Global average: 1.1296

User average: 1.0651

Movie average: 1.0533
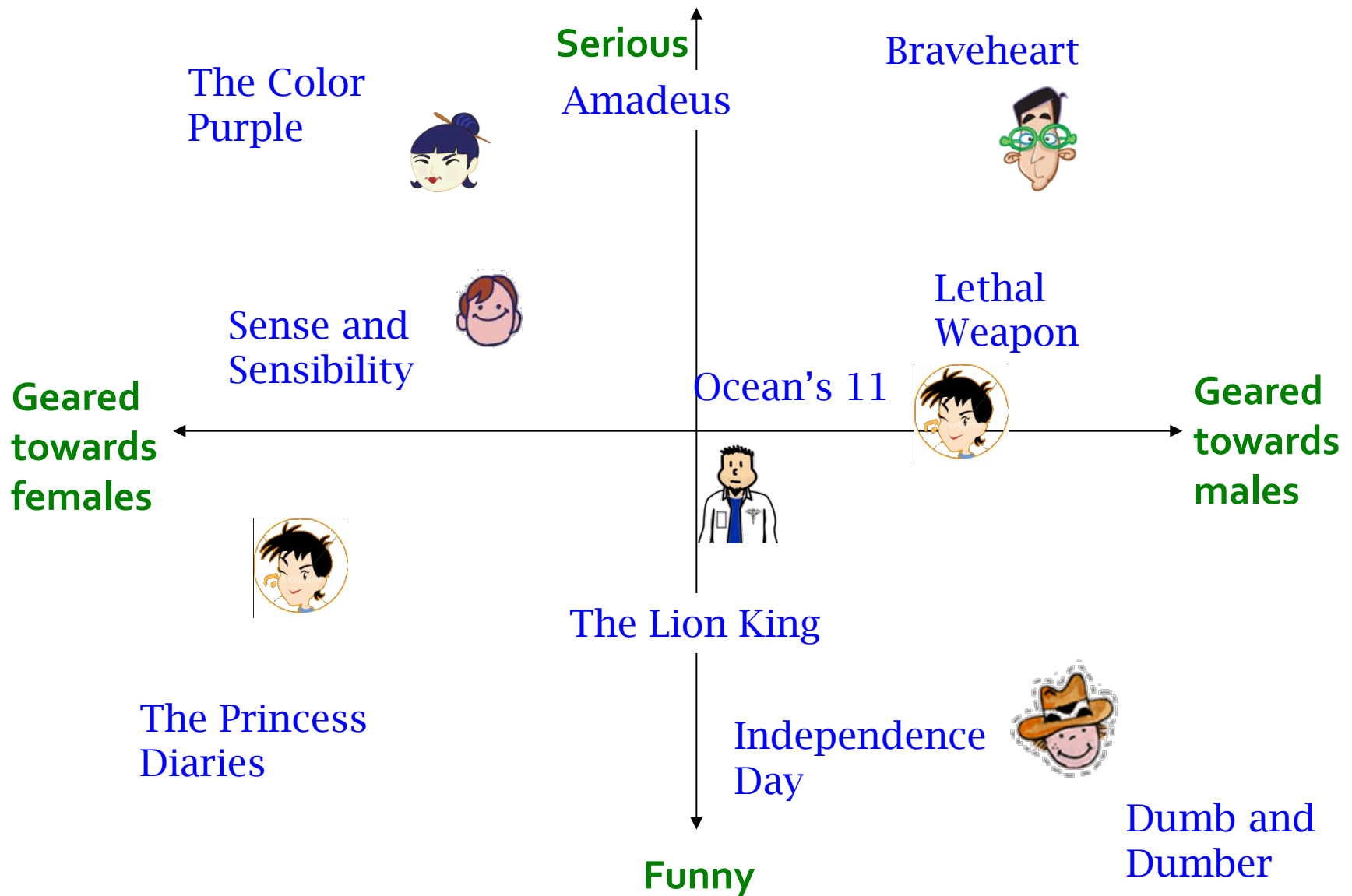
Netflix: 0.9514

Basic Collaborative filtering: 0.94

CF+Biases+learned weights: 0.91

Grand Prize: 0.8563

# Latent Factor Models (e.g., SVD)

# Latent Factor Models

- ■ **"SVD" on Netflix data: R ≈ Q · P$^T$**



- ■ For now let's assume we can approximate the rating matrix ***R*** as a product of "thin" ***Q · P***$^T$

  - ■ ***R*** has missing entries but let's ignore that for now!

