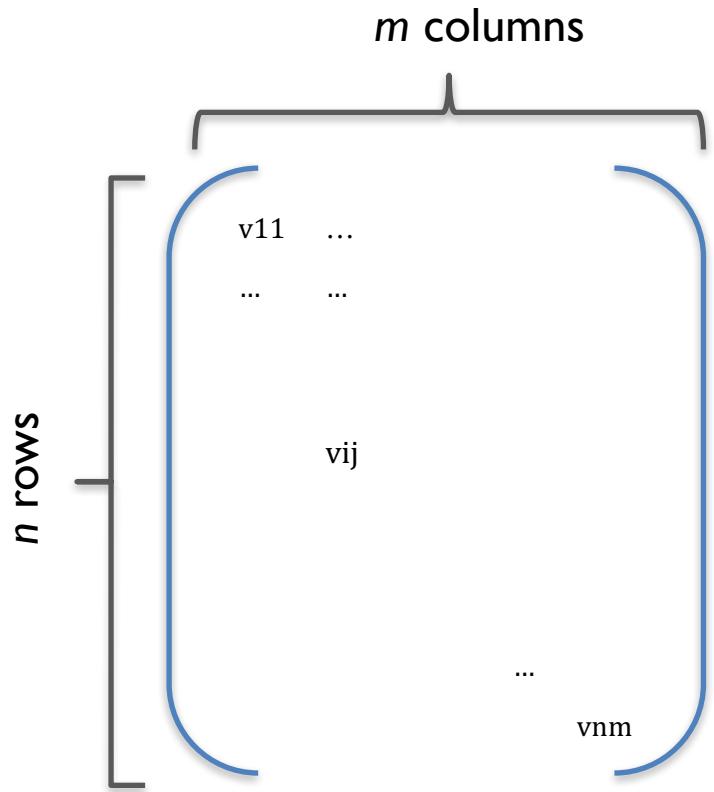
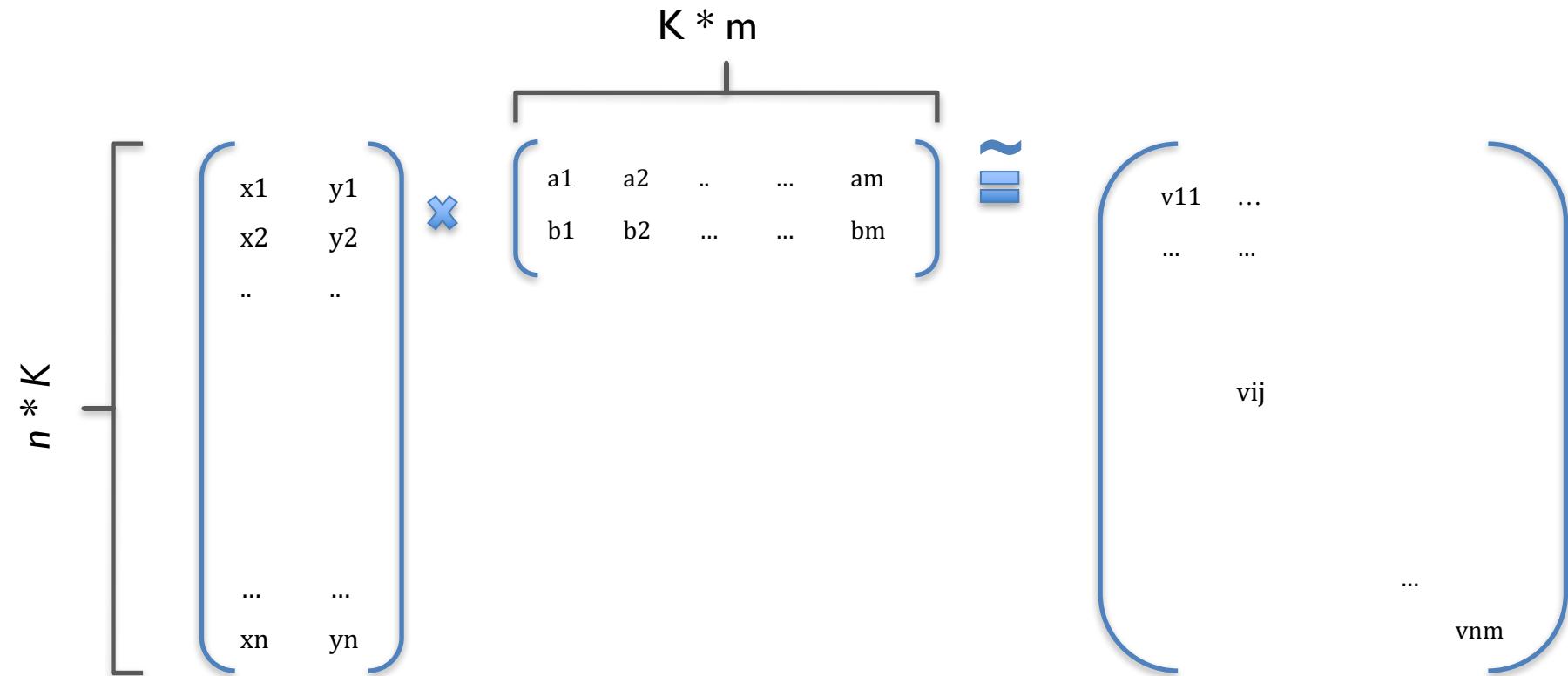


Matrix Factorization

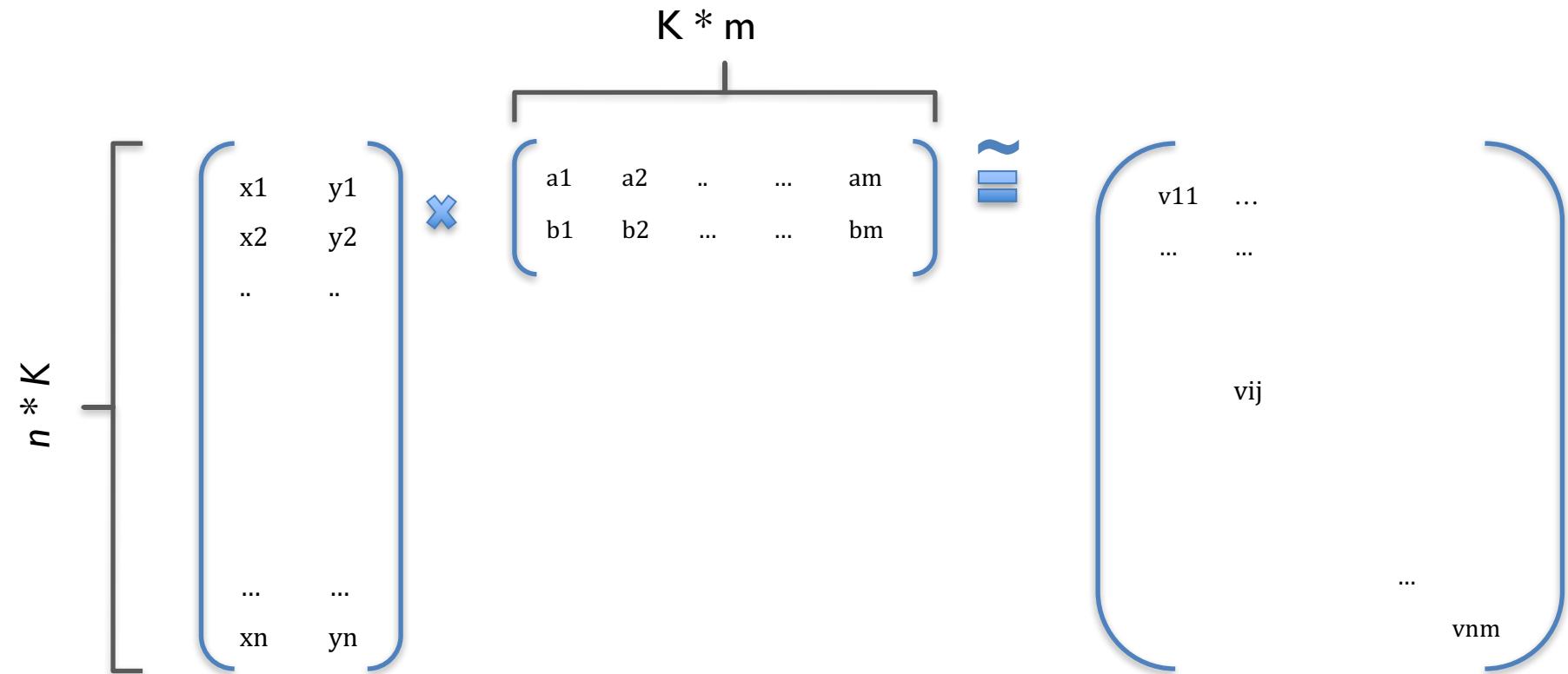
Recovering latent factors in a matrix



Recovering latent factors in a matrix



What is this for?



MF for collaborative filtering

What is collaborative filtering?

Your Amazon.com

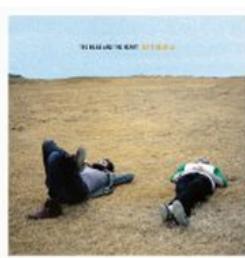
Featured Recommendations MP3 Albums Kindle eBooks Books Health & Personal Care Apparel Sports & Outdoors See All Recommendations

MP3 Albums

Page 1 of 20



New Release
Build Me Up From ...
Sarah Jarosz
★★★★★ (28)
\$9.49
[Why recommended?](#)



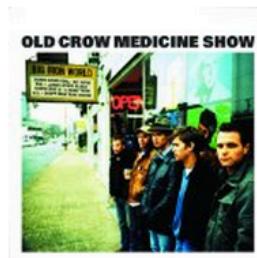
New Release
Let's Be Still
The Head And The Heart
★★★★★ (21)
\$9.49
[Why recommended?](#)



Leaving Eden
Carolina Chocolate Drops
★★★★★ (66)
\$10.49
[Why recommended?](#)



Who's Feeling Young ...
Punch Brothers
★★★★★ (60)
\$10.49
[Why recommended?](#)



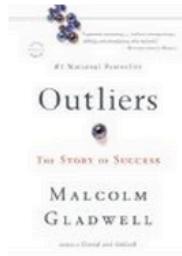
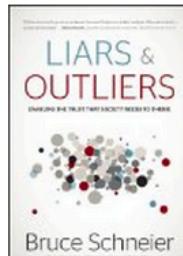
Big Iron World
Old Crow Medicine Show
★★★★★ (39)
\$9.49
[Why recommended?](#)

[See all recommendations in MP3 Albums](#)



Kindle eBooks

Page 1 of 20



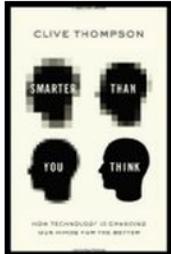
6



What is collaborative filtering?

Books

Page 1 of 20



New Release

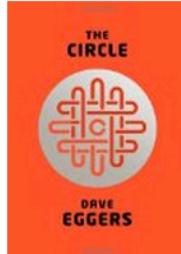
Smarter Than You ...

› Clive Thompson

★★★★★ (26)

\$27.95 \$20.82

Why recommended?



New Release

The Circle

› Dave Eggers

★★★★★ (77)

\$27.95 \$16.77

Why recommended?



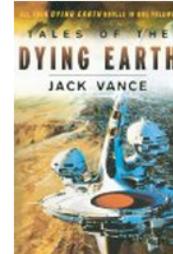
Lord of Light

› Roger Zelazny

★★★★★ (186)

\$13.99 \$10.68

Why recommended?



Tales of the Dying ...

› Jack Vance

★★★★★ (81)

\$22.99 \$15.94

Why recommended?



Latro in the Mist

› Gene Wolfe

★★★★★ (24)

\$21.99 \$15.25

Why recommended?

› See all recommendations in Books

Sports & Outdoors

Page 1 of 17



Halo-V Velcro ...

★★★★★ (30)

\$6.45 - \$19.64

Why recommended?



Halo Headband

★★★★★ (101)

\$3.40 - \$18.34

Why recommended?



Halo Super Wide ...

★★★★★ (15)

\$7.95 - \$14.95

Why recommended?



Headsweats ...

★★★★★ (126)

\$12.06 - \$28.99

Why recommended?



Sweat Gtr Headband

★★★★★ (180)

\$15.77 - \$53.17

Why recommended?

What is collaborative filtering?

Your Amazon.com > Improve Your Recommendations

(If you're not William Cohen, click here.)

Help us make better recommendations. You can refine your recommendations by rating items or adjusting the checkboxes.

EDIT YOUR COLLECTION

Items you've purchased

1.  **LOOK INSIDE!**
Love Is Strange (A Paranormal Romance)
by Bruce Sterling
Your tags: **Add** ([What's this?](#))
Click to Add: [paranormal romance](#), [nerd](#), [futurist](#), [science fiction romance](#), [science fiction](#), [technology](#), [scifi](#), [literature](#)

Your Rating: 
 This was a gift
 Don't use for recommendations

2. 
MAD Magazine #1
by Harvey Kurtzman
Your tags: **Add** ([What's this?](#))
Click to Add: [harvey kurtzman](#), [dc](#)

Your Rating: 
 This was a gift
 Don't use for recommendations

3. 
Ahoy!
Punch Brothers | Format: MP3 Music
Your tags: **Add** ([What's this?](#))
Click to Add: [bluegrass](#), [music](#), [punch brothers](#), [singer-songwriters](#)

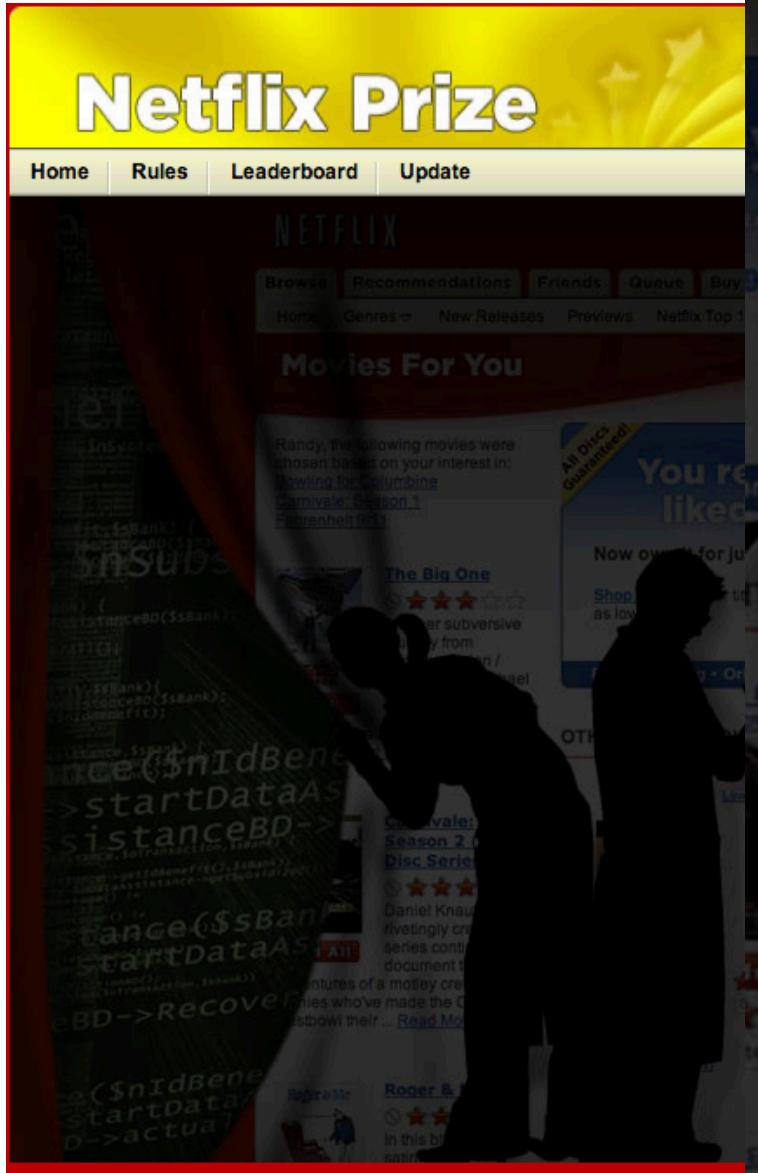
Your Rating: 
 This was a gift
 Don't use for recommendations

EDIT YOUR PREFERENCES

Show Amazon book recommendations as Kindle editions when possible.

Need Help?
Visit our [help](#) area to learn more.

What is collaboration?



Congratulations!

The Netflix Prize sought to substantially improve the accuracy of predictions about how much someone is going to enjoy a movie based on their movie preferences.

On September 21, 2009 we awarded the \$1M Grand Prize to team “BellKor’s Pragmatic Chaos”. Read about [their algorithm](#), checkout team scores on the [Leaderboard](#), and join the discussions on the [Forum](#).

We applaud all the contributors to this quest, which improves our ability to connect people to the movies they love.

