8 | .7 | .5 | 1.4 | .3 | -1 | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| 2.1 | -.4 | .6 | 1.7 | 2.4 | .9 | -.3 | .4 | .8 | .7 | -.6 | .1 |

~

**A rank-3 SVD approximation**

# Estimate unknown ratings as inner-products of factors:

users

| 1 |   | 3 |   |   | 5 |   |   | 5 |   | 4 |
|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 5 | ? |   | 4 |   |   | 2 | 1 | 3 |
| 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |
|   | 2 | 4 |   | 5 |   |   | 4 |   | 2 |   |
|   |   | 4 | 3 | 4 | 2 |   |   |   | 2 | 5 |
| 1 |   | 3 |   | 3 |   |   | 2 |   | 4 |   |

items  ~

$\sim$

| .1  | -.4 | .2 |
|-----|-----|-----|
| -.5 | .6  | .5 |
| -.2 | .3  | .5 |
| 1.1 | 2.1 | .3 |
| -.7 | 2.1 | -2 |
| -1  | .7  | .3 |

items  ~

●

users

| 1.1 | -.2 | .3 | .5  | -2  | -.5 | .8  | -.4 | .3  | 1.4 | 2.4 | -.9 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| -.8 | .7  | .5 | 1.4 | .3  | -1  | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| 2.1 | -.4 | .6 | 1.7 | 2.4 | .9  | -.3 | .4  | .8  | .7  | -.6 | .1  |

A rank-3 SVD approximation

# Estimate unknown ratings as inner-products of factors:

users

items

Top matrix (ratings):

| 1 |   | 3 |   |   | 5 |   | 5 |   | 4 |
|---|---|---|---|---|---|---|---|---|---|
|   |   | 5 | **?** |   | 4 |   |   | 2 | 1 | 3 |
| 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |
|   | 2 | 4 |   | 5 |   |   | 4 |   |   | 2 |
|   |   |   | 4 | 3 | 4 | 2 |   |   | 2 | 5 |
| 1 |   | 3 |   |   | 3 |   |   | 2 |   | 4 |

~

items

| .1 | -.4 | .2 |
|----|-----|----|
| **-.5** | **.6** | **.5** |
| -.2 | .3 | .5 |
| 1.1 | 2.1 | .3 |
| -.7 | 2.1 | -2 |
| -1 | .7 | .3 |

●

users

| 1.1 | -.2 | .3 | .5 | **-2** | -.5 | .8 | -.4 | .3 | 1.4 | 2.4 | -.9 |
|-----|-----|----|----|-------|-----|----|-----|----|-----|-----|-----|
| -.8 | .7 | .5 | 1.4 | **.3** | -1 | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| 2.1 | -.4 | .6 | 1.7 | **2.4** | .9 | -.3 | .4 | .8 | .7 | -.6 | .1 |

A rank-3 SVD approximation

# Estimate unknown ratings as inner-products of factors:

users

| 1 |   | 3 |   |   | 5 |   |   | 5 |   | 4 |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 5 | 4 | 2.4 | 4 |   |   |   | 2 | 1 | 3 |
| 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |   |
|   | 2 | 4 |   | 5 |   |   | 4 |   |   | 2 |   |
|   |   | 4 | 3 | 4 | 2 |   |   |   |   | 2 | 5 |
| 1 |   | 3 |   | 3 |   |   | 2 |   | 4 |   |   |

items

~

~

items

| .1  | -.4 | .2 |
|-----|-----|----|
| -.5 | .6  | .5 |
| -.2 | .3  | .5 |
| 1.1 | 2.1 | .3 |
| -.7 | 2.1 | -2 |
| -1  | .7  | .3 |

●

users

| 1.1 | -.2 | .3 | .5  | -2  | -.5 | .8  | -.4 | .3  | 1.4 | 2.4 | -.9 |
|-----|-----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| -.8 | .7  | .5 | 1.4 | .3  | -1  | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| 2.1 | -.4 | .6 | 1.7 | 2.4 | .9  | -.3 | .4  | .8  | .7  | -.6 | .1  |

A rank-3 SVD approximation

# Matrix factorization model

| 1 | | 3 | | | 5 | | | 5 | | 4 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 5 | 4 | | | 4 | | | 2 | 1 | 3 |
| 2 | 4 | | 1 | 2 | | 3 | | 4 | 3 | 5 | |
| | 2 | 4 | | 5 | | | 4 | | | 2 | |
| | | 4 | 3 | 4 | 2 | | | | | 2 | 5 |
| 1 | | 3 | | 3 | | | 2 | | | 4 | |