    - ■ Basically, we will want the reconstruction error to be small on known ratings and we don't care about the values on the missing ones

# Ratings as Products of Factors

- **How to estimate the missing rating of user *x* for item *i*?**

$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

$q_i$ = row *i* of **Q**
$p_x$ = column *x* of $P^T$

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

≈

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

items — factors

**Q**

●

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 |

factors

$P^T$

# Ratings as Products of Factors

- **How to estimate the missing rating of user *x* for item *i*?**

$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

$q_i$ = row *i* of $Q$
$p_x$ = column *x* of $P^\mathrm{T}$



≈

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  /  factors

$Q$

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

•

users

factors

$P^\mathrm{T}$

| 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 |

# Ratings as Products of Factors

- **How to estimate the missing rating of user *x* for item *i*?**



$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

$q_i$ = row *i* of **Q**
$p_x$ = column *x* of $P^{\mathrm{T}}$

# Latent Factor Models

Serious

The Color Purple

Braveheart

Amadeus

Lethal Weapon

Sense and Sensibility

Ocean's 11

Geared towards females

Factor 1

Geared towards males

The Lion King

Factor 2

The Princess Diaries
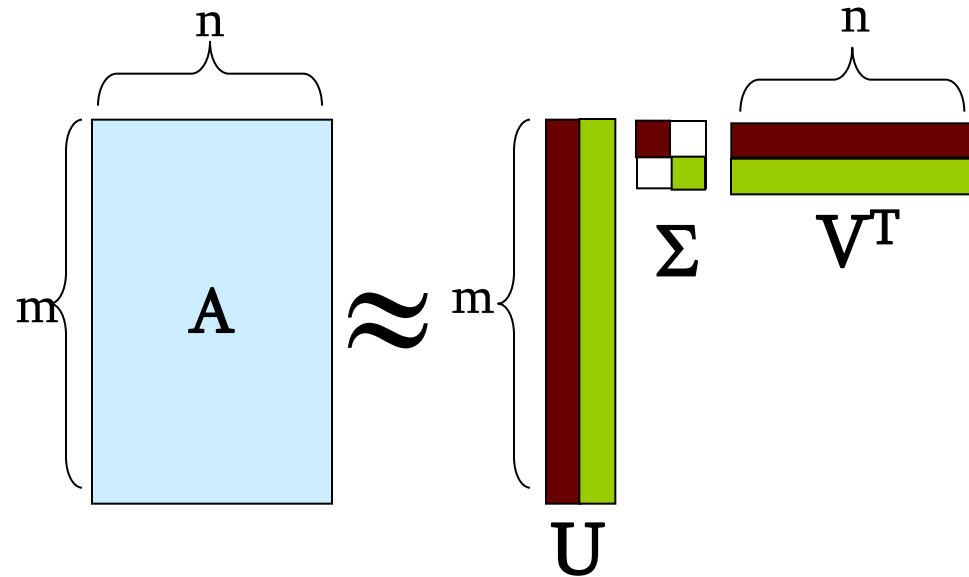
Independence Day

Dumb and Dumber

Funny

# Latent Factor Models

# Recap: SVD

- **Remember SVD:**
  - **A**: Input data matrix
  - **U**: Left singular vecs
  - **V**: Right singular vecs
  - $\Sigma$: Singular values



- **So in our case:**
  **"SVD" on Netflix data: $R \approx Q \cdot P^T$**
  $A = R, \quad Q = U, \quad P^T = \Sigma\, V^T$

$$\hat{r}_{xi} = q_i \cdot p_x$$

# SVD: More good stuff

- **We already know that SVD gives minimum reconstruction error** (Sum of Squared Errors):

$$\min_{U,V,\Sigma} \sum_{ij \in A} \left( A_{ij} - \left[ U\Sigma V^{\mathrm{T}} \right]_{ij} \right)^2$$

- **Note two things:**

  - **SSE** and **RMSE** are monotonically related:

    - $RMSE = \frac{1}{c}\sqrt{SSE}$   **Great news: SVD is minimizing RMSE**

  - **Complication:** The sum in SVD error term is over all entries (no-rating in interpreted as zero-rating). But our **R** has missing entries!