Leaderboard

Showing Test Score. [Click here to show quiz score](#)

Display top 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

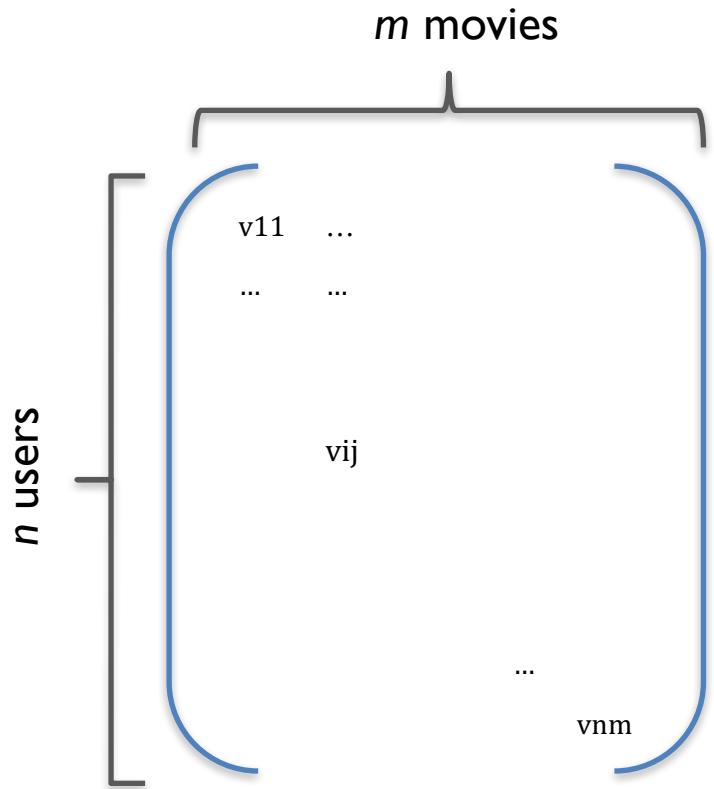
Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos

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

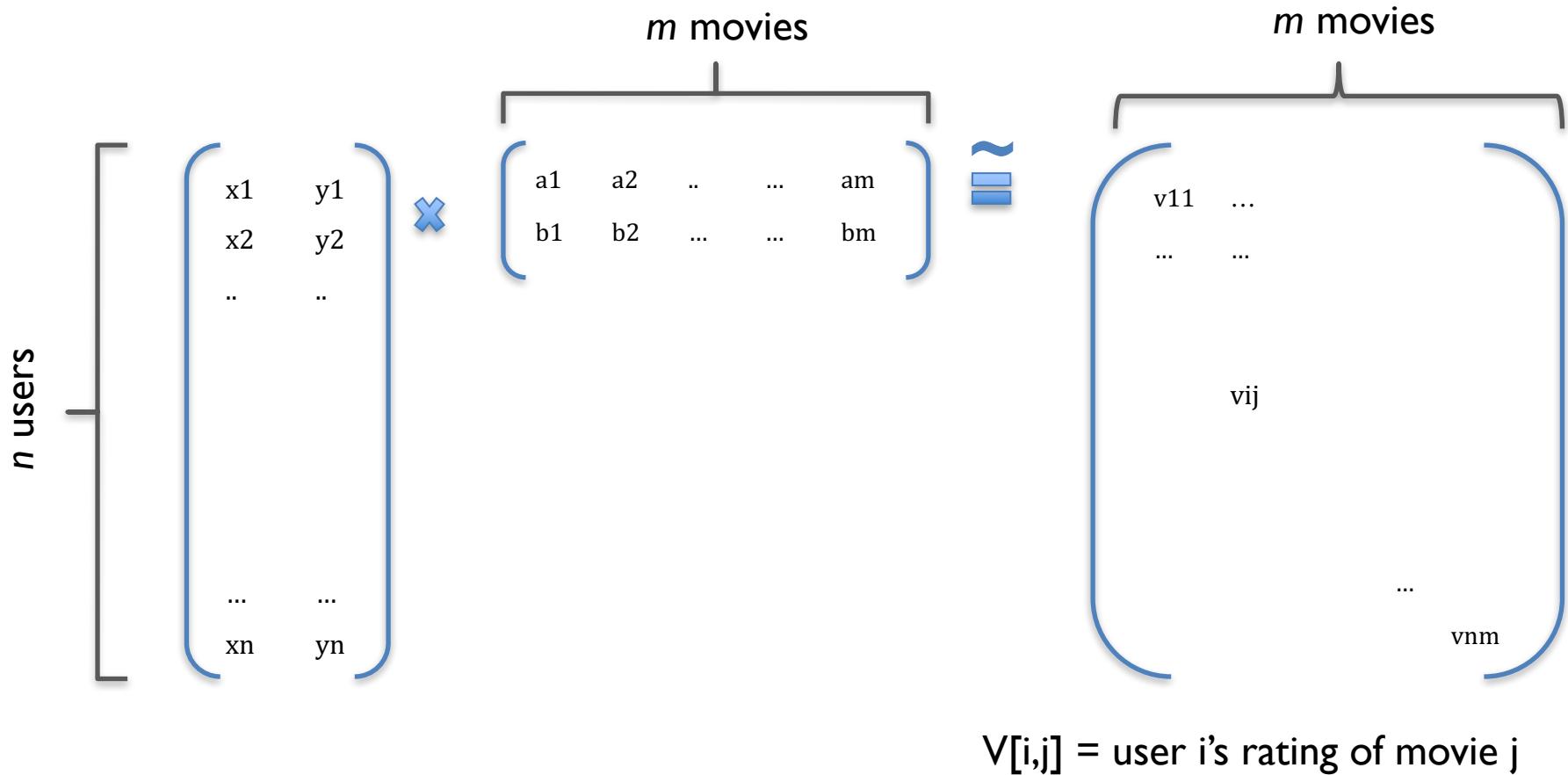
Cinematch score - RMSE = 0.9525

Recovering latent factors in a matrix

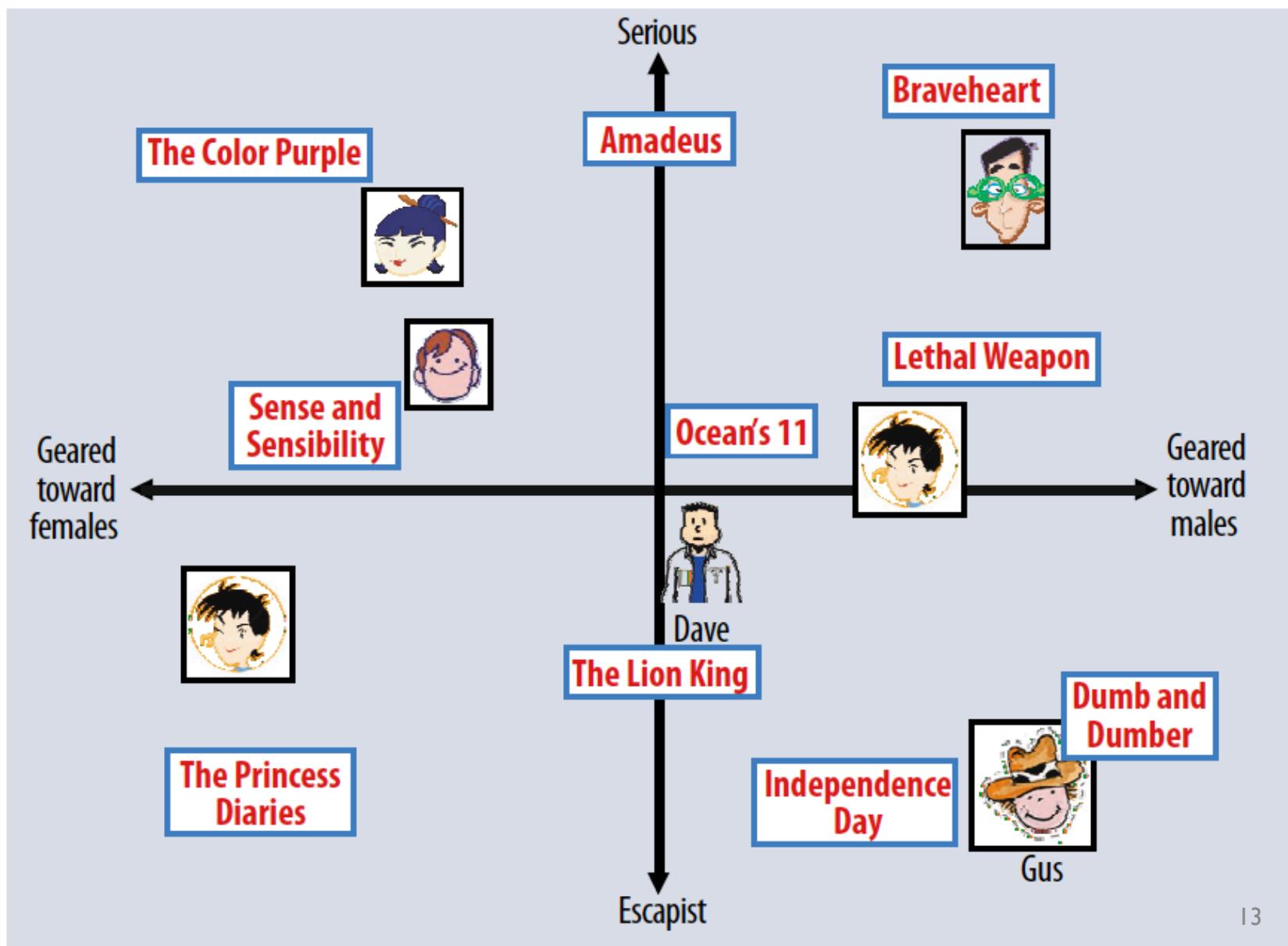


$V[i,j] =$ user i's rating of movie j

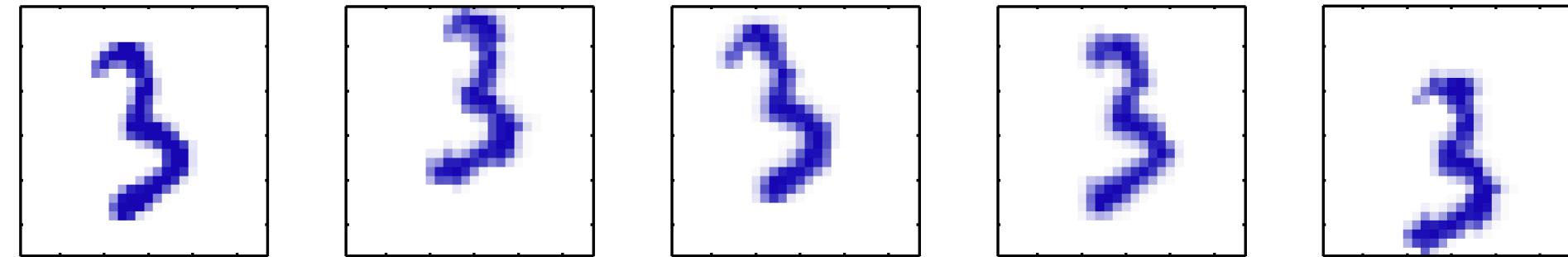
Recovering latent factors in a matrix



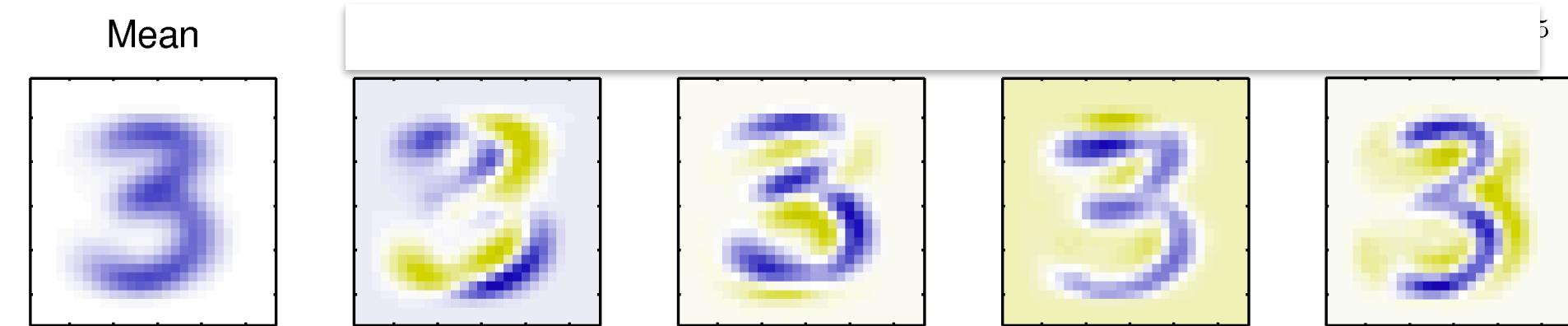
Semantic Factors (Koren et al., 2009)



