

Collaborative Filtering: Models and Algorithms

Andrea Montanari

Jose Bento, Yash Deshpande, Adel Javanmard, Raghunandan Keshavan, Sewoong Oh,
Stratis Ioannidis, Nadia Fawaz, Amy Zhang
Stanford University, Technicolor

September 15, 2012

Problem statement

Given data on the activity of a set of users, provide personalized recommendations to users X, Y, Z, \dots

Example

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to [see all recommendations](#).

Page 1 of 35



[Probabilistic Graphical Models...](#) (Hardcover) by
Daphne Koller
★★★★★ (4) \$74.90
[Fix this recommendation](#)



[Elements of Information Theory...](#) (Hardcover) by
Thomas M. Cover
★★★★★ (27) \$80.51
[Fix this recommendation](#)



[Networks: An Introduction](#) (Hardcover) by Mark Newman
★★★★★ (3) \$70.10
[Fix this recommendation](#)



[The Elements of Statistical Learning...](#) (Hardcover) by Trevor Hastie
★★★★★ (45) \$62.32
[Fix this recommendation](#)



[Bayesian Data Analysis, Second...](#) (Hardcover) by Andrew Gelman
★★★★★ (16) \$62.41
[Fix this recommendation](#)

Andrea, Welcome to Your Amazon.com (If you're not Andrea Montanari, [click here](#).)

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to [see all recommendations](#).

Page 5 of 35 (Start over)



[Large-Scale Inference: Empirical...](#) (Hardcover) by Bradley Efron
★★★★★ (2) \$59.31
[Fix this recommendation](#)



[Sesame Street - Fiesta! DVD](#) ~
Celia Cruz
★★★★★ (58) \$8.49
[Fix this recommendation](#)



[Introducing Monte Carlo M... \(Paperback\)](#) by Christian P. Robert
\$52.25
[Fix this recommendation](#)

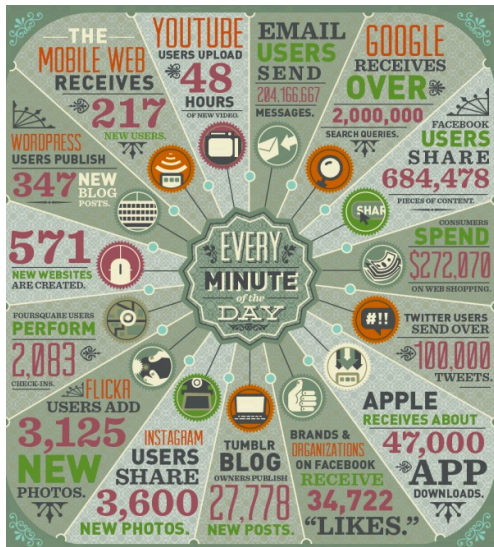


[Maple Teethers](#)
★★★★★ (9) \$13.45
[Fix this recommendation](#)



[Data Manipulation with R \(Use R!\) \(Paperback\)](#) by Phil Spector
★★★★★ (15) \$46.83
[Fix this recommendation](#)

An obviously useful technology



An obviously useful technology

The screenshot shows the Amazon.com search results for the author 'montanari'. The page layout includes the Amazon logo, navigation links, a search bar with the query 'montanari', and a list of search results. On the left, there are category filters for 'Books' and 'Kindle Store'. A 'Listmania!' section highlights a mystery list by RosieVeit. The search results list four books: 'Echo Man', 'The Echo Man: A Novel of Suspense', 'Killing Room', and 'Violet Hour', each with its price, number of offers, and user ratings.

amazon Your Amazon.com Today's Deals Gift Cards Help

The All-New **kindle fire HD**

Shop by Department Hello, Sign In Your Account

Department

Books

- Contemporary Literature & Fiction
- Action & Adventure Fiction
- Thrillers
- Mysteries
- Police Procedurals

Kindle Store

- Contemporary Fiction
- Fiction
- Mystery & Thrillers
- Police Procedurals
- Suspense Thrillers

+ See All 14 Departments

Shipping Option (books.usa)
Free Super Saver Shipping

Listmania!

RosieVeit's Mystery List: A list by RosieVeit

"montanari"

Showing 1 - 16 of 1,755 Results

Echo Man by Richard Montanari (Sep 1, 2011)

\$6.99 new (15 offers)
\$0.44 used (37 offers)

★★★★☆ (9)
Books: See all 743 items

The Echo Man: A Novel of Suspense (Jessica Balzano and Kevin Byrne) by Richard Montanari (Sep 19, 2011)

\$4.99 Kindle Edition
Auto-delivered wirelessly

★★★★☆ (9)
Kindle Store: See all 33 items

Killing Room by Richard Montanari (Feb 1, 2012)

\$14.72 used (9 offers)

★★★★☆ (2)
Sell this back for an Amazon.com Gift Card
Books: See all 743 items

Violet Hour by Richard Montanari (Dec 1, 2010)

\$0.02 used (53 offers)

★★★★☆ (21)
Books: See all 743 items

Outline

- 1 A model
- 2 Algorithms and accuracy
- 3 Challenge #1: Privacy
- 4 Challenge #2: Interactivity
- 5 Conclusion

A model

Setting

Users : $i \in \{1, 2, \dots, m\}$

Movies : $j \in \{1, 2, \dots, n\}$

When user i watches movie j , she enters her rating R_{ij} .

Want to predict ratings for missing pairs.

Setting

Users : $i \in \{1, 2, \dots, m\}$

Movies : $j \in \{1, 2, \dots, n\}$

When user i watches movie j , she enters her rating R_{ij} .

Want to predict ratings for missing pairs.