~

| .1 | -.4 | .2 |
|---|---|---|
| -.5 | .6 | .5 |
| -.2 | .3 | .5 |
| 1.1 | 2.1 | .3 |
| -.7 | 2.1 | -2 |
| -1 | .7 | .3 |

| 1.1 | -.2 | .3 | .5 | -2 | -.5 | .8 | -.4 | .3 | 1.4 | 2.4 | -.9 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| -.8 | .7 | .5 | 1.4 | .3 | -1 | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| 2.1 | -.4 | .6 | 1.7 | 2.4 | .9 | -.3 | .4 | .8 | .7 | -.6 | .1 |

Why can't use standard SVD R = M $\Sigma$ V?

- SVD isn't defined when entries are missing

- Regularization is necessary:
  Estimate as much signal as possible where there are sufficient data, without over fitting where data are scarce

# A regularized model

- Limit the values the factors can take

  Unlimited values will produce an overfit

  Reduce the optimization space

- User factors:

  Model a user u as a vector $p_u \sim N_k(\mu, \Sigma)$

- Item factors:

  Model an item i as a vector $q_i \sim N_k(\gamma, \Lambda)$

- Ratings:

  Measure "agreement" between u and i: $r_{ui} \sim N(p_u^T q_i, \varepsilon^2)$

- Simplifying assumptions:

  $\mu = \gamma = 0, \ \Sigma = \Lambda = \lambda I$

# Matrix factorization as a cost function

$$\text{Min}_{p_*,q_*} \sum_{\text{known } r_{ui}} \left( r_{ui} - p_u^T q_i \right)^2 + \underbrace{\lambda \left( \| p_u \|^2 + \| q_i \|^2 \right)}_{\text{regularization}}$$

$p_u$ - user-factor of u

$q_i$ - item-factor of i

$r_{ui}$ - rating by u for i

- Optimize by either stochastic gradient-descent or alternating least squares

# Stochastic gradient descent optimization

Perform till convergence:

- For each training example $r_{ui}$ :
  - Compute prediction error: $e_{ui} = r_{ui} - p_u^T q_i$
  - Update item factor: $q_i \leftarrow q_i + \gamma(p_u e_{ui} - \lambda q_i)$
  - Update user factor: $p_u \leftarrow p_u + \gamma(q_i e_{ui} - \lambda p_u)$

- Two constants to tune: $\gamma$ (step size) and $\lambda$ (regularization)
- Find values that minimize error on **validation** set

See, e.g., *Simon Funk, "Netflix Update: Try This at Home"*,
http://sifter.org/~simon/journal/20061211.html

**R**

| | Inglourious Basterds | 2012 | LOST 5 | Matrix | Monty Python |
|---|---|---|---|---|---|
| (person 1) | 1 | 4 | | 3 | |
| (person 2) | | | 4 | 4 | |
| (person 3) | 4 | | 2 | | 4 |

**P**

| | |
|---|---|
| 1.2 | -0.4 |
| 1.2 | 0.8 |
| 0.4 | -0.4 |

**Q**

| | | | | |
|---|---|---|---|---|
| 1.4 | 0.8 | -1.3 | -0.0 | 0.6 |
| -0.0 | 0.4 | -0.4 | 1.6 | 0.3 |

**R**

| | Inglourious Basterds | 2012 | LOST | Matrix | Monty Python |
|---|---|---|---|---|---|
| (person 1) | 1 | 4 | 3.3 | 3 | 2.4 |
| (person 2) | -0.5 | 3.5 | 4 | 4 | 1.5 |
| (person 3) | 4 | 4.9 | 2 | 1.1 | 4 |

**P**

| | |
|---|---|
| 1.4 | 1.1 |
| 0.9 | 1.9 |
| 2.5 | -0.3 |

**Q**

| | | | | |
|---|---|---|---|---|
| 1.5 | 2.1 | 1.0 | 0.7 | 1.6 |
| -1.0 | 0.8 | 1.6 | 1.8 | 0.0 |

# Matrix factorization with biases

$$\hat{r}_{ui} = \underbrace{\mu + b_u + b_i} + p_u^T q_i$$

Baseline predictors:
μ – global average
$b_u$ – bias of u
$b_i$ – bias of i

➔ Minimization problem:

$$\text{Min}_{p_*, p_*, b_*} \sum_{\text{known } r_{ui}} \left( r_{ui} - (\mu + b_u + b_i + p_u^T q_i) \right)^2 + \underbrace{\lambda \left( \|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2 \right)}_{\text{regularization}}$$

See, e.g., **A. Paterek, "Improving regularized singular value decomposition for collaborative filtering",** Proc. KDD Cup and Workshop 2007

# Ratings values vs rating occurrences

- There is information in the fact that user has rated a movie
- The user chose to see the movie
- The user chose to rate the movie
- The choice depends on many factors
- Can we use this information to improve the factorization?