MF for image modeling



Mean



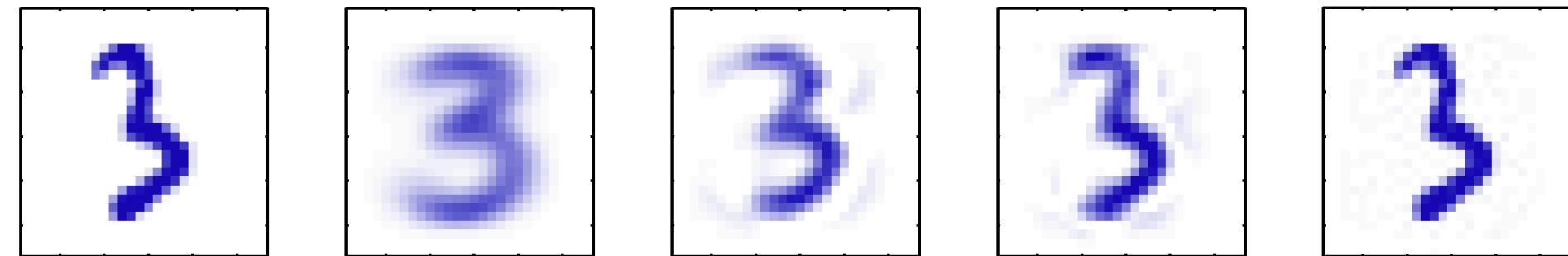
Original

$M = 1$

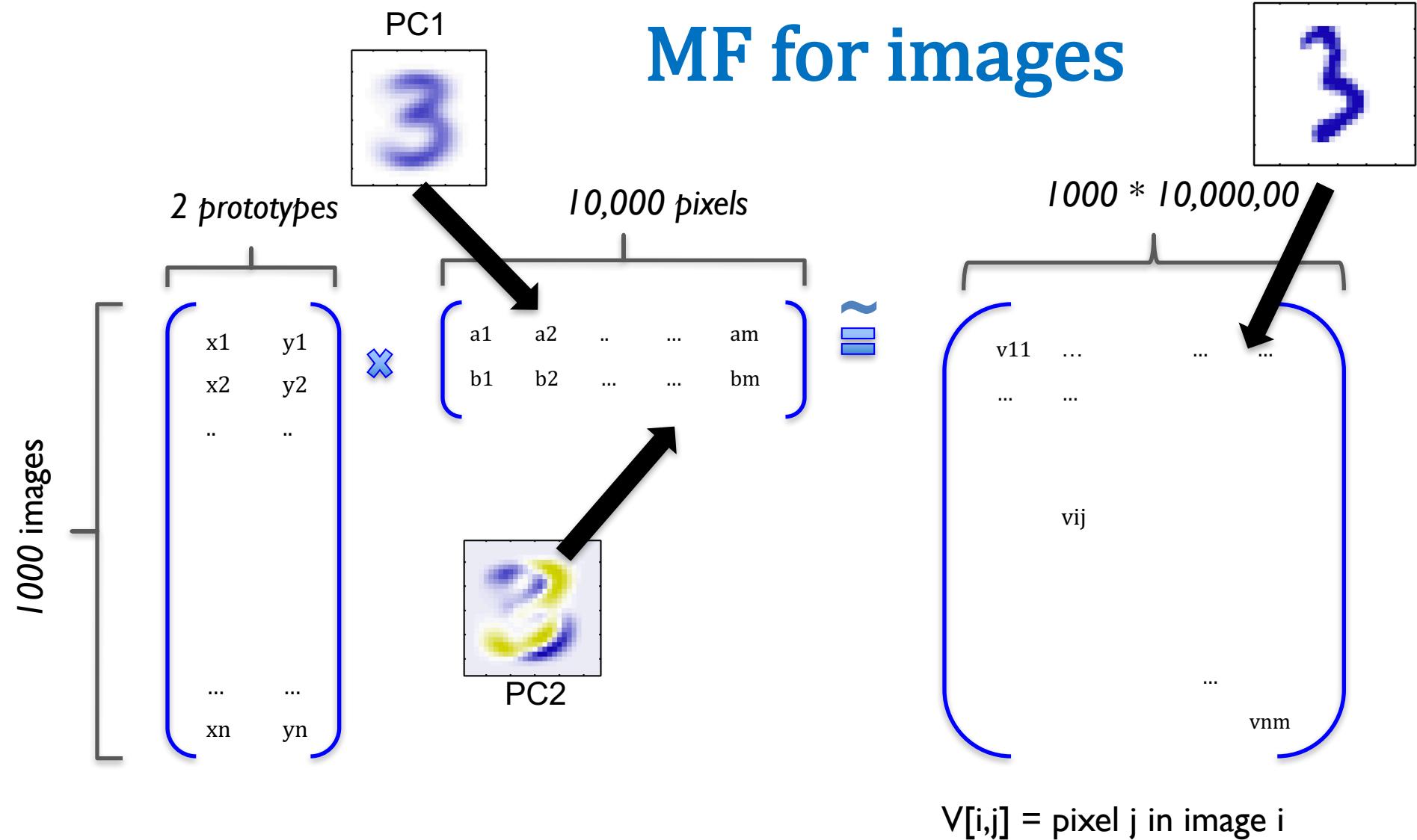
$M = 10$

$M = 50$

$M = 250$



MF for images



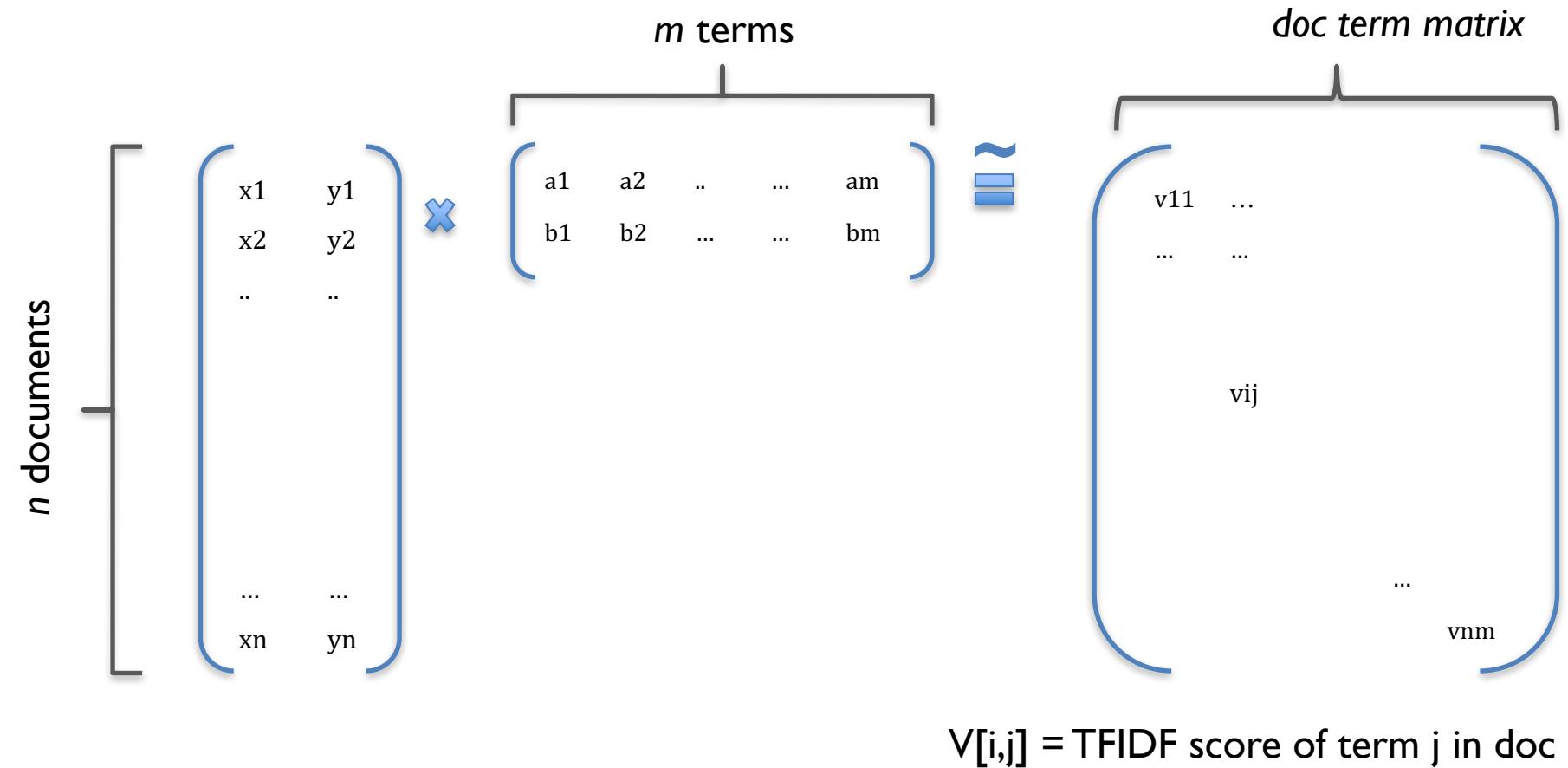
MF for modeling text

- The Neatest Little Guide to Stock Market Investing
- Investing For Dummies, 4th Edition
- The Little Book of Common Sense Investing: The Only Way to Guarantee Your Fair Share of Stock Market Returns
- The Little Book of Value Investing
- Value Investing: From Graham to Buffett and Beyond
- Rich Dad's Guide to Investing: What the Rich Invest in, That the Poor and the Middle Class Do Not!
- Investing in Real Estate, 5th Edition
- Stock Investing For Dummies
- Rich Dad's Advisors: The ABC's of Real Estate Investing: The Secrets of Finding Hidden Profits Most Investors Miss

TFIDF counts would be better

Index Words	Titles								
	T1	T2	T3	T4	T5	T6	T7	T8	T9
book			1	1					
dads						1			1
dummies		1						1	
estate							1		1
guide	1					1			
investing	1	1	1	1	1	1	1	1	1
market	1		1						
real							1		1
rich						2			1
stock	1		1					1	
value				1	1				

Recovering latent factors in a matrix

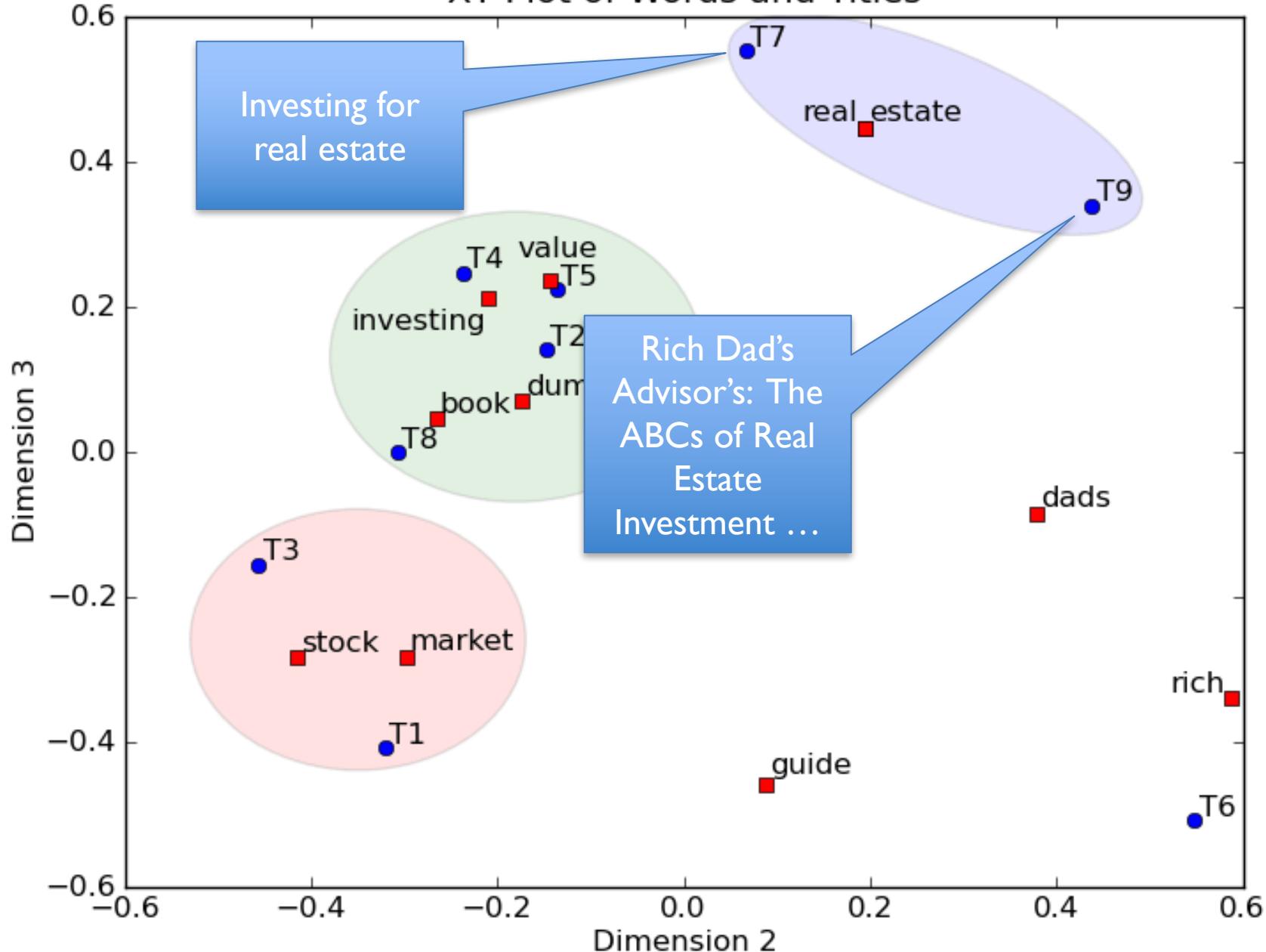


$$\begin{array}{|c|c|c|} \hline 3.91 & 0 & 0 \\ \hline 0 & 2.61 & 0 \\ \hline 0 & 0 & 2 \\ \hline \end{array} * \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline \text{T1} & \text{T2} & \text{T3} & \text{T4} & \text{T5} & \text{T6} & \text{T7} & \text{T8} & \text{T9} \\ \hline 0.35 & 0.22 & 0.34 & 0.26 & 0.22 & 0.49 & 0.28 & 0.29 & 0.44 \\ \hline -0.32 & -0.15 & -0.46 & -0.24 & -0.14 & 0.55 & 0.07 & -0.31 & 0.44 \\ \hline -0.41 & 0.14 & -0.16 & 0.25 & 0.22 & -0.51 & 0.55 & 0 & 0.34 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline \end{array}$$

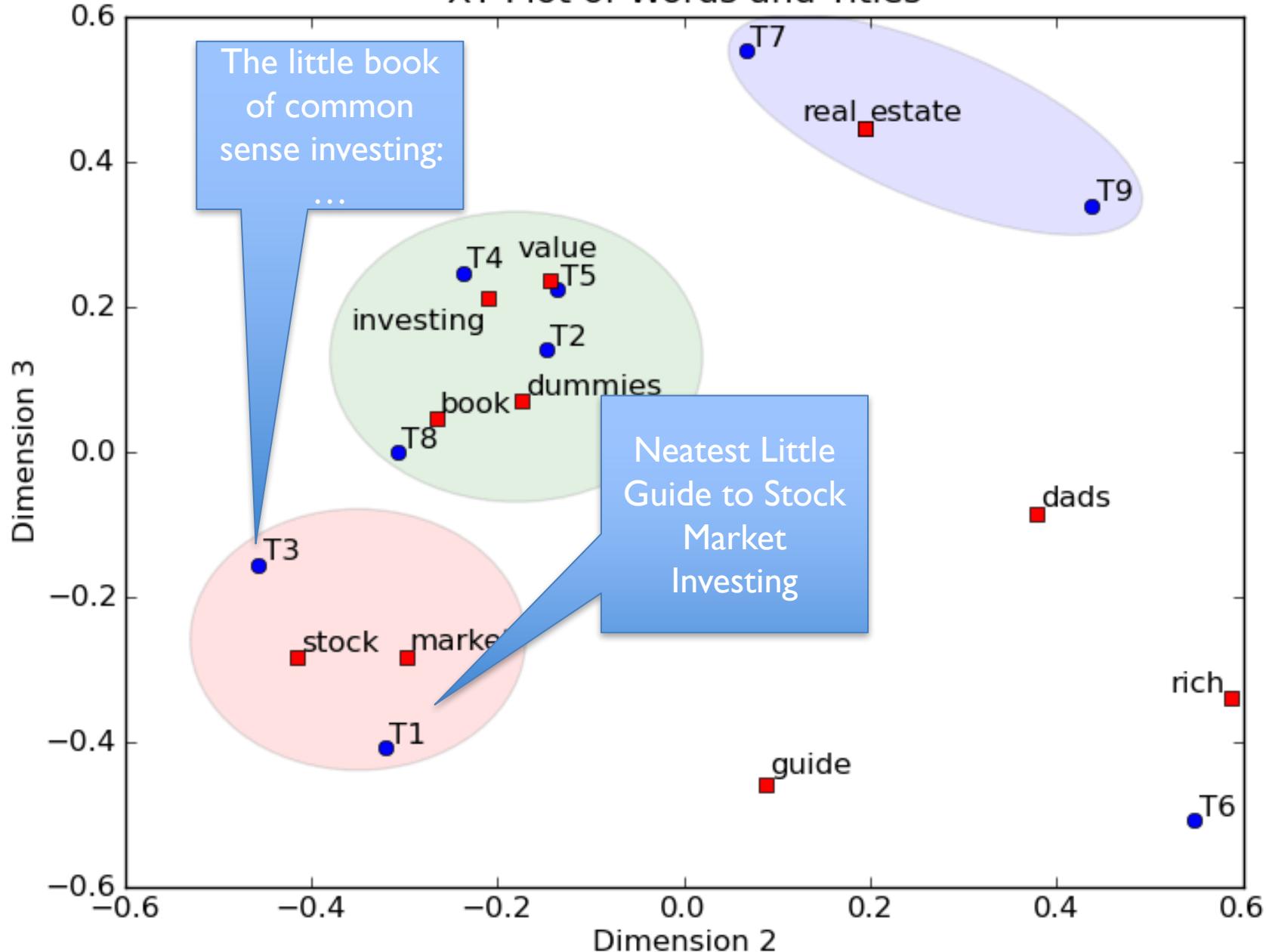
book	0.15	-0.27	0.04
dads	0.24	0.38	-0.09
dummies	0.13	-0.17	0.07
estate	0.18	0.19	0.45
guide	0.22	0.09	-0.46
investing	0.74	-0.21	0.21
market	0.18	-0.3	-0.28
real	0.18	0.19	0.45
rich	0.36	0.59	-0.34
stock	0.25	-0.42	-0.28
value	0.12	-0.14	0.23

*

XY Plot of Words and Titles



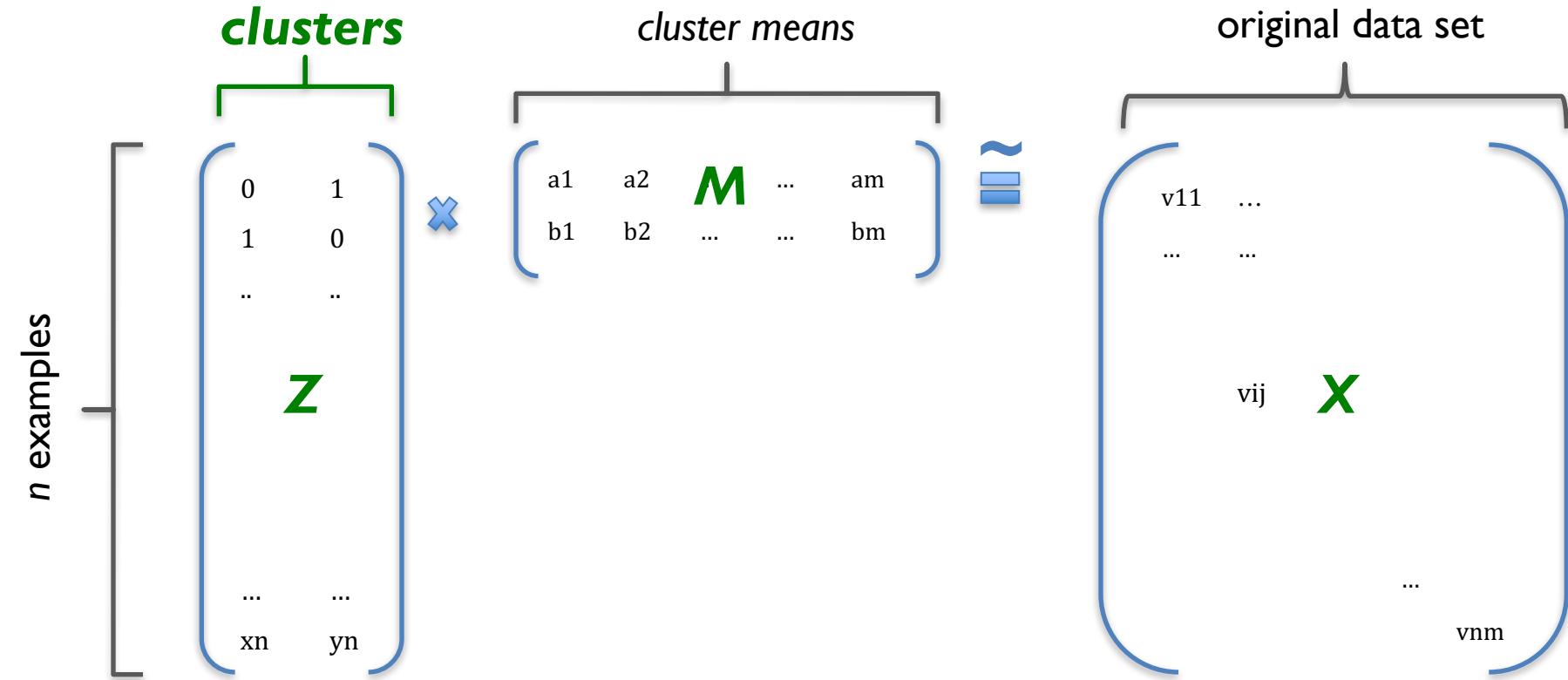
XY Plot of Words and Titles



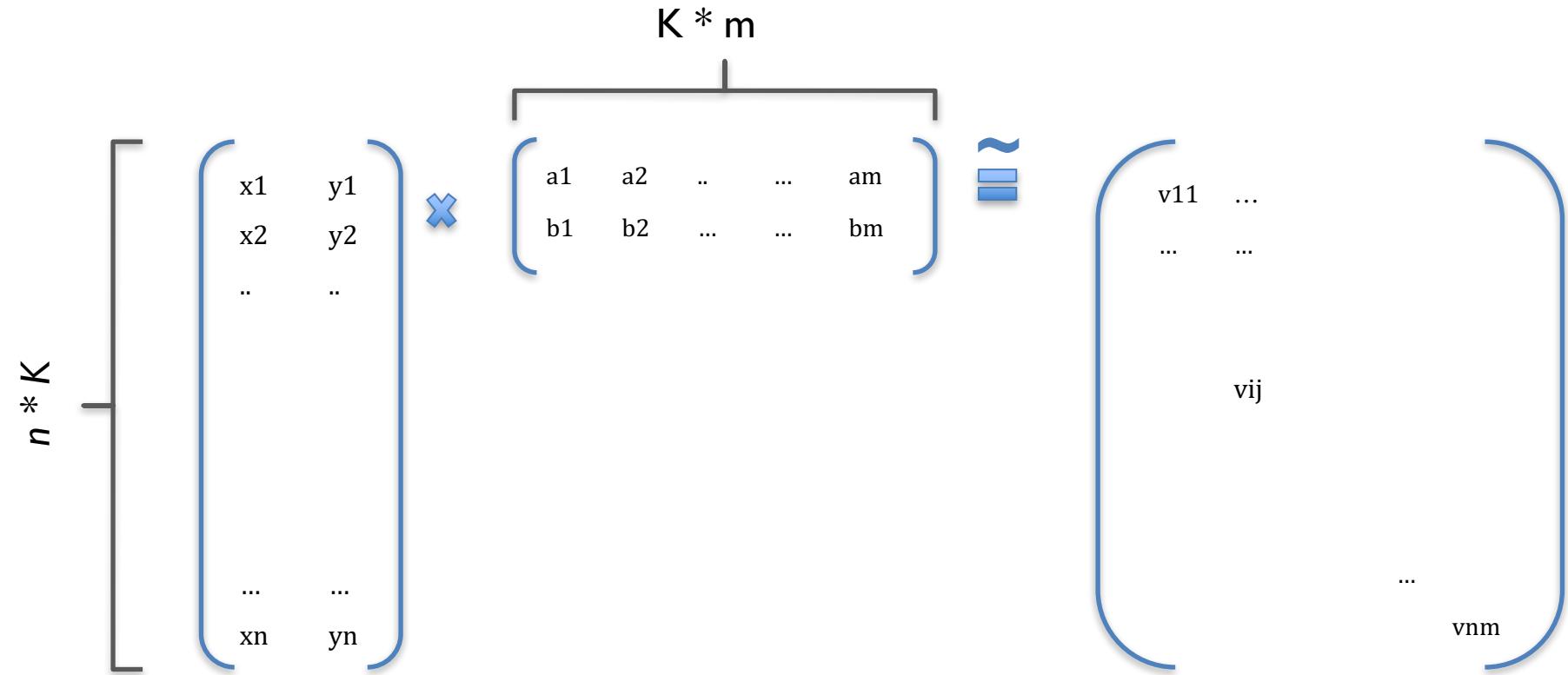
MF is like clustering

k-means as MF

indicators for r



How do you do it?