Linear regression model

Movie j

$$v_j = (\text{genre; main actor; supporting actor; year; } \dots) \in \mathbb{R}^r$$

$$R_{ij} \sim c_i + \langle u_i, v_j \rangle + \varepsilon_{ij}$$

Want: User parameters vector u_i

Linear regression model

Movie j

$$v_j = (\text{genre; main actor; supporting actor; year; } \dots) \in \mathbb{R}^r$$

$$R_{ij} \sim c_i + \langle u_i, v_j \rangle + \epsilon_{ij}$$

Want: User parameters vector u_i

Linear regression model

Movie j

$$v_j = (\text{genre; main actor; supporting actor; year; } \dots) \in \mathbb{R}^r$$

$$R_{ij} \sim c_i + \langle u_i, v_j \rangle + \epsilon_{ij}$$

Want: User parameters vector u_i

Linear regression model

Movie j

$$v_j = (\text{genre; main actor; supporting actor; year; } \dots) \in \mathbb{R}^r$$

$$R_{ij} \sim \cancel{\mu_i} + \langle u_i, v_j \rangle + \epsilon_{ij}$$

Want: User parameters vector u_i

Least squares

$$u_i = \arg \min_{x_i \in \mathbb{R}^r} \left\{ \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle x_i, v_j \rangle \right)^2 \right\}$$

Ridge regression

$$u_i = \arg \min_{x_i \in \mathbb{R}^r} \left\{ \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle x_i, v_j \rangle \right)^2 + \lambda \|x_i\|^2 \right\}$$

- ▶ How do we construct the v_j 's?
- ▶ Ad hoc definitions are not suited to recommendation!

Ridge regression

$$u_i = \arg \min_{x_i \in \mathbb{R}^r} \left\{ \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle x_i, v_j \rangle \right)^2 + \lambda \|x_i\|^2 \right\}$$

- ▶ How do we construct the v_j 's?
- ▶ Ad hoc definitions are not suited to recommendation!

If I knew the users' feature vectors

$$R_{ij} \sim \langle u_i, v_j \rangle + \varepsilon_{ij}$$

$$v_j = \arg \min_{y_j \in \mathbb{R}^r} \left\{ \sum_{i \in \text{Watched}(j)} \left(R_{ij} - \langle u_i, y_j \rangle \right)^2 + \lambda \|y_j\|^2 \right\}$$

If I knew the users' feature vectors

$$R_{ij} \sim \langle \mathbf{u}_i, \mathbf{v}_j \rangle + \varepsilon_{ij}$$

$$\mathbf{v}_j = \arg \min_{\mathbf{y}_j \in \mathbb{R}^r} \left\{ \sum_{i \in \text{Watched}(j)} \left(R_{ij} - \langle \mathbf{u}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \|\mathbf{y}_j\|^2 \right\}$$

Everything together

$$\begin{aligned}u_i &= \arg \min_{x_i \in \mathbb{R}^r} \left\{ \sum_{j \in \text{WatchedBy}(i)} (R_{ij} - \langle x_i, y_j \rangle)^2 + \lambda \|x_i\|^2 \right\} \\v_j &= \arg \min_{y_j \in \mathbb{R}^r} \left\{ \sum_{i \in \text{Watched}(j)} (R_{ij} - \langle u_i, y_j \rangle)^2 + \lambda \|y_j\|^2 \right\}\end{aligned}$$

Minimize ($E = \text{Watched}$)

$$F(X, Y) = \sum_{(i,j) \in E} (R_{ij} - \langle x_i, y_j \rangle)^2 + \lambda \sum_{i=1}^m \|x_i\|_2^2 + \lambda \sum_{j=1}^n \|y_j\|_2^2$$

Everything together

$$\begin{aligned} \mathbf{u}_i &= \arg \min_{\mathbf{x}_i \in \mathbb{R}^r} \left\{ \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \|\mathbf{x}_i\|^2 \right\} \\ \mathbf{v}_j &= \arg \min_{\mathbf{y}_j \in \mathbb{R}^r} \left\{ \sum_{i \in \text{Watched}(j)} \left(R_{ij} - \langle \mathbf{u}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \|\mathbf{y}_j\|^2 \right\} \end{aligned}$$

Minimize ($E = \text{Watched}$)

$$F(X, Y) = \sum_{(i,j) \in E} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 + \lambda \sum_{j=1}^n \|\mathbf{y}_j\|_2^2$$

Objective function

Minimize ($E = \text{Watched}$)

$$\begin{aligned} F(X, Y) &= \sum_{(i,j) \in E} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 + \lambda \sum_{j=1}^n \|\mathbf{y}_j\|_2^2 \\ &\equiv \|\mathcal{P}_E(R - XY^T)\|_F^2 + \lambda \|X\|_F^2 + \lambda \|Y\|_F^2 \end{aligned}$$

$$\mathcal{P}_E(A) = \begin{cases} A_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise} \end{cases}$$

$$X^T = \left[\mathbf{x}_1 \mid \mathbf{x}_2 \mid \cdots \mid \mathbf{x}_m \right]$$

$$Y^T = \left[\mathbf{y}_1 \mid \mathbf{y}_2 \mid \cdots \mid \mathbf{y}_n \right]$$

Objective function

Minimize ($E = \text{Watched}$)

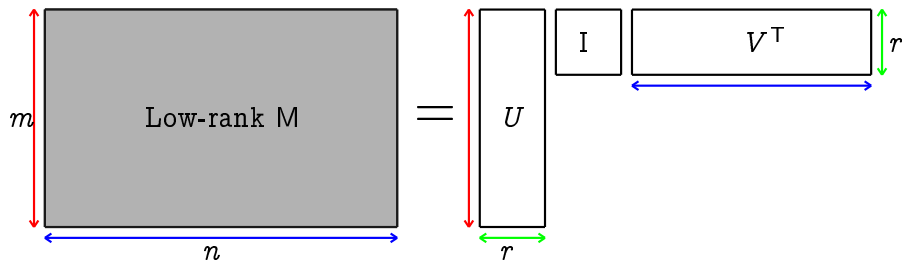
$$\begin{aligned} F(X, Y) &= \sum_{(i,j) \in E} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 + \lambda \sum_{j=1}^n \|\mathbf{y}_j\|_2^2 \\ &\equiv \|\mathcal{P}_E(\mathbf{R} - \mathbf{X}\mathbf{Y}^\top)\|_F^2 + \lambda \|\mathbf{X}\|_F^2 + \lambda \|\mathbf{Y}\|_F^2 \end{aligned}$$

$$\mathcal{P}_E(A) = \begin{cases} A_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{X}^\top = \left[\mathbf{x}_1 \mid \mathbf{x}_2 \mid \cdots \mid \mathbf{x}_m \right]$$

$$\mathbf{Y}^\top = \left[\mathbf{y}_1 \mid \mathbf{y}_2 \mid \cdots \mid \mathbf{y}_n \right]$$