# Latent Factor Models



- **SVD isn't defined when entries are missing!**
- **Use specialized methods to find *P, Q***

$$\min_{P,Q} \sum_{(i,x)\in \mathrm{R}}\left(r_{xi} - q_i \cdot p_x\right)^2 \qquad \hat{r}_{xi} = q_i \cdot p_x$$

- **Note:**
  - We don't require cols of ***P, Q*** to be orthogonal/unit length
  - ***P, Q*** map users/movies to a latent space
  - The most popular model among Netflix contestants

# Finding the Latent Factors

# Latent Factor Models

- **Our goal is to find P and Q such tat:**

$$\min_{P,Q} \sum_{(i,x)\in R} \left(r_{xi} - q_i \cdot p_x\right)^2$$

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 |  |

≈

factors

Q

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

users

$P^{\mathrm{T}}$

factors

| 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 |

# Back to Our Problem

- **Want to minimize SSE for unseen test data**
- **Idea: Minimize SSE on training data**
  - Want large $k$ (# of factors) to capture all the signals
  - But, **SSE** on test data begins to rise for $k > 2$

- This is a classical example of **overfitting:**
  - With too much freedom (too many free parameters) the model starts fitting noise
    - That is it fits too well the training data and thus **not generalizing** well to unseen test data

# Dealing with Missing Entries

- **To solve overfitting we introduce regularization:**
  - Allow rich model where there are sufficient data
  - Shrink aggressively where data are scarce

$$\min_{P,Q} \underbrace{\sum_{training} (r_{xi} - q_i p_x)^2}_{\text{"error"}} + \left[ \underbrace{\lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2}_{\text{"length"}} \right]$$

$\lambda_1, \lambda_2$ … user set regularization parameters

**Note:** We do not care about the "raw" value of the objective function, but we care in P,Q that achieve the minimum of the objective

# The Effect of Regularization



serious

The Color Purple

Amadeus

Braveheart

Sense and Sensibility

Lethal Weapon

Ocean's 11

Geared towards females

**Factor 1**

Geared towards males

The Princess Diaries

The Lion King

Dumb and Dumber

**Factor 2**

Independence Day

funny

$$\min_{P,Q} \sum_{training} (r_{xi} - q_i p_x)^2 + \lambda \left[ \sum_x \|p_x\|^2 + \sum_i \|q_i\|^2 \right]$$

$\min_{factors}$ "error" + $\lambda$ "length"

# The Effect of Regularization

serious

Braveheart

The Color Purple

Amadeus

Lethal Weapon

Sense and Sensibility

Ocean's 11

**Geared towards females**

**Factor 1**

**Geared towards males**

The Princess Diaries

The Lion King

Dumb and Dumber

**Factor 2**

Independence Day

$$\min_{P,Q} \sum_{training} (r_{xi} - q_i p_x)^2 + \lambda \left[ \sum_x \|p_x\|^2 + \sum_i \|q_i\|^2 \right]$$

$\min_{factors}$ "error" + $\lambda$ "length"

funny

# The Effect of Regularization

serious

Braveheart

The Color Purple

Amadeus

Lethal Weapon

Sense and Sensibility

Ocean's 11

**Geared towards females**

**Factor 1**

**Geared towards males**

The Princess Diaries

The Lion King

Dumb and Dumber

**Factor 2**

Independence Day

$$\min_{P,Q} \sum_{training} (r_{xi} - q_i p_x)^2 + \lambda \left[ \sum_x \|p_x\|^2 + \sum_i \|q_i\|^2 \right]$$

$\min_{factors}$ "error" + $\lambda$ "length"

funny

# The Effect of Regularization

serious

Amadeus

Braveheart

The Color
Purple

Sense and
Sensibility

Lethal
Weapon

Ocean's 11

**Geared
towards
females**

**Factor 1**

**Geared
towards
males**

The Princess
Diaries

The Lion King

Dumb and
Dumber

**Factor 2**

Independence
Day

funny

$$\min_{P,Q} \sum_{training} (r_{xi} - q_i p_x)^2 + \lambda \left[ \sum_x \|p_x\|^2 + \sum_i \|q_i\|^2 \right]$$

$\min_{factors}$ "error" + $\lambda$ "length"

# Stochastic Gradient Descent

- **Want to find matrices _P_ and _Q_:**

$$\min_{P,Q} \sum_{training} (r_{xi} - q_i p_x)^2 + \left[ \lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]$$

- **Gradient decent:**

  - Initialize **_P_** and **_Q_** (using SVD, pretend missing ratings are 0)

  - Do gradient descent:

    How to compute gradient of a matrix?
    Compute gradient of every element independently!