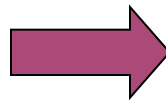# Ratings are not given at random!

**Distribution of ratings**



Netflix ratings     Yahoo! music ratings     Yahoo! survey answers

**B. Marlin et al., "Collaborative Filtering and the Missing at Random Assumption"** UAI 2007

# Which items users rate?

- A powerful source of information:
  Characterize **users** by **which** items they rated, rather than **how** they rated

- ➔ A dense binary representation of the data:

users

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 3 | | | 5 | | | 5 | | 4 |
| | | 5 | 4 | | | 4 | | | 2 | 1 | 3 |
| 2 | 4 | | 1 | 2 | | 3 | | 4 | 3 | 5 |
| | 2 | 4 | | 5 | | | 4 | | | 2 |
| | | 4 | 3 | 4 | 2 | | | | | 2 | 5 |
| 1 | | 3 | | 3 | | | 2 | | | 4 |

items

$$R = \left\{ r_{ui} \right\}_{u,i}$$

➡

users

| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 |
| 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |

items

$$B = \left\{ b_{ui} \right\}_{u,i}$$

# Factoring the binary view

- Describe both R and B using a factor model:

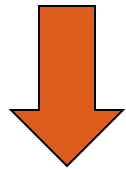$$r_{ui} = \mu + b_u + b_i + p_u^T q_i$$

$$b_{ui} = p_u^T x_i$$

- User factors are shared across models
- Each item $i$ is associated with two factor vectors: $q_i$ and $x_i$

# Factoring the binary view
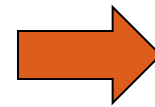
$$\forall i:\ b_{ui} = p_u^T x_i$$

$$X = \left( x_1, x_2, \cdots, x_n \right)$$

$$p_u = (XX^T)^{-1} X B_u$$

$$B_u = \left( b_{u1}, b_{u2}, \cdots, b_{un} \right)$$

$$p_u \propto X B_u = \sum_j b_{uj} x_j$$

User factor is indirectly defined by item factors:

sum of item factors for items rated by u
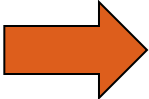
# Integrating ratings and binary views

So far:

- Each user $u$ is associated with a factor vector $p_u$

- Each item $i$ is associated with two factor vectors: $q_i$ and $x_i$

- The pure rating model:

$$r_{ui} = \mu + b_u + b_i + q_i^T p_u$$
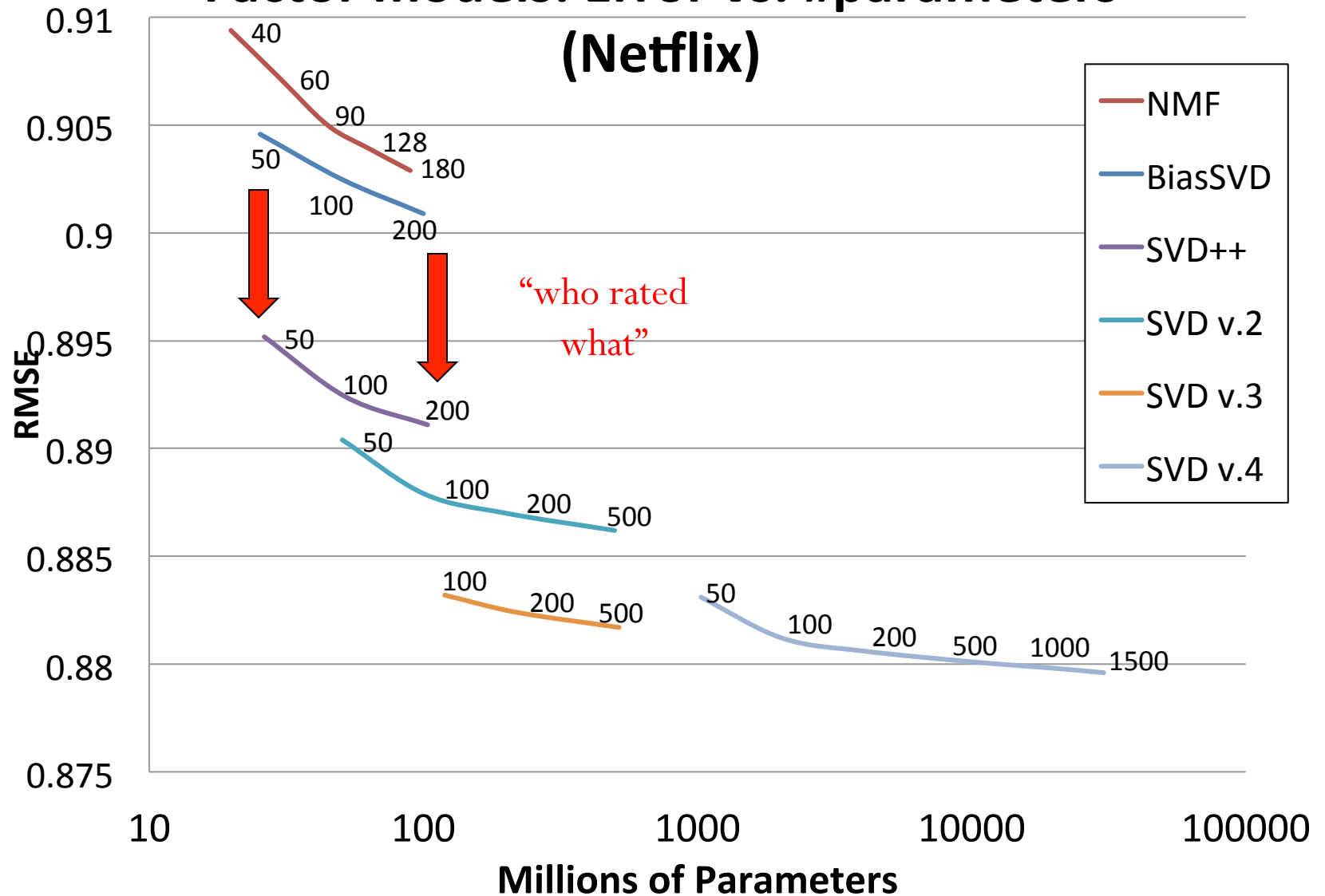
- The binary view of the user factors:

$$p_u \propto \sum_j b_{uj} x_j = \sum_{j\ rated\ by\ u} x_j$$

# Integrating ratings and binary views

So far:

- **Each user $u$ is associated with a factor vector $p_u$**

- Each item $i$ is associated with two factor vectors: $q_i$ and $x_i$

- The pure rating model:

$$r_{ui} = \mu + b_u + b_i + q_i^T p_u$$

- The binary view of the user factors:

$$p_u \propto \sum_j b_{uj} x_j = \sum_{j \; rated \; by \; u} x_j$$

$$r_{ui} = \mu + b_u + b_i + q_i^T \left( p_u + \sum_j b_{uj} x_j \right)$$

**"Factorization Meets the Neighborhood…",** KDD'08
R. Salakhutdinov and A. Mnih, **"Probabilistic Matrix Factorization"**, NIPS'07

Factor models: Error vs. #parameters (Netflix)