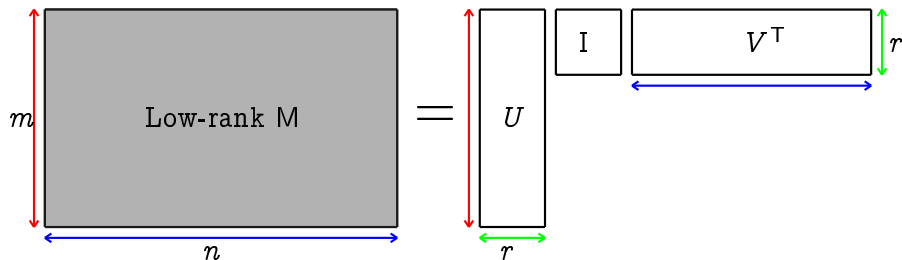
Low rank structure



1. Low-rank matrix M
2. $R = M + Z$
3. Observed subset E

$$\mathcal{P}_E(R)_{ij} = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

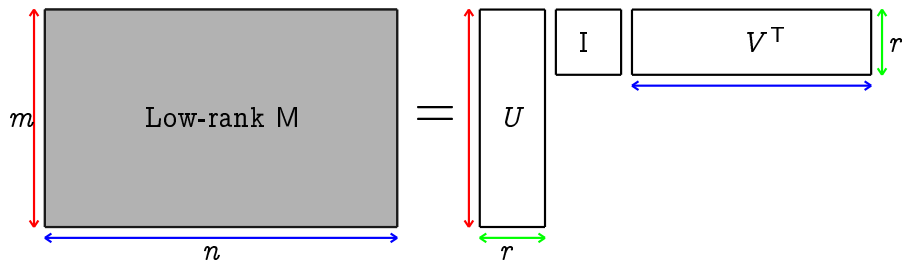
Low rank structure



1. Low-rank matrix M
2. $R = M + Z$
3. Observed subset E

$$\mathcal{P}_E(R)_{ij} = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

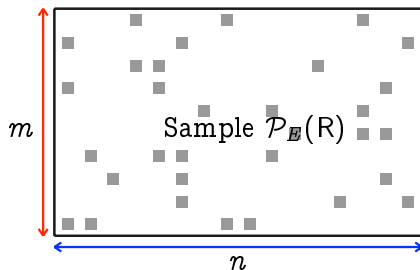
Low rank structure



1. Low-rank matrix M
2. $R = M + Z$
3. Observed subset E

$$\mathcal{P}_E(R)_{ij} = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Low rank structure



1. Low-rank matrix M
2. $R = M + Z$
3. Observed subset E

$$\mathcal{P}_E(R)_{ij} = \begin{cases} M_{ij} + Z_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

Algorithms and accuracy

Questions

▶ How do we minimize $F(X, Y)$?

▶ What prediction accuracy?

$$(\text{RMSE} = \|M - \hat{M}\|_F / \sqrt{mn})$$

How do we minimize $F(X, Y)$?

▶ Spectral methods.

▶ Gradient method.

[Srebro, Rennie, Jaakkola, 2003]

▶ Convex relaxations.

[Fazel, Hindi, Boyd, 2001]

Spectral method (for simplicity $\lambda = 0$)

Replace

$$\begin{aligned} F(X, Y) &= \|\mathcal{P}_E(R - XY^\top)\|_F^2 \\ &= \|\mathcal{P}_E(XY^\top)\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const.} \end{aligned}$$

with ($p = |E|/mn$ fraction of observed entries)

$$\begin{aligned} \tilde{F}(X, Y) &= p \|XY^\top\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const} \\ &= p \left\| XY^\top - \frac{1}{p} \mathcal{P}_E(R) \right\|_F^2 \end{aligned}$$

Spectral method (for simplicity $\lambda = 0$)

Replace

$$\begin{aligned} F(X, Y) &= \|\mathcal{P}_E(R - XY^\top)\|_F^2 \\ &= \|\mathcal{P}_E(XY^\top)\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const.} \end{aligned}$$

with ($p = |E|/mn$ fraction of observed entries)

$$\begin{aligned} \tilde{F}(X, Y) &= p \|XY^\top\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const} \\ &= p \left\| XY^\top - \frac{1}{p} \mathcal{P}_E(R) \right\|_F^2 \end{aligned}$$

Spectral method (for simplicity $\lambda = 0$)

Replace

$$\begin{aligned} F(X, Y) &= \|\mathcal{P}_E(R - XY^\top)\|_F^2 \\ &= \|\mathcal{P}_E(XY^\top)\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const.} \end{aligned}$$

with ($p = |E|/mn$ fraction of observed entries)

$$\begin{aligned} \tilde{F}(X, Y) &= p \|XY^\top\|_F^2 - 2\langle \mathcal{P}_E(R), XY^\top \rangle + \text{const} \\ &= p \left\| XY^\top - \frac{1}{p} \mathcal{P}_E(R) \right\|_F^2 \end{aligned}$$

Spectral method (for simplicity $\lambda = 0$)

Minimize

$$\tilde{F}(X, Y) = \left\| XY^T - \frac{1}{p} \mathcal{P}_E(R) \right\|_F^2$$

Solved by SVD

$$\mathcal{P}_E(R) = XS Y^T \Rightarrow \hat{M} = \frac{1}{p} X_{m \times r} (S)_{r \times r} Y_{n \times r}^T$$

Spectral method (for simplicity $\lambda = 0$)