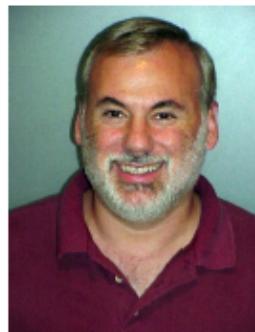
Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla

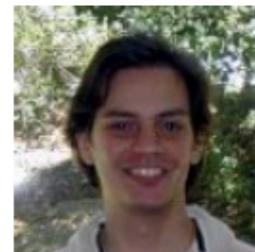


talk pilfered from →

Peter J. Haas



Yannis Sismanis



Erik Nijkamp



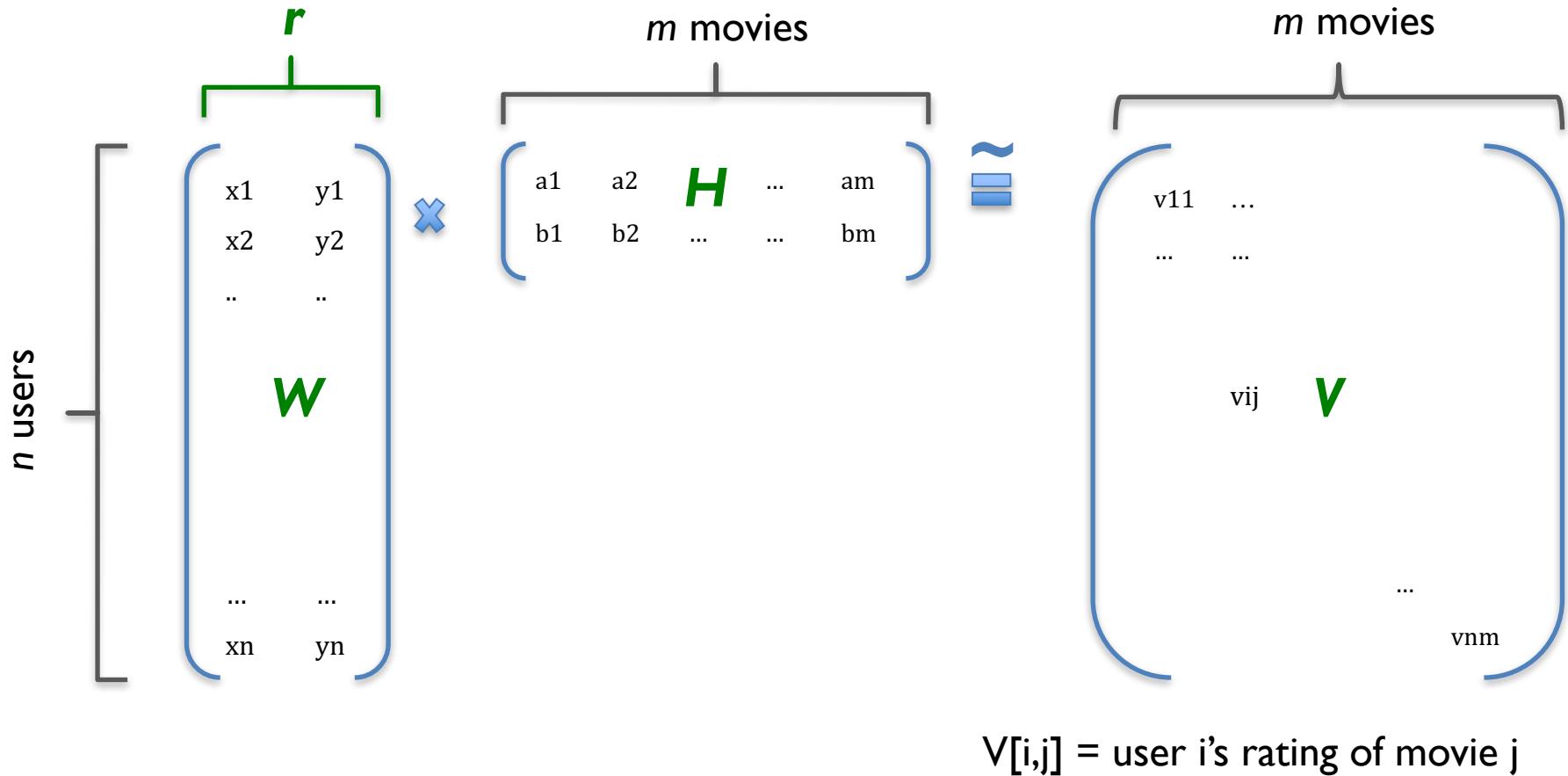
Collaborative Filtering

- ▶ Problem
 - ▶ Set of users
 - ▶ Set of items (movies, books, jokes, products, stories, ...)
 - ▶ Feedback (ratings, purchase, click-through, tags, ...)
- ▶ Predict additional items a user may like
 - ▶ Assumption: Similar feedback \implies Similar taste
- ▶ Example

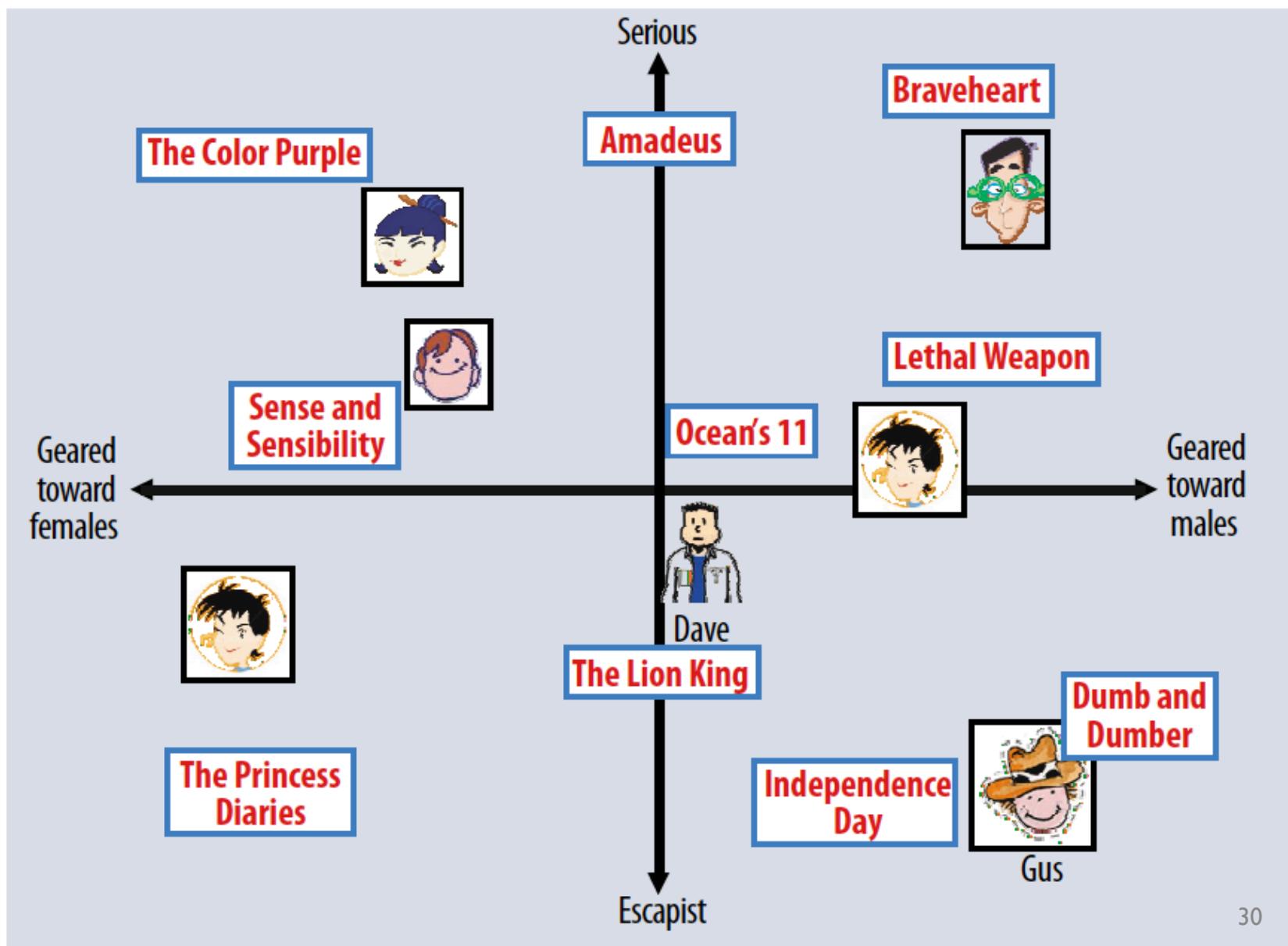
	<i>Avatar</i>	<i>The Matrix</i>	<i>Up</i>
<i>Alice</i>	?	4	2
<i>Bob</i>	3	2	?
<i>Charlie</i>	5	?	3

- ▶ Netflix competition: 500k users, 20k movies, 100M movie ratings, 3M question marks

Recovering latent factors in a matrix



Semantic Factors (Koren et al., 2009)



Latent Factor Models

- ▶ Discover latent factors ($r = 1$)

	Avatar (2.24)	The Matrix (1.92)	Up (1.18)
Alice (1.98)		4 (3.8)	2 (2.3)
Bob (1.21)	3 (2.7)	2 (2.3)	
Charlie (2.30)	5 (5.2)		3 (2.7)

- ▶ Minimum loss

$$\min_{\mathbf{W}, \mathbf{H}} \sum_{(i,j) \in Z} (\mathbf{v}_{ij} - [\mathbf{WH}]_{ij})^2$$

Latent Factor Models

- ▶ Discover latent factors ($r = 1$)

	Avatar (2.24)	The Matrix (1.92)	Up (1.18)
Alice (1.98)	?	4	2
Bob (1.21)	3 (2.7)	2 (2.3)	?
Charlie (2.30)	5 (5.2)	?	3 (2.7)

- ▶ Minimum loss

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{u}, \mathbf{m}} \sum_{(i,j) \in Z} (\mathbf{v}_{ij} - \mu - \mathbf{u}_i - \mathbf{m}_j - [\mathbf{WH}]_{ij})^2 + \lambda (\|\mathbf{W}\| + \|\mathbf{H}\| + \|\mathbf{u}\| + \|\mathbf{m}\|)$$

- ▶ Bias, regularization

Matrix completion for image denoising



Matrix factorization as SGD

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

Algorithm 1 SGD for Matrix Factorization

Require: A training set Z , initial values \mathbf{W}_0 and \mathbf{H}_0

while not converged **do** {step}

Select a training point $(i, j) \in Z$ uniformly at random.

$$\mathbf{W}'_{i*} \leftarrow \mathbf{W}_{i*} - \epsilon_n N \frac{\partial}{\partial \mathbf{W}_{i*}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

$$\mathbf{H}_{*j} \leftarrow \mathbf{H}_{*j} - \epsilon_n N \frac{\partial}{\partial \mathbf{H}_{*j}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$

$$\mathbf{W}_{i*} \leftarrow \mathbf{W}'_{i*}$$

end while

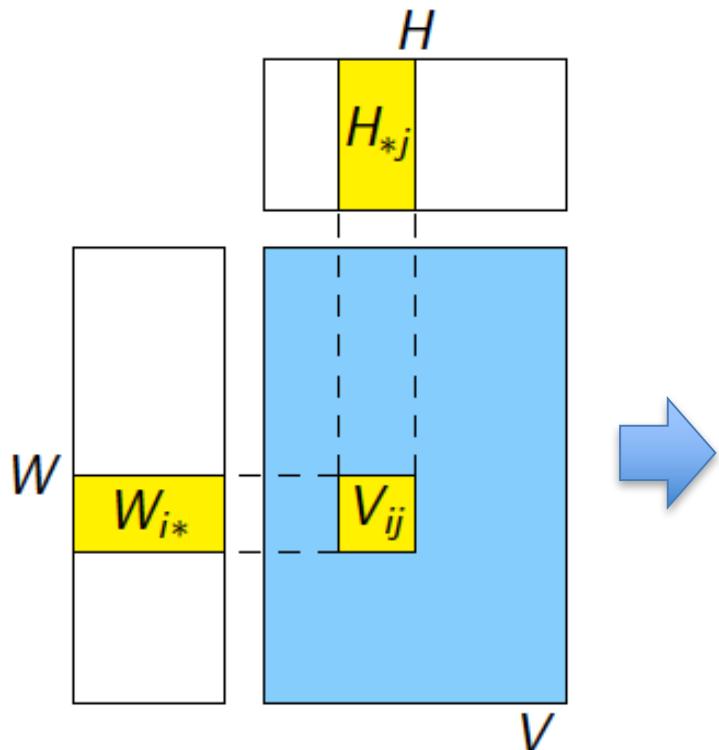
step size

why does this work?

Matrix factorization as SGD - why does this work? Here's the key claim:

require that the loss can be written as

$$L = \sum_{(i,j) \in Z} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j})$$



$$\frac{\partial}{\partial \mathbf{W}_{i'k}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } i \neq i' \\ \frac{\partial}{\partial \mathbf{W}_{ik}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial \mathbf{H}_{kj'}} L_{ij}(\mathbf{W}, \mathbf{H}) = \begin{cases} 0 & \text{if } j \neq j' \\ \frac{\partial}{\partial \mathbf{H}_{kj}} l(\mathbf{V}_{ij}, \mathbf{W}_{i*}, \mathbf{H}_{*j}) & \text{otherwise} \end{cases}$$

Checking the claim

$$\frac{\partial}{\partial \mathbf{W}_{i*}} L(\mathbf{W}, \mathbf{H}) = \frac{\partial}{\partial \mathbf{W}_{i*}} \sum_{(i',j) \in Z} L_{i'j}(\mathbf{W}_{i'*}, \mathbf{H}_{*j}) = \sum_{j \in Z_{i*}} \frac{\partial}{\partial \mathbf{W}_{i*}} L_{ij}(\mathbf{W}_{i*}, \mathbf{H}_{*j}),$$

where $Z_{i*} = \{ j : (i, j) \in Z \}$.

$$\frac{\partial}{\partial \mathbf{H}_{*j}} L(\mathbf{W}, \mathbf{H}) = \sum_{i \in Z_{*j}} \frac{\partial}{\partial \mathbf{W}_{*j}} L_{ij}(\mathbf{W}_{i*}, \mathbf{H}_{*j}),$$

where $Z_{*j} = \{ i : (i, j) \in Z \}$.

Think for SGD for logistic regression

- LR loss = compare y and $\hat{y} = \text{dot}(w, x)$
- similar but now update w (user weights) and x (movie weight)

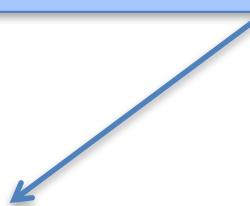
What loss functions are possible?

N1, N2 - diagonal matrixes, sort of like IDF factors for the users/movies

$$L_{\text{NZSL}} = \sum_{(i,j) \in Z} (\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2$$

$$L_{\text{L2}} = L_{\text{NZSL}} + \lambda(\|\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H}\|_{\text{F}}^2)$$

$$L_{\text{NZL2}} = L_{\text{NZSL}} + \lambda(\|\mathbf{N}_1\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H}\mathbf{N}_2\|_{\text{F}}^2)$$



What loss functions are possible?