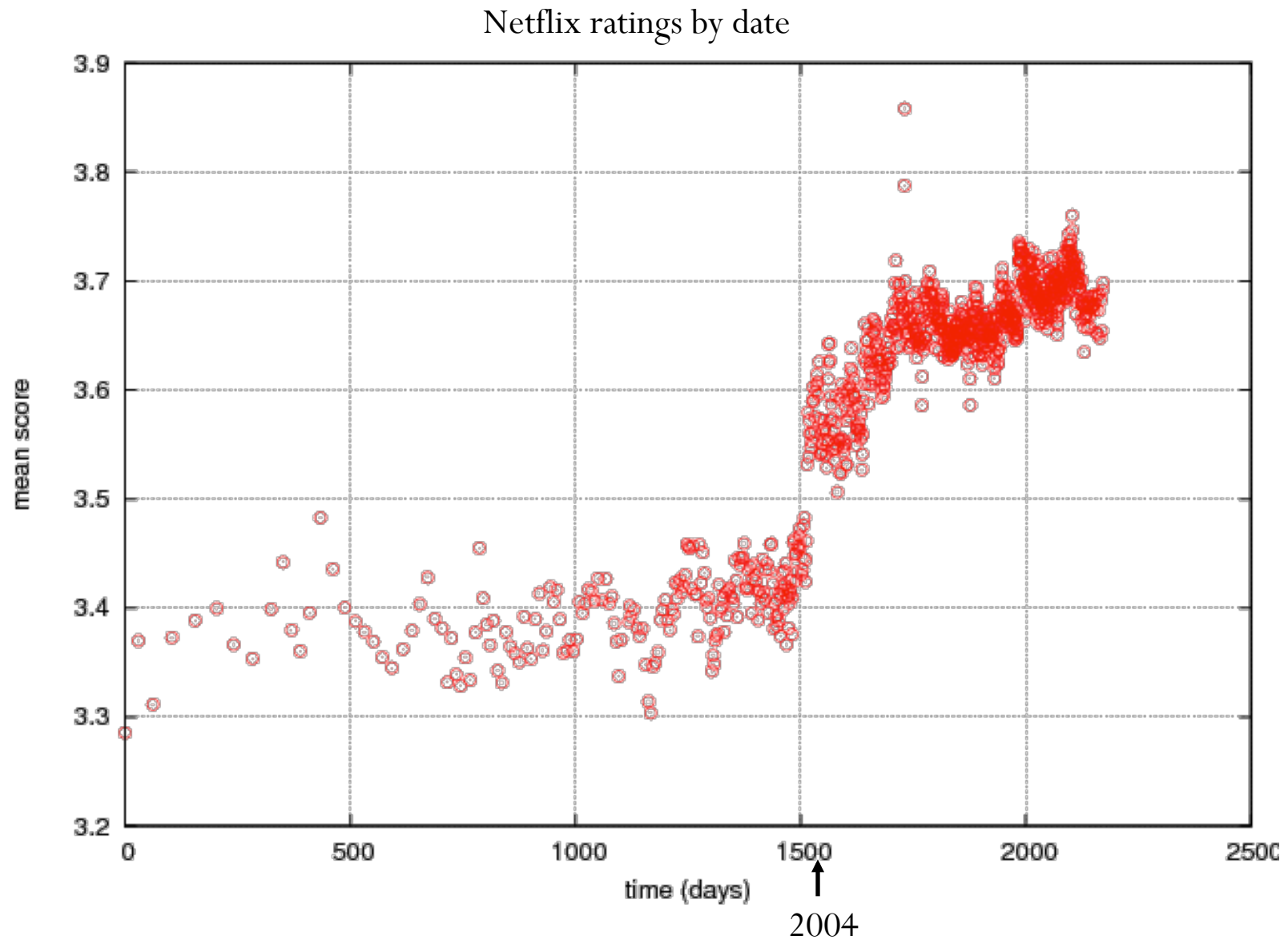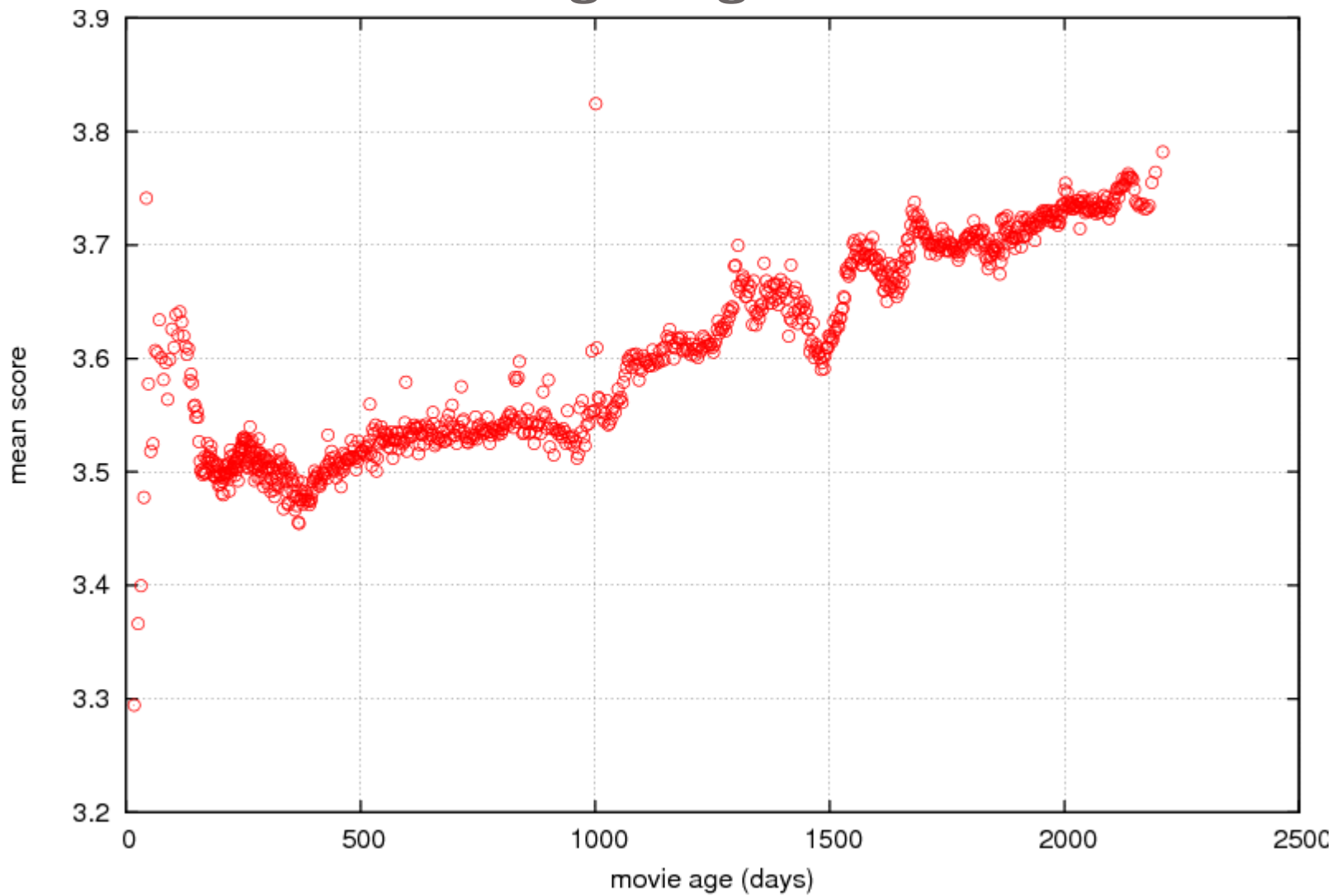# Temporal dynamics