    - $P \leftarrow P - \eta \cdot \nabla P$
    - $Q \leftarrow Q - \eta \cdot \nabla Q$
    - where $\nabla Q$ is gradient/derivative of matrix **_Q_**:
      $$\nabla Q = [\nabla q_{if}] \text{ and } \nabla q_{if} = \sum_{x,i} -2(r_{xi} - q_i p_x)p_{xf} + 2\lambda_2 q_{if}$$
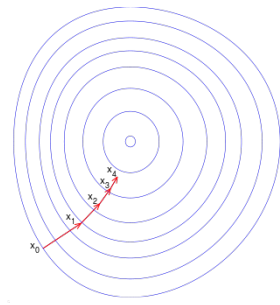      - Here $q_{if}$ is entry **_f_** of row **_q_i_** of matrix **_Q_**

  - **Observation: Computing gradients is slow!**

# Stochastic Gradient Descent

- **Gradient Descent (GD) vs. Stochastic GD**
    - **Observation:** $\nabla Q = [\nabla q_{if}]$ where
$$\nabla q_{if} = \sum_{x,i} -2(r_{xi} - q_{if}p_{xf})p_{xf} + 2\lambda q_{if} = \sum_{x,i} \nabla \boldsymbol{Q}(r_{xi})$$
        - Here $\boldsymbol{q_{if}}$ is entry $\boldsymbol{f}$ of row $\boldsymbol{q_i}$ of matrix $\boldsymbol{Q}$
    - $\boldsymbol{Q} = \boldsymbol{Q} - \square \quad \boldsymbol{Q} = \boldsymbol{Q} - \eta\left[\sum_{x,i} \nabla \boldsymbol{Q}(r_{xi})\right]$
    - **Idea:** Instead of evaluating gradient over all ratings evaluate it for each individual rating and make a step
- **GD:** $\boldsymbol{Q} \leftarrow \boldsymbol{Q} - \eta\left[\sum_{r_{xi}} \nabla \boldsymbol{Q}(r_{xi})\right]$
- **SGD:** $\boldsymbol{Q} \leftarrow \boldsymbol{Q} - \mu\nabla \boldsymbol{Q}(r_{xi})$
    - **Faster convergence!**
        - Need more steps but each step is computed much faster

# GD vs SGD on linear regression

$$h_\theta(x) = \sum_{j=0}^{n} \theta_j x_j$$

## GD

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

Repeat{

$$\theta_j := \theta_j - \eta \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

$$(\text{for every } j = 0, \dots, n)$$

}

## SGD

New cost function

$$cost\left(\theta, (x^{(i)}, y^{(i)})\right) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^{m} cost\left(\theta, (x^{(i)}, y^{(i)})\right)$$

1. Randomly shuffle dataset;

2. Repeat{
   for $i = 1, \dots, m${

   $$\theta_j := \theta_j - \eta\,(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$
   $$(\text{for every } j = 0, \dots, n)$$
   }
   }

# SGD vs. GD

- **Convergence of <span style="color:red">GD</span> vs. <span style="color:blue">SGD</span>**



**Value of the objective function**

**Iteration/step**

**GD** improves the value of the objective function at every step.
**SGD** improves the value but in a "noisy" way.
**GD** takes fewer steps to converge but each step takes much longer to compute.
In practice, **SGD** is much faster!

# Stochastic Gradient Descent

## ■ Stochastic gradient decent:

- Initialize **$P$** and **$Q$** (using SVD, pretend missing ratings are 0)

- Then iterate over the ratings (multiple times if necessary) and update factors:

  **For each $r_{xi}$:**

  - $\varepsilon_{xi} = 2(r_{xi} - q_i \cdot p_x)$      (derivative of the "error")
  - $q_i \leftarrow q_i - \mu_1 (\varepsilon_{xi} p_x - \lambda_2 q_i)$      (update equation)
  - $p_x \leftarrow p_x - \mu_2 (\varepsilon_{xi} q_i - \lambda_1 p_x)$      (update equation)

  $\mu$ … learning rate

## ■ 2 for loops:

- For until convergence:
  - For each **$r_{xi}$**
    - Compute gradient, do a "step"

# Extending Latent Factor Model to Include Biases

# Modeling Biases and Interactions

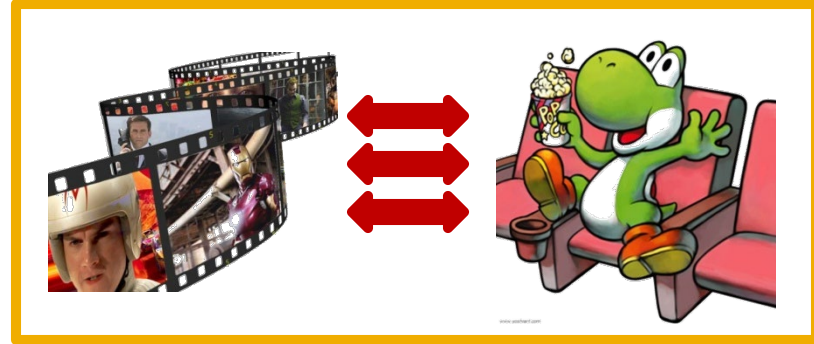**user bias**                    **movie bias**                    **user-movie interaction**