Loss Function Definition and Derivatives

$$L_{\text{NZSL}} \quad L_{\text{NZSL}} = \sum_{(i,j) \in Z} (\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})^2$$

$$\frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})\mathbf{H}_{kj}$$

$$\frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} = -2(\mathbf{V}_{ij} - [\mathbf{W}\mathbf{H}]_{ij})\mathbf{W}_{ik}$$

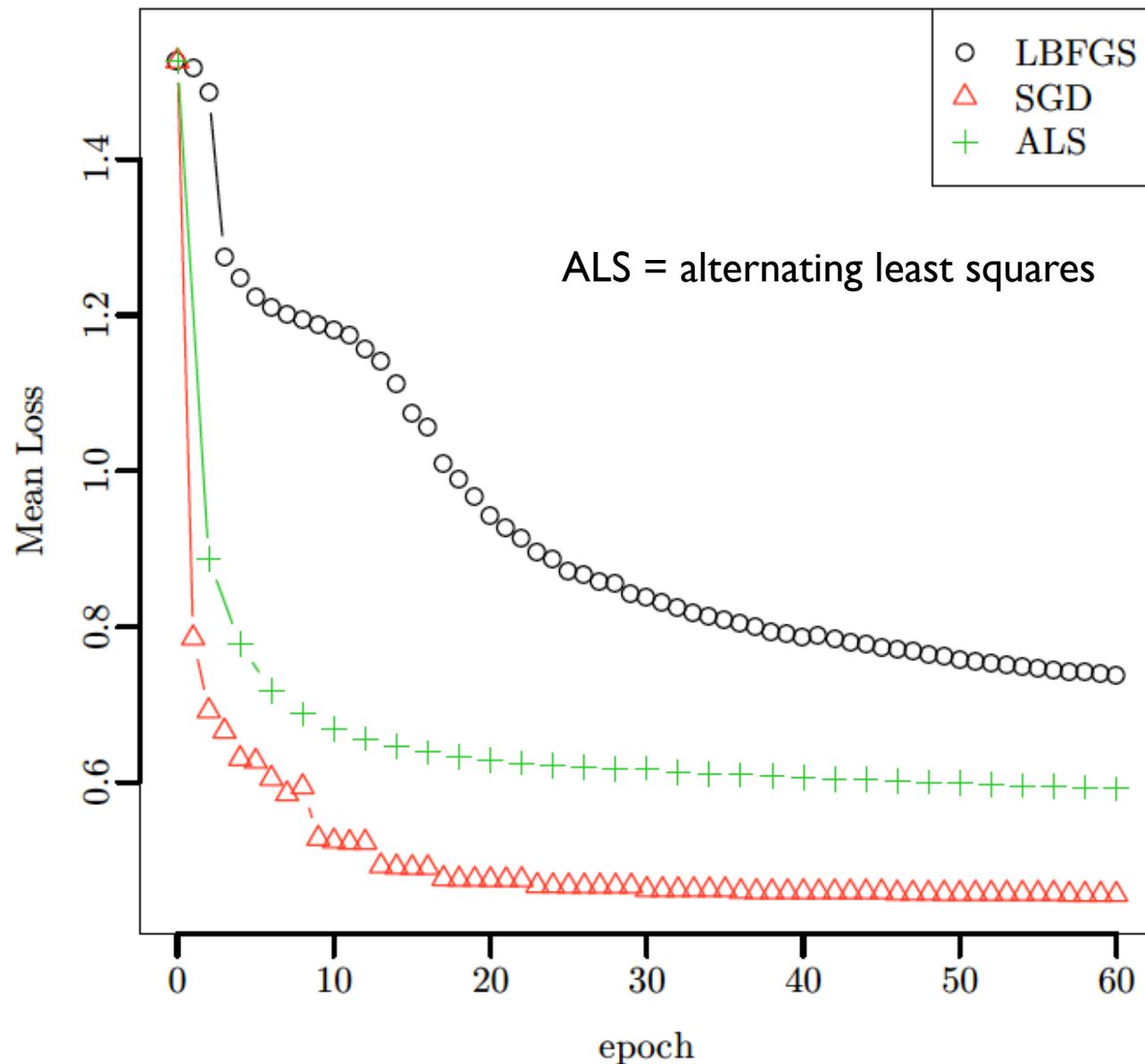
What loss functions are possible?

Loss Function Definition and Derivatives

$$\begin{aligned} L_{\text{L2}} &= L_{\text{NZSL}} + \lambda (\|\mathbf{W}\|_{\text{F}}^2 + \|\mathbf{H}\|_{\text{F}}^2) \\ &= \sum_{(i,j) \in Z} \left[(\mathbf{V}_{ij} - [\mathbf{WH}]_{ij})^2 + \lambda \left(\frac{\|\mathbf{W}_{i*}\|_{\text{F}}^2}{N_{i*}} + \frac{\|\mathbf{H}_{*j}\|_{\text{F}}^2}{N_{*j}} \right) \right] \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{W}_{ik}} L_{ij} &= -2(\mathbf{V}_{ij} - [\mathbf{WH}]_{ij}) \mathbf{H}_{kj} + 2\lambda \frac{\mathbf{W}_{ik}}{N_{i*}} \\ \frac{\partial}{\partial \mathbf{H}_{kj}} L_{ij} &= -2(\mathbf{V}_{ij} - [\mathbf{WH}]_{ij}) \mathbf{W}_{ik} + 2\lambda \frac{\mathbf{H}_{kj}}{N_{*j}} \end{aligned}$$

Stochastic Gradient Descent on Netflix Data



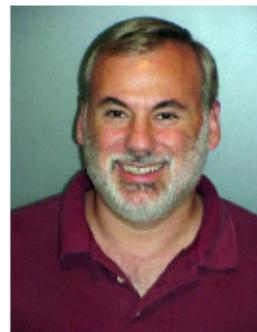
Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent

Rainer Gemulla

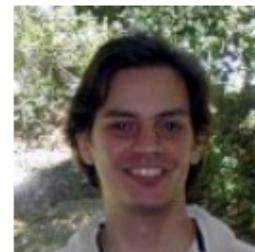


talk pilfered from →

Peter J. Haas



Yannis Sismanis



Erik Nijkamp



Outline

Matrix Factorization

Stochastic Gradient Descent

Distributed SGD with MapReduce

Experiments

Summary

Averaging Techniques

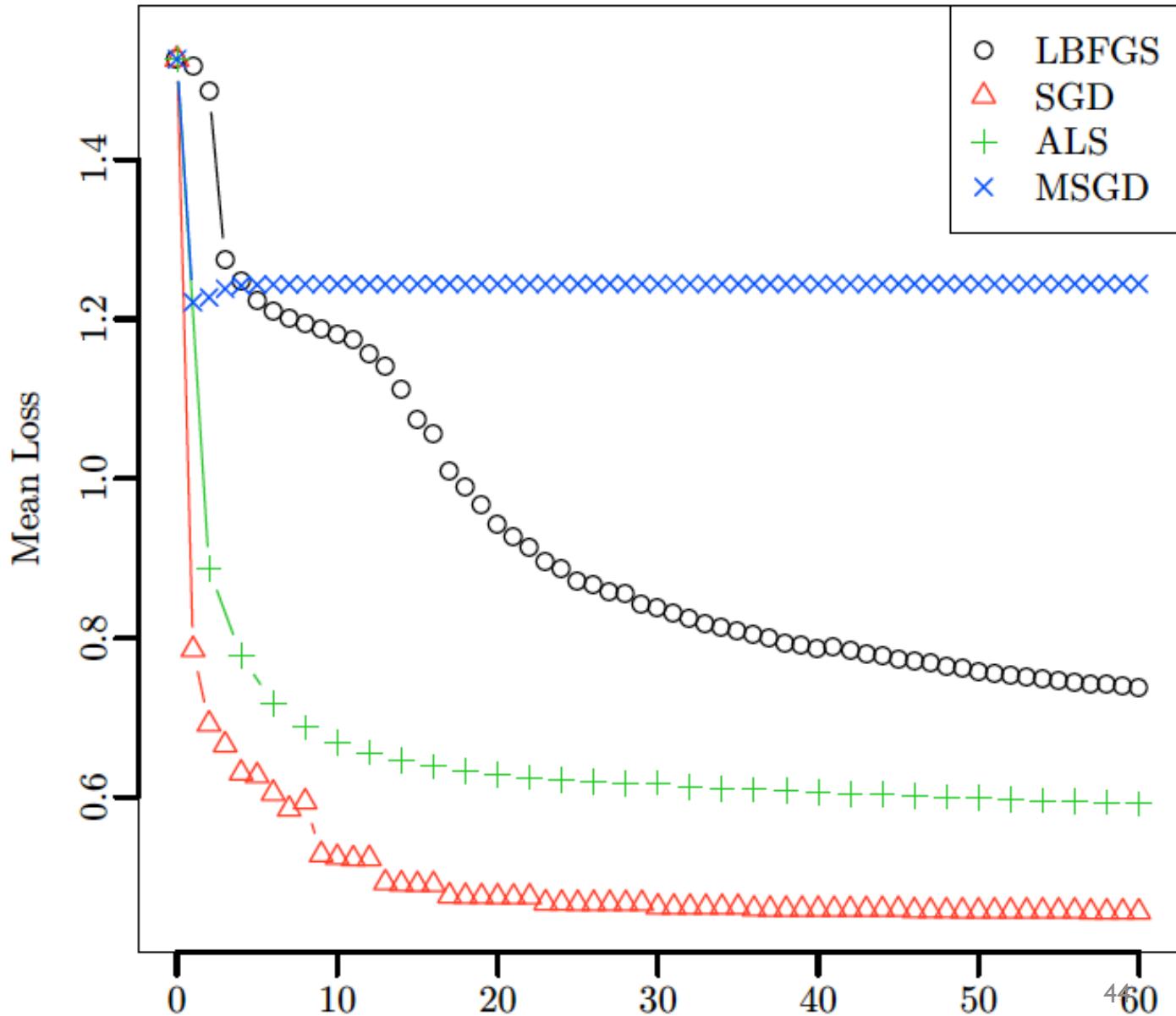
- ▶ SGD steps depend on each other

$$\theta_{n+1} = \theta_n - \epsilon_n \hat{L}'(\theta_n)$$

How to distribute?

- ▶ Parameter mixing (MSGD)
 - ▶ *Map*: Run independent instances of SGD on subsets of the data (until convergence)
 - ▶ *Reduce*: Average results

Averaging Techniques



Averaging Techniques

- ▶ SGD steps depend on each other

$$\theta_{n+1} = \theta_n - \epsilon_n \hat{L}'(\theta_n)$$

How to distribute?

- ▶ Parameter mixing (MSGD)

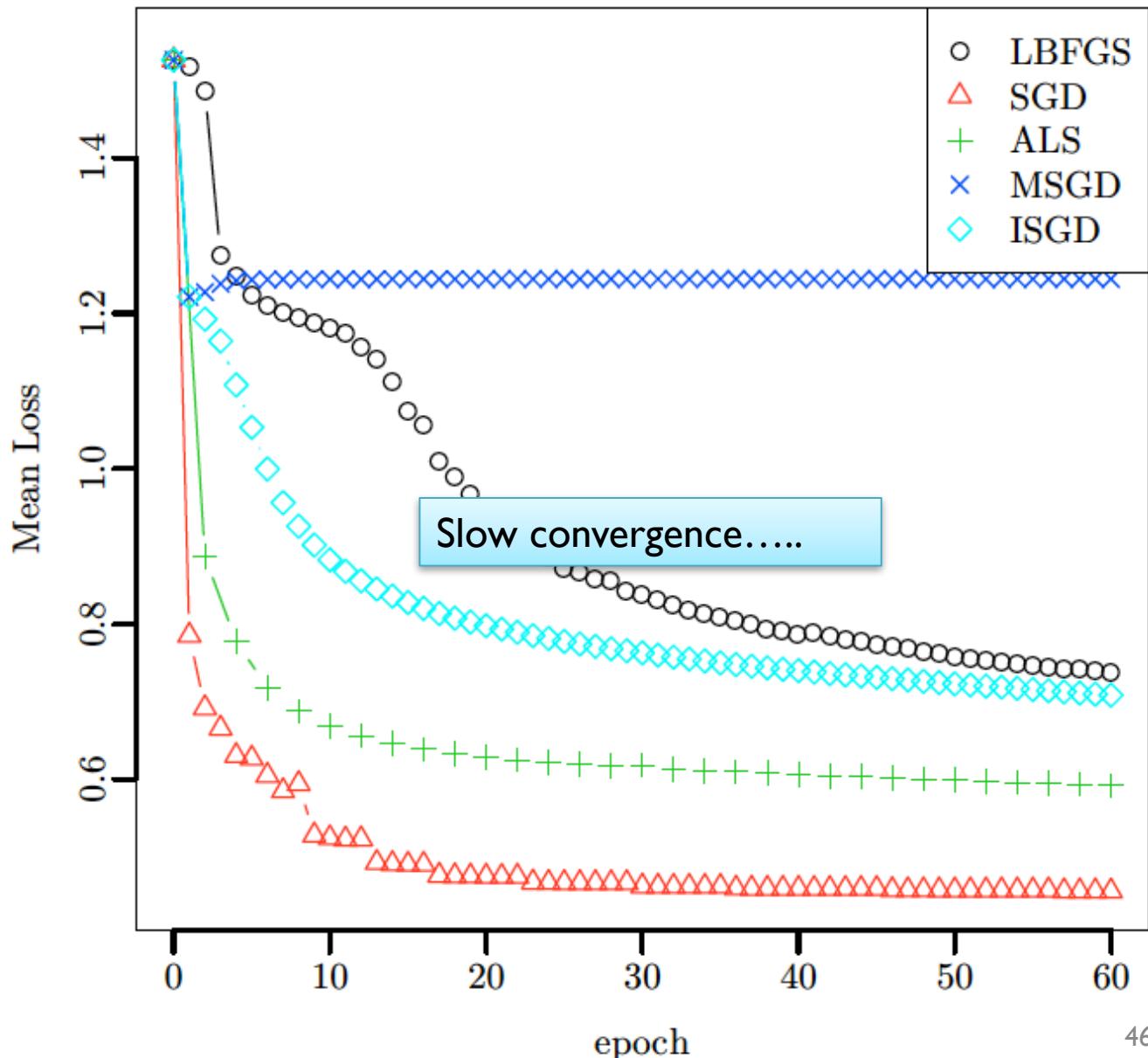
- ▶ *Map*: Run independent instances of SGD on subsets of the data (until convergence)
- ▶ *Reduce*: Average results
- ▶ Does not converge to correct solution!

Similar to McDonnell et al
with perceptron learning

- ▶ Iterative Parameter mixing (ISGD)

- ▶ *Map*: Run independent instances of SGD on subsets of the data (for some time)
- ▶ *Reduce*: Average results
- ▶ Repeat

Averaging Techniques



Problem Structure

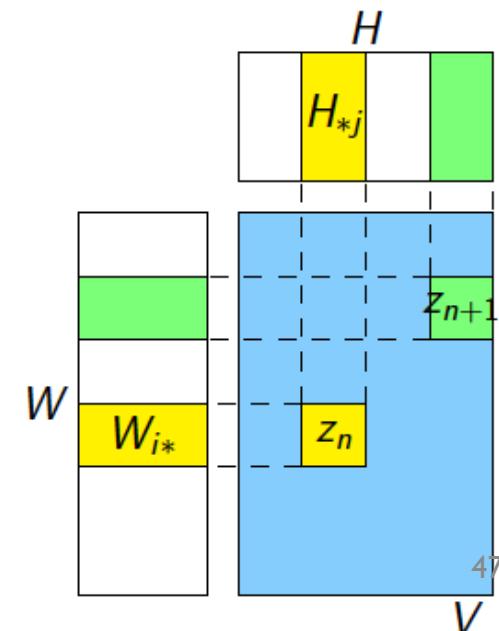
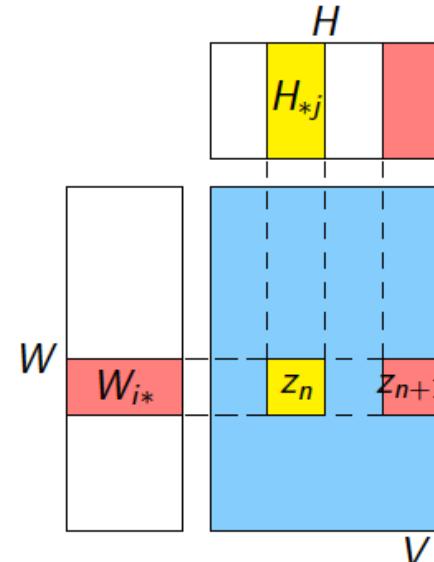
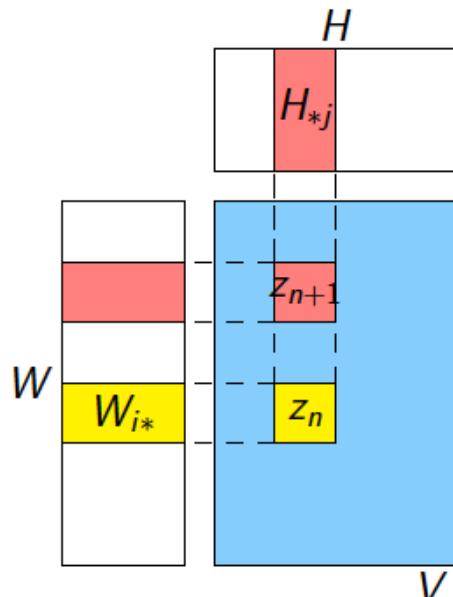
- ▶ SGD steps depend on each other

$$\theta_{n+1} = \theta_n - \epsilon_n \hat{L}'(\theta_n)$$

- ▶ An SGD step on example $z \in Z \dots$

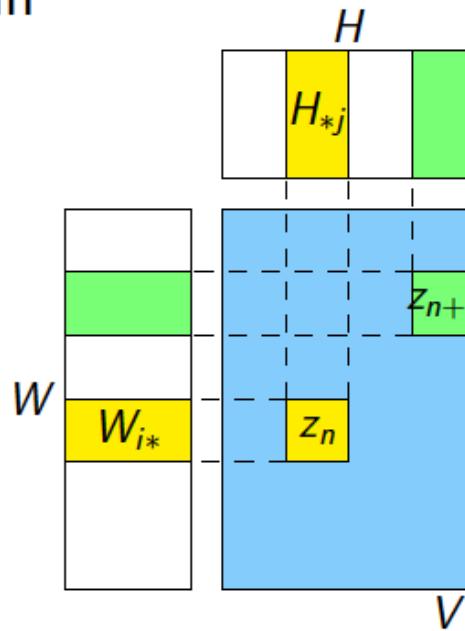
1. Reads W_{i_z*} and H_{*j_z}
2. Performs gradient computation $L'_{ij}(W_{i_z*}, H_{*j_z})$
3. Updates W_{i_z*} and H_{*j_z}

- ▶ Not all steps are dependent



Interchangeability

- ▶ Two elements $z_1, z_2 \in Z$ are *interchangeable* if they share neither row nor column



- ▶ When z_n and z_{n+1} are interchangeable, the SGD steps

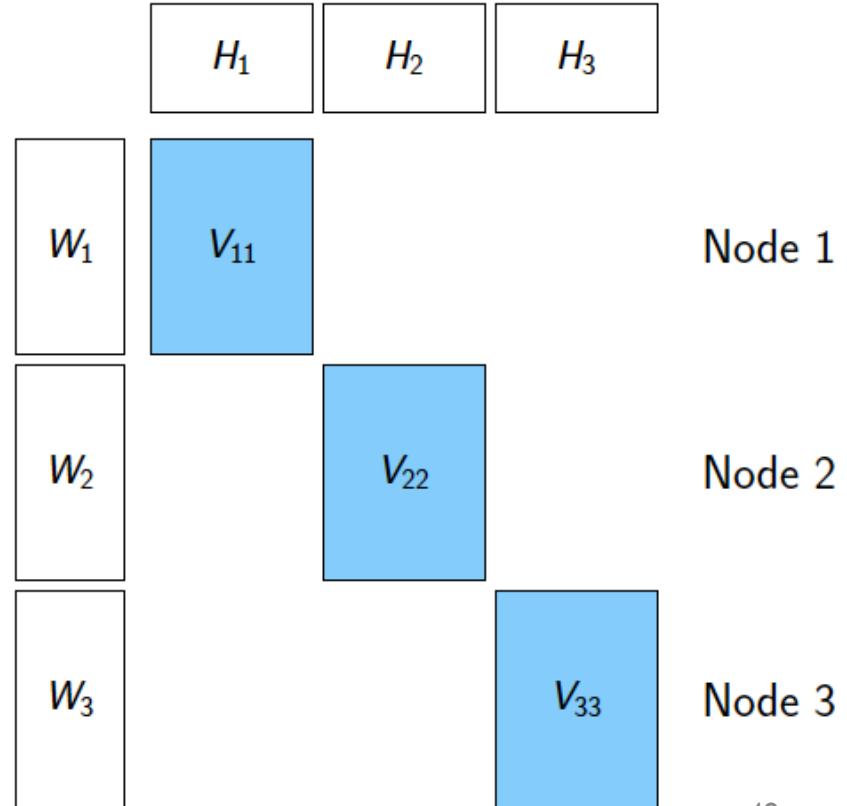
$$\begin{aligned}\theta_{n+2} &= \theta_n - \epsilon \hat{L}'(\theta_n, z_n) - \epsilon \hat{L}'(\theta_{n+1}, z_{n+1}) \\ &= \theta_n - \epsilon \hat{L}'(\theta_n, z_n) - \epsilon \hat{L}'(\theta_n, z_{n+1}),\end{aligned}$$

become parallelizable!