Panta rhei

# Something Happened in Early 2004...

Netflix ratings by date
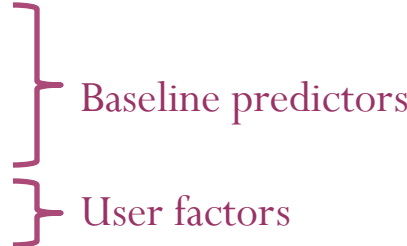
Are movies getting better with time?

# Multiple sources of temporal dynamics

- Item-side effects:
  - Product perception and popularity are constantly changing
  - Seasonal patterns influence items' popularity
- User-side effects:
  - Customers redefine their taste
  - Transient, short-term bias; anchoring
  - Drifting rating scale
  - Change of rater within household

# Temporal dynamics - challenges

- Multiple effects: Both items and users are changing over time
  → Scarce data per target

- Inter-related targets: Signal needs to be shared among users — foundation of **collaborative** filtering
  → cannot isolate multiple problems

→ Common "concept drift" methodologies won't hold.
E.g., **underweighting older instances is unappealing**

# Addressing temporal dynamics

- Factor model conveniently allows separately treating different aspects

- We observe changes in:
  1. Rating scale of individual users
  2. Popularity of individual items
  3. User preferences

  Baseline predictors

  User factors

$$r_{ui}(t) = \mu + b_u(t) + b_i(t) + q_i^T p_u(t)$$

# Parameterizing the model

$$r_{ui}(t) = \mu + b_u(t) + b_i(t) + q_i^T p_u(t)$$

- Use functional forms: $b_u(t) = f(u,t)$, $b_i(t) = g(i,t)$, $p_u(t) = h(u,t)$
- Need to find adequate $f()$, $g()$, $h()$
- General guidelines:
  - Items show slower temporal changes
  - Users exhibit frequent and sudden changes
  - Factors $-p_u(t)-$ are expensive to model
  - Gain flexibility by heavily parameterizing the functions

**"Collaborative Filtering with Temporal Dynamics", KDD'09**

Factor models: Error vs. #parameters