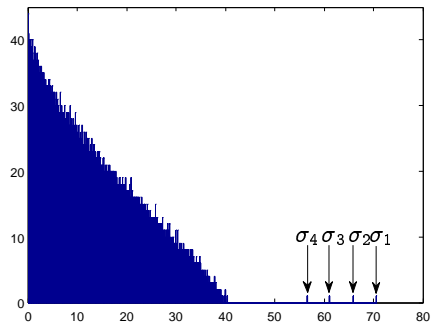
Minimize

$$\tilde{F}(X, Y) = \left\| XY^T - \frac{1}{p} \mathcal{P}_E(R) \right\|_F^2$$

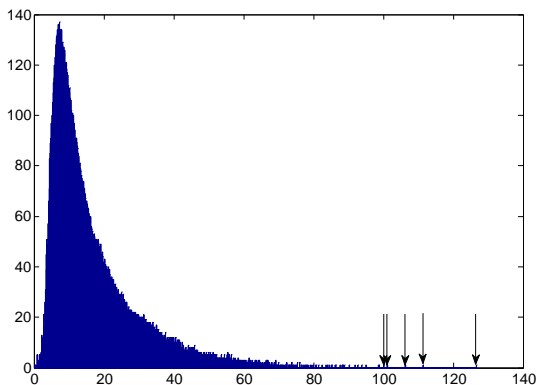
Solved by SVD

$$\mathcal{P}_E(R) = XS Y^T \Rightarrow \hat{M} = \frac{1}{p} X_{m \times r} (S)_{r \times r} Y_{n \times r}^T$$

Random matrix $r = 4$, $m = n = 10000$, $p = 0.0012$



Netflix data (trimmed)



RMSE \approx 0.99

Accuracy guarantees

Theorem (Keshavan, M, Oh, 2009)

Assume $|M_{ij}| \leq M_{\max}$. Then, w.h.p., rank- r projection achieves

$$\text{RMSE} \leq CM_{\max} \sqrt{nr/|\mathbf{E}|} + C' \|Z^{\mathbf{E}}\|_2 n \sqrt{r/|\mathbf{E}|}.$$

E.g. Gaussian noise: $C''(1 + \sigma_z) \sqrt{rn/|\mathbf{E}|}$
[Improves over Achlioptas-McSherry 2003]

Gradient descent

$$F(X, Y) = \sum_{(i,j) \in E} \left(R_{ij} - \langle x_i, y_j \rangle \right)^2 + \lambda \sum_{i=1}^m \|x_i\|_2^2 + \lambda \sum_{j=1}^n \|y_j\|_2^2$$

Update rule: (γ = 'learning rate')

$$x_i \leftarrow (1 - \lambda\gamma)x_i + \gamma \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle x_i, y_j \rangle \right) y_j$$

$$y_j \leftarrow (1 - \lambda\gamma)y_j + \gamma \sum_{i \in \text{Watched}(j)} \left(R_{ij} - \langle x_i, y_j \rangle \right) x_i$$

Gradient descent

$$F(X, Y) = \sum_{(i,j) \in E} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right)^2 + \lambda \sum_{i=1}^m \|\mathbf{x}_i\|_2^2 + \lambda \sum_{j=1}^n \|\mathbf{y}_j\|_2^2$$

Update rule: (γ = 'learning rate')

$$\mathbf{x}_i \leftarrow (1 - \lambda\gamma)\mathbf{x}_i + \gamma \sum_{j \in \text{WatchedBy}(i)} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right) \mathbf{y}_j$$

$$\mathbf{y}_j \leftarrow (1 - \lambda\gamma)\mathbf{y}_j + \gamma \sum_{i \in \text{Watched}(j)} \left(R_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle \right) \mathbf{x}_i$$

Interpretation

$$x_i \leftarrow (1 - \lambda\gamma)x_i + \gamma \sum_{j \in \text{WatchedBy}(i)} w_{ij} y_j$$
$$y_j \leftarrow (1 - \lambda\gamma)y_j + \gamma \sum_{i \in \text{Watched}(j)} w_{ij} x_i$$

user \leftarrow avg of movies she liked
movie \leftarrow avg of users that liked it

Interpretation

$$x_i \leftarrow (1 - \lambda\gamma)x_i + \gamma \sum_{j \in \text{WatchedBy}(i)} w_{ij} y_j$$
$$y_j \leftarrow (1 - \lambda\gamma)y_j + \gamma \sum_{i \in \text{Watched}(j)} w_{ij} x_i$$

user \leftarrow avg of movies she liked
movie \leftarrow avg of users that liked it

A variant (stochastic gradient)

Pick $(i, j) \in E$:

$$x_i \leftarrow (1 - \lambda\gamma)x_i + \gamma w_{ij} y_j$$

$$y_j \leftarrow (1 - \lambda\gamma)y_j + \gamma w_{ij} x_i$$

[Srebro, Rennie, Jaakkola, 2003]

[Srebro, Jaakkola, 2005]

A variant (stochastic gradient)

Pick $(i, j) \in E$:

$$x_i \leftarrow (1 - \lambda\gamma)x_i + \gamma w_{ij} y_j$$

$$y_j \leftarrow (1 - \lambda\gamma)y_j + \gamma w_{ij} x_i$$

[Srebro, Rennie, Jaakkola, 2003]

[Srebro, Jaakkola, 2005]

In the words of SIMON FUNK

Only problem is, we don't have 8.5B entries, we have 100M entries and 8.4B empty cells. Ok, there's another problem too, which is that computing the SVD of ginormous matrices is... well, no fun. Unless you're into that sort of thing.

But, just because there are five hundred really complicated ways of computing singular value decompositions in the literature doesn't mean there isn't a really simple way too: Just take the derivative of the approximation error and follow it. This has the added bonus that we can choose to simply ignore the unknown error on the 8.4B empty slots.

So, yeah, you mathy guys are rolling your eyes right now as it dawns on you how short the path was.