**Baseline predictor**

- Separates users and movies
- Benefits from insights into user's behavior
- Among the main practical contributions of the competition

**User-Movie interaction**

- Characterizes the matching between users and movies
- Attracts most research in the field
- Benefits from algorithmic and mathematical innovations

- $\mu$ = overall mean rating
- $b_x$ = bias of user $x$
- $b_i$ = bias of movie $i$

# Baseline Predictor

- We have expectations on the rating by user $x$ of movie $i$, even without estimating $x$'s attitude towards movies like $i$



- Rating scale of user $x$
- Values of other ratings user gave recently (day-specific mood, anchoring, multi-user accounts)

- (Recent) popularity of movie $i$
- Selection bias; related to number of ratings user gave on the same day ("frequency")

# Putting It All Together

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

Overall mean rating    Bias for user *x*    Bias for movie *i*    User-Movie interaction

- **Example:**
  - Mean rating: $\mu$ **= 3.7**
  - You are a critical reviewer: your ratings are 1 star lower than the mean: **$b_x$ = -1**
  - Star Wars gets a mean rating of *0.5* higher than average movie: **$b_i$ = + 0.5**
  - Predicted rating for you on Star Wars:
    **= 3.7 - 1 + 0.5 = 3.2**

# Fitting the New Model

- **Solve:**

$$\min_{Q,P} \sum_{(x,i)\in R} \left( r_{xi} - (\mu + b_x + b_i + q_i\, p_x) \right)^2$$

goodness of fit

$$+ \left( \lambda_1 \sum_i \|q_i\|^2 + \lambda_2 \sum_x \|p_x\|^2 + \lambda_3 \sum_x \|b_x\|^2 + \lambda_4 \sum_i \|b_i\|^2 \right)$$
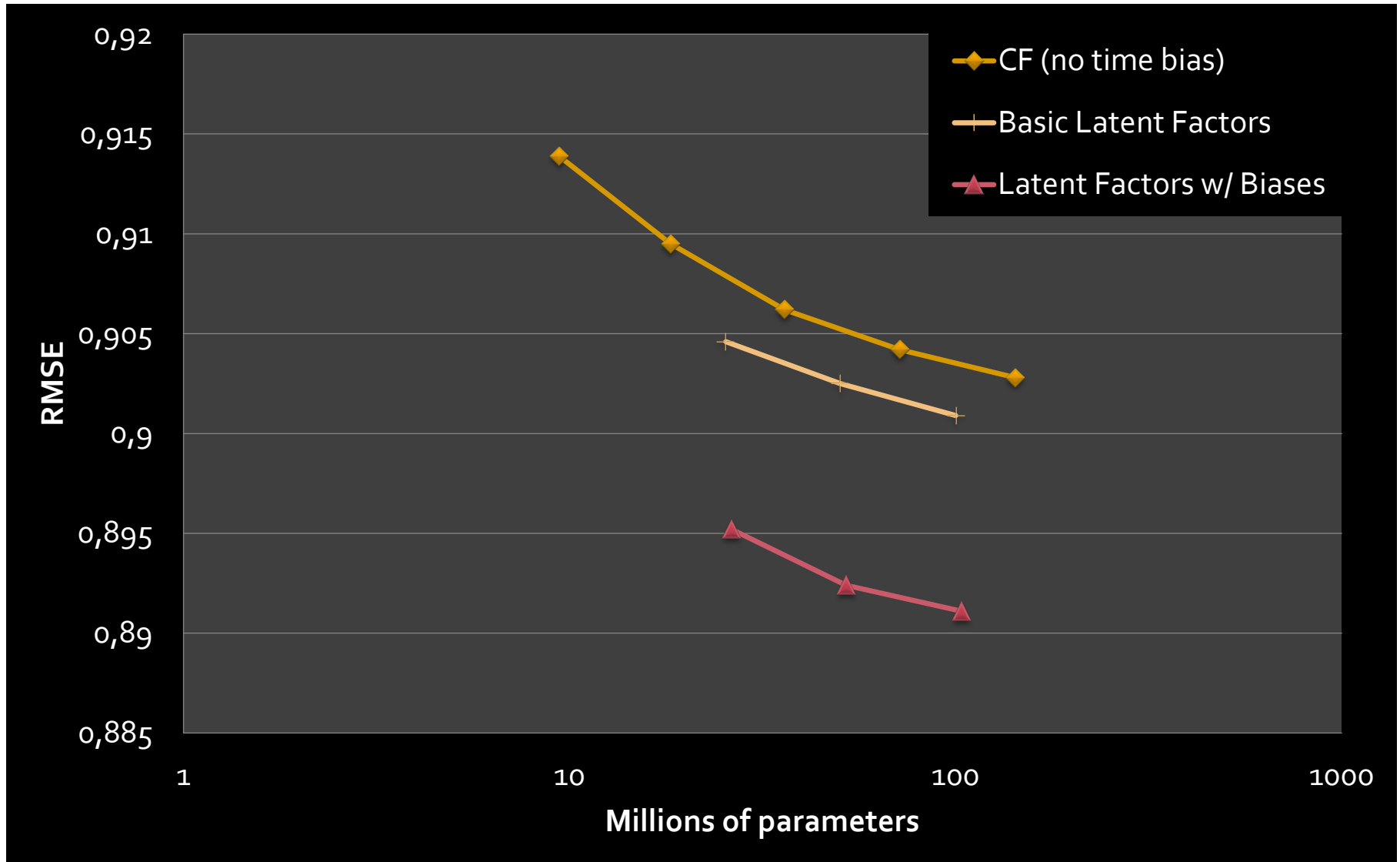
regularization