Exploitation

- ▶ Block and distribute the input matrix \mathbf{V}
- ▶ High-level approach (Map only)
 1. Pick a “diagonal”
 2. Run SGD on the diagonal (in parallel)
 3. Merge the results
 4. Move on to next “diagonal”

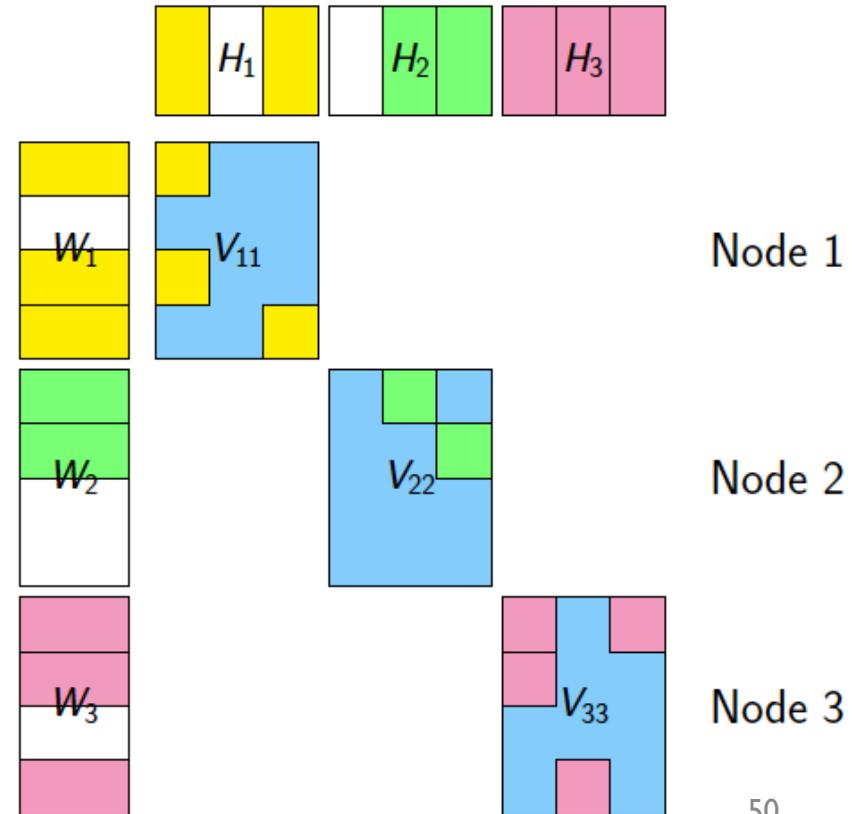
▶ Steps 1–3 form a *cycle*



Exploitation

- ▶ Block and distribute the input matrix \mathbf{V}
- ▶ High-level approach (Map only)
 1. Pick a “diagonal”
 2. Run SGD on the diagonal (in parallel)
 3. Merge the results
 4. Move on to next “diagonal”
 - ▶ Steps 1–3 form a *cycle*

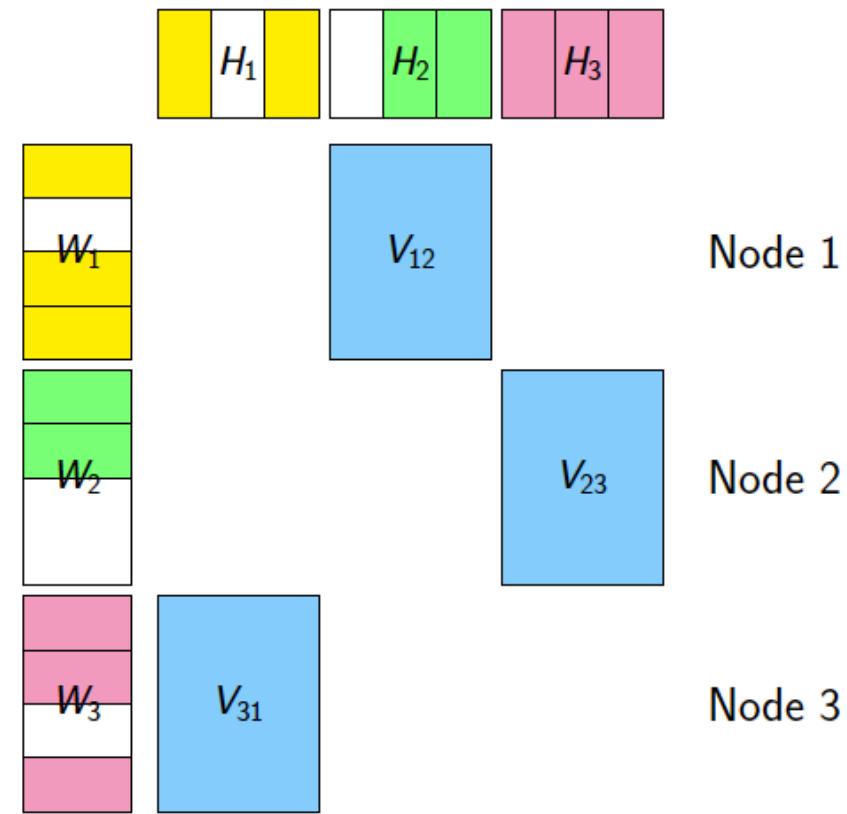
- ▶ Step 2:
Simulate sequential SGD
 - ▶ Interchangeable blocks
 - ▶ Throw dice of how many iterations per block
 - ▶ Throw dice of which step sizes per block



Exploitation

- ▶ Block and distribute the input matrix \mathbf{V}
- ▶ High-level approach (Map only)
 1. Pick a “diagonal”
 2. Run SGD on the diagonal (in parallel)
 3. Merge the results
 4. Move on to next “diagonal”
 - ▶ Steps 1–3 form a *cycle*

- ▶ Step 2:
Simulate sequential SGD
 - ▶ Interchangeable blocks
 - ▶ Throw dice of how many iterations per block
 - ▶ Throw dice of which step sizes per block
 - ▶ Instance of “stratified SGD”
 - ▶ Provably correct



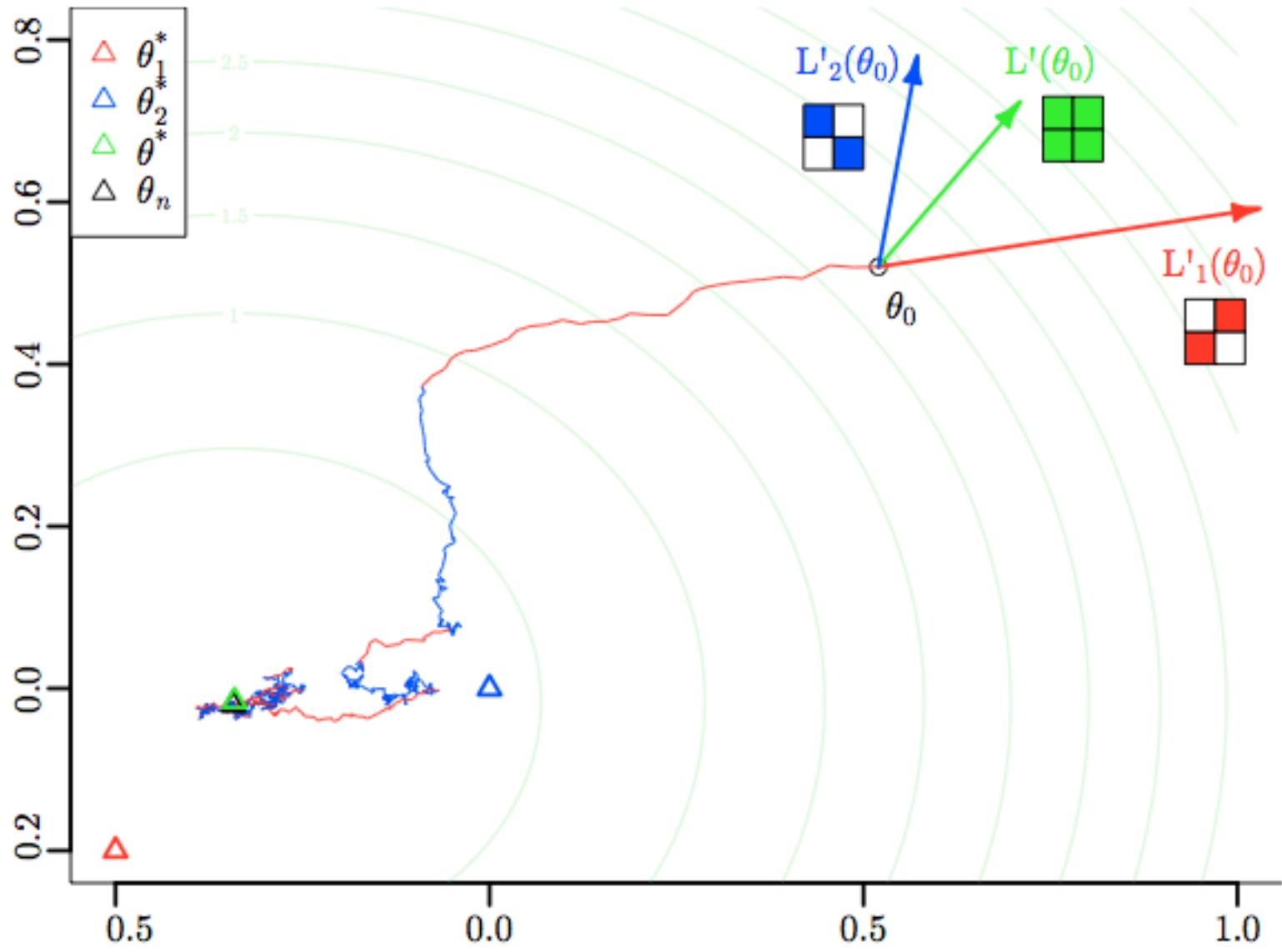


Figure 2: Example of stratified SGD

More detail....

- Randomly permute rows/cols of matrix
- Chop V, W, H into blocks of size $d \times d$
 - m/d blocks in W , n/d blocks in H
- Group the data:
 - Pick a set of blocks with no overlapping rows or columns (a *stratum*)
 - Repeat until all blocks in V are covered
- Train the SGD
 - Process strata in series
 - Process blocks within a stratum in parallel

More detail....

Algorithm 2 DSGD for Matrix Factorization

Require: Z, W_0, H_0 , cluster size d

Z was V

$W \leftarrow W_0$

$H \leftarrow H_0$

Block $Z / W / H$ into $d \times d / d \times 1 / 1 \times d$ blocks

while not converged **do** /* epoch */

 Pick step size ϵ

for $s = 1, \dots, d$ **do** /* subepoch */

 Pick d blocks $\{Z^{1j_1}, \dots, Z^{dj_d}\}$ to form a stratum

for $b = 1, \dots, d$ **do** /* in parallel */

 Run SGD on the training points in Z^{bj_b} (step size = ϵ)

end for

end for

end while

More detail....

- Initialize W, H randomly
 - not at zero ☺
- Choose a random ordering (random sort) of the points in a stratum in each “sub-epoch”
- Pick strata sequence by permuting rows and columns of M , and using $M'[k,i]$ as column index of row i in subepoch k
- Use “bold driver” to set step size:
 - increase step size when loss decreases (in an epoch)
 - decrease step size when loss increases
- Implemented in Hadoop and R/Snowfall

$$M = \begin{pmatrix} 1 & 2 & \cdots & d \\ 2 & 3 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ d & 1 & \cdots & d-1 \end{pmatrix}.$$