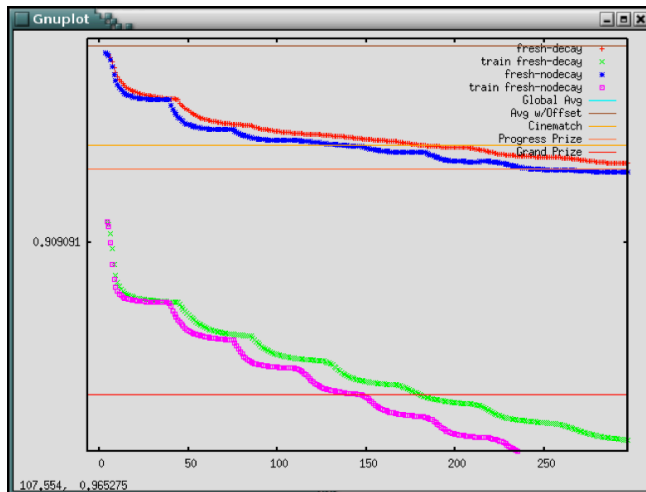
If you write out the equations for the error between the SVD-like model and the original data—just the given values, not the empties—and then take the derivative with respect to the parameters we're trying to infer, you get a rather simple result which I'll give here in C code to save myself the trouble of formatting the math:

```
userValue[user] += lrate * err * movieValue[movie];  
movieValue[movie] += lrate * err * userValue[user];
```

This is kind of like the scene in the Wizard of Oz where Toto pulls back the curtain, isn't it. But wait... let me fluff it up some and make it sound more impressive.

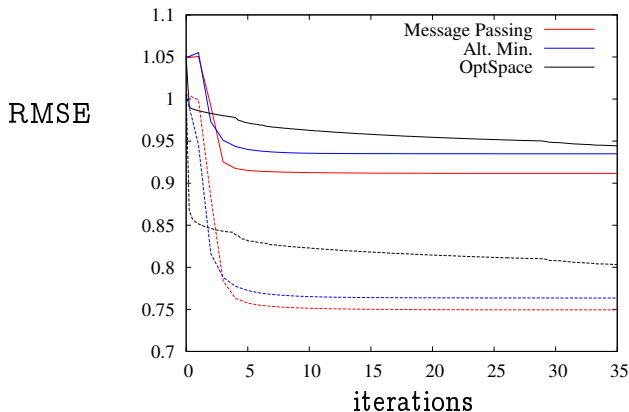
[Netflix challenge, 2006-2009]

And his results



Target $\text{RMSE} \leq 0.8564$

Three variants



OPTSPACE ~ Gradient descent on Grassmannian

ALTERNATING LEAST SQUARES

MESSAGE PASSING

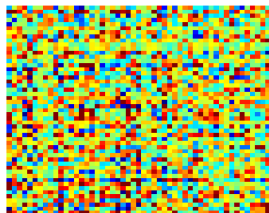
[Keshavan-M.-Oh 2009]

[Koren-Bell 2008]

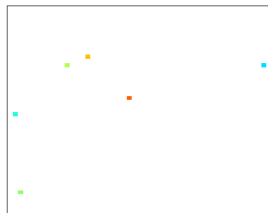
[Keshavan-M. 2011]

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

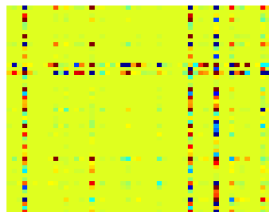
low-rank matrix M



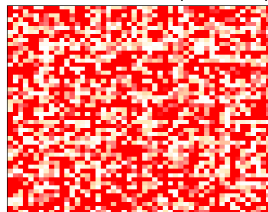
sampled matrix M^E



OPTSPACE output \hat{M}



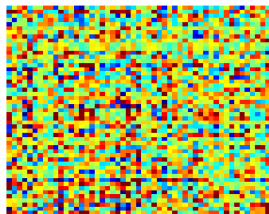
squared error $(M - \hat{M})^2$



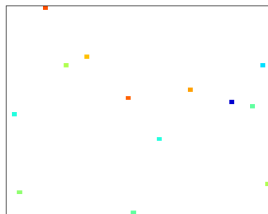
0.25% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

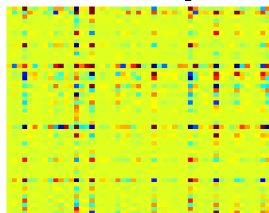
low-rank matrix M



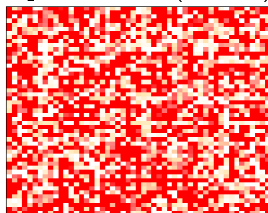
sampled matrix M^E



OPTSPACE output \hat{M}



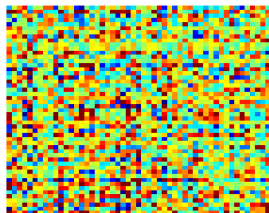
squared error $(M - \hat{M})^2$



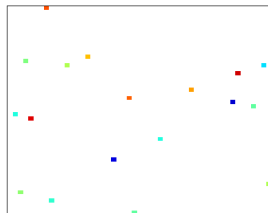
0.50% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

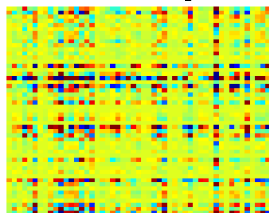
low-rank matrix M



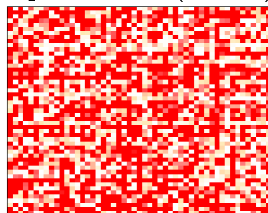
sampled matrix M^E



OPTSPACE output \hat{M}



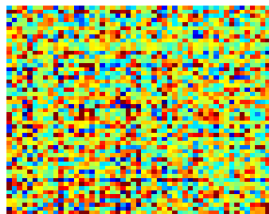
squared error $(M - \hat{M})^2$



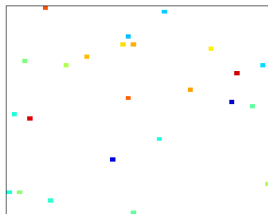
0.75% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