$\lambda$ is selected via grid-search on a validation set

- **Stochastic gradient decent to find parameters**
  - **Note:** Both biases $b_x$, $b_i$ as well as interactions $q_i$, $p_x$ are treated as parameters (we estimate them)

# Performance of Various Methods

Global average: 1.1296

User average: 1.0651

Movie average: 1.0533

Netflix: 0.9514

Basic Collaborative filtering: 0.94

Collaborative filtering++: 0.91
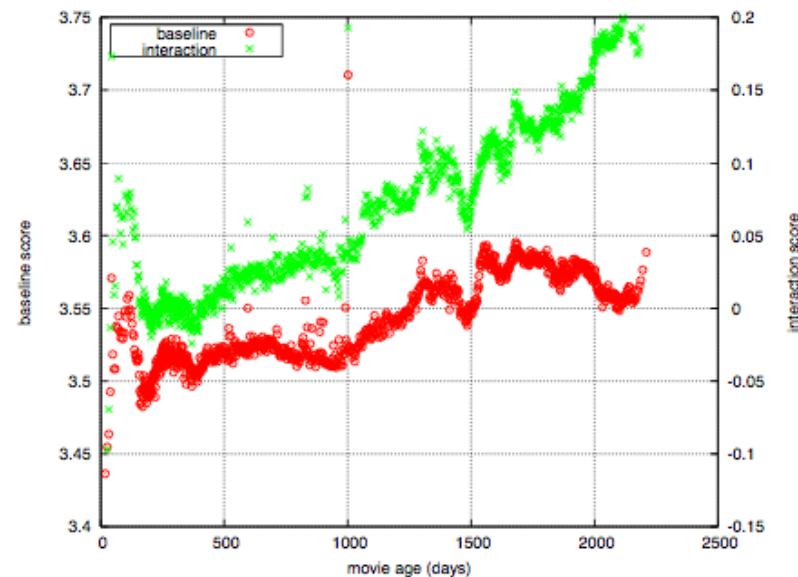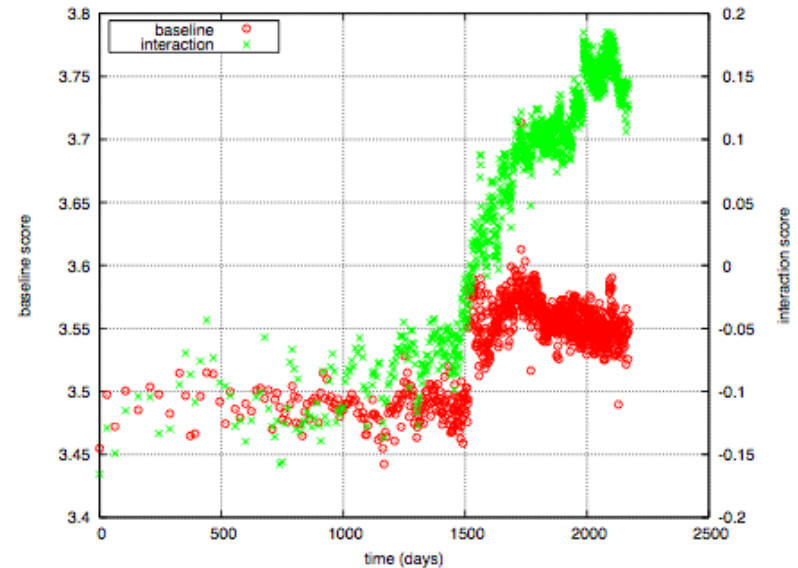
Latent factors: 0.90