Outline

Matrix Factorization

Stochastic Gradient Descent

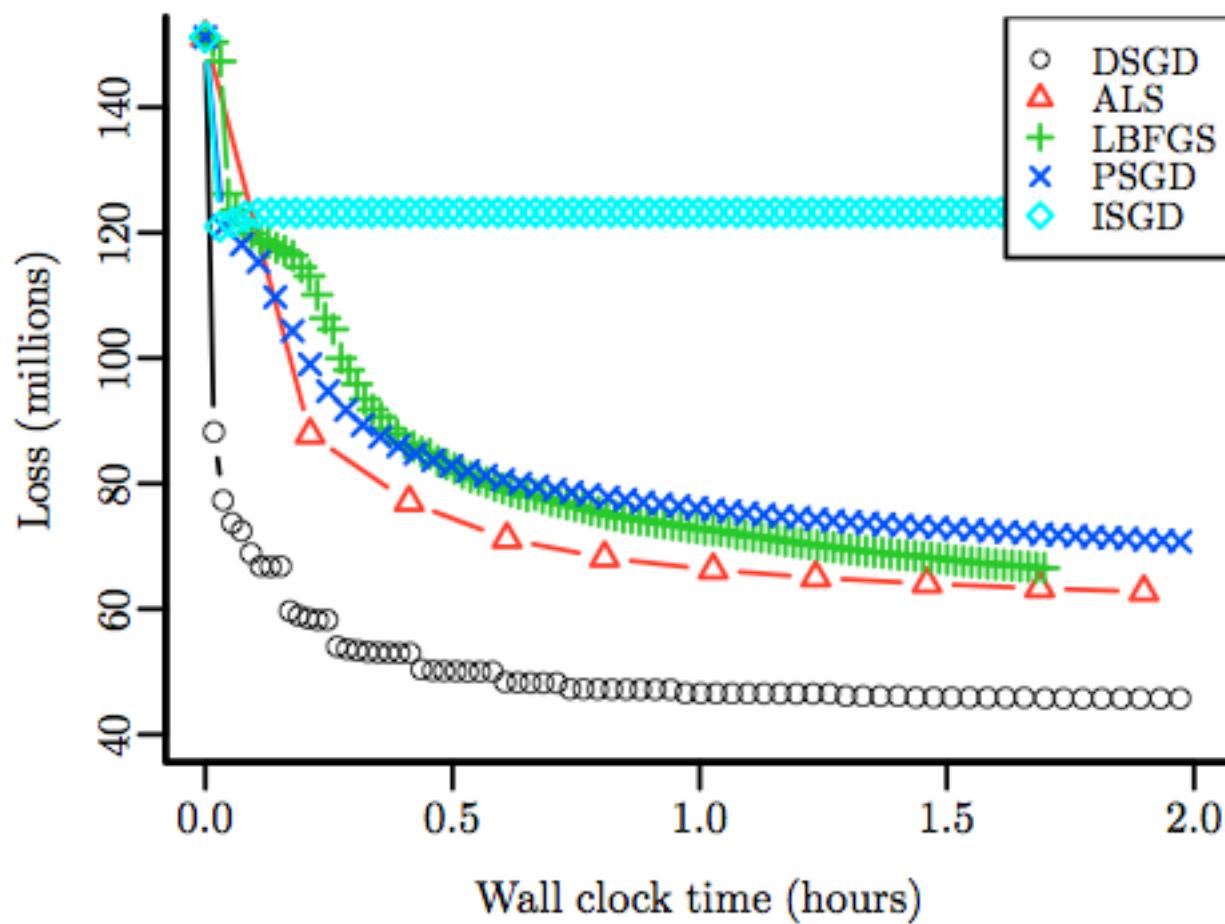
Distributed SGD with MapReduce

Experiments

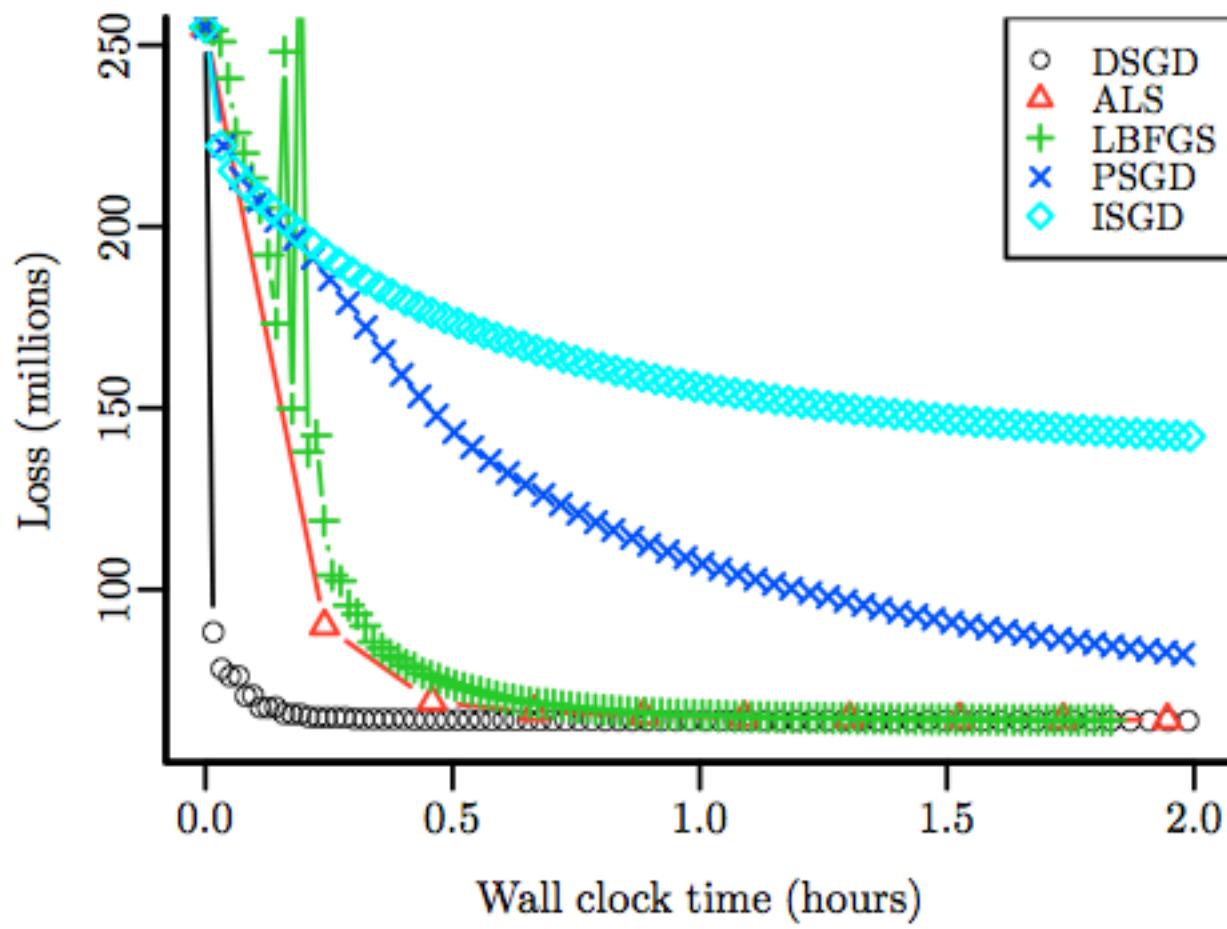
Summary

Wall Clock Time

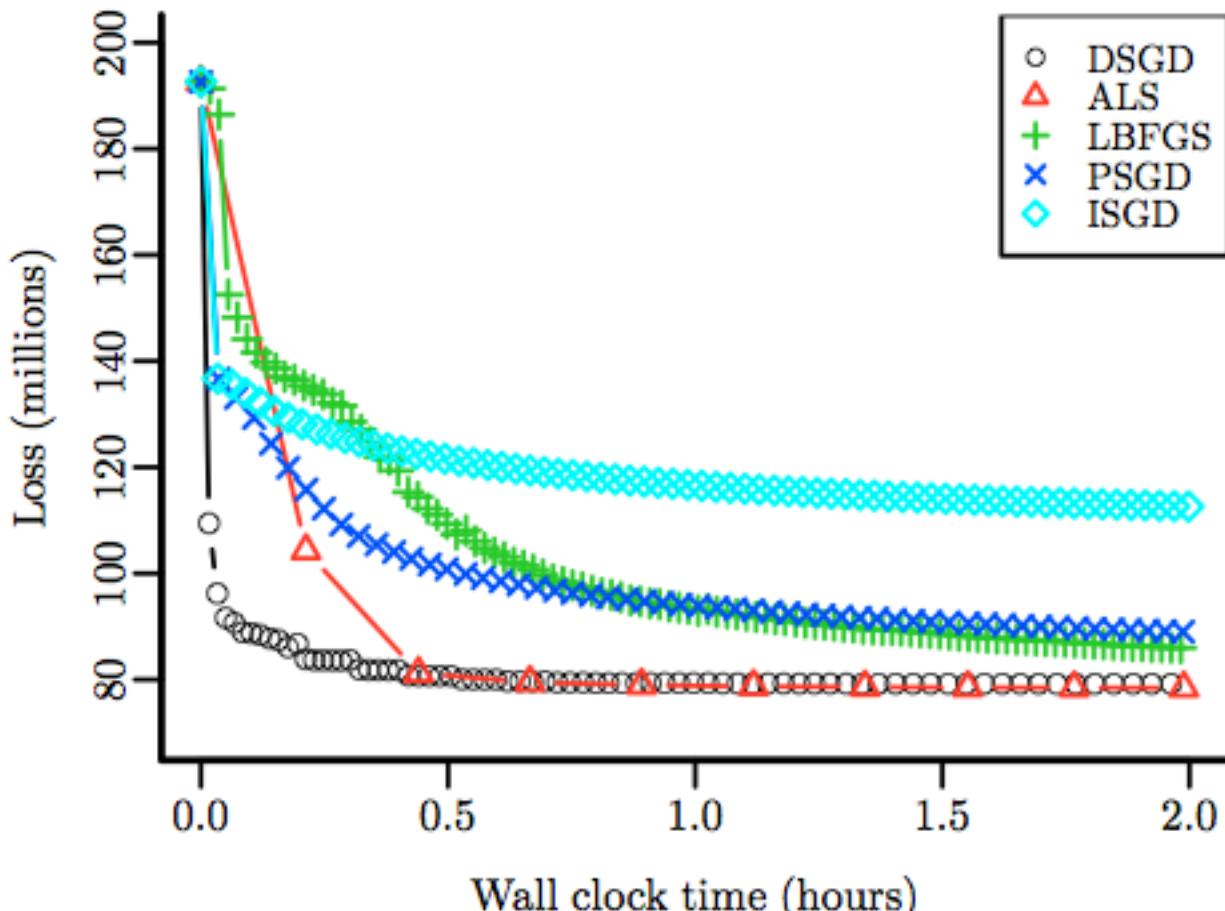
8 nodes, 64 cores, R/snow



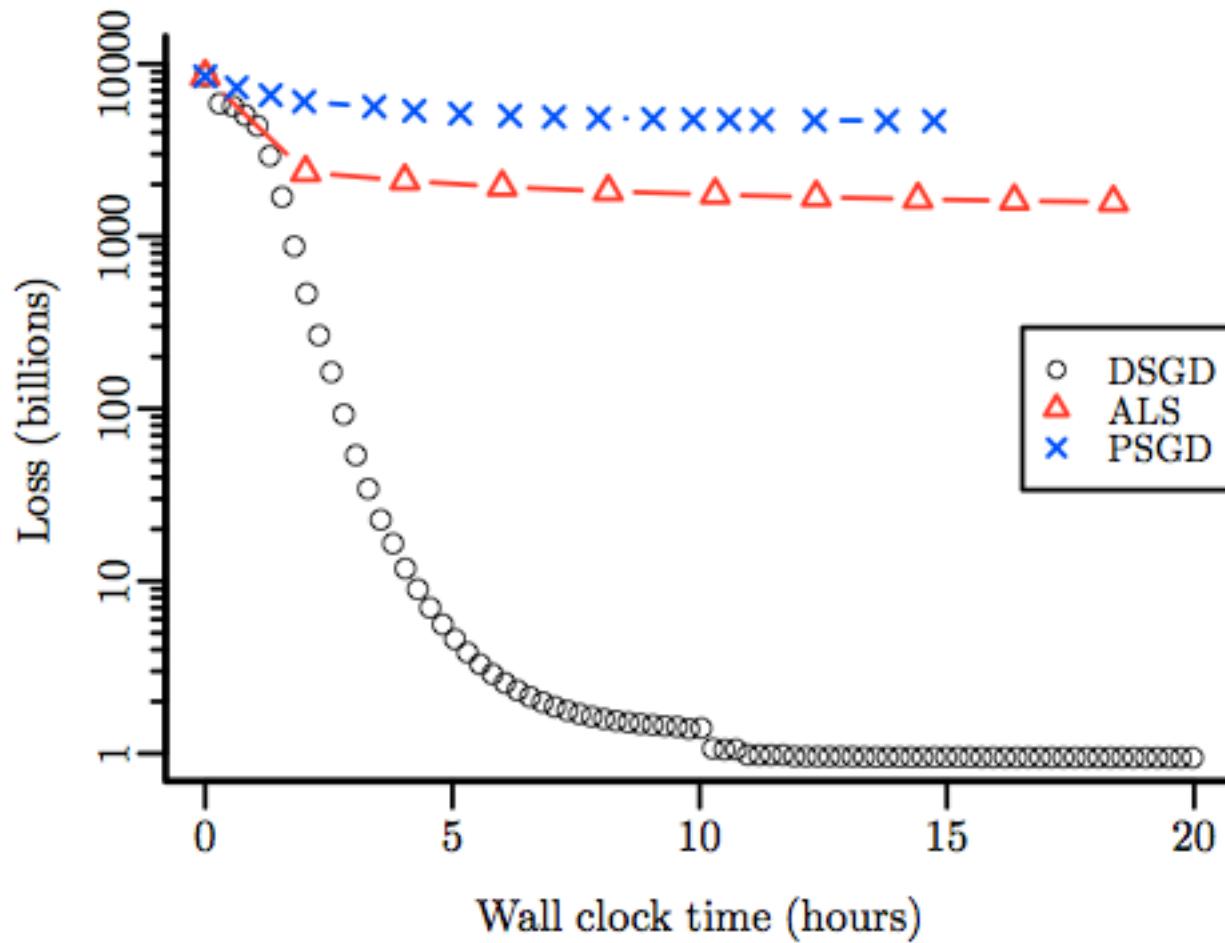
(a) Netflix, NZSL



(b) Netflix, L2, $\lambda = 50$

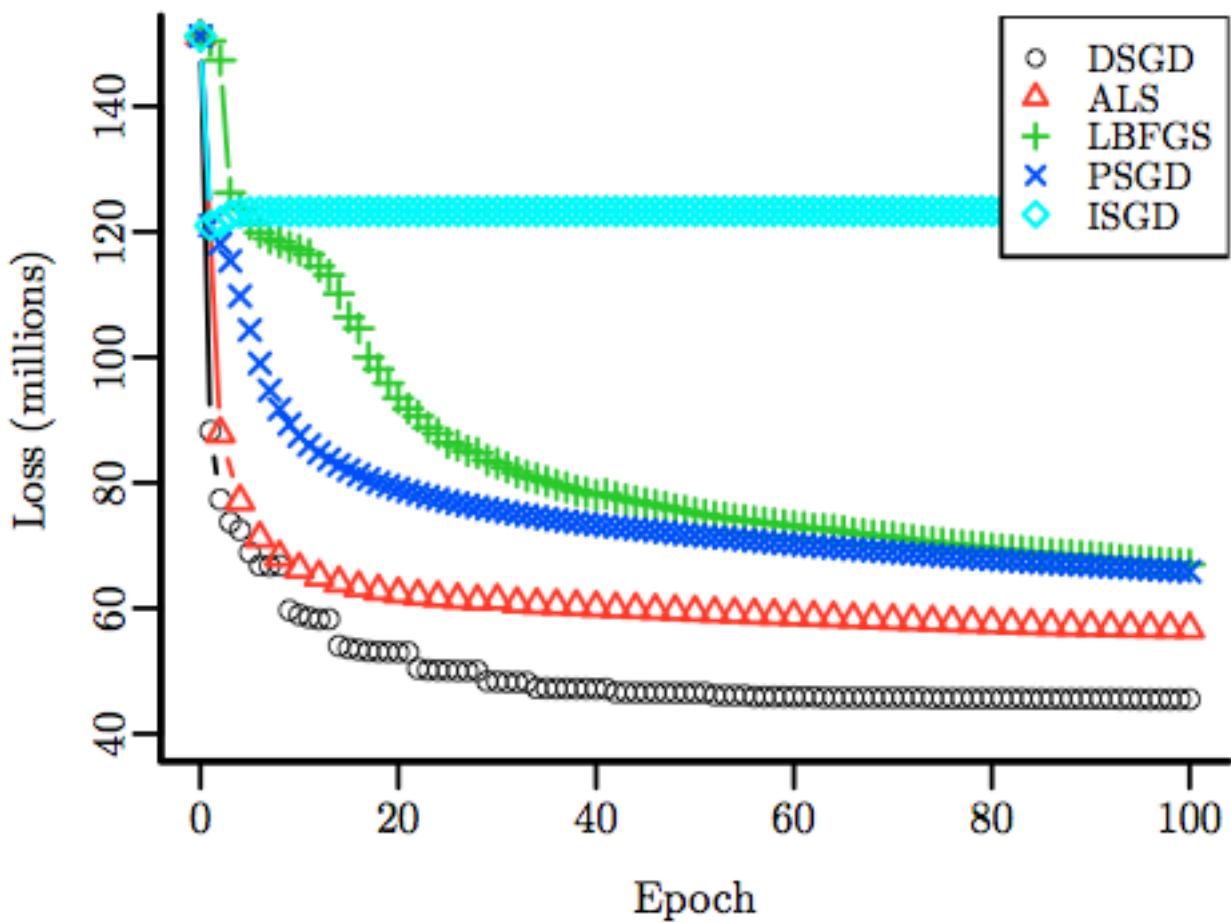


(c) Netflix, NZL2, $\lambda = 0.05$

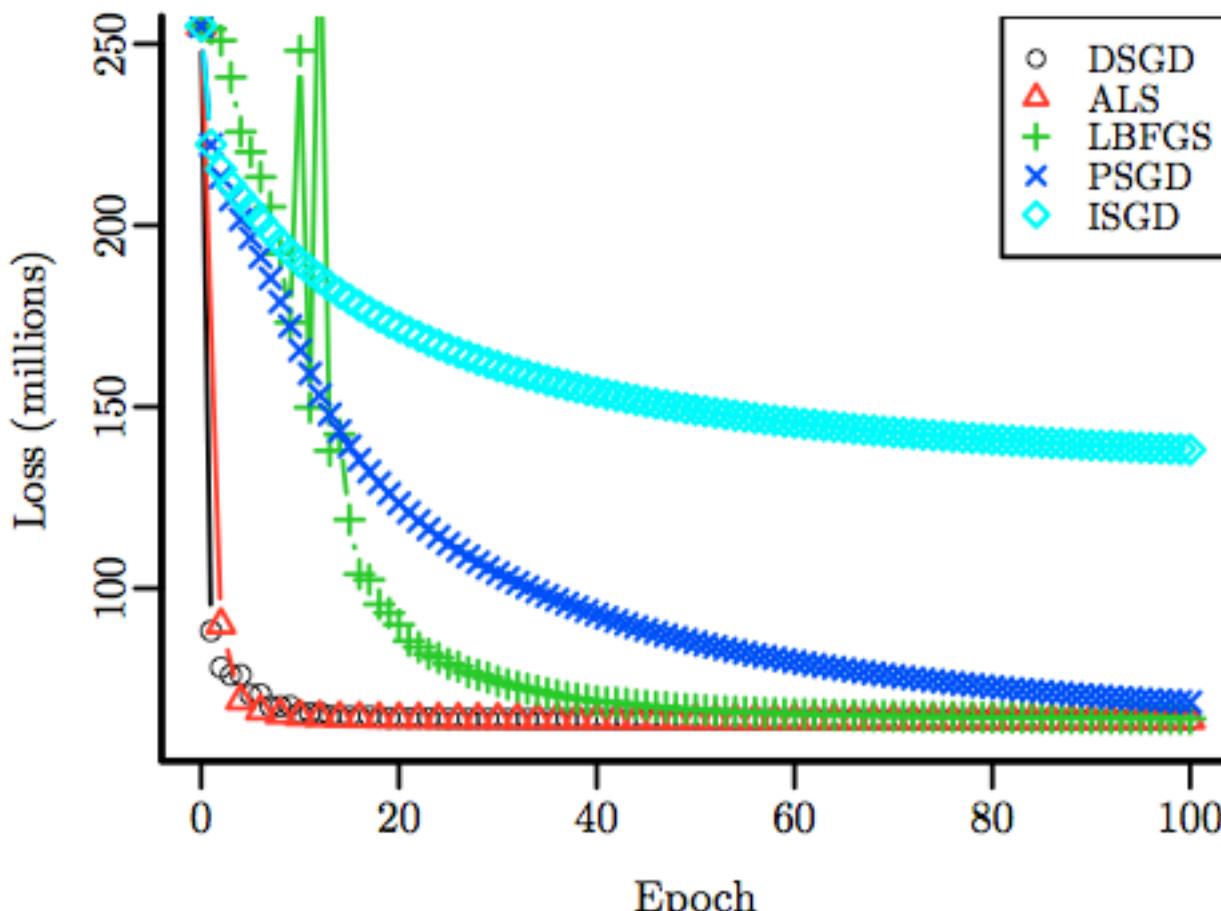


(d) Synthetic data, L2, $\lambda = 0.1$

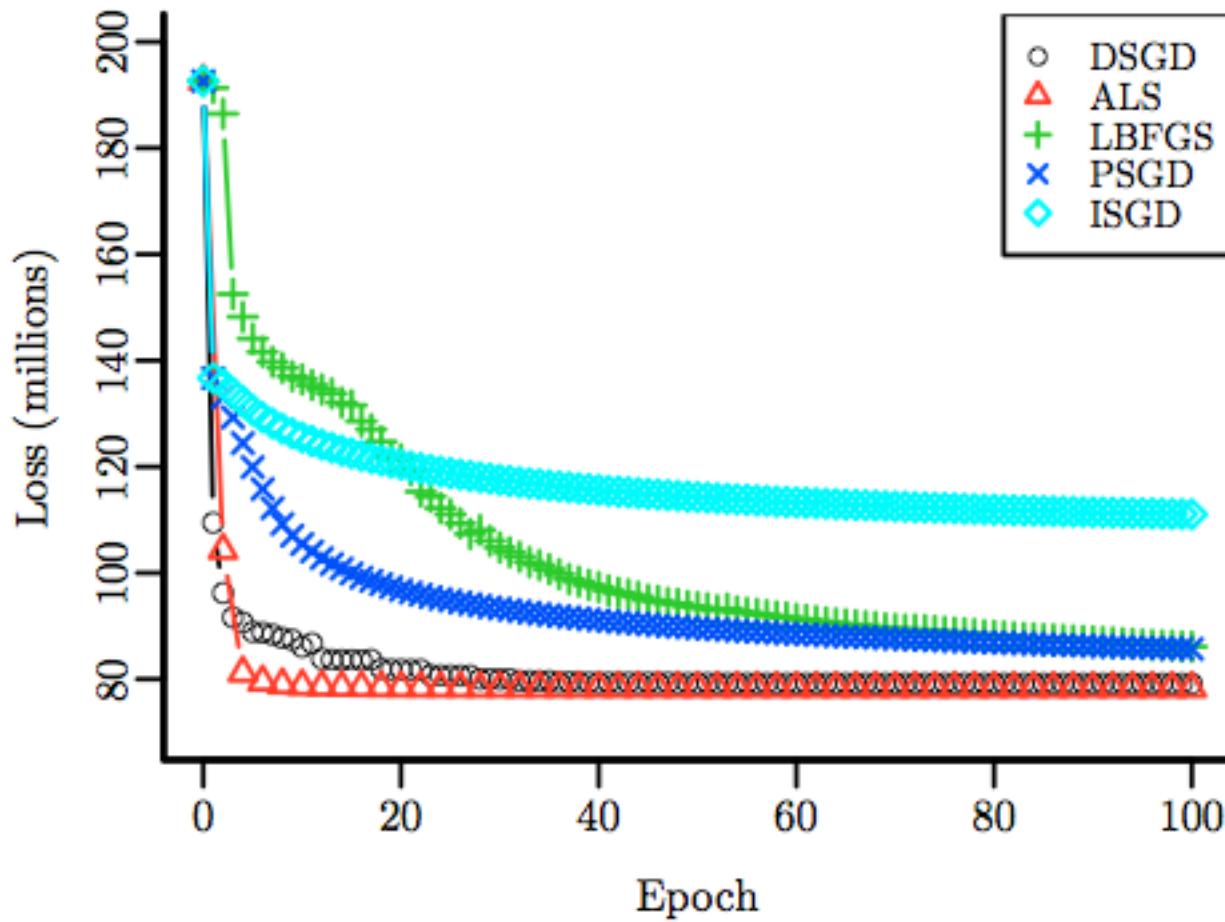
Number of Epochs



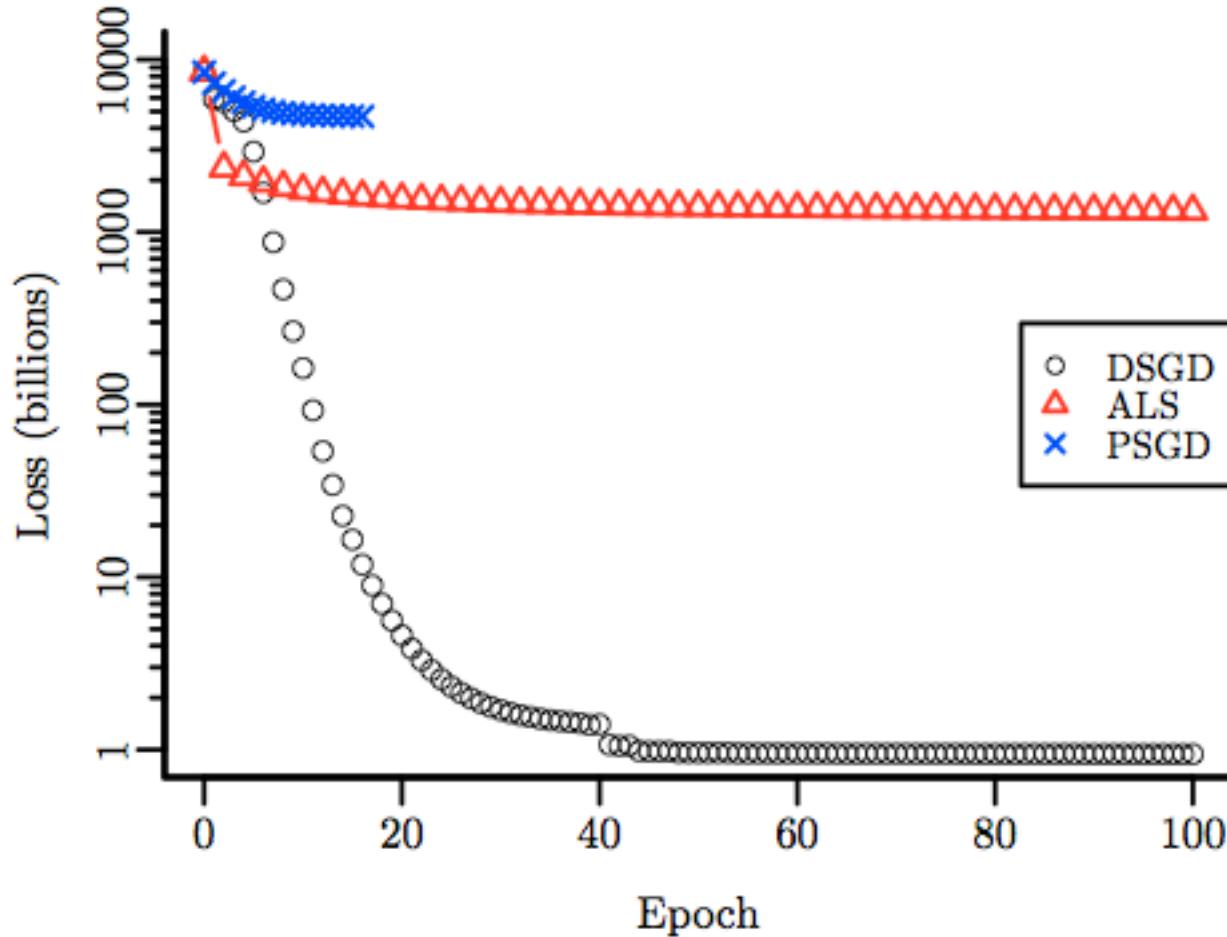
(a) Netflix, NZSL



(b) Netflix, L2, $\lambda = 50$

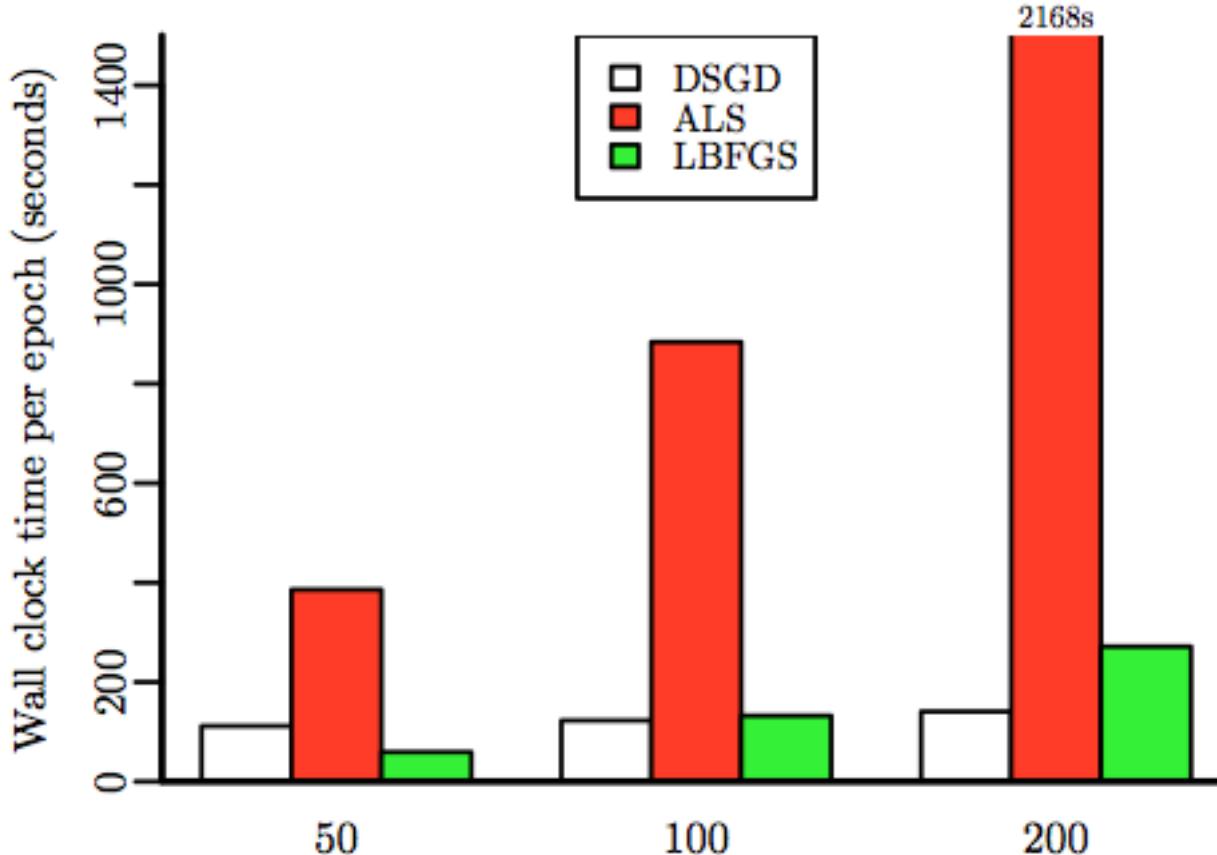


(c) Netflix, NZL2, $\lambda = 0.05$

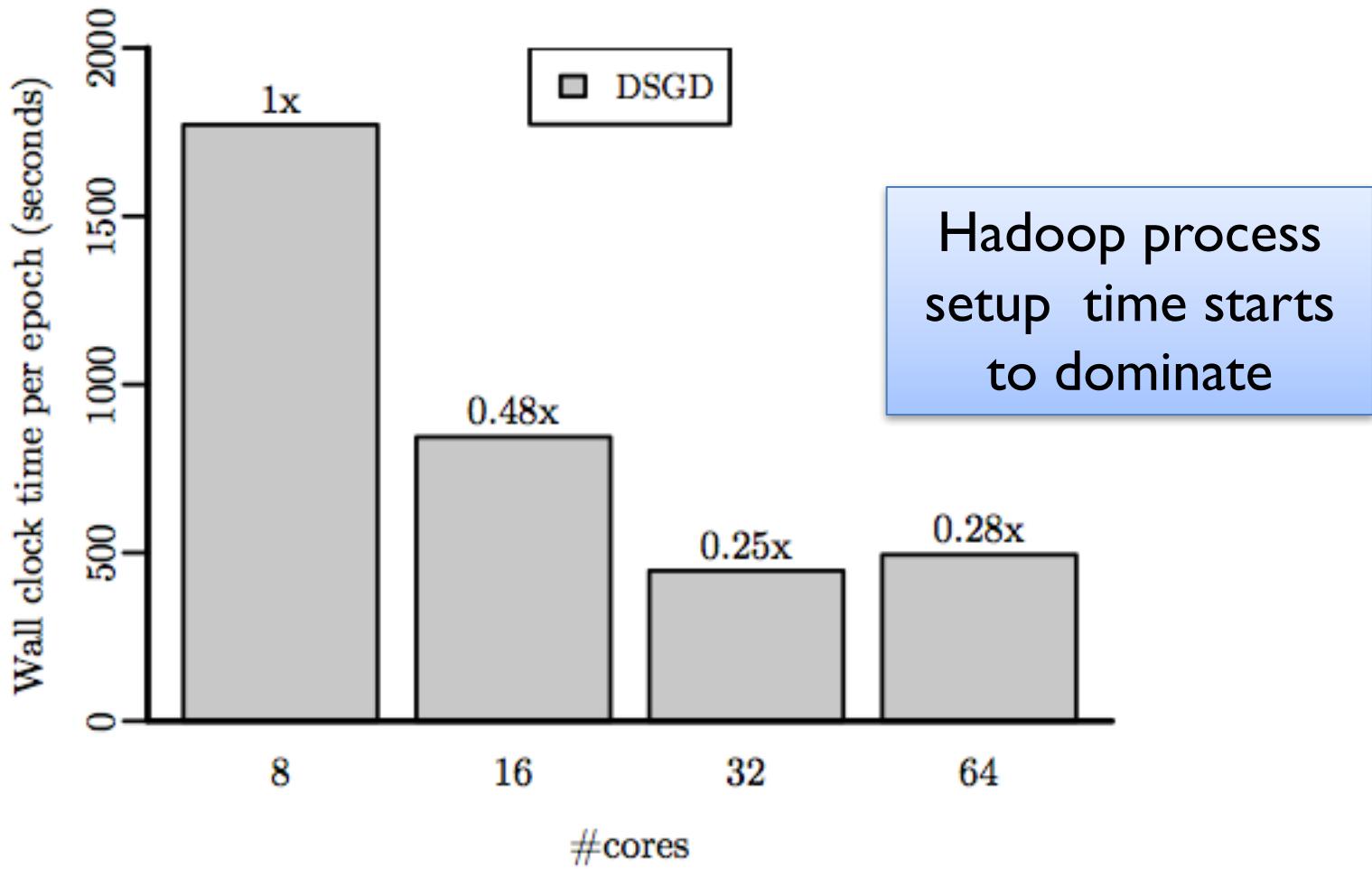