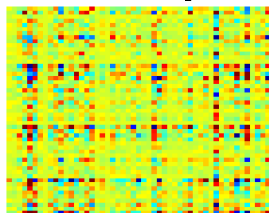
low-rank matrix M



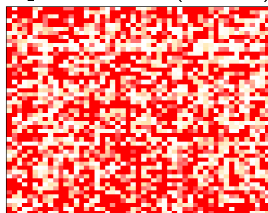
sampled matrix M^E



OPTSPACE output \hat{M}



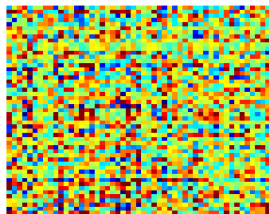
squared error $(M - \hat{M})^2$



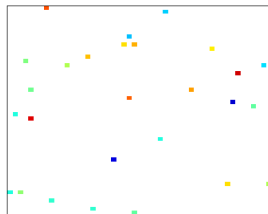
1.00% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

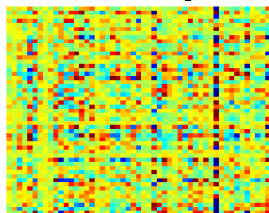
low-rank matrix M



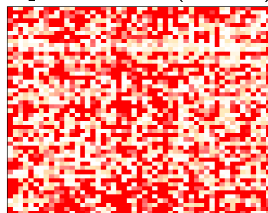
sampled matrix M^E



OPTSPACE output \hat{M}



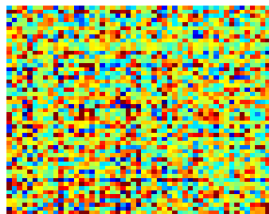
squared error $(M - \hat{M})^2$



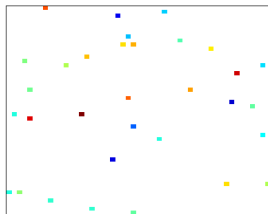
1.25% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

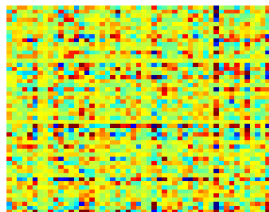
low-rank matrix M



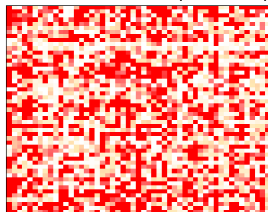
sampled matrix M^E



OPTSPACE output \hat{M}



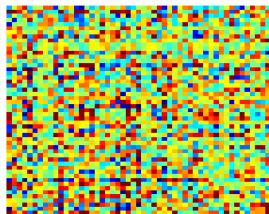
squared error $(M - \hat{M})^2$



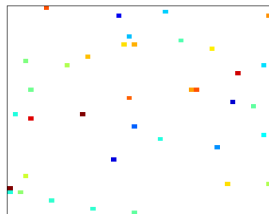
1.50% sampled

Accuracy guarantees: Vignette ($m = n = 2000$ rank-8)

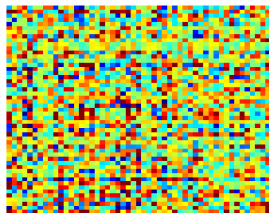
low-rank matrix M



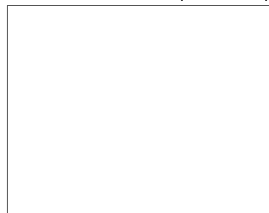
sampled matrix M^E



OPTSPACE output \hat{M}



squared error $(M - \hat{M})^2$



1.75% sampled

Accuracy guarantees: Unstructured factors

Incoherence

$$\|u_i\|^2 \leq \mu \langle \|u\|^2 \rangle_{\text{av}}, \quad \|v_j\|_2^2 \leq \mu \langle \|v\|^2 \rangle_{\text{av}}.$$

[Candés, Recht 2008]

Accuracy guarantees

Theorem (Keshavan, M, Oh, 2009)

Assume $|M_{ij}| \leq M_{\max}$. Then, w.h.p., rank- r projection achieves

$$\text{RMSE} \leq CM_{\max} \sqrt{nr/|E|} + C' \|Z^E\|_2 n \sqrt{r/|E|}.$$

Theorem (Keshavan, M, Oh, 2009)

Let M be *incoherent* with $\sigma_1(M)/\sigma_r(M) = O(1)$. If $|E| \geq Cn \min\{r(\log n)^2, r^2 \log n\}$ then, w.h.p., OPTSPACE achieves

$$\text{RMSE} \leq C'' \frac{n\sqrt{r}}{|E|} \|Z^E\|_2,$$

with complexity $O(nr^3(\log n)^2)$.

E.g. Gaussian noise: $C''\sigma_z \sqrt{rn/|E|}$

Accuracy guarantees

Theorem (Keshavan, M, Oh, 2009)

Assume $|M_{ij}| \leq M_{\max}$. Then, w.h.p., rank- r projection achieves

$$\text{RMSE} \leq CM_{\max} \sqrt{nr/|E|} + C' \|Z^E\|_2 n\sqrt{r/|E|}.$$

Theorem (Keshavan, M, Oh, 2009)

Let M be *incoherent* with $\sigma_1(M)/\sigma_r(M) = O(1)$. If $|E| \geq Cn \min\{r(\log n)^2, r^2 \log n\}$ then, w.h.p., OPTSPACE achieves

$$\text{RMSE} \leq C'' \frac{n\sqrt{r}}{|E|} \|Z^E\|_2,$$

with complexity $O(nr^3(\log n)^2)$.

E.g. Gaussian noise: $C''\sigma_z \sqrt{rn/|E|}$

Two surprises

Can do **much better** than SVD !

Error = Noise / Sampling factor

Analogous guarantees for convex relaxations

Candés, Recht, 2008

Candés, Plan, 2009

Candés, Tao, 2009

Gross, 2010

Negahbani, Wainwright, 2010

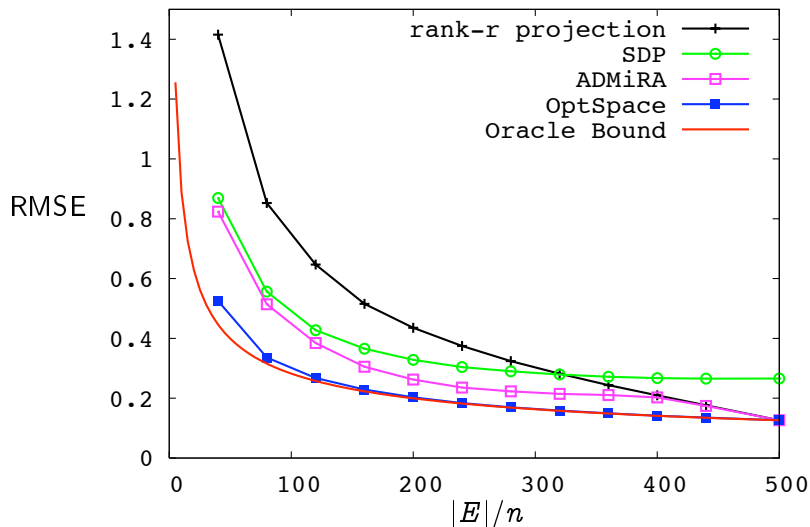
Koltchinskii, Lounici, Tsybakov, 2011

...

*

A noisy example

- $n = 500, r = 4, \sigma_z = 1$, example from [Candés, Plan, 2009]



Challenge #1: Privacy

Research question

User ratings \rightarrow User feature vector

User ratings \rightarrow User private attributes ???

Let us try to do it!

Research question

User ratings \rightarrow User feature vector

User ratings \rightarrow User private attributes ???

Let us try to do it!

Research question

User ratings \rightarrow User feature vector

User ratings \rightarrow User private attributes ???

Let us try to do it!

Which attribute?

Number of persons in the household

- ▶ Non-obvious
- ▶ Recommender has incentive
- ▶ 2011 CAMRA CHALLENGE

[no prize :-(we won it :-)]

Which attribute?

Number of persons in the household

- ▶ Non-obvious
- ▶ Recommender has incentive
- ▶ 2011 CAMRA CHALLENGE

[no prize :-(we won it :-)]

CAMRA dataset

- ▶ $m \approx 2 \cdot 10^5$ users, $n \approx 2 \cdot 10^4$ movies
- ▶ $|E| \approx 4.5 \cdot 10^6$ ratings ($p \approx 0.001$)
- ▶ 272 households of size 2
- ▶ 14 households of size 3
- ▶ 4 households of size 4

CAMRA dataset

- ▶ $m \approx 2 \cdot 10^5$ users, $n \approx 2 \cdot 10^4$ movies
- ▶ $|E| \approx 4.5 \cdot 10^6$ ratings ($p \approx 0.001$)
- ▶ 272 households of size 2
- ~~▶ 14 households of size 3~~
- ~~▶ 4 households of size 4~~

CAMRA dataset

- ▶ $m \approx 2 \cdot 10^5$ users, $n \approx 2 \cdot 10^4$ movies
- ▶ $|E| \approx 4.5 \cdot 10^6$ ratings ($p \approx 0.001$)
- ▶ 272 households of size 2

- ▶ Can you identify whether two users shared an account?
- ▶ Can you identify which user watched a movie?

Short answer

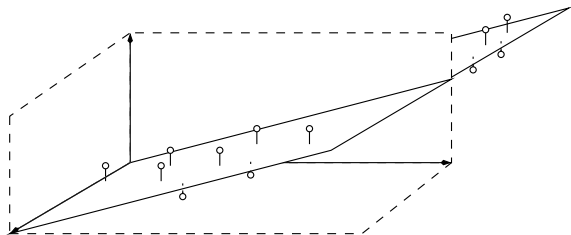
Yes: For a significant fraction of the accounts.

Two examples (Netflix dataset, no ground truth)

User 1	User 2
TLOTR: The Fellowship of the Ring [†] (5), TLOTR: The Return of the King [†] (5), TLOTR: The Two Towers [†] (5), The Whole Nine Yards(4), Immortal [†] (1), The Deep End(2), Toys [†] (4), The Addams Family(5)	H.R. Pufnstuf(5), Sex and the City: Season 5 [♡] (1), Me Myself & Irene(1), All the Real Girls ^{♡△} (5), Titanic [♡] (5), George Washington [△] (5), The Siege(1), In the Bedroom [△] (5)
User 1	User 2
Monsters Inc. [◇] (5), Finding Nemo [◇] (5), Whale Rider(5), Con Air(4), Lilo and Stitch [◇] (4), Ice Age [◇] (5), Ring of Fire(4), Star Trek: Nemesis(3),	In America [♣] (2), Super Size Me(2), A Very Long Engagement [♣] (1), Bend It Like Beckham(2), 21 Grams [♣] (1), Airplane II: The Sequel(4), Spun [♣] (1), Fahrenheit 9/11(1)

No movie info used!

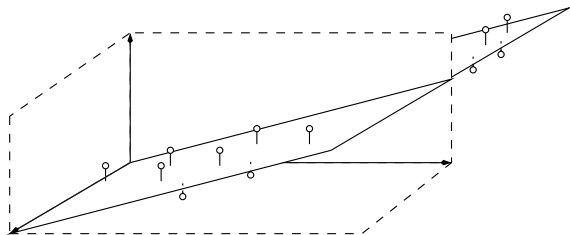
How did you do it?



$$R_{ij} \sim \langle \mathbf{u}_i, \mathbf{v}_j \rangle + \epsilon_{ij}$$

$$\left\{ (R_{ij}, -v_j) \in \mathbb{R}^{r+1}, \quad j \in \text{WatchedBy}(i) \right\} \subseteq \text{Hyperplane}$$

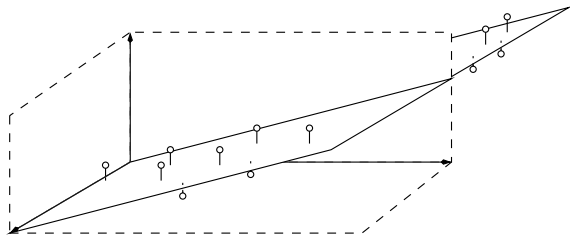
How did you do it?



$$R_{ij} \sim \langle \mathbf{u}_i, \mathbf{v}_j \rangle + \epsilon_{ij}$$

$$\left\{ (R_{ij}, -v_j) \in \mathbb{R}^{r+1}, \quad j \in \text{WatchedBy}(i) \right\} \subseteq \text{Hyperplane}$$

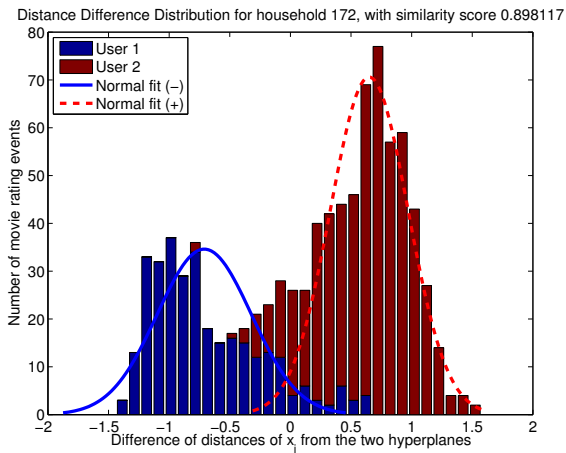
How did you do it? Subspace clustering



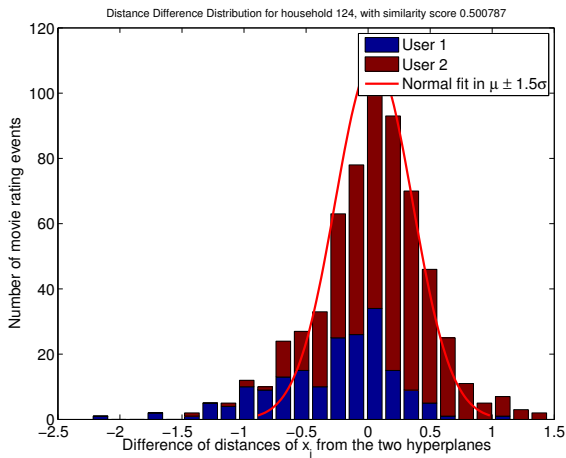
One user \rightarrow One hyperplane

Two users \rightarrow Two hyperplanes

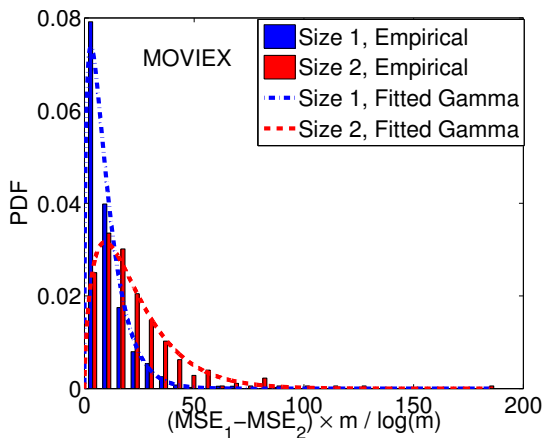
Example: A two-user household (CAMRA)



Example: Another two-user household (CAMRA)



Classifying households

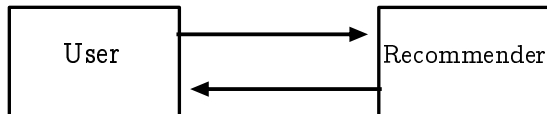


MSE_1 → MSE using one hyperplane

MSE_2 → MSE using two hyperplanes

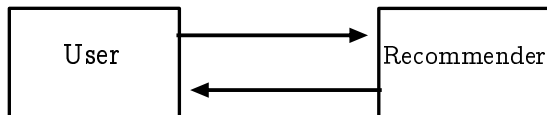
Challenge #2: Interactivity

We want to design the system



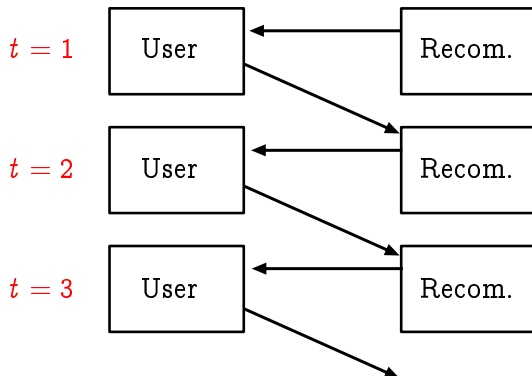
We took time out of the picture.

We want to design the system



We took time out of the picture.

Interactive system



Abstract from other users

At time t

- ▶ Recommender suggests movie v_t ;
- ▶ User gives feedback R_t

Focus on user u

$$R_t \sim \langle u, v_t \rangle + \epsilon_t$$

Abstract from other users

At time t

- ▶ Recommender suggests movie v_t ;
- ▶ User gives feedback R_t

Focus on user u

$$R_t \sim \langle u, v_t \rangle + \epsilon_t$$

Linear bandits

Decision:

$$v_t$$

Observations:

$$R_t \sim \langle u, v_t \rangle + \varepsilon_t$$

Reward:

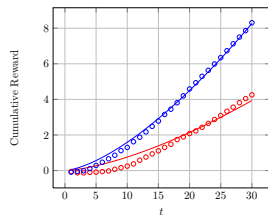
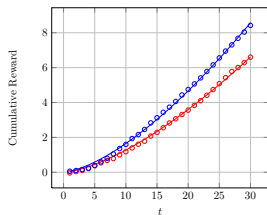
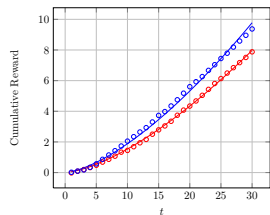
$$O_t = \sum_{\ell=1}^t \langle u, v_\ell \rangle$$

[Rusmevichientong, Tsitsiklis, 2008]

Important differences

- ▶ Number of observations \sim Dimensions
- ▶ Cannot explore completely at random.

Simulating an interactive system (Netflix data)



- ▶ 3 'typical' users
- ▶ Constant-optimal policy

Conclusion

Conclusion

- ▶ A crucial technology for modern information networks.
- ▶ Only scratched the surface.

Thanks!

Conclusion

- ▶ A crucial technology for modern information networks.
- ▶ Only scratched the surface.

Thanks!