**Latent factors+Biases: 0.89**

Grand Prize: 0.8563

# The Netflix Challenge: 2006-09

# Temporal Biases Of Users

- **Sudden rise in the average movie rating** (early 2004)
  - Improvements in Netflix
  - GUI improvements
  - Meaning of rating changed
- **Movie age**
  - Users prefer new movies without any reasons
  - Older movies are just inherently better than newer ones

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09

# Temporal Biases & Factors

- **Original model:**

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

- **Add time dependence to biases:**

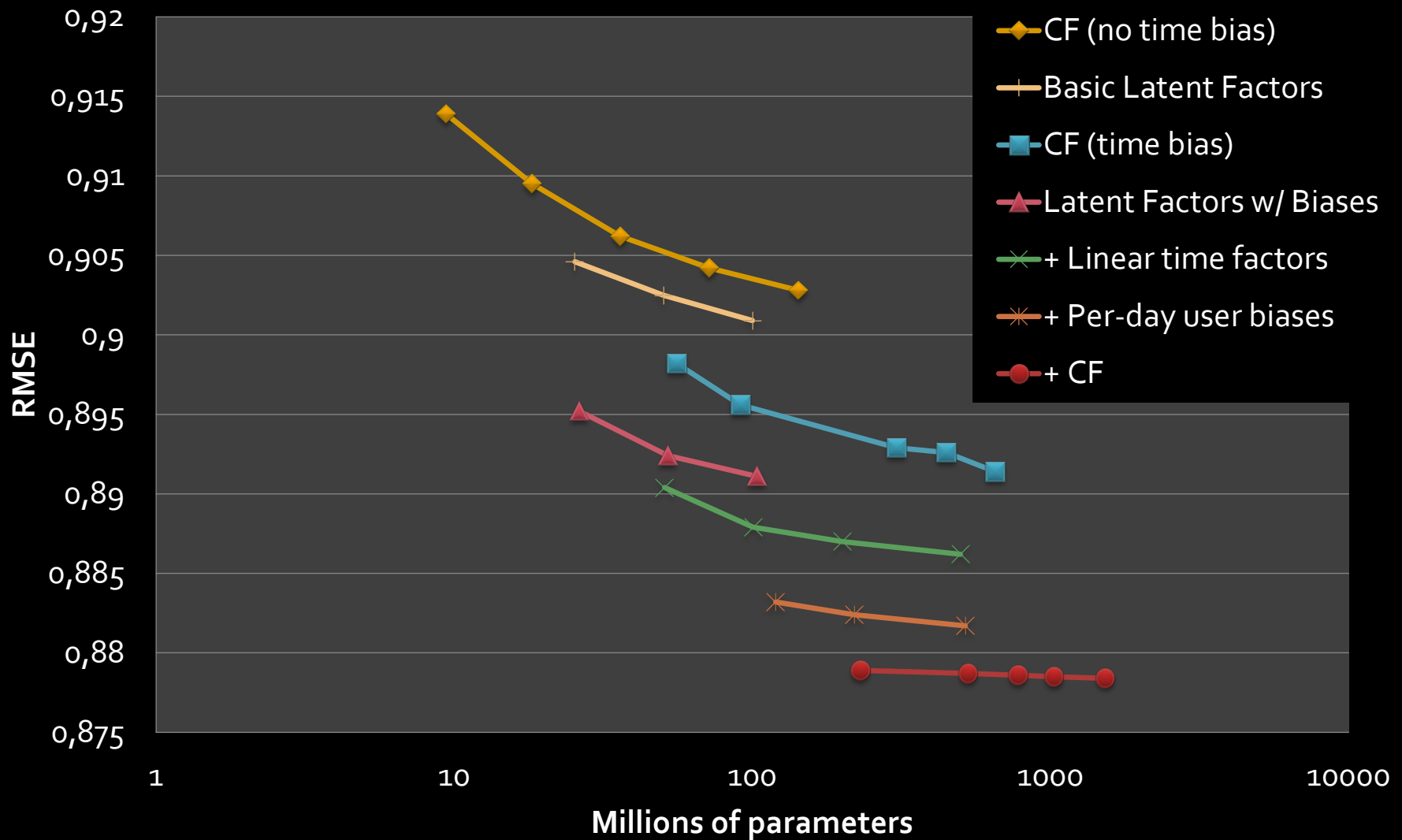$$r_{xi} = \mu + b_x(t) + b_i(t) + q_i \cdot p_x$$