(d) Synthetic data, L2, $\lambda = 0.1$

Varying rank 100 epochs for all

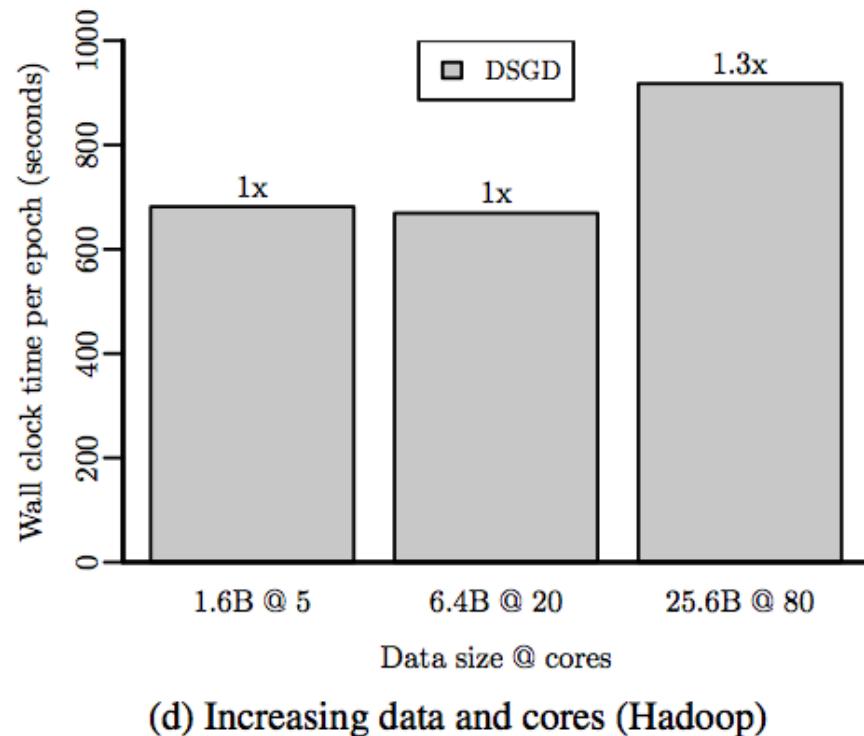
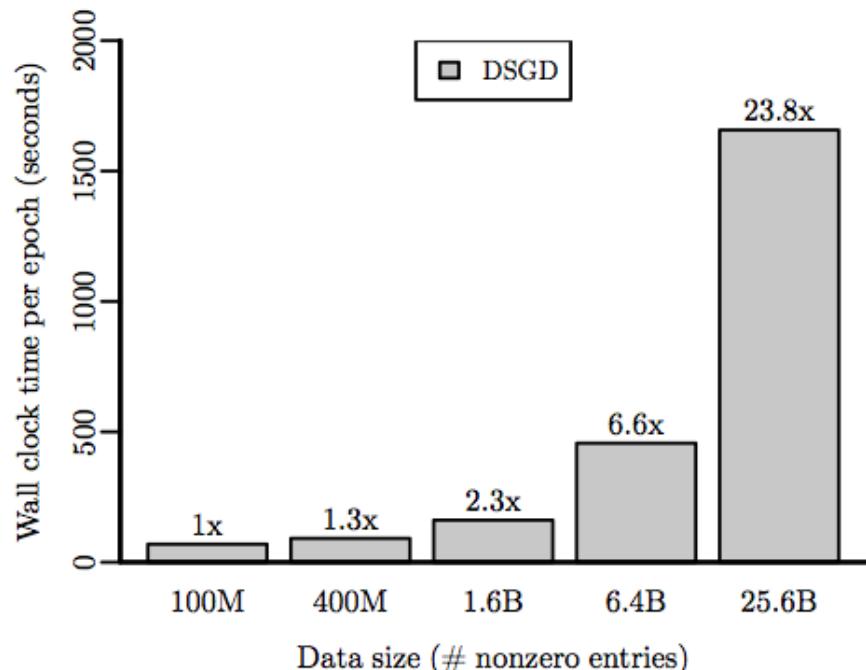


Hadoop scalability



(b) Increasing cores (Hadoop, 6.4B entries)

Hadoop scalability



Summary

- ▶ Matrix factorization
 - ▶ Widely applicable via customized loss functions
 - ▶ Large instances (millions × millions with billions of entries)
- ▶ Distributed Stochastic Gradient Descent
 - ▶ Simple and versatile
 - ▶ Avoids averaging via novel “stratified SGD” variant
 - ▶ Achieves
 - ▶ Fully distributed data/model
 - ▶ Fully distributed processing
 - ▶ Competitive to alternative algorithms
 - ▶ Fast, scalable
- ▶ Future Directions
 - ▶ Improved stratification
 - ▶ Simultaneous computation & communication
 - ▶ Stratification for other models
 - ▶ ...