  - Make parameters $b_x$ and $b_i$ to depend on time

  - **(1)** Parameterize time-dependence by linear trends
    **(2)** Each bin corresponds to 10 consecutive weeks

$$b_i(t) = b_i + b_{i,\text{Bin}(t)}$$

- **Add temporal dependence to factors**

  - $p_x(t)$… user preference vector on day $t$

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09

# Adding Temporal Effects

# The big picture
## Solution of BellKor's Pragmatic Chaos

All developed CF models

BRISMF  SBRAMF  SVD-Time  Split RBM  BK3  BK2
MF1  NSVDD  RBM day  GE  FRBM  3K4  3K1  BK5-SVD++
Movie KNN V.  Baseline 1/2/3  DRBM  SVD++  NSVD2  GTE
KNN+time  Integrated M.  RBM  MF2
NSVD1
SVD-AUF  Movie KNN  CTD/MTD  SVDNN
User KNN  Classif. Mode  KNN 1...5  Asym. 1/2/3

Latent User and Movie Features

approx. 500 predictors

Probe Blending

Probe Blending

200 blends

30 blends

Linear Blend    10.09 % improvement

Michael Jahrer / Andreas Töscher  –  Team BigChaos  –  September 21, 2009

# Standing on June 26[th] 2009



**June 26[th] submission triggers 30-day "last call"**

# NETFLIX

Home    Rules    Leaderboard    Update    Download

# Leaderboard

Showing Test Score. Click here to show quiz score

Display top [ 20 ▲▼ ] leaders.

| Rank | Team Name | Best Test Score | % Improvement | Best Submit Time |
|---|---|---|---|---|
| **Grand Prize** – RMSE = 0.8567 – Winning Team: BellKor's Pragmatic Chaos | | | | |
| 1 | BellKor's Pragmatic Chaos | 0.8567 | 10.06 | 2009-07-26 18:18:28 |
| 2 | The Ensemble | 0.8567 | 10.06 | 2009-07-26 18:38:22 |
| 3 | Grand Prize Team | 0.8582 | 9.90 | 2009-07-10 21:24:40 |
| 4 | Opera Solutions and Vandelay United | 0.8588 | 9.84 | 2009-07-10 01:12:31 |
| 5 | Vandelay Industries ! | 0.8591 | 9.81 | 2009-07-10 00:32:20 |
| 6 | PragmaticTheory | 0.8594 | 9.77 | 2009-06-24 12:06:56 |
| 7 | BellKor in BigChaos | 0.8601 | 9.70 | 2009-05-13 08:14:09 |
| 8 | Dace_ | 0.8612 | 9.59 | 2009-07-24 17:18:43 |
| 9 | Feeds2 | 0.8622 | 9.48 | 2009-07-12 13:11:51 |
| 10 | BigChaos | 0.8623 | 9.47 | 2009-04-07 12:33:59 |
| 11 | Opera Solutions | 0.8623 | 9.47 | 2009-07-24 00:34:07 |
| 12 | BellKor | 0.8624 | 9.46 | 2009-07-26 17:19:11 |
| **Progress Prize 2008** – RMSE = 0.8627 – Winning Team: BellKor in BigChaos | | | | |
| 13 | xiangliang | 0.8642 | 9.27 | 2009-07-15 14:53:22 |
| 14 | Gravity | 0.8643 | 9.26 | 2009-04-22 18:31:32 |
| 15 | Ces | 0.8651 | 9.18 | 2009-06-21 19:24:53 |
| 16 | Invisible Ideas | 0.8653 | 9.15 | 2009-07-15 15:53:04 |
| 17 | Just a guy in a garage | 0.8662 | 9.06 | 2009-05-24 10:02:54 |
| 18 | J Dennis Su | 0.8666 | 9.02 | 2009-03-07 17:16:17 |
| 19 | Craig Carmichael | 0.8666 | 9.02 | 2009-07-25 16:00:54 |
| 20 | acmehill | 0.8668 | 9.00 | 2009-03-21 16:20:50 |
| **Progress Prize 2007** – RMSE = 0.8723 – Winning Team: KorBell | | | | |

# Million $ Awarded Sept 21ˢᵗ 2009

# Acknowledgments

- Some slides and plots borrowed from Yehuda Koren, Robert Bell and Padhraic Smyth
- **Further reading:**
  - Y. Koren, Collaborative filtering with temporal dynamics, KDD '09
  - http://www2.research.att.com/~volinsky/netflix/bpc.html
  - http://www.the-ensemble